

Contents

1	Notation and conventions	2
1.0.1	Background Information	3
1.1	Some Useful Mathematical Facts	4
1.2	Acknowledgements	4
2	Some Preliminaries	5
2.1	Notation and conventions	5
2.1.1	Background Information	6
2.2	Some Useful Mathematical Facts	7
2.3	Acknowledgements	7
2.4	The Curse of Dimension	7
2.4.1	The Curse: Data isn't Where You Think it is	7
2.4.2	Minor Banes of Dimension	9
3	Learning to Classify	10
3.1	Classification, Error, and Loss	10
3.1.1	Loss and the Cost of Misclassification	11
3.1.2	Building a Classifier from Probabilities	11
3.1.3	Building a Classifier using Decision Boundaries	12
3.1.4	What will happen on Test Data?	12
3.1.5	The Class Confusion Matrix	14
3.1.6	Statistical Learning Theory and Generalization	15
3.2	Classifying with Naive Bayes	16
3.3	The Support Vector Machine	20
3.3.1	Choosing a Classifier with the Hinge Loss	21
3.3.2	Finding a Minimum: General Points	23
3.3.3	Finding a Minimum: Stochastic Gradient Descent	24
3.3.4	Example: Training an SVM with Stochastic Gradient Descent	26
3.3.5	Multi-Class Classifiers	28
3.4	Classifying with Random Forests	29
3.4.1	Building a Decision Tree	29
3.4.2	Entropy and Information Gain	32
3.4.3	Entropy and Splits	34
3.4.4	Choosing a Split with Information Gain	35
3.4.5	Forests	37
3.4.6	Building and Evaluating a Decision Forest	37
3.4.7	Classifying Data Items with a Decision Forest	38
3.5	Classifying with Nearest Neighbors	40
3.6	You should	43
3.6.1	be able to:	43
3.6.2	remember:	43

4	Extracting Important Relationships in High Dimensions	49
4.1	Some Plots of High Dimensional Data	49
4.1.1	Understanding Blobs with Scatterplot Matrices - CLEANUP	49
4.1.2	Parallel Plots	49
4.1.3	Scatterplot Matrices	50
4.2	Summaries of High Dimensional Data	59
4.2.1	The Mean	59
4.2.2	Using Covariance to encode Variance and Correlation	59
4.3	Blob Analysis of High-Dimensional Data	63
4.3.1	Transforming High Dimensional Data	63
4.3.2	Transforming Blobs	64
4.3.3	Whitening Data	67
4.4	Principal Components Analysis	69
4.4.1	The Blob Coordinate System and Smoothing	70
4.4.2	The Low-Dimensional Representation of a Blob	72
4.4.3	Smoothing Data with a Low-Dimensional Representation	74
4.4.4	The Error of the Low-Dimensional Representation	76
4.4.5	Example: Representing Spectral Reflectances	78
4.4.6	Example: Representing Faces with Principal Components	79
4.5	High Dimensions, SVD and NIPALS	81
4.5.1	Principal Components by SVD	81
4.5.2	Just a few Principal Components with NIPALS	82
4.5.3	Projection and Discriminative Problems	84
4.5.4	Just a few Discriminative Directions with PLS1	85
4.6	Multi-Dimensional Scaling	86
4.6.1	Principal Coordinate Analysis	86
4.6.2	Example: Mapping with Multidimensional Scaling	88
4.7	Example: Understanding Height and Weight	90
4.8	What you should remember - NEED SOMETHING	93
5	Clustering: Models of High Dimensional Data	97
5.1	Agglomerative and Divisive Clustering	97
5.1.1	Clustering and Distance	99
5.2	The K-Means Algorithm and Variants	104
5.2.1	How to choose K	106
5.2.2	Soft Assignment	108
5.2.3	General Comments on K-Means	109
5.3	Describing Repetition with Vector Quantization	111
5.3.1	Vector Quantization	112
5.3.2	Example: Groceries in Portugal	113
5.3.3	Efficient Clustering and Hierarchical K Means	115
5.3.4	Example: Activity from Accelerometer Data	116
5.4	You should	120
5.4.1	remember:	120
6	Clustering using Probability Models	122
6.1	The Multivariate Normal Distribution	122

6.1.1	Affine Transformations and Gaussians	123
6.1.2	Plotting a 2D Gaussian: Covariance Ellipses	123
6.2	Mixture Models and Clustering	124
6.2.1	A Finite Mixture of Blobs	125
6.2.2	Topics and Topic Models	126
6.3	The EM Algorithm	129
6.3.1	Example: Mixture of Normals: The E-step	130
6.3.2	Example: Mixture of Normals: The M-step	132
6.3.3	Example: Topic Model: The E-Step	133
6.3.4	Example: Topic Model: The M-step	134
6.3.5	EM in Practice	134
6.4	You should	136
6.4.1	remember:	136
7	Regression	138
7.1	Overview	138
7.1.1	Regression to Spot Trends	140
7.2	Linear Regression and Least Squares	142
7.2.1	Linear Regression	142
7.2.2	Choosing β	143
7.2.3	Residuals	145
7.2.4	R-squared	146
7.2.5	Transforming Variables	148
7.2.6	Can you Trust Your Regression?	151
7.3	Problem Data Points	152
7.3.1	Problem Data Points have Significant Impact	153
7.3.2	The Hat Matrix and Leverage	155
7.3.3	Cook's Distance	156
7.3.4	Standardized Residuals	157
7.4	Many Explanatory Variables	159
7.4.1	Functions of One Explanatory Variable	160
7.4.2	Regularizing Linear Regressions	161
7.4.3	Example: Weight against Body Measurements	165
7.5	You should	169
7.5.1	remember:	169
8	Regression: Choosing and Managing Models	173
8.1	Model Selection: Which Model is Best?	173
8.1.1	Bias and Variance	173
8.1.2	Choosing a Model using Penalties: AIC and BIC	175
8.1.3	Choosing a Model using Cross-Validation	177
8.1.4	A Search Process: Forward and Backward Stagewise Regression	177
8.1.5	Significance: What Variables are Important?	178
8.2	Robust Regression	179
8.2.1	M-Estimators and Iteratively Reweighted Least Squares	180
8.2.2	Scale for M-Estimators	182
8.3	Generalized Linear Models	183

8.3.1	Logistic Regression	183
8.3.2	Multiclass Logistic Regression	185
8.3.3	Regressing Count Data	185
8.3.4	Deviance	186
8.4	L1 Regularization and Sparse Models	187
8.4.1	Dropping Variables with L1 Regularization	187
8.4.2	Wide Datasets	194
8.4.3	Using Sparsity Penalties with Other Models	196
8.5	You should	198
8.5.1	remember:	198
9	Bumps and Non-Parametric Regression	199
9.1	Linear Models using Bumps	199
9.1.1	Scattered Data: Smoothing and Interpolation	199
9.1.2	Model Selection	204
9.2	Non Parametric Bump Models	205
9.2.1	Kernel Smoothing	207
9.2.2	Density Estimation	210
9.3	Exploiting Your Neighbors for Regression	213
9.3.1	Local Polynomial Regression	215
9.3.2	Using your Neighbors to Predict More than a Number	216
10	Neural Networks	219
10.1	Units and Classification	219
10.1.1	Building a Classifier out of Units: The Cost Function	219
10.1.2	Building a Classifier out of Units: Strategy	220
10.1.3	Building a Classifier out of Units: Training	221
10.2	Layers and Networks	224
10.2.1	Notation	225
10.2.2	Training, Gradients and Backpropagation	226
10.2.3	Training Multiple Layers	230
10.2.4	Gradient Scaling Tricks	231
10.2.5	Dropout	234
10.2.6	It's Still Difficult..	234
10.3	Convolutional Neural Networks	236
10.3.1	Images and Convolutional Layers	237
10.3.2	Convolutional Layers upon Convolutional Layers	239
10.3.3	Pooling	239
10.4	Building an Image Classifier	239
10.4.1	An Image Classification Architecture	240
10.4.2	Useful Tricks - 1: Preprocessing Data	241
10.4.3	Useful Tricks - 2: Enhancing Training Data	243
10.4.4	Useful Tricks - 3: Batch Normalization	243
10.4.5	Useful Tricks - 4: Residual Networks	245
11	Structured Models, Learning and Inference	247
11.1	HMM's, Dynamic Programming, and Tree-Structured Models	247

11.1.1	Hidden Markov Models	247
11.1.2	Inference for an HMM	249
11.1.3	Fitting an HMM with EM	254
11.1.4	Tree-Structured Energy Models	256
11.2	Examples	259
11.2.1	Transcribing Speech	259
11.2.2	Tracking a Blob	259
11.2.3	Reading Text	259
11.2.4	Human Parsing	259
11.3	Simple Conditional Random Field Models	259
11.3.1	Drawing a Model	260
11.3.2	MEMM's and Label Bias	260
11.3.3	Conditional Random Field Models	262
11.4	Discriminative Learning of CRFs	263
11.4.1	Representing the Model	263
11.4.2	Setting Up the Learning Problem	265
12	Background: First Tools for Looking at Data	266
12.1	Datasets	266
12.2	What's Happening? - Plotting Data	268
12.2.1	Bar Charts	269
12.2.2	Histograms	269
12.2.3	How to Make Histograms	270
12.2.4	Conditional Histograms	272
12.3	Summarizing 1D Data	272
12.3.1	The Mean	272
12.3.2	Standard Deviation and Variance	275
12.3.3	Variance	279
12.3.4	The Median	280
12.3.5	Interquartile Range	282
12.3.6	Using Summaries Sensibly	283
12.4	Plots and Summaries	284
12.4.1	Some Properties of Histograms	284
12.4.2	Standard Coordinates and Normal Data	287
12.4.3	Boxplots	290
12.5	Whose is bigger? Investigating Australian Pizzas	291
12.6	You should	296
12.6.1	be able to:	296
12.6.2	remember:	296
13	Background: Looking at Relationships	297
13.1	Plotting 2D Data	297
13.1.1	Categorical Data, Counts, and Charts	297
13.1.2	Series	301
13.1.3	Scatter Plots for Spatial Data	303
13.1.4	Exposing Relationships with Scatter Plots	306
13.2	Correlation	309

13.2.1	The Correlation Coefficient	311
13.2.2	Using Correlation to Predict	316
13.2.3	Confusion caused by correlation	320
13.3	Sterile Males in Wild Horse Herds	321
13.4	You should	324
13.4.1	be able to:	324
13.4.2	remember:	324
14	Background: Useful Probability Distributions	325
14.1	Discrete Distributions	326
14.1.1	The Discrete Uniform Distribution	326
14.1.2	Bernoulli Random Variables	326
14.1.3	The Geometric Distribution	327
14.1.4	The Binomial Probability Distribution	328
14.1.5	Multinomial probabilities	330
14.1.6	The Poisson Distribution	331
14.2	Continuous Distributions	333
14.2.1	The Continuous Uniform Distribution	333
14.2.2	The Beta Distribution	333
14.2.3	The Gamma Distribution	335
14.2.4	The Exponential Distribution	335
14.3	The Normal Distribution	336
14.3.1	The Standard Normal Distribution	336
14.3.2	The Normal Distribution	337
14.3.3	Properties of the Normal Distribution	338
14.4	Approximating Binomials with Large N	340
14.4.1	Large N	341
14.4.2	Getting Normal	343
14.4.3	So What?	345
14.5	You should	346
14.5.1	remember:	346
15	Background: Inference: Making Point Estimates	347
15.1	Estimating Model Parameters with Maximum Likelihood	348
15.1.1	The Maximum Likelihood Principle	349
15.1.2	Cautions about Maximum Likelihood	357
15.2	Incorporating Priors with Bayesian Inference	358
15.2.1	Constructing the Posterior	359
15.2.2	Normal Prior and Normal Likelihood	363
15.2.3	MAP Inference	366
15.2.4	Filtering	367
15.2.5	Cautions about Bayesian Inference	370
15.3	Samples, Urns and Populations	370
15.3.1	Estimating the Population Mean from a Sample	371
15.3.2	The Variance of the Sample Mean	372
15.3.3	The Probability Distribution of the Sample Mean	376
15.3.4	When The Urn Model Works	376

15.4	You should	378
15.4.1	be able to:	378
15.4.2	remember:	378
16	Markov Chains and Simulation	379
16.1	Markov Chains	379
16.1.1	Motivating Example: Multiple Coin Flips	379
16.1.2	Motivating Example: The Gambler's Ruin	381
16.1.3	Motivating Example: A Virus	383
16.1.4	Markov Chains	383
16.1.5	Example: Particle Motion as a Markov Chain	387
16.2	Simulation	388
16.2.1	Obtaining Random Numbers	388
16.2.2	Computing Expectations with Simulations	389
16.2.3	Computing Probabilities with Simulations	390
16.2.4	Simulation Results as Random Variables	391
16.3	Simulation Examples	393
16.3.1	Simulating Experiments	393
16.3.2	Simulating Markov Chains	394
16.3.3	Example: Ranking the Web by Simulating a Markov Chain	396
16.3.4	Example: Simulating a Complicated Game	399
16.4	You should	403
16.4.1	remember	403
17	Math Resources	409
17.1	Useful Material about Matrices	409
17.1.1	The Singular Value Decomposition	410
17.1.2	Approximating A Symmetric Matrix	411

CHAPTER 1

Notation and conventions

A dataset as a collection of d -tuples (a d -tuple is an ordered list of d elements). Tuples differ from vectors, because we can always add and subtract vectors, but we cannot necessarily add or subtract tuples. There are always N items in any dataset. There are always d elements in each tuple in a dataset. The number of elements will be the same for every tuple in any given tuple. Sometimes we may not know the value of some elements in some tuples.

We use the same notation for a tuple and for a vector. Most of our data will be vectors. We write a vector in bold, so \mathbf{x} could represent a vector or a tuple (the context will make it obvious which is intended).

The entire data set is $\{\mathbf{x}\}$. When we need to refer to the i 'th data item, we write \mathbf{x}_i . Assume we have N data items, and we wish to make a new dataset out of them; we write the dataset made out of these items as $\{\mathbf{x}_i\}$ (the i is to suggest you are taking a set of items and making a dataset out of them). If we need to refer to the j 'th component of a vector \mathbf{x}_i , we will write $x_i^{(j)}$ (notice this isn't in bold, because it is a component not a vector, and the j is in parentheses because it isn't a power). Vectors are always column vectors.

When I write $\{kx\}$, I mean the dataset created by taking each element of the dataset $\{x\}$ and multiplying by k ; and when I write $\{x + c\}$, I mean the dataset created by taking each element of the dataset $\{x\}$ and adding c .

Terms:

- $\text{mean}(\{x\})$ is the mean of the dataset $\{x\}$ (definition 12.1, page 272).
- $\text{std}(\{x\})$ is the standard deviation of the dataset $\{x\}$ (definition 12.2, page 275).
- $\text{var}(\{x\})$ is the standard deviation of the dataset $\{x\}$ (definition 12.3, page 279).
- $\text{median}(\{x\})$ is the standard deviation of the dataset $\{x\}$ (definition 12.4, page 280).
- $\text{percentile}(\{x\}, k)$ is the $k\%$ percentile of the dataset $\{x\}$ (definition 12.5, page 282).
- $\text{iqr}\{x\}$ is the interquartile range of the dataset $\{x\}$ (definition 12.7, page 283).
- $\{\hat{x}\}$ is the dataset $\{x\}$, transformed to standard coordinates (definition 12.8, page 287).
- Standard normal data is defined in definition 12.9, (page 288).
- Normal data is defined in definition 12.10, (page 288).
- $\text{corr}(\{(x, y)\})$ is the correlation between two components x and y of a dataset (definition 13.1, page 311).

- \emptyset is the empty set.
- Ω is the set of all possible outcomes of an experiment.
- Sets are written as \mathcal{A} .
- \mathcal{A}^c is the complement of the set \mathcal{A} (i.e. $\Omega - \mathcal{A}$).
- \mathcal{E} is an event (page 402).
- $P(\{\mathcal{E}\})$ is the probability of event \mathcal{E} (page 402).
- $P(\{\mathcal{E}\}|\{\mathcal{F}\})$ is the probability of event \mathcal{E} , conditioned on event \mathcal{F} (page 402).
- $p(x)$ is the probability that random variable X will take the value x ; also written $P(\{X = x\})$ (page 402).
- $p(x, y)$ is the probability that random variable X will take the value x and random variable Y will take the value y ; also written $P(\{X = x\} \cap \{Y = y\})$ (page 402).
- $\operatorname{argmax}_x f(x)$ means the value of x that maximises $f(x)$.
- $\operatorname{argmin}_x f(x)$ means the value of x that minimises $f(x)$.
- $\max_i(f(x_i))$ means the largest value that f takes on the different elements of the dataset $\{x_i\}$.
- $\hat{\theta}$ is an estimated value of a parameter θ .

1.0.1 Background Information

Cards: A standard deck of playing cards contains 52 cards. These cards are divided into four suits. The suits are: spades and clubs (which are black); and hearts and diamonds (which are red). Each suit contains 13 cards: Ace, 2, 3, 4, 5, 6, 7, 8, 9, 10, Jack (sometimes called Knave), Queen and King. It is common to call Jack, Queen and King *court cards*.

Dice: If you look hard enough, you can obtain dice with many different numbers of sides (though I've never seen a three sided die). We adopt the convention that the sides of an N sided die are labeled with the numbers $1 \dots N$, and that no number is used twice. Most dice are like this.

Fairness: Each face of a fair coin or die has the same probability of landing upmost in a flip or roll.

Roulette: A roulette wheel has a collection of slots. There are 36 slots numbered with the digits $1 \dots 36$, and then one, two or even three slots numbered with zero. There are no other slots. A ball is thrown at the wheel when it is spinning, and it bounces around and eventually falls into a slot. If the wheel is properly balanced, the ball has the same probability of falling into each slot. The number of the slot the ball falls into is said to "come up". There are a variety of bets available.

1.1 SOME USEFUL MATHEMATICAL FACTS

The gamma function $\Gamma(x)$ is defined by a series of steps. First, we have that for n an integer,

$$\Gamma(n) = (n - 1)!$$

and then for z a complex number with positive real part (which includes positive real numbers), we have

$$\Gamma(z) = \int_0^\infty t^z \frac{e^{-t}}{t} dt.$$

By doing this, we get a function on positive real numbers that is a smooth interpolate of the factorial function. We won't do any real work with this function, so won't expand on this definition. In practice, we'll either look up a value in tables or require a software environment to produce it.

1.2 ACKNOWLEDGEMENTS

Typos spotted by: Han Chen (numerous!), Henry Lin (numerous!), Eric Huber, Brian Lunt, Yusuf Sobh, Scott Walters, — Your Name Here — TA's for this course have helped improve the notes. Thanks to Zicheng Liao, Michael Sittig, Nikita Spirin, Saurabh Singh, Daphne Tsatsoulis, Henry Lin, Karthik Ramaswamy.

CHAPTER 2

Some Preliminaries

2.1 NOTATION AND CONVENTIONS

A dataset as a collection of d -tuples (a d -tuple is an ordered list of d elements). Tuples differ from vectors, because we can always add and subtract vectors, but we cannot necessarily add or subtract tuples. There are always N items in any dataset. There are always d elements in each tuple in a dataset. The number of elements will be the same for every tuple in any given tuple. Sometimes we may not know the value of some elements in some tuples.

We use the same notation for a tuple and for a vector. Most of our data will be vectors. We write a vector in bold, so \mathbf{x} could represent a vector or a tuple (the context will make it obvious which is intended).

The entire data set is $\{\mathbf{x}\}$. When we need to refer to the i 'th data item, we write \mathbf{x}_i . Assume we have N data items, and we wish to make a new dataset out of them; we write the dataset made out of these items as $\{\mathbf{x}_i\}$ (the i is to suggest you are taking a set of items and making a dataset out of them). If we need to refer to the j 'th component of a vector \mathbf{x}_i , we will write $x_i^{(j)}$ (notice this isn't in bold, because it is a component not a vector, and the j is in parentheses because it isn't a power). Vectors are always column vectors.

When I write $\{kx\}$, I mean the dataset created by taking each element of the dataset $\{x\}$ and multiplying by k ; and when I write $\{x + c\}$, I mean the dataset created by taking each element of the dataset $\{x\}$ and adding c .

Terms:

- $\text{mean}(\{x\})$ is the mean of the dataset $\{x\}$ (definition 12.1, page 272).
- $\text{std}(\{x\})$ is the standard deviation of the dataset $\{x\}$ (definition 12.2, page 275).
- $\text{var}(\{x\})$ is the variance of the dataset $\{x\}$ (definition 12.3, page 279).
- $\text{median}(\{x\})$ is the standard deviation of the dataset $\{x\}$ (definition 12.4, page 280).
- $\text{percentile}(\{x\}, k)$ is the $k\%$ percentile of the dataset $\{x\}$ (definition 12.5, page 282).
- $\text{iqr}\{x\}$ is the interquartile range of the dataset $\{x\}$ (definition 12.7, page 283).
- $\{\hat{x}\}$ is the dataset $\{x\}$, transformed to standard coordinates (definition 12.8, page 287).
- Standard normal data is defined in definition 12.9, (page 288).
- Normal data is defined in definition 12.10, (page 288).

- $\text{corr}(\{(x, y)\})$ is the correlation between two components x and y of a dataset (definition 13.1, page 311).
- \emptyset is the empty set.
- Ω is the set of all possible outcomes of an experiment.
- Sets are written as \mathcal{A} .
- \mathcal{A}^c is the complement of the set \mathcal{A} (i.e. $\Omega - \mathcal{A}$).
- \mathcal{E} is an event (page 402).
- $P(\{\mathcal{E}\})$ is the probability of event \mathcal{E} (page 402).
- $P(\{\mathcal{E}\}|\{\mathcal{F}\})$ is the probability of event \mathcal{E} , conditioned on event \mathcal{F} (page 402).
- $p(x)$ is the probability that random variable X will take the value x ; also written $P(\{X = x\})$ (page 402).
- $p(x, y)$ is the probability that random variable X will take the value x and random variable Y will take the value y ; also written $P(\{X = x\} \cap \{Y = y\})$ (page 402).
- $\operatorname{argmax}_x f(x)$ means the value of x that maximises $f(x)$.
- $\operatorname{argmin}_x f(x)$ means the value of x that minimises $f(x)$.
- $\max_i(f(x_i))$ means the largest value that f takes on the different elements of the dataset $\{x_i\}$.
- $\hat{\theta}$ is an estimated value of a parameter θ .

2.1.1 Background Information

Cards: A standard deck of playing cards contains 52 cards. These cards are divided into four suits. The suits are: spades and clubs (which are black); and hearts and diamonds (which are red). Each suit contains 13 cards: Ace, 2, 3, 4, 5, 6, 7, 8, 9, 10, Jack (sometimes called Knave), Queen and King. It is common to call Jack, Queen and King *court cards*.

Dice: If you look hard enough, you can obtain dice with many different numbers of sides (though I've never seen a three sided die). We adopt the convention that the sides of an N sided die are labeled with the numbers $1 \dots N$, and that no number is used twice. Most dice are like this.

Fairness: Each face of a fair coin or die has the same probability of landing upmost in a flip or roll.

Roulette: A roulette wheel has a collection of slots. There are 36 slots numbered with the digits $1 \dots 36$, and then one, two or even three slots numbered with zero. There are no other slots. A ball is thrown at the wheel when it is spinning, and it bounces around and eventually falls into a slot. If the wheel is properly balanced, the ball has the same probability of falling into each slot. The number of the slot the ball falls into is said to “come up”. There are a variety of bets available.

2.2 SOME USEFUL MATHEMATICAL FACTS

The gamma function $\Gamma(x)$ is defined by a series of steps. First, we have that for n an integer,

$$\Gamma(n) = (n - 1)!$$

and then for z a complex number with positive real part (which includes positive real numbers), we have

$$\Gamma(z) = \int_0^\infty t^z \frac{e^{-t}}{t} dt.$$

By doing this, we get a function on positive real numbers that is a smooth interpolate of the factorial function. We won't do any real work with this function, so won't expand on this definition. In practice, we'll either look up a value in tables or require a software environment to produce it.

2.3 ACKNOWLEDGEMENTS

Typos spotted by: Han Chen (numerous!), Henry Lin (numerous!), Paris Smaragdis (numerous!), Johnny Chang, Eric Huber, Brian Lunt, Yusuf Sobh, Scott Walters, — Your Name Here — TA's for this course have helped improve the notes. Thanks to Zicheng Liao, Michael Sittig, Nikita Spirin, Saurabh Singh, Daphne Tsatsoulis, Henry Lin, Karthik Ramaswamy.

2.4 THE CURSE OF DIMENSION

High dimensional models display unintuitive behavior (or, rather, it can take years to make your intuition see the true behavior of high-dimensional models as natural). In these models, most data lies in places you don't expect. We will do several simple calculations with an easy high-dimensional distribution to build some intuition.

2.4.1 The Curse: Data isn't Where You Think it is

Assume our data lies within a cube, with edge length two, centered on the origin. This means that each component of \mathbf{x}_i lies in the range $[-1, 1]$. One simple model for such data is to assume that each dimension has uniform probability density in this range. In turn, this means that $P(x) = \frac{1}{2a}$. The mean of this model is at the origin, which we write as $\mathbf{0}$.

The first surprising fact about high dimensional data is that most of the data can lie quite far away from the mean. For example, we can divide our dataset into two pieces. $\mathcal{A}(\epsilon)$ consists of all data items where *every* component of the data has a value in the range $[-(1 - \epsilon), (1 - \epsilon)]$. $\mathcal{B}(\epsilon)$ consists of all the rest of the data. If you think of the data set as forming a cubical orange, then $\mathcal{B}(\epsilon)$ is the rind (which has thickness ϵ) and $\mathcal{A}(\epsilon)$ is the fruit.

Your intuition will tell you that there is more fruit than rind. This is true, for three dimensional oranges, but not true in high dimensions. The fact that the orange is cubical just simplifies the calculations, but has nothing to do with the real problem.

We can compute $P(\{\mathbf{x} \in \mathcal{A}(\epsilon)\})$ and $P(\{\mathbf{x} \in \mathcal{B}(\epsilon)\})$. These probabilities tell us the probability a data item lies in the fruit (resp. rind). $P(\{\mathbf{x} \in \mathcal{A}(\epsilon)\})$ is easy

to compute as

$$P(\{\mathbf{x} \in \mathcal{A}(\epsilon)\}) = (2(1 - \epsilon))^d \left(\frac{1}{2^d}\right) = (1 - \epsilon)^d$$

and

$$P(\{\mathbf{x} \in \mathcal{B}(\epsilon)\}) = 1 - P(\{\mathbf{x} \in \mathcal{A}(\epsilon)\}) = 1 - (1 - \epsilon)^d.$$

But notice that, as $d \rightarrow \infty$,

$$P(\{\mathbf{x} \in \mathcal{A}(\epsilon)\}) \rightarrow 0.$$

This means that, for large d , we expect most of the data to be in $\mathcal{B}(\epsilon)$. Equivalently, for large d , we expect that at least one component of each data item is close to either 1 or -1 .

This suggests (correctly) that much data is quite far from the origin. It is easy to compute the average of the squared distance of data from the origin. We want

$$\mathbb{E}[\mathbf{x}^T \mathbf{x}] = \int_{\text{box}} \left(\sum_i x_i^2 \right) P(\mathbf{x}) d\mathbf{x}$$

but we can rearrange, so that

$$\mathbb{E}[\mathbf{x}^T \mathbf{x}] = \sum_i \mathbb{E}[x_i^2] = \sum_i \int_{\text{box}} x_i^2 P(\mathbf{x}) d\mathbf{x}.$$

Now each component of \mathbf{x} is independent, so that $P(\mathbf{x}) = P(x_1)P(x_2)\dots P(x_d)$. Now we substitute, to get

$$\mathbb{E}[\mathbf{x}^T \mathbf{x}] = \sum_i \mathbb{E}[x_i^2] = \sum_i \int_{-1}^1 x_i^2 P(x_i) dx_i = \sum_i \frac{1}{2} \int_{-1}^1 x_i^2 dx_i = \frac{d}{3},$$

so as d gets bigger, most data points will be further and further from the origin. Worse, as d gets bigger, data points tend to get further and further from one another. We can see this by computing the average of the squared distance of data points from one another. Write \mathbf{u} for one data point and \mathbf{v} ; we can compute

$$\mathbb{E}[d(\mathbf{u}, \mathbf{v})^2] = \int_{\text{box}} \int_{\text{box}} \sum_i (u_i - v_i)^2 d\mathbf{u} d\mathbf{v} = \mathbb{E}[\mathbf{u}^T \mathbf{u}] + \mathbb{E}[\mathbf{v}^T \mathbf{v}] - \mathbb{E}[\mathbf{u}^T \mathbf{v}]$$

but since \mathbf{u} and \mathbf{v} are independent, we have $\mathbb{E}[\mathbf{u}^T \mathbf{v}] = \mathbb{E}[\mathbf{u}]^T \mathbb{E}[\mathbf{v}] = 0$. This yields

$$\mathbb{E}[d(\mathbf{u}, \mathbf{v})^2] = 2\frac{d}{3}.$$

This means that, for large d , we expect our data points to be quite far apart.

2.4.2 Minor Banes of Dimension

High dimensional data presents a variety of important practical nuisances which follow from the curse of dimension. It is hard to estimate covariance matrices, and it is hard to build histograms.

Covariance matrices are hard to work with for two reasons. The number of entries in the matrix grows as the square of the dimension, so the matrix can get big and so difficult to store. More important, the amount of data we need to get an accurate estimate of all the entries in the matrix grows fast. As we are estimating more numbers, we need more data to be confident that our estimates are reasonable. There are a variety of straightforward work-arounds for this effect. In some cases, we have so much data there is no need to worry. In other cases, we assume that the covariance matrix has a particular form, and just estimate those parameters. There are two strategies that are usual. In one, we assume that the covariance matrix is diagonal, and estimate only the diagonal entries. In the other, we assume that the covariance matrix is a scaled version of the identity, and just estimate this scale. You should see these strategies as acts of desperation, to be used only when computing the full covariance matrix seems to produce more problems than using these approaches.

It is difficult to build histogram representations for high dimensional data. The strategy of dividing the domain into boxes, then counting data into them, fails miserably because there are too many boxes. In the case of our cube, imagine we wish to divide each dimension in half (i.e. between $[-1, 0]$ and between $[0, 1]$). Then we must have 2^d boxes. This presents two problems. First, we will have difficulty representing this number of boxes. Second, unless we are exceptionally lucky, most boxes must be empty because we will not have 2^d data items.

However, one representation is extremely effective. We can represent data as a collection of **clusters** — coherent blobs of similar datapoints that could, under appropriate circumstances, be regarded as the same. We could then represent the dataset by, for example, the center of each cluster and the number of data items in each cluster. Since each cluster is a blob, we could also report the covariance of each cluster, if we can compute it.

CHAPTER 3

Learning to Classify

A **classifier** is a procedure that accepts a set of features and produces a class label for them. There could be two, or many, classes. Classifiers are immensely useful, and find wide application, because many problems are naturally classification problems. For example, if you wish to determine whether to place an advert on a web-page or not, you would use a classifier (i.e. look at the page, and say yes or no according to some rule). As another example, if you have a program that you found for free on the web, you would use a classifier to decide whether it was safe to run it (i.e. look at the program, and say yes or no according to some rule). As yet another example, credit card companies must decide whether a transaction is good or fraudulent.

All these examples are two class classifiers, but in many cases it is natural to have more classes. You can think of sorting laundry as applying a multi-class classifier. You can think of doctors as complex multi-class classifiers. In this (crude) model, the doctor accepts a set of features, which might be your complaints, answers to questions, and so on, and then produces a response which we can describe as a class. The grading procedure for any class is a multi-class classifier: it accepts a set of features — performance in tests, homeworks, and so on — and produces a class label (the letter grade).

Classifiers are built by taking a set of labeled examples and using them to come up with a procedure that assigns a label to any new example. In the general problem, we have a training dataset (\mathbf{x}_i, y_i) ; each of the **feature vectors** \mathbf{x}_i consists of measurements of the properties of different types of object, and the y_i are labels giving the type of the object that generated the example. We will then use the training dataset to find a procedure that will predict an accurate label (y) for any new object (\mathbf{x}).

Definition: 3.1 *Classifier*

A classifier is a procedure that accepts a feature vector and produces a label.

3.1 CLASSIFICATION, ERROR, AND LOSS

You should think of a classifier as a procedure — we pass in a feature vector, and get a class label in return. We want to use the training data to find the procedure that is “best” when used on the test data. This problem has two tricky features.

First, we need to be clear on what a good procedure is. Second, we really want the procedure to be good on test data, which we haven't seen and won't see; we only get to see the training data. These two considerations shape much of what we do.

3.1.1 Loss and the Cost of Misclassification

The choice of procedure must depend on the cost of making a mistake. This cost can be represented with a **loss function**, which specifies the cost of making each type of mistake. I will write $L(j \rightarrow k)$ for the loss incurred when classifying an example of class j as having class k .

A two-class classifier can make two kinds of mistake. Because two-class classifiers are so common, there is a special name for each kind of mistake. A **false positive** occurs when a negative example is classified positive (which we can write $L(- \rightarrow +)$ and avoid having to remember which index refers to which class); a **false negative** occurs when a positive example is classified negative (similarly $L(+ \rightarrow -)$). By convention, the loss of getting the right answer is zero, and the loss for any wrong answer is non-negative.

The choice of procedure should depend quite strongly on the cost of each mistake. For example, pretend there is only one disease; then doctors would be classifiers, deciding whether a patient had it or not. If this disease is dangerous, but is safely and easily treated, false negatives are expensive errors, but false positives are cheap. In this case, procedures that tend to make more false positives than false negatives are better. Similarly, if the disease is not dangerous, but the treatment is difficult and unpleasant, then false positives are expensive errors and false negatives are cheap, and so we prefer false negatives to false positives.

You might argue that the best choice of classifier makes no mistake. But for most practical cases, the best choice of classifier is guaranteed to make mistakes. As an example, consider an alien who tries to classify humans into male and female, using only height as a feature. However the alien's classifier uses that feature, it will make mistakes. This is because the classifier must choose, for each value of height, whether to label the humans with that height male or female. But for the vast majority of heights, there are some males and some females with that height, and so the alien's classifier must make some mistakes whatever gender it chooses for that height.

For many practical problems, it is difficult to know what loss function to use. There is seldom an obvious choice. One common choice is to assume that all errors are equally bad. This gives the **0-1 loss** — every error has loss 1, and all right answers have loss zero.

3.1.2 Building a Classifier from Probabilities

Assume that we have a reliable model of $p(y|\mathbf{x})$. This case occurs less often than you might think for practical data, because building such a model is often very difficult. However, when we do have a model and a loss function, it is easy to determine the best classifier. We should choose the rule that gives minimum expected loss.

We start with a two-class classifier. At \mathbf{x} , the expected loss of saying $-$ is $L(+ \rightarrow -)p(+|\mathbf{x})$ (remember, $L(- \rightarrow -) = 0$); similarly, the expected loss of saying $+$ is $L(- \rightarrow +)p(-|\mathbf{x})$. At most points, one of $L(- \rightarrow +)p(-|\mathbf{x})$ and

$L(+ \rightarrow -)p(+|\mathbf{x})$ is larger than the other, and so the choice is clear. The remaining set of points (where $L(- \rightarrow +)p(-|\mathbf{x}) = L(+ \rightarrow -)p(+|\mathbf{x})$) is “small” (formally, it has zero measure) for most models and problems, and so it doesn’t matter what we choose at these points. This means that the rule

$$\text{say } \begin{cases} + & \text{if } L(+ \rightarrow -)p(+|\mathbf{x}) > L(- \rightarrow +)p(-|\mathbf{x}) \\ - & \text{if } L(+ \rightarrow -)p(+|\mathbf{x}) < L(- \rightarrow +)p(-|\mathbf{x}) \\ \text{random choice} & \text{otherwise} \end{cases}$$

is the best available. Because it doesn’t matter what we do when $L(+ \rightarrow -)p(+|\mathbf{x}) = L(- \rightarrow +)p(-|\mathbf{x})$, it is fine to use

$$\text{say } \begin{cases} + & \text{if } L(+ \rightarrow -)p(+|\mathbf{x}) > L(- \rightarrow +)p(-|\mathbf{x}) \\ - & \text{otherwise} \end{cases}$$

The same reasoning applies in the multi-class case. We choose the class where the expected loss from that choice is smallest. In the case of 0-1 loss, this boils down to:

$$\text{choose } k \text{ such that } p(k|\mathbf{x}) \text{ is largest.}$$

3.1.3 Building a Classifier using Decision Boundaries

Building a classifier out of posterior probabilities is less common than you might think, for two reasons. First, it’s often very difficult to get a good posterior probability model. Second, most of the model doesn’t matter to the choice of classifier. What is important is knowing which class has the lowest expected loss, not the exact values of the expected losses, so we should be able to get away without an exact posterior model.

Look at the rules in section 3.1.2. Each of them carves up the domain of \mathbf{x} into pieces, and then attaches a class – the one with the lowest expected loss – to each piece. There isn’t necessarily one piece per class, (though there’s always one class per piece). The important factor here is the boundaries between the pieces, which are known as **decision boundaries**. A powerful strategy for building classifiers is to choose some way of building decision boundaries, then adjust it to perform well on the data one has. This involves modelling considerably less detail than modelling the whole posterior.

For example, in the two-class case, we will spend some time discussing the decision boundary given by

$$\text{choose } \begin{cases} - & \text{if } \mathbf{x}^T \mathbf{a} + b < 0 \\ + & \text{otherwise} \end{cases}$$

often written as $\text{sign} \mathbf{x}^T \mathbf{a} + b$ (section 16.20). In this case we choose \mathbf{a} and b to obtain low loss.

3.1.4 What will happen on Test Data?

What we really want from a classifier is to have small loss on test data. But this is difficult to measure or achieve directly. For example, think about the case of classifying credit-card transactions as “good” or “bad”. We could certainly obtain

a set of examples that have been labelled for training, because the card owner often complains some time after a fraudulent use of their card. But what is important here is to see a new transaction and label it without holding it up for a few months to see what the card owner says. The classifier may never know if the label is right or not.

Generally, we will assume that the training data is “like” the test data, and so we will try to make the classifier perform well on the training data. Classifiers that have small training error might not have small test error. One example of this problem is the (silly) classifier that takes any data point and, if it is the same as a point in the training set, emits the class of that point and otherwise chooses randomly between the classes. This classifier has been learned from data, and has a zero error rate on the training dataset; it is likely to be unhelpful on any other dataset, however.

Test error is usually worse than training error, because of an effect that is sometimes called **overfitting**, so called because the classification procedure fits the training data better than it fits the test data. Other names include **selection bias**, because the training data has been selected and so isn’t exactly like the test data, and **generalizing badly**, because the classifier fails to generalize. The effect occurs because the classifier has been trained to perform well *on the training dataset*, and the training dataset is not the same as the test dataset. First, it is quite likely smaller. Second, it might be biased through a variety of accidents. This means that small training error may have to do with quirks of the training dataset that don’t occur in other sets of examples. One consequence of overfitting is that classifiers should always be evaluated on data that was not used in training.

Remember this: *Classifiers should always be evaluated on data that was not used in training.*

Now assume that we are using the 0-1 loss, so that the loss of using a classifier is the same as the **error rate**, that is, the percentage of classification attempts on a test set that result in the wrong answer. We could also use the **accuracy**, which is the percentage of classification attempts that result in the right answer. We cannot estimate the error rate of the classifier using training data, because the classifier has been trained to do well on that data, which will mean our error rate estimate will be too low. An alternative is to separate out some training data to form a **validation set** (confusingly, this is often called a **test set**), then train the classifier on the rest of the data, and evaluate on the validation set. This has the difficulty that the classifier will not be the best estimate possible, because we have left out some training data when we trained it. This issue can become a significant nuisance when we are trying to tell which of a set of classifiers to use—did the classifier perform poorly on validation data because it is not suited to the problem representation or because it was trained on too little data?

We can resolve this problem with **cross-validation**, which involves repeatedly: splitting data into training and validation sets uniformly and at random,

training a classifier on the training set, evaluating it on the validation set, and then averaging the error over all splits. This allows an estimate of the likely future performance of a classifier, at the expense of substantial computation. You should notice that cross-validation, in some sense, looks at the sensitivity of the classifier to a change in the training set. The most usual form of this algorithm involves omitting single items from the dataset and is known as **leave-one-out cross-validation**.

You should usually compare the error rate of a classifier to two important references. The first is the error rate if you assign classes to examples uniformly at random, which for a two class classifier is 50%. A two class classifier should never have an error rate higher than 50%. If you have one that does, all you need to do is swap its class assignment, and the resulting error rate would be lower than 50%. The second is the error rate if you assign all data to the most common class. If one class is uncommon and the other is common, this error rate can be hard to beat. Data where some classes occur very seldom requires careful, and quite specialized, handling.

3.1.5 The Class Confusion Matrix

Evaluating a multi-class classifier is more complex than evaluating a binary classifier. The error rate if you assign classes to examples uniformly at random can be rather high. If each class has about the same frequency, then this error rate is $(1 - 100/\text{number of classes})\%$. A multi-class classifier can make many more kinds of mistake than a binary classifier can. It is useful to know the total error rate of the classifier (percentage of classification attempts that produce the wrong answer) or the accuracy, (the percentage of classification attempts that produce the *right* answer). If the error rate is low enough, or the accuracy is high enough, there's not much to worry about. But if it's not, you can look at the **class confusion matrix** to see what's going on.

	Predict 0	Predict 1	Predict 2	Predict 3	Predict 4	Class error
True 0	151	7	2	3	1	7.9%
True 1	32	5	9	9	0	91%
True 2	10	9	7	9	1	81%
True 3	6	13	9	5	2	86%
True 4	2	3	2	6	0	100%

TABLE 3.1: The class confusion matrix for a multiclass classifier. Further details about the dataset and this example appear in worked example 3.3.

Table 3.1 gives an example. This is a class confusion matrix from a classifier built on a dataset where one tries to predict the degree of heart disease from a collection of physiological and physical measurements. There are five classes (0...4). The i, j 'th cell of the table shows the number of data points of true class i that were classified to have class j . As I find it hard to recall whether rows or columns represent true or predicted classes, I have marked this on the table. For each row, there is a class error rate, which is the percentage of data points of that class that

were misclassified. The first thing to look at in a table like this is the diagonal; if the largest values appear there, then the classifier is working well. This clearly isn't what is happening for table 3.1. Instead, you can see that the method is very good at telling whether a data point is in class 0 or not (the class error rate is rather small), but cannot distinguish between the other classes. This is a strong hint that the data can't be used to draw the distinctions that we want. It might be a lot better to work with a different set of classes.

3.1.6 Statistical Learning Theory and Generalization

What is required in a classifier is an ability to *predict*—we should like to be confident that the classifier chosen on a particular data set has a low risk on future data items. The family of decision boundaries from which a classifier is chosen is an important component of the problem. Some decision boundaries are more flexible than others (in a sense we don't intend to make precise). This has nothing to do with the number of parameters in the decision boundary. For example, if we were to use a point to separate points on the line, there are very small sets of points that are not linearly separable (the smallest set has three points in it). This means that relatively few sets of points on the line are linearly separable, so that if our dataset is sufficiently large and linearly separable, the resulting classifier is likely to behave well in future. However, using the sign of $\sin \lambda x$ to separate points on the line produces a completely different qualitative phenomenon; for *any* labeling of distinct points on the line into two classes, we can choose a value of λ to achieve this labeling. This flexibility means that the classifier is wholly unreliable—it can be made to fit any set of examples, meaning the fact that it fits the examples is uninformative.

There is a body of theory that treats this question, which rests on two important points.

- **A large enough dataset yields a good representation of the source of the data:** this means that if the dataset used to train the classifier is very large, there is a reasonable prospect that the performance on the training set will represent the future performance. However, for this to be helpful, we need the data set to be large with respect to the “flexibility” of the family
- **The “flexibility” of a family of decision boundaries can be formalized:** yielding the **Vapnik-Chervonenkis dimension** (or **V-C dimension**) of the family. This dimension is independent of the number of parameters of the family. Families with finite V-C dimension can yield classifiers whose future performance can be bounded using the number of training elements; families with infinite V-C dimension (like the $\sin \lambda x$ example above) cannot be used to produce reliable classifiers.

The essence of the theory is as follows: if one chooses a decision boundary from an inflexible family, and the resulting classifier performs well on a large data set, there is strong reason to believe that it will perform well on future items drawn from the same source. This statement can be expressed precisely in terms of bounds on total risk to be expected for particular classifiers as a function of the size of the data set used to train the classifier. These bounds hold in probability. These bounds

tend not to be used in practice, because they appear to be extremely pessimistic. Space doesn't admit an exposition of this theory—which is somewhat technical—but interested readers can look it up in (? , ? , ?).

3.2 CLASSIFYING WITH NAIVE BAYES

One reason it is difficult to build a posterior probability model is the dependencies between features. However, if we *assume* that features are conditionally independent conditioned on the class of the data item, we can get a simple expression for the posterior. This assumption is hardly ever true in practice. Remarkably, this doesn't matter very much, and the classifier we build from the assumption often works extremely well. It is the classifier of choice for very high dimensional data.

Recall bayes' rule. If we have $p(\mathbf{x}|y)$ (often called either a **likelihood** or **class conditional probability**), and $p(y)$ (often called a **prior**) then we can form

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

(the **posterior**). We write x_j for the j 'th component of \mathbf{x} . Our assumption is

$$p(\mathbf{x}|y) = \prod_i p(x_i|y)$$

(again, this isn't usually the case; it just turns out to be fruitful to *assume* that it is true). This assumption means that

$$\begin{aligned} p(y|\mathbf{x}) &= \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} \\ &= \frac{\prod_i p(x_i|y)p(y)}{p(\mathbf{x})} \\ &\propto \prod_i p(x_i|y)p(y). \end{aligned}$$

Now because we need only to know the posterior values up to scale at \mathbf{x} to make a decision (check the rules above if you're unsure), we don't need to estimate $p(\mathbf{x})$. In the case of 0-1 loss, this yields the rule

choose y such that $\prod_i p(x_i|y)p(y)$ is largest.

Naive bayes is particularly good when there are a large number of features, but there are some things to be careful about. You can't actually multiply a large number of probabilities and expect to get an answer that a floating point system thinks is different from zero. Instead, you should add the log probabilities. A model with many different features is likely to have many strongly negative log probabilities, so you should not just add up all the log probabilities then exponentiate, or else you will find that each class has a posterior probability of zero. Instead, subtract the largest log from all the others, then exponentiate; you will obtain a vector proportional to the class probabilities, where the largest element has the value 1.

We still need models for $p(x_i|y)$ for each x_i . It turns out that simple parametric models work really well here. For example, one could fit a normal distribution to each x_i in turn, for each possible value of y , using the training data. The logic of the measurements might suggest other distributions, too. If one of the x_i 's was a count, we might fit a Poisson distribution. If it was a 0-1 variable, we might fit a Bernoulli distribution. If it was a numeric variable that took one of several values, then we might use either a multinomial model.

Many effects cause missing values: measuring equipment might fail; a record could be damaged; it might be too hard to get information in some cases; survey respondents might not want to answer a question; and so on. As a result, missing values are quite common in practical datasets. A nice feature of naive bayes classifiers is that they can handle missing values for particular features rather well.

Dealing with missing data during learning is easy. For example, assume for some i , we wish to fit $p(x_i|y)$ with a normal distribution. We need to estimate the mean and standard deviation of that normal distribution (which we do with maximum likelihood, as one should). If not every example has a known value of x_i , this really doesn't matter; we simply omit the unknown number from the estimate. Write $x_{i,j}$ for the value of x_i for the j 'th example. To estimate the mean, we form

$$\frac{\sum_{j \in \text{cases with known values}} x_{i,j}}{\text{number of cases with known values}}$$

and so on.

Dealing with missing data during classification is easy, too. We need to look for the y that produces the largest value of $\sum_i \log p(x_i|y)$. We can't evaluate $p(x_i|y)$ if the value of that feature is missing - but it is missing for each class. We can just leave that term out of the sum, and proceed. This procedure is fine if data is missing as a result of "noise" (meaning that the missing terms are independent of class). If the missing terms depend on the class, there is much more we could do — for example, we might build a model of the class-conditional density of missing terms.

Notice that if some values of a discrete feature x_i don't appear for some class, you could end up with a model of $p(x_i|y)$ that had zeros for some values. This almost inevitably leads to serious trouble, because it means your model states you cannot ever observe that value for a data item of that class. This isn't a safe property: it is hardly ever the case that not observing something means you cannot observe it. A simple, but useful, fix is to add one to all small counts.

The usual way to find a model of $p(y)$ is to count the number of training examples in each class, then divide by the number of classes. If there are some classes with very little data, then the classifier is likely to work poorly, because you will have trouble getting reasonable estimates of the parameters for the $p(x_i|y)$.

Worked example 3.1 *Classifying breast tissue samples*

The “breast tissue” dataset at <https://archive.ics.uci.edu/ml/datasets/Breast+Tissue> contains measurements of a variety of properties of six different classes of breast tissue. Build and evaluate a naive bayes classifier to distinguish between the classes automatically from the measurements.

Solution: The main difficulty here is finding appropriate packages, understanding their documentation, and checking they’re right, unless you want to write the source yourself (which really isn’t all that hard). I used the R package `caret` to do train-test splits, cross-validation, etc. on the naive bayes classifier in the R package `klaR`. I separated out a test set randomly (approx 20% of the cases for each class, chosen at random), then trained with cross-validation on the remainder. The class-confusion matrix on the test set was:

Prediction	adi	car	con	fad	gla	mas
adi	2	0	0	0	0	0
car	0	3	0	0	0	1
con	2	0	2	0	0	0
fad	0	0	0	0	1	0
gla	0	0	0	0	2	1
mas	0	1	0	3	0	1

which is fairly good. The accuracy is 52%. In the training data, the classes are nearly balanced and there are six classes, meaning that chance is about 16%. The κ is 4.34. These numbers, and the class-confusion matrix, will vary with test-train split. I have not averaged over splits, which would be the next thing.

Worked example 3.2 *Classifying mouse protein expression*

Build a naive bayes classifier to classify the “mouse protein” dataset from the UC Irvine machine learning repository. The dataset is at <http://archive.ics.uci.edu/ml/datasets/Mice+Protein+Expression>.

Solution: There’s only one significant difficulty here; many of the data items are incomplete. I dropped all incomplete data items, which is about half of the dataset. One can do somewhat more sophisticated things, but we don’t have the tools yet. I used the R package `caret` to do train-test splits, cross-validation, etc. on the naive bayes classifier in the R package `klaR`. I separated out a test set, then trained with cross-validation on the remainder. The class-confusion matrix on the test set was:

Pred'n	c-CS-m	c-CS-s	c-SC-m	c-SC-s	t-CS-m	t-CS-s	t-SC-m	t-SC-s
c-CS-m	9	0	0	0	0	0	0	0
c-CS-s	0	15	0	0	0	0	0	0
c-SC-m	0	0	12	0	0	0	0	0
c-SC-s	0	0	0	15	0	0	0	0
t-CS-m	0	0	0	0	18	0	0	0
t-CS-s	0	0	0	0	0	15	0	0
t-SC-m	0	0	0	0	0	0	12	0
t-SC-s	0	0	0	0	0	0	0	14

which is as accurate as you can get. Again, I have not averaged over splits, which would be the next thing.

Naive bayes with normal class-conditional distributions takes an interesting and suggestive form. Assume we have two classes. Recall our decision rule is

$$\text{say } \begin{cases} + & \text{if } L(+ \rightarrow -)p(+|\mathbf{x}) > L(- \rightarrow +)p(-|\mathbf{x}) \\ - & \text{otherwise} \end{cases}$$

Now as p gets larger, so does $\log p$ (logarithm is a monotonically increasing function), and the rule isn’t affected by adding the same constant to both sides, so we can rewrite as:

$$\text{say } \begin{cases} + & \text{if } \log L(+ \rightarrow -) + \log p(\mathbf{x}|+) + \log p(+) > \log L(- \rightarrow +) + \log p(\mathbf{x}|-) + \log p(-) \\ - & \text{otherwise} \end{cases}$$

Write μ_j^+ , σ_j^+ respectively for the mean and standard deviation for the class-conditional density for the j ’th component of \mathbf{x} for class + (and so on); the comparison becomes $\log L(+ \rightarrow -) - \sum_j \frac{(x_j - \mu_j^+)^2}{2(\sigma_j^+)^2} - \sum_j \log \sigma_j^+ + \log p(+)$ $>$ $\log L(- \rightarrow +) - \sum_j \frac{(x_j - \mu_j^-)^2}{2(\sigma_j^-)^2} - \sum_j \log \sigma_j^- + \log p(-)$ Now we can expand and collect terms really aggressively to get

$$\left(\sum_j c_j x_j^2 - d_j x_j \right) - e > 0$$

(where c_j , d_j , e are functions of the means and standard deviations and losses and priors). Rather than forming these by estimating the means, etc., we could directly search for good values of c_j , d_j and e .

3.3 THE SUPPORT VECTOR MACHINE

Assume we have a set of N example points \mathbf{x}_i that belong to two classes, which we indicate by 1 and -1 . These points come with their class labels, which we write as y_i ; thus, our dataset can be written as

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}.$$

We wish to predict the sign of y for any point \mathbf{x} . We will use a linear classifier, so that for a new data item \mathbf{x} , we will predict

$$\text{sign}(\mathbf{a} \cdot \mathbf{x} + b)$$

and the particular classifier we use is given by our choice of \mathbf{a} and b .

You should think of \mathbf{a} and b as representing a hyperplane, given by the points where $\mathbf{a} \cdot \mathbf{x} + b = 0$. This hyperplane separates the positive data from the negative data, and is known as the **decision boundary**. Notice that the magnitude of $\mathbf{a} \cdot \mathbf{x} + b$ grows as the point \mathbf{x} moves further away from the hyperplane.

Example: 3.1 *A linear model with a single feature*

Assume we use a linear model with one feature. Then the model has the form $y_i^{(p)} = \text{sign}(ax_i + b)$. For any particular example which has the feature value x^* , this means we will test whether x^* is larger than, or smaller than, $-b/a$.

Example: 3.2 *A linear model with two features*

Assume we use a linear model with two features. Then the model has the form $y_i^{(p)} = \text{sign}(\mathbf{a}^T \mathbf{x}_i + b)$. The sign changes along the line $\mathbf{a}^T \mathbf{x} + b = 0$. You should check that this is, indeed, a line. On one side of this line, the model makes positive predictions; on the other, negative. Which side is which can be swapped by multiplying \mathbf{a} and b by -1 .

This family of classifiers may look bad to you, and it is easy to come up with examples that it misclassifies badly. In fact, the family is extremely strong. First, it is easy to estimate the best choice of rule for very large datasets. Second, linear

classifiers have a long history of working very well in practice on real data. Third, linear classifiers are fast to evaluate.

In fact, examples that are classified badly by the linear rule usually are classified badly because there are two few features. Remember the case of the alien who classified humans into male and female by looking at their heights; if that alien had looked at their chromosomes as well, the error rate would be extremely small. In practical examples, experience shows that the error rate of a poorly performing linear classifier can usually be improved by adding features to the vector \mathbf{x} .

Recall that using naive bayes with a normal model for the class conditional distributions boiled down to testing $(\sum_j c_j x_j^2 - d_j x_j) - e > 0$ for some values of c_j , d_j , and e . This may not look to you like a linear classifier, but it is. Imagine that, for an example \mathbf{u}_i , you form the feature vector

$$\mathbf{x} = (u_{i,1}^2, u_{i,1}, u_{i,2}^2, u_{i,2}, \dots, u_{i,d})^T.$$

Then we can interpret testing $\mathbf{a}^T \mathbf{x} + b > 0$ as testing $a_1 u_{i,1}^2 - (-a_2) u_{i,1} + a_3 u_{i,2}^2 - (-a_4) u_{i,2} + \dots - (-b) > 0$, and pattern matching to the expression for naive bayes suggests that the two cases are equivalent (i.e. for any choice of \mathbf{a} , b , there is a corresponding naive bayes case and vice versa; exercises).

3.3.1 Choosing a Classifier with the Hinge Loss

We will choose \mathbf{a} and b by choosing values that minimize a cost function. We will adopt a cost function of the form:

Training error cost + penalty term.

For the moment, we will ignore the penalty term and focus on the training error cost. Write

$$\gamma_i = \mathbf{a}^T \mathbf{x}_i + b$$

for the value that the linear function takes on example i . Write $C(\gamma_i, y_i)$ for a function that compares γ_i with y_i . The training error cost will be of the form

$$(1/N) \sum_{i=1}^N C(\gamma_i, y_i).$$

A good choice of C should have some important properties. If γ_i and y_i have different signs, then C should be large, because the classifier will make the wrong prediction for this training example. Furthermore, if γ_i and y_i have different signs and γ_i has large magnitude, then the classifier will very likely make the wrong prediction for test examples that are close to \mathbf{x}_i . This is because the magnitude of $(\mathbf{a} \cdot \mathbf{x} + b)$ grows as \mathbf{x} gets further from the decision boundary. So C should get larger as the magnitude of γ_i gets larger in this case.

If γ_i and y_i have the same signs, but γ_i has small magnitude, then the classifier will classify \mathbf{x}_i correctly, but might not classify points that are nearby correctly. This is because a small magnitude of γ_i means that \mathbf{x}_i is close to the decision boundary. So C should not be zero in this case. Finally, if γ_i and y_i have the same

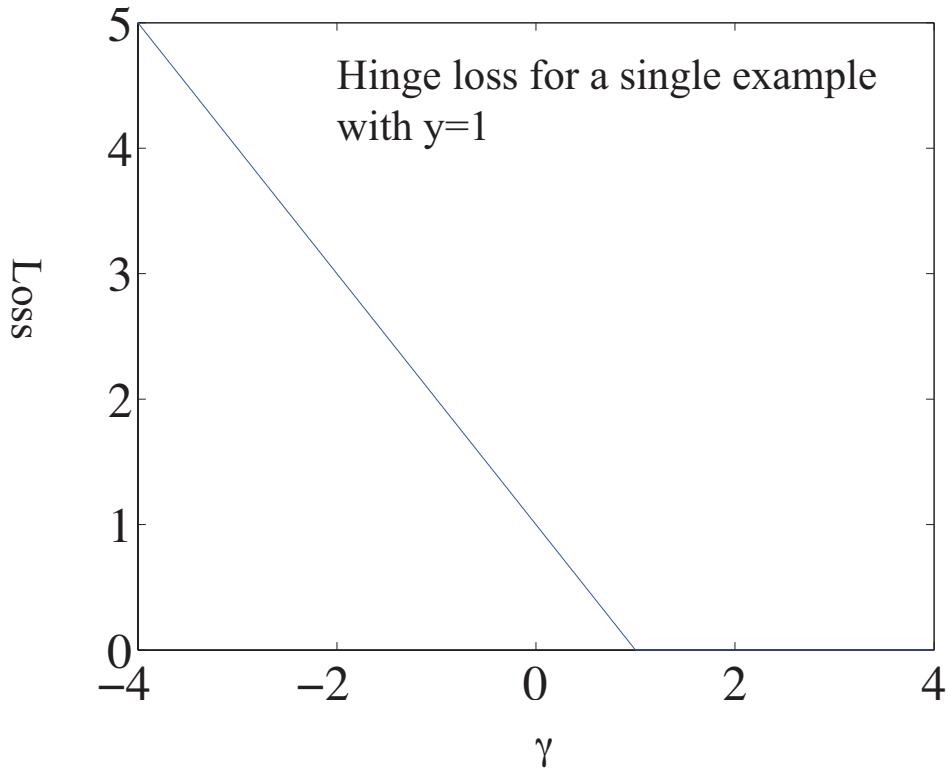


FIGURE 3.1: The hinge loss, plotted for the case $y_i = 1$. The horizontal variable is the $\gamma_i = \mathbf{a} \cdot \mathbf{x}_i + b$ of the text. Notice that giving a strong negative response to this positive example causes a loss that grows linearly as the magnitude of the response grows. Notice also that giving an insufficiently positive response also causes a loss. Giving a strongly positive response is free.

signs and γ_i has large magnitude, then C can be zero because \mathbf{x}_i is on the right side of the decision boundary and so are all the points near to \mathbf{x}_i .

The choice

$$C(y_i, \gamma_i) = \max(0, 1 - y_i \gamma_i)$$

has these properties. If $y_i \gamma_i > 1$ (so the classifier predicts the sign correctly *and* \mathbf{x}_i is far from the boundary) there is no cost. But in any other case, there is a cost. The cost rises if \mathbf{x}_i moves toward the decision boundary from the correct side, and grows linearly as \mathbf{x}_i moves further away from the boundary on the wrong side (Figure 3.1). This means that minimizing the loss will encourage the classifier to (a) make strong positive (or negative) predictions for positive (or negative) examples and (b) for examples it gets wrong, make the most positive (negative) prediction that it can. This choice is known as the **hinge loss**.

Now we think about future examples. We don't know what their feature values will be, and we don't know their labels. But we do know that an example

with feature vector \mathbf{x} will be classified with the rule $\text{sign}(\mathbf{a} \cdot \mathbf{x} + b)$. If we classify this example wrongly, we should like $|\mathbf{a} \cdot \mathbf{x} + b|$ to be small. Achieving this would mean that at least some nearby examples will have the right sign. The way to achieve this is to ensure that $\|\mathbf{a}\|$ is small. By this argument, we would like to achieve a small value of the hinge loss using a small value of $\|\mathbf{a}\|$. Thus, we add a penalty term to the loss so that pairs (\mathbf{a}, b) that have small values of the hinge loss and large values of $\|\mathbf{a}\|$ are expensive. We minimize

$$S(\mathbf{a}, b; \lambda) = \left[(1/N) \sum_{i=1}^N \max(0, 1 - y_i (\mathbf{a} \cdot \mathbf{x}_i + b)) \right] + \frac{\lambda}{2} \mathbf{a}^T \mathbf{a}$$

(hinge loss) (penalty)

where λ is some weight that balances the importance of a small hinge loss against the importance of a small $\|\mathbf{a}\|$. There are now two problems to solve. First, assume we know λ ; we will need to find \mathbf{a} and b that minimize $S(\mathbf{a}, b; \lambda)$. Second, we will need to estimate λ .

3.3.2 Finding a Minimum: General Points

I will first summarize general recipes for finding a minimum. Write $\mathbf{u} = [\mathbf{a}, b]$ for the vector obtained by stacking the vector \mathbf{a} together with b . We have a function $g(\mathbf{u})$, and we wish to obtain a value of \mathbf{u} that achieves the minimum for that function. Sometimes we can solve this problem in closed form by constructing the gradient and finding a value of \mathbf{u} that makes the gradient zero. This happens mainly for specially chosen problems that occur in textbooks. For practical problems, we tend to need a numerical method.

Typical methods take a point $\mathbf{u}^{(i)}$, update it to $\mathbf{u}^{(i+1)}$, then check to see whether the result is a minimum. This process is started from a start point. The choice of start point may or may not matter for general problems, but for our problem it won't matter. The update is usually obtained by computing a direction $\mathbf{p}^{(i)}$ such that for small values of h , $g(\mathbf{u}^{(i)} + h\mathbf{p}^{(i)})$ is smaller than $g(\mathbf{u}^{(i)})$. Such a direction is known as a **descent direction**. We must then determine how far to go along the descent direction, a process known as **line search**.

One method to choose a descent direction is **gradient descent**, which uses the negative gradient of the function. Recall our notation that

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \dots \\ u_d \end{pmatrix}$$

and that

$$\nabla g = \begin{pmatrix} \frac{\partial g}{\partial u_1} \\ \frac{\partial g}{\partial u_2} \\ \dots \\ \frac{\partial g}{\partial u_d} \end{pmatrix}.$$

We can write a Taylor series expansion for the function $g(\mathbf{u}^{(i)} + h\mathbf{p}^{(i)})$. We have that

$$g(\mathbf{u}^{(i)} + h\mathbf{p}^{(i)}) = g(\mathbf{u}^{(i)}) + h(\nabla g)^T \mathbf{p}^{(i)} + O(h^2)$$

This means that we can expect that if

$$\mathbf{p}^{(i)} = -\nabla g(\mathbf{u}^{(i)}),$$

we expect that, at least for small values of h , $g(\mathbf{u}^{(i)} + h\mathbf{p}^{(i)})$ will be less than $g(\mathbf{u}^{(i)})$. This works (as long as g is differentiable, and quite often when it isn't) because g must go down for at least small steps in this direction.

3.3.3 Finding a Minimum: Stochastic Gradient Descent

Assume we wish to minimize some function $g(\mathbf{u}) = g_0(\mathbf{u}) + (1/N) \sum_{i=1}^N g_i(\mathbf{u})$, as a function of \mathbf{u} . Gradient descent would require us to form

$$-\nabla g(\mathbf{u}) = -\left(\nabla g_0(\mathbf{u}) + (1/N) \sum_{i=1}^N \nabla g_i(\mathbf{u})\right)$$

and then take a small step in this direction. But if N is large, this is unattractive, as we might have to sum a lot of terms. This happens a lot in building classifiers, where you might quite reasonably expect to deal with millions of examples. For some cases, there might be trillions of examples. Touching each example at each step really is impractical.

Instead, assume that, at each step, we choose a number k in the range $1 \dots N$ uniformly and at random, and form

$$\mathbf{p}_k = -(\nabla g_0(\mathbf{u}) + \nabla g_k(\mathbf{u}))$$

and then take a small step along \mathbf{p}_k . Our new point becomes

$$\mathbf{a}^{(i+1)} = \mathbf{a}^{(i)} + \eta \mathbf{p}_k^{(i)},$$

where η is called the **steplength** (even though it very often isn't the length of the step we take!). It is easy to show that

$$\mathbb{E}[\mathbf{p}_k] = \nabla g(\mathbf{u})$$

(where the expectation is over the random choice of k). This implies that if we take many small steps along \mathbf{p}_k , they should average out to a step backwards along the gradient. This approach is known as **stochastic gradient descent** (because we're not going along the gradient, but along a random vector which is the gradient only in expectation). It isn't obvious that stochastic gradient descent is a good idea. Although each step is easy to take, we may need to take more steps. The question is then whether we gain in the increased speed of the step what we lose by having to take more steps. Not much is known theoretically, but in practice the approach is hugely successful for training classifiers.

Choosing a steplength η takes some work. We can't search for the step that gives us the best value of g , because we don't want to evaluate the function g (doing

so involves looking at each of the g_i terms). Instead, we use a steplength that is large at the start — so that the method can explore large changes in the values of the classifier parameters — and small steps later — so that it settles down. One useful strategy is to divide training into **epochs**. Each epoch is a block of a fixed number of iterations. Each iteration is one of the steps given above, with fixed steplength. However, the steplength changes from epoch to epoch. In particular, in the r 'th epoch, the steplength is

$$\eta^{(r)} = \frac{a}{r + b}$$

where a and b are constants chosen by experiment with small subsets of the dataset.

One cannot really test whether stochastic gradient descent has converged to the right answer. A better approach is to plot the error as a function of epoch on a validation set. This should vary randomly, but generally go down as the epochs proceed. I have summarized this discussion in box 3.1. You should be aware that the recipe there admits many useful variations, though.

Procedure: 3.1 *Stochastic Gradient Descent*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each x_i is a d -dimensional feature vector, and each y_i is a label, either 1 or -1 . Choose a set of possible values of the regularization weight λ . We wish to train a model that minimizes a cost function of the form $g(\mathbf{u}) = \frac{\lambda}{2} \mathbf{u}^T \mathbf{u} + (\frac{1}{N}) \sum_{i=1}^N g_i(\mathbf{u})$. Separate the data into three sets: test, training and validation. For each value of the regularization weight, train a model, and evaluate the model on validation data. Keep the model that produces the lowest error rate on the validation data, and report its performance on the test data.

Train a model by choosing a fixed number of epochs N_e , and the number of steps per epoch N_s . Choose a random start point, $\mathbf{u}_0 = [\mathbf{a}, b]$. For each epoch, first compute the steplength. In the e 'th epoch, the steplength is typically $\eta = \frac{1}{ae+b}$ for constants a and b chosen by small-scale experiments (you try training a model with different values and see what happens). For the e 'th epoch, choose a subset of the training set for validation for that epoch. Now repeat until the model has been updated N_s times:

- Take k steps. Each step is taken by selecting a single data item uniformly and at random. Assume we select the i 'th data item. We then compute $\mathbf{p} = -\nabla g_i(\mathbf{u}) - \lambda \mathbf{u}$, and update the model by computing

$$\mathbf{u}_{n+1} = \mathbf{u}_n + \eta \mathbf{p}$$

- Evaluate the current model by computing the accuracy on the validation set for that epoch. Plot the accuracy as a function of step number.

3.3.4 Example: Training an SVM with Stochastic Gradient Descent

I have summarized stochastic gradient descent in algorithm 3.1, but here is an example in more detail. We need to choose \mathbf{a} and b to minimize

$$C(\mathbf{a}, b) = (1/N) \sum_{i=1}^N \max(0, 1 - y_i (\mathbf{a} \cdot \mathbf{x}_i + b)) + \frac{\lambda}{2} \mathbf{a}^T \mathbf{a}.$$

This is a support vector machine, because it uses hinge loss. For a support vector machine, stochastic gradient descent is particularly easy. We have estimates $\mathbf{a}^{(n)}$ and $b^{(n)}$ of the classifier parameters, and we want to improve the estimates. We pick the k 'th example at random. We must now compute

$$\nabla \left(\max(0, 1 - y_k (\mathbf{a} \cdot \mathbf{x}_k + b)) + \frac{\lambda}{2} \mathbf{a}^T \mathbf{a} \right).$$

Assume that $y_k (\mathbf{a} \cdot \mathbf{x}_k + b) > 1$. In this case, the classifier predicts a score with the right sign, and a magnitude that is greater than one. Then the first term is zero, and the gradient of the second term is easy. Now if $y_k (\mathbf{a} \cdot \mathbf{x}_k + b) < 1$, we can ignore the max, and the first term is $1 - y_k (\mathbf{a} \cdot \mathbf{x}_k + b)$; the gradient is again easy. But what if $y_k (\mathbf{a} \cdot \mathbf{x}_k + b) = 1$? there are two distinct values we could choose for the gradient, because the max term isn't differentiable. It turns out not to matter which term we choose (Figure ??), so we can write the gradient as

$$p_k = \begin{cases} \begin{bmatrix} \lambda \mathbf{a} \\ 0 \end{bmatrix} & \text{if } y_k (\mathbf{a} \cdot \mathbf{x}_k + b) \geq 1 \\ \begin{bmatrix} \lambda \mathbf{a} - y_k \mathbf{x} \\ -y_k \end{bmatrix} & \text{otherwise} \end{cases}$$

We choose a steplength η , and update our estimates using this gradient. This yields:

$$\mathbf{a}^{(n+1)} = \mathbf{a}^{(n)} - \eta \begin{cases} \lambda \mathbf{a} & \text{if } y_k (\mathbf{a} \cdot \mathbf{x}_k + b) \geq 1 \\ \lambda \mathbf{a} - y_k \mathbf{x} & \text{otherwise} \end{cases}$$

and

$$b^{(n+1)} = b^{(n)} - \eta \begin{cases} 0 & \text{if } y_k (\mathbf{a} \cdot \mathbf{x}_k + b) \geq 1 \\ -y_k & \text{otherwise} \end{cases}.$$

To construct figures, I downloaded the dataset at <http://archive.ics.uci.edu/ml/datasets/Adult>. This dataset apparently contains 48, 842 data items, but I worked with only the first 32, 000. Each consists of a set of numeric and categorical features describing a person, together with whether their annual income is larger than or smaller than 50K\$. I ignored the categorical features to prepare these figures. This isn't wise if you want a good classifier, but it's fine for an example. I used these features to predict whether income is over or under 50K\$. I split the data into 5, 000 test examples, and 27,000 training examples. It's important to do so at random. There are 6 numerical features. I subtracted the mean (which doesn't usually make much difference) and rescaled each so that the variance was 1 (which is often very important). I used two different training regimes.

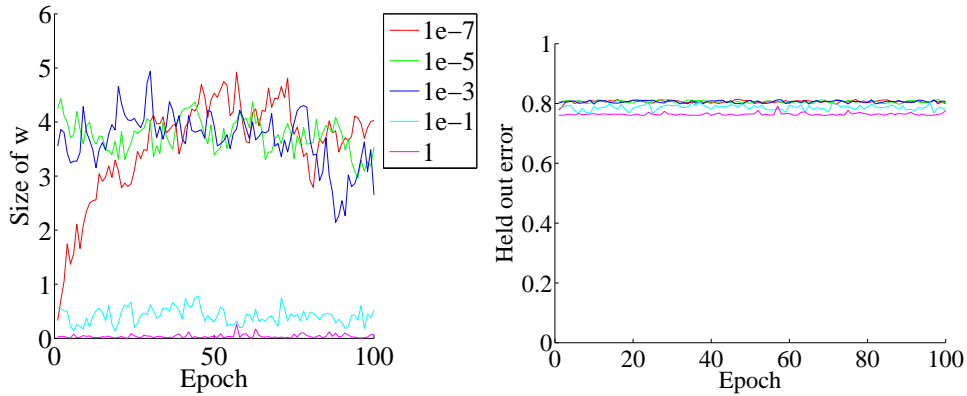


FIGURE 3.2: On the **left**, the magnitude of the weight vector \mathbf{a} at the end of each epoch for the first training regime described in the text. On the **right**, the accuracy on held out data at the end of each epoch. Notice how different choices of regularization parameter lead to different magnitudes of \mathbf{a} ; how the method isn't particularly sensitive to choice of regularization parameter (they change by factors of 100); how the accuracy settles down fairly quickly; and how overlarge values of the regularization parameter do lead to a loss of accuracy.

In the first training regime, there were 100 epochs. In each epoch, I applied 426 steps. For each step, I selected one data item uniformly at random (sampling with replacement), then stepped down the gradient. This means the method sees a total of 42,600 data items. This means that there is a high probability it has touched each data item once (27,000 isn't enough, because we are sampling with replacement, so some items get seen more than once). I chose 5 different values for the regularization parameter and trained with a steplength of $1/(0.01 * e + 50)$, where e is the epoch. At the end of each epoch, I computed $\mathbf{a}^T \mathbf{a}$ and the accuracy (fraction of examples correctly classified) of the current classifier on the held out test examples. Figure 3.2 shows the results. You should notice that the accuracy changes slightly each epoch; that for larger regularizer values $\mathbf{a}^T \mathbf{a}$ is smaller; and that the accuracy settles down to about 0.8 very quickly.

In the second training regime, there were 100 epochs. In each epoch, I applied 50 steps. For each step, I selected one data item uniformly at random (sampling with replacement), then stepped down the gradient. This means the method sees a total of 5,000 data items, and about 3,216 unique data items — it hasn't seen the whole training set. I chose 5 different values for the regularization parameter and trained with a steplength of $1/(0.01 * e + 50)$, where e is the epoch. At the end of each epoch, I computed $\mathbf{a}^T \mathbf{a}$ and the accuracy (fraction of examples correctly classified) of the current classifier on the held out test examples. Figure 3.3 shows the results. You should notice that the accuracy changes slightly each epoch; that for larger regularizer values $\mathbf{a}^T \mathbf{a}$ is smaller; and that the accuracy settles down to about 0.8 very quickly; and that there isn't much difference between the two training regimes. All of these points are relatively typical of stochastic gradient

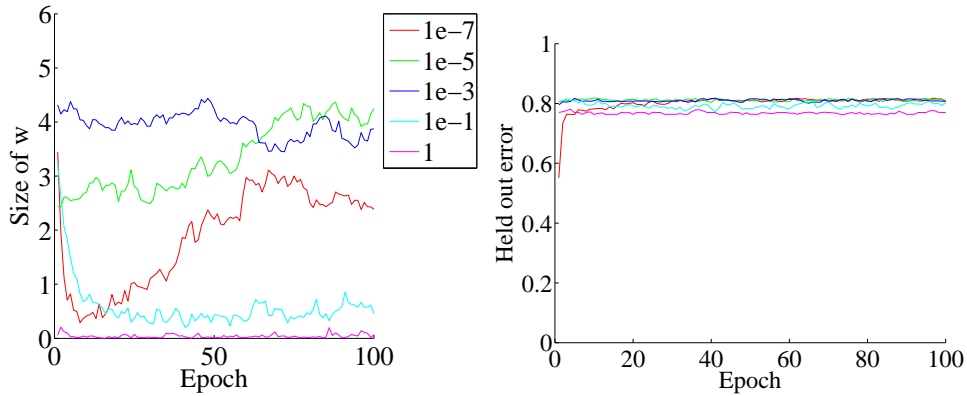


FIGURE 3.3: On the **left**, the magnitude of the weight vector \mathbf{a} at the end of each epoch for the second training regime described in the text. On the **right**, the accuracy on held out data at the end of each epoch. Notice how different choices of regularization parameter lead to different magnitudes of \mathbf{a} ; how the method isn't particularly sensitive to choice of regularization parameter (they change by factors of 100); how the accuracy settles down fairly quickly; and how overlarge values of the regularization parameter do lead to a loss of accuracy.

descent with very large datasets.

Remember this: *Linear SVM's are a go-to classifier. When you have a binary classification problem, the first step should be to try a linear SVM. There is an immense quantity of good software available.*

3.3.5 Multi-Class Classifiers

I have shown how one trains a linear SVM to make a binary prediction (i.e. predict one of two outcomes). But what if there are three, or more, labels? In principle, you could write a binary code for each label, then use a different SVM to predict each bit of the code. It turns out that this doesn't work terribly well, because an error by one of the SVM's is usually catastrophic.

There are two methods that are widely used. In the **all-vs-all** approach, we train a binary classifier for each pair of classes. To classify an example, we present it to each of these classifiers. Each classifier decides which of two classes the example belongs to, then records a vote for that class. The example gets the class label with the most votes. This approach is simple, but scales very badly with the number of classes (you have to build $O(N^2)$ different SVM's for N classes).

In the **one-vs-all** approach, we build a binary classifier for each class. This classifier must distinguish its class from all the other classes. We then take the class

with the largest classifier score. One can think up quite good reasons this approach shouldn't work. For one thing, the classifier isn't told that you intend to use the score to tell similarity between classes. In practice, the approach works rather well and is quite widely used. This approach scales a bit better with the number of classes ($O(N)$).

Remember this: *It is straightforward to build a multi-class classifier out of binary classifiers. Any decent SVM package will do this for you.*

3.4 CLASSIFYING WITH RANDOM FORESTS

I described a classifier as a rule that takes a feature, and produces a class. One way to build such a rule is with a sequence of simple tests, where each test is allowed to use the results of all previous tests. This class of rule can be drawn as a tree (Figure ??), where each node represents a test, and the edges represent the possible outcomes of the test. To classify a test item with such a tree, you present it to the first node; the outcome of the test determines which node it goes to next; and so on, until the example arrives at a leaf. When it does arrive at a leaf, we label the test item with the most common label in the leaf. This object is known as a **decision tree**. Notice one attractive feature of this decision tree: it deals with multiple class labels quite easily, because you just label the test item with the most common label in the leaf that it arrives at when you pass it down the tree.

Figure 3.5 shows a simple 2D dataset with four classes, next to a decision tree that will correctly classify at least the training data. Actually classifying data with a tree like this is straightforward. We take the data item, and pass it down the tree. Notice it can't go both left and right, because of the way the tests work. This means each data item arrives at a single leaf. We take the most common label at the leaf, and give that to the test item. In turn, this means we can build a geometric structure on the feature space that corresponds to the decision tree. I have illustrated that structure in figure 3.5, where the first decision splits the feature space in half (which is why the term split is used so often), and then the next decisions split each of those halves into two.

The important question is how to get the tree from data. It turns out that the best approach for building a tree incorporates a great deal of randomness. As a result, we will get a different tree each time we train a tree on a dataset. None of the individual trees will be particularly good (they are often referred to as “weak learners”). The natural thing to do is to produce many such trees (a **decision forest**), and allow each to vote; the class that gets the most votes, wins. This strategy is extremely effective.

3.4.1 Building a Decision Tree

There are many algorithms for building decision trees. We will use an approach chosen for simplicity and effectiveness; be aware there are others. We will always

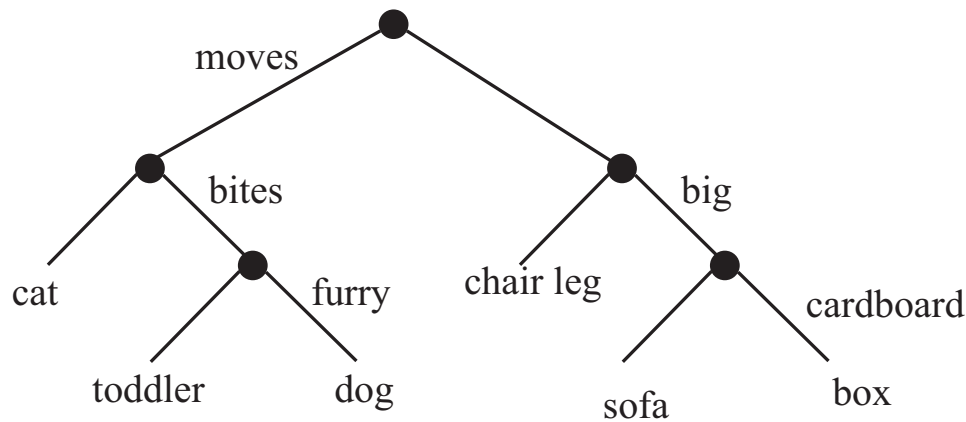


FIGURE 3.4: *This — the household robot’s guide to obstacles — is a typical decision tree. I have labelled only one of the outgoing branches, because the other is the negation. So if the obstacle moves, bites, but isn’t furry, then it’s a toddler. In general, an item is passed down the tree until it hits a leaf. It is then labelled with the leaf’s label.*

use a binary tree, because it’s easier to describe and because that’s usual (it doesn’t change anything important, though). Each node has a **decision function**, which takes data items and returns either 1 or -1.

We train the tree by thinking about its effect on the training data. We pass the whole pool of training data into the root. Any node splits its incoming data into two pools, left (all the data that the decision function labels 1) and right (ditto, -1). Finally, each leaf contains a pool of data, which it can’t split because it is a leaf.

Training the tree uses a straightforward algorithm. First, we choose a class of decision functions to use at each node. It turns out that a very effective algorithm is to choose a single feature at random, then test whether its value is larger than, or smaller than a threshold. For this approach to work, one needs to be quite careful about the choice of threshold, which is what we describe in the next section. Some minor adjustments, described below, are required if the feature chosen isn’t ordinal. Surprisingly, being clever about the choice of *feature* doesn’t seem add a great deal of value. We won’t spend more time on other kinds of decision function, though there are lots.

Now assume we use a decision function as described, and we know how to choose a threshold. We start with the root node, then recursively either split the pool of data at that node, passing the left pool left and the right pool right, or stop splitting and return. Splitting involves choosing a decision function from the class to give the “best” split for a leaf. The main questions are how to choose the best split (next section), and when to stop.

Stopping is relatively straightforward. Quite simple strategies for stopping are very good. It is hard to choose a decision function with very little data, so we must stop splitting when there is too little data at a node. We can tell this is the

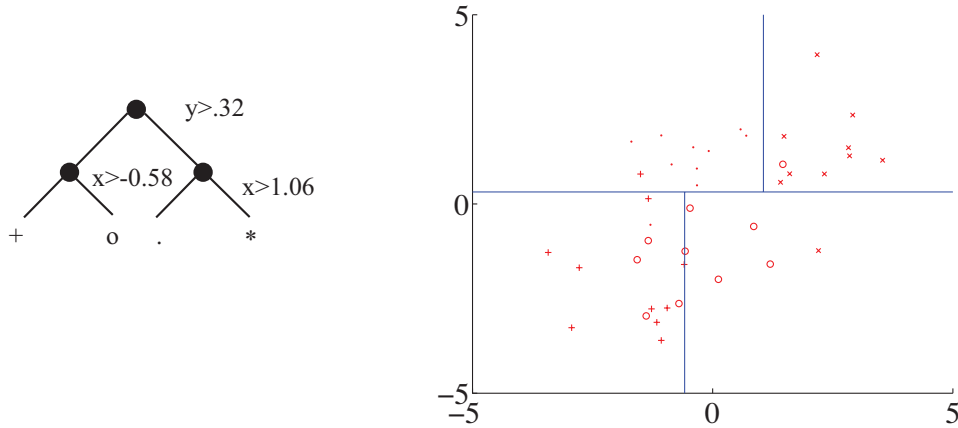


FIGURE 3.5: A straightforward decision tree, illustrated in two ways. On the left, I have given the rules at each split; on the right, I have shown the data points in two dimensions, and the structure that the tree produces in the feature space.

case by testing the amount of data against a threshold, chosen by experiment. If all the data at a node belongs to a single class, there is no point in splitting. Finally, constructing a tree that is too deep tends to result in generalization problems, so we usually allow no more than a fixed depth D of splits. Choosing the best splitting threshold is more complicated.

Figure 3.6 shows two possible splits of a pool of training data. One is quite obviously a lot better than the other. In the good case, the split separates the pool into positives and negatives. In the bad case, each side of the split has the same number of positives and negatives. We cannot usually produce splits as good as the good case here. What we are looking for is a split that will make the proper label more certain.

Figure 3.7 shows a more subtle case to illustrate this. The splits in this figure are obtained by testing the horizontal feature against a threshold. In one case, the left and the right pools contain about the same fraction of positive ('x') and negative ('o') examples. In the other, the left pool is all positive, and the right pool is mostly negative. This is the better choice of threshold. If we were to label any item on the left side positive and any item on the right side negative, the error rate would be fairly small. If you count, the best error rate for the informative split is 20% on the training data, and for the uninformative split it is 40% on the training data.

But we need some way to score the splits, so we can tell which threshold is best. Notice that, in the uninformative case, knowing that a data item is on the left (or the right) does not tell me much more about the data than I already knew. We have that $p(1|\text{left pool, uninformative}) = 2/3 \approx 3/5 = p(1|\text{parent pool})$ and $p(1|\text{right pool, uninformative}) = 1/2 \approx 3/5 = p(1|\text{parent pool})$. For the informative pool, knowing a data item is on the left classifies it completely, and knowing that it is on the right allows us to classify it an error rate of $1/3$. The informative

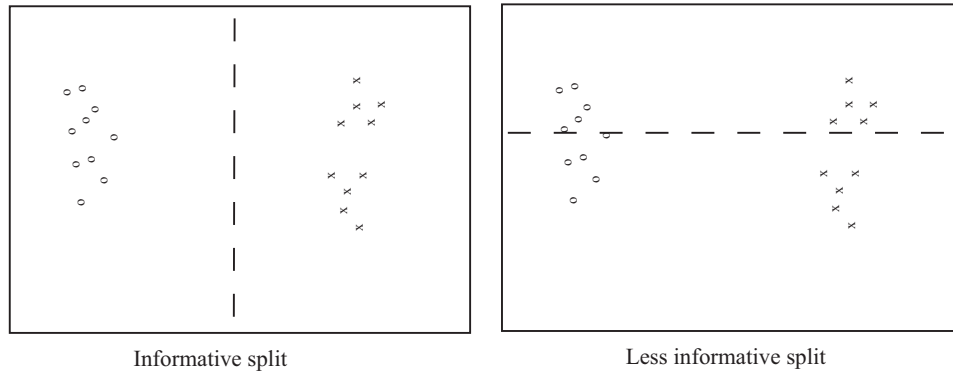


FIGURE 3.6: *Two possible splits of a pool of training data. Positive data is represented with an 'x', negative data with a 'o'. Notice that if we split this pool with the informative line, all the points on the left are 'o's, and all the points on the right are 'x's. This is an excellent choice of split — once we have arrived in a leaf, everything has the same label. Compare this with the less informative split. We started with a node that was half 'x' and half 'o', and now have two nodes each of which is half 'x' and half 'o' — this isn't an improvement, because we do not know more about the label as a result of the split.*

split means that my uncertainty about what class the data item belongs to is significantly reduced if I know whether it goes left or right. To choose a good threshold, we need to keep track of how informative the split is.

3.4.2 Entropy and Information Gain

It turns out to be straightforward to keep track of information, in simple cases. We will start with an example. Assume I have 4 classes. There are 8 examples in class 1, 4 in class 2, 2 in class 3, and 2 in class 4. How much information *on average* will you need to send me to tell me the class of a given example? Clearly, this depends on how you communicate the information. You could send me the complete works of Edward Gibbon to communicate class 1; the Encyclopaedia for class 2; and so on. But this would be redundant. The question is how little can you send me. Keeping track of the amount of information is easier if we encode it with bits (i.e. you can send me sequences of '0's and '1's).

Imagine the following scheme. If an example is in class 1, you send me a '1'. If it is in class 2, you send me '01'; if it is in class 3, you send me '001'; and in class 4, you send me '101'. Then the expected number of bits you will send me is

$$p(\text{class} = 1)1 + p(2)2 + p(3)3 + p(4)3 = \frac{1}{2}1 + \frac{1}{4}2 + \frac{1}{8}3 + \frac{1}{8}3$$

which is 1.75 bits. This number doesn't have to be an integer, because it's an expectation.

Notice that for the i 'th class, you have sent me $-\log_2 p(i)$ bits. We can write



FIGURE 3.7: *Two possible splits of a pool of training data. Positive data is represented with an 'x', negative data with a 'o'. Notice that if we split this pool with the informative line, all the points on the left are 'x's, and two-thirds of the points on the right are 'o's. This means that knowing which side of the split a point lies would give us a good basis for estimating the label. In the less informative case, about two-thirds of the points on the left are 'x's and about half on the right are 'x's — knowing which side of the split a point lies is much less useful in deciding what the label is.*

the expected number of bits you need to send me as

$$-\sum_i p(i) \log_2 p(i).$$

This expression handles other simple cases correctly, too. You should notice that it isn't really important *how many* objects appear in each class. Instead, the *fraction* of all examples that appear in the class is what matters. This fraction is the prior probability that an item will belong to the class. You should try what happens if you have two classes, with an even number of examples in each; 256 classes, with an even number of examples in each; and 5 classes, with $p(1) = 1/2$, $p(2) = 1/4$, $p(3) = 1/8$, $p(4) = 1/16$ and $p(5) = 1/16$. If you try other examples, you may find it hard to construct a scheme where you can send as few bits *on average* as this expression predicts. It turns out that, in general, the smallest number of bits you will need to send me is given by the expression

$$-\sum_i p(i) \log_2 p(i)$$

under all conditions, though it may be hard or impossible to determine what representation is required to achieve this number.

The **entropy** of a probability distribution is a number that scores how many bits, on average, would need to be known to identify an item sampled from that probability distribution. For a discrete probability distribution, the entropy is computed as

$$-\sum_i p(i) \log_2 p(i)$$

where i ranges over all the numbers where $p(i)$ is not zero. For example, if we have two classes and $p(1) = 0.99$, then the entropy is 0.0808, meaning you need very little information to tell which class an object belongs to. This makes sense, because there is a very high probability it belongs to class 1; you need very little information to tell you when it is in class 2. If you are worried by the prospect of having to send 0.0808 bits, remember this is an average, so you can interpret the number as meaning that, if you want to tell which class each of 10^4 independent objects belong to, you could do so in principle with only 808 bits.

Generally, the entropy is larger if the class of an item is more uncertain. Imagine we have two classes and $p(1) = 0.5$, then the entropy is 1, and this is the largest possible value for a probability distribution on two classes. You can always tell which of two classes an object belongs to with just one bit (though you might be able to tell with even less than one bit).

3.4.3 Entropy and Splits

Now we return to the splits. Write \mathcal{P} for the set of all data at the node. Write \mathcal{P}_l for the left pool, and \mathcal{P}_r for the right pool. The **entropy** of a pool \mathcal{C} that scores how many bits would be required to represent the class of an item in that pool, on average. Write $n(i; \mathcal{C})$ for the number of items of class i in the pool, and $N(\mathcal{C})$ for the number of items in the pool. Then the entropy is $H(\mathcal{C})$ of the pool \mathcal{C} is

$$-\sum_i \frac{n(i; \mathcal{C})}{N(\mathcal{C})} \log_2 \frac{n(i; \mathcal{C})}{N(\mathcal{C})}.$$

It is straightforward that $H(\mathcal{P})$ bits are required to classify an item in the parent pool \mathcal{P} . For an item in the left pool, we need $H(\mathcal{P}_l)$ bits; for an item in the right pool, we need $H(\mathcal{P}_r)$ bits. If we split the parent pool, we expect to encounter items in the left pool with probability

$$\frac{N(\mathcal{P}_l)}{N(\mathcal{P})}$$

and items in the right pool with probability

$$\frac{N(\mathcal{P}_r)}{N(\mathcal{P})}.$$

This means that, on average, we must supply

$$\frac{N(\mathcal{P}_l)}{N(\mathcal{P})} H(\mathcal{P}_l) + \frac{N(\mathcal{P}_r)}{N(\mathcal{P})} H(\mathcal{P}_r)$$

bits to classify data items if we split the parent pool. Now a good split is one that results in left and right pools that are informative. In turn, we should need fewer bits to classify once we have split than we need before the split. You can see the difference

$$I(\mathcal{P}_l, \mathcal{P}_r; \mathcal{P}) = H(\mathcal{P}) - \left(\frac{N(\mathcal{P}_l)}{N(\mathcal{P})} H(\mathcal{P}_l) + \frac{N(\mathcal{P}_r)}{N(\mathcal{P})} H(\mathcal{P}_r) \right)$$

as the **information gain** caused by the split. This is the average number of bits that you *don't* have to supply if you know which side of the split an example lies. Better splits have larger information gain.

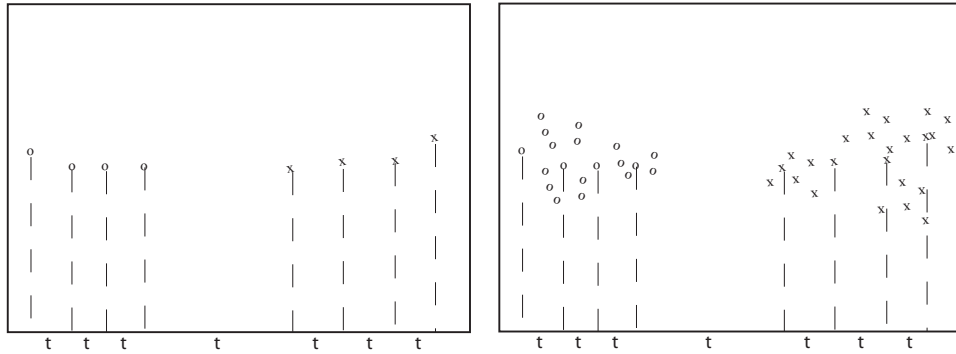


FIGURE 3.8: We search for a good splitting threshold by looking at values of the chosen component that yield different splits. On the **left**, I show a small dataset and its projection onto the chosen splitting component (the horizontal axis). For the 8 data points here, there are only 7 threshold values that produce interesting splits, and these are shown as 't's on the axis. On the **right**, I show a larger dataset; in this case, I have projected only a subset of the data, which results in a small set of thresholds to search.

3.4.4 Choosing a Split with Information Gain

Recall that our decision function is to choose a feature at random, then test its value against a threshold. Any data point where the value is larger goes to the left pool; where the value is smaller goes to the right. This may sound much too simple to work, but it is actually effective and popular. Assume that we are at a node, which we will label k . We have the pool of training examples that have reached that node. The i 'th example has a feature vector \mathbf{x}_i , and each of these feature vectors is a d dimensional vector.

We choose an integer j in the range $1 \dots d$ uniformly and at random. We will split on this feature, and we store j in the node. Recall we write $x_i^{(j)}$ for the value of the j 'th component of the i 'th feature vector. We will choose a threshold t_k , and split by testing the sign of $x_i^{(j)} - t_k$. Choosing the value of t_k is easy. Assume there are N_k examples in the pool. Then there are $N_k - 1$ possible values of t_k that lead to different splits. To see this, sort the N_k examples by $x_i^{(j)}$, then choose values of t_k halfway between example values (Figure 3.8). For each of these values, we compute the information gain of the split. We then keep the threshold with the best information gain.

We can elaborate this procedure in a useful way, by choosing m features at random, finding the best split for each, then keeping the feature and threshold value that is best. It is important that m is a lot smaller than the total number of features — a usual root of thumb is that m is about the square root of the total number of features. It is usual to choose a single m , and choose that for all the splits.

Now assume we happen to have chosen to work with a feature that isn't ordinal, and so can't be tested against a threshold. A natural, and effective, strategy

is as follows. We can split such a feature into two pools by flipping an unbiased coin for each value — if the coin comes up H , any data point with that value goes left, and if it comes up T , any data point with that value goes right. We chose this split at random, so it might not be any good. We can come up with a good split by repeating this procedure F times, computing the information gain for each split, then keeping the one that has the best information gain. We choose F in advance, and it usually depends on the number of values the categorical variable can take.

We now have a relatively straightforward blueprint for an algorithm, which I have put in a box. It's a blueprint, because there are a variety of ways in which it can be revised and changed.

Procedure: 3.2 *Building a decision tree*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each x_i is a d -dimensional feature vector, and each y_i is a label. Call this dataset a **pool**. Now recursively apply the following procedure:

- If the pool is too small, or if all items in the pool have the same label, or if the depth of the recursion has reached a limit, stop.
- Otherwise, search the features for a good split that divides the pool into two, then apply this procedure to each child.

We search for a good split by the following procedure:

- Choose a subset of the feature components at random. Typically, one uses a subset whose size is about the square root of the feature dimension.
- For each component of this subset, search for the best splitting threshold. Do so by selecting a set of possible values for the threshold, then for each value splitting the dataset (every data item with a value of the component below the threshold goes left, others go right), and computing the information gain for the split. Keep the threshold that has the largest information gain.

A good set of possible values for the threshold will contain values that separate the data “reasonably”. If the pool of data is small, you can project the data onto the feature component (i.e. look at the values of that component alone), then choose the $N - 1$ distinct values that lie between two data points. If it is big, you can randomly select a subset of the data, then project that subset on the feature component and choose from the values between data points.

3.4.5 Forests

A single decision tree tends to yield poor classifications. One reason is because the tree is not chosen to give the best classification of its training data. We used a random selection of splitting variables at each node, so the tree can't be the "best possible". Obtaining the best possible tree presents significant technical difficulties. It turns out that the tree that gives the best possible results on the training data can perform rather poorly on test data. The training data is a small subset of possible examples, and so must differ from the test data. The best possible tree on the training data might have a large number of small leaves, built using carefully chosen splits. But the choices that are best for training data might not be best for test data.

Rather than build the best possible tree, we have built a tree efficiently, but with number of random choices. If we were to rebuild the tree, we would obtain a different result. This suggests the following extremely effective strategy: build many trees, and classify by merging their results.

3.4.6 Building and Evaluating a Decision Forest

There are two important strategies for building and evaluating decision forests. I am not aware of evidence strongly favoring one over the other, but different software packages use different strategies, and you should be aware of the options. In one strategy, we separate labelled data into a training and a test set. We then build multiple decision trees, training each using the whole training set. Finally, we evaluate the forest on the test set. In this approach, the forest has not seen some fraction of the available labelled data, because we used it to test. However, each tree has seen every training data item.

Procedure: 3.3 *Building a decision forest*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each \mathbf{x}_i is a d -dimensional feature vector, and each y_i is a label. Separate the dataset into a test set and a training set. Train multiple distinct decision trees on the training set, recalling that the use of a random set of components to find a good split means you will obtain a distinct tree each time.

In the other strategy, sometimes called **bagging**, each time we train a tree we randomly subsample the labelled data with replacement, to yield a training set the same size as the original set of labelled data. Notice that there will be duplicates in this training set, which is like a bootstrap replicate. This training set is often called a **bag**. We keep a record of the examples that do not appear in the bag (the "out of bag" examples). Now to evaluate the forest, we evaluate each tree on its out of bag examples, and average these error terms. In this approach, the entire forest has seen all labelled data, and we also get an estimate of error, but no tree has seen all the training data.

Procedure: 3.4 *Building a decision forest using bagging*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each \mathbf{x}_i is a d -dimensional feature vector, and each y_i is a label. Now build k bootstrap replicates of the training data set. Train one decision tree on each replicate.

3.4.7 Classifying Data Items with a Decision Forest

Once we have a forest, we must classify test data items. There are two major strategies. The simplest is to classify the item with each tree in the forest, then take the class with the most votes. This is effective, but discounts some evidence that might be important. For example, imagine one of the trees in the forest has a leaf with many data items with the same class label; another tree has a leaf with exactly one data item in it. One might not want each leaf to have the same vote.

Procedure: 3.5 *Classification with a decision forest*

Given a test example \mathbf{x} , pass it down each tree of the forest. Now choose one of the following strategies.

- Each time the example arrives at a leaf, record one vote for the label that occurs most often at the leaf. Now choose the label with the most votes.
- Each time the example arrives at a leaf, record N_l votes for each of the labels that occur at the leaf, where N_l is the number of times the label appears in the training data at the leaf. Now choose the label with the most votes.

An alternative strategy that takes this observation into account is to pass the test data item down each tree. When it arrives at a leaf, we record one vote for each of the training data items in that leaf. The vote goes to the class of the training data item. Finally, we take the class with the most votes. This approach allows big, accurate leaves to dominate the voting process. Both strategies are in use, and I am not aware of compelling evidence that one is always better than the other. This may be because the randomness in the training process makes big, accurate leaves uncommon in practice.

Worked example 3.3 *Classifying heart disease data*

Build a random forest classifier to classify the “heart” dataset from the UC Irvine machine learning repository. The dataset is at <http://archive.ics.uci.edu/ml/datasets/Heart+Disease>. There are several versions. You should look at the processed Cleveland data, which is in the file “processed.cleveland.data.txt”.

Solution: I used the R random forest package. This uses a bagging strategy. There is sample code in listing ???. This package makes it quite simple to fit a random forest, as you can see. In this dataset, variable 14 (V14) takes the value 0, 1, 2, 3 or 4 depending on the severity of the narrowing of the arteries. Other variables are physiological and physical measurements pertaining to the patient (read the details on the website). I tried to predict all five levels of variable 14, using the random forest as a multivariate classifier. This works rather poorly, as the out-of-bag class confusion matrix below shows. The total out-of-bag error rate was 45%.

	Predict 0	Predict 1	Predict 2	Predict 3	Predict 4	Class error
True 0	151	7	2	3	1	7.9%
True 1	32	5	9	9	0	91%
True 2	10	9	7	9	1	81%
True 3	6	13	9	5	2	86%
True 4	2	3	2	6	0	100%

This is the example of a class confusion matrix from table 3.1. Fairly clearly, one can predict narrowing or no narrowing from the features, but not the degree of narrowing (at least, not with a random forest). So it is natural to quantize variable 14 to two levels, 0 (meaning no narrowing), and 1 (meaning any narrowing, so the original value could have been 1, 2, or 3). I then built a random forest to predict this from the other variables. The total out-of-bag error rate was 19%, and I obtained the following out-of-bag class confusion matrix

	Predict 0	Predict 1	Class error
True 0	138	26	16%
True 1	31	108	22%

Notice that the false positive rate (16%, from 26/164) is rather better than the false negative rate (22%). Looking at these class confusion matrices, you might wonder whether it is better to predict 0, . . . , 4, then quantize. But this is not a particularly good idea. While the false positive rate is 7.9%, the false negative rate is much higher (36%, from 50/139). In this application, a false negative is likely more of a problem than a false positive, so the tradeoff is unattractive.

Listing 3.1: R code used for the random forests of worked example 3.3

```

setwd('/users/daf/Current/courses/Probcourse/Trees/RCode');
install.packages('randomForest')
library(randomForest)
heart<-read.csv('processed.cleveland.data.txt', header=FALSE)
heart$levels<-as.factor(heart$V14)
heartforest.allvals<-
  randomForest(formula=levels~V1+V2+V3+V4+V5+V6
               +V7+V8+V9+V10+V11+V12+V13,
               data=heart, type='classification', mtry=5)
# this fits to all levels
# I got the OCM by typing
heartforest.allvals
heart$yesno<-cut(heart$V14, c(-Inf, 0.1, Inf))
heartforest<-
  randomForest(formula=yesno~V1+V2+V3+V4+V5+V6
               +V7+V8+V9+V10+V11+V12+V13,
               data=heart, type='classification', mtry=5)
# this fits to the quantized case
# I got the OCM by typing
heartforest

```

Remember this: *Random forests are straightforward to build, and very effective. They can predict any kind of label. Good software implementations are easily available.*

3.5 CLASSIFYING WITH NEAREST NEIGHBORS

Generally we expect that, if two points \mathbf{x}_i and \mathbf{x}_j are close, then their labels will be the same most of the time. If we were unlucky, the two points could lie on either side of a decision boundary, but most pairs of nearby points will have the same labels.

This observation suggests the extremely useful and general strategy of exploiting a data item's neighbors. If you want to classify a data item, find the closest example, and report the class of that example. Alternatively, you could find the closest k examples, and vote. Imagine we have a data point \mathbf{x} that we wish to classify (a query point). Our strategy will be to find the closest training example, and report its class.

How well can we expect this strategy to work? Remember that any classifier will slice up the space of examples into cells (which might be quite complicated) where every point in a cell has the same label. The boundaries between cells are decision boundaries — when a point passes over the decision boundary, the label changes. Now assume we have a large number of labelled training examples, and we know the best possible set of decision boundaries. If there are many training examples, there should be at least one training example that is close to the query

point. If there are enough training examples, then the closest point should be inside the same cell as the query point.

You may be worried that, if the query point is close to a decision boundary, the closest point might be on the other side of that boundary. But if it were, we could improve things by simply having more training points. All this suggests that, *with enough training points*, our classifier should work about as well as the best possible classifier. This intuition turns out to be correct, though the number of training points required is wholly impractical, particularly for high-dimensional feature vectors.

One important generalization is to find the k nearest neighbors, then choose a label from those. A (k, l) nearest neighbor classifier finds the k example points closest to the point being considered, and classifies this point with the class that has the highest number of votes, as long as this class has more than l votes (otherwise, the point is classified as unknown). A $(k, 0)$ -nearest neighbor classifier is usually known as a **k-nearest neighbor classifier**, and a $(1, 0)$ -nearest neighbor classifier is usually known as a **nearest neighbor classifier**. In practice, one seldom uses more than three nearest neighbors. Finding the k nearest points for a particular query can be difficult, and Section ?? reviews this point.

There are three practical difficulties in building nearest neighbor classifiers. You need a lot of labelled examples. You need to be able to find the nearest neighbors for your query point. And you need to use a sensible choice of distance. For features that are obviously of the same type, such as lengths, the usual metric may be good enough. But what if one feature is a length, one is a color, and one is an angle? One possibility is to whiten the features (section 10.1). This may be hard if the dimension is so large that the covariance matrix is hard to estimate. It is almost always a good idea to scale each feature independently so that the variance of each feature is the same, or at least consistent; this prevents features with very large scales dominating those with very small scales. Notice that nearest neighbors (fairly obviously) doesn't like categorical data. If you can't give a clear account of how far apart two things are, you shouldn't be doing nearest neighbors. It is possible to fudge this point a little, by (say) putting together a distance between the levels of each factor, but it's probably unwise.

Nearest neighbors is wonderfully flexible about the labels the classifier predicts. Nothing changes when you go from a two-class classifier to a multi-class classifier.

Worked example 3.4 *Classifying using nearest neighbors*

Build a nearest neighbor classifier to classify the digit data originally constructed by Yann Lecun. You can find it at several places. The original dataset is at <http://yann.lecun.com/exdb/mnist/>. The version I used was used for a Kaggle competition (so I didn't have to decompress Lecun's original format). I found it at <http://www.kaggle.com/c/digit-recognizer>.

Solution: As you'd expect, R has nearest neighbor code that seems quite good (I haven't had any real problems with it, at least). There isn't really all that much to say about the code. I used the R FNN package. This uses a bagging strategy. There is sample code in listing ???. I trained on 1000 of the 42000 examples, so you could see how in the code. I tested on the next 200 examples. For this (rather small) case, I found the following class confusion matrix

	P									
	0	1	2	3	4	5	6	7	8	9
0	12	0	0	0	0	0	0	0	0	0
1	0	20	4	1	0	1	0	2	2	1
2	0	0	20	1	0	0	0	0	0	0
3	0	0	0	12	0	0	0	0	4	0
4	0	0	0	0	18	0	0	0	1	1
5	0	0	0	0	0	19	0	0	1	0
6	1	0	0	0	0	0	18	0	0	0
7	0	0	1	0	0	0	0	19	0	2
8	0	0	1	0	0	0	0	0	16	0
9	0	0	0	2	3	1	0	1	1	14

There are no class error rates here, because I was in a rush and couldn't recall the magic line of R to get them. However, you can see the classifier works rather well for this case.

Remember this: *Nearest neighbor classifiers are often very effective. They can predict any kind of label. You do need to be careful to have enough data, and to have a meaningful distance function.*

3.6 YOU SHOULD

3.6.1 be able to:

- construct a naive bayes classifier for continuous and discrete features
- train a linear SVM with stochastic gradient descent and evaluate the resulting classifier
- train a random forest using a package and evaluate the resulting classifier
- train and evaluate a nearest neighbors classifier

3.6.2 remember:

New term: classifier	10
New term: feature vectors	10
Definition: Classifier	10
New term: decision boundaries	12
New term: overfitting	13
New term: selection bias	13
New term: generalizing badly	13
Do not evaluate a classifier on training data.	13
New term: validation set	13
New term: test set	13
New term: class confusion matrix	14
New term: Vapnik-Chervonenkis dimension	15
New term: V-C dimension	15
New term: likelihood	16
New term: class conditional probability	16
New term: prior	16
New term: posterior	16
New term: decision boundary	20
New term: descent direction	23
New term: line search	23
New term: gradient descent	23
New term: steplength	24
New term: stochastic gradient descent	24
Linear SVM's are a go-to classifier.	28
Any SVM package should build a multi-class classifier for you.	29
New term: decision tree	29
New term: decision forest	29
New term: decision function	30
New term: entropy	33
New term: entropy	34
New term: information gain	34
New term: bagging	37
New term: bag	37
Random forests are good and easy.	40
Nearest neighbors are good and easy.	42

PROBLEMS

PROGRAMMING EXERCISES

- 3.1.** The UC Irvine machine learning data repository hosts a famous collection of data on whether a patient has diabetes (the Pima Indians dataset), originally owned by the National Institute of Diabetes and Digestive and Kidney Diseases and donated by Vincent Sigillito. This can be found at <http://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>. This data has a set of attributes of patients, and a categorical variable telling whether the patient is diabetic or not. For several attributes in this data set, a value of 0 may indicate a missing value of the variable.
- Build a simple naive Bayes classifier to classify this data set. You should hold out 20% of the data for evaluation, and use the other 80% for training. You should use a normal distribution to model each of the class-conditional distributions. You should write this classifier yourself (it's quite straightforward), but you may find the function `createDataPartition` in the R package `caret` helpful to get the random partition.
 - Now adjust your code so that, for attribute 3 (Diastolic blood pressure), attribute 4 (Triceps skin fold thickness), attribute 6 (Body mass index), and attribute 8 (Age), it regards a value of 0 as a missing value when estimating the class-conditional distributions, and the posterior. R uses a special number `NA` to flag a missing value. Most functions handle this number in special, but sensible, ways; but you'll need to do a bit of looking at manuals to check. Does this affect the accuracy of your classifier?
 - Now use the `caret` and `klaR` packages to build a naive bayes classifier for this data, assuming that no attribute has a missing value. The `caret` package does cross-validation (look at `train`) and can be used to hold out data. The `klaR` package can estimate class-conditional densities using a density estimation procedure that I will describe much later in the course. Use the cross-validation mechanisms in `caret` to estimate the accuracy of your classifier. I have not been able to persuade the combination of `caret` and `klaR` to handle missing values the way I'd like them to, but that may be ignorance (look at the `na.action` argument).
 - Now install `SVMLight`, which you can find at <http://svmlight.joachims.org>, via the interface in `klaR` (look for `svmlight` in the manual) to train and evaluate an SVM to classify this data. You don't need to understand much about SVM's to do this — we'll do that in following exercises. You should hold out 20% of the data for evaluation, and use the other 80% for training. You should NOT substitute `NA` values for zeros for attributes 3, 4, 6, and 8.
- 3.2.** The UC Irvine machine learning data repository hosts a collection of data on student performance in Portugal, donated by Paulo Cortez, University of Minho, in Portugal. You can find this data at <https://archive.ics.uci.edu/ml/datasets/Student+Performance>. It is described in P. Cortez and A. Silva. Using Data Mining to Predict Secondary School Student Performance. In A. Brito and J. Teixeira Eds., Proceedings of 5th Future Business Technology Conference (FUBUTEC 2008) pp. 5-12, Porto, Portugal, April, 2008, EUROISIS, ISBN 978-9077381-39-7.
- There are two datasets (for grades in mathematics and for grades in Portuguese). There are 30 attributes each for 649 students, and 3 values that can

be predicted ($G1$, $G2$ and $G3$). Of these, ignore $G1$ and $G2$.

- (a) Use the mathematics dataset. Take the $G3$ attribute, and quantize this into two classes, $G3 > 12$ and $G3 \leq 12$. Build and evaluate a naive bayes classifier that predicts $G3$ from all attributes except $G1$ and $G2$. You should build this classifier from scratch (i.e. DON'T use the packages described in the code snippets). For binary attributes, you should use a binomial model. For the attributes described as “numeric”, which take a small set of values, you should use a multinomial model. For the attributes described as “nominal”, which take a small set of values, you should again use a multinomial model. Ignore the “absence” attribute. Estimate accuracy by cross-validation. You should use at least 10 folds, excluding 15% of the data at random to serve as test data, and average the accuracy over those folds. Report the mean and standard deviation of the accuracy over the folds.
 - (b) Now revise your classifier of the previous part so that, for the attributes described as “numeric”, which take a small set of values, you use a multinomial model. For the attributes described as “nominal”, which take a small set of values, you should still use a multinomial model. Ignore the “absence” attribute. Estimate accuracy by cross-validation. You should use at least 10 folds, excluding 15% of the data at random to serve as test data, and average the accuracy over those folds. Report the mean and standard deviation of the accuracy over the folds.
 - (c) Which classifier do you believe is more accurate and why?
- 3.3.** The UC Irvine machine learning data repository hosts a collection of data on heart disease. The data was collected and supplied by Andras Janosi, M.D., of the Hungarian Institute of Cardiology, Budapest; William Steinbrunn, M.D., of the University Hospital, Zurich, Switzerland; Matthias Pfisterer, M.D., of the University Hospital, Basel, Switzerland; and Robert Detrano, M.D., Ph.D., of the V.A. Medical Center, Long Beach and Cleveland Clinic Foundation. You can find this data at <https://archive.ics.uci.edu/ml/datasets/Heart+Disease>. Use the processed Cleveland dataset, where there are a total of 303 instances with 14 attributes each. The irrelevant attributes described in the text have been removed in these. The 14'th attribute is the disease diagnosis. There are records with missing attributes, and you should drop these.
- (a) Take the disease attribute, and quantize this into two classes, $\text{num} = 0$ and $\text{num} > 0$. Build and evaluate a naive bayes classifier that predicts the class from all other attributes. Estimate accuracy by cross-validation. You should use at least 10 folds, excluding 15% of the data at random to serve as test data, and average the accuracy over those folds. Report the mean and standard deviation of the accuracy over the folds.
 - (b) Now revise your classifier to predict each of the possible values of the disease attribute (0-4 as I recall). Estimate accuracy by cross-validation. You should use at least 10 folds, excluding 15% of the data at random to serve as test data, and average the accuracy over those folds. Report the mean and standard deviation of the accuracy over the folds.
- 3.4.** The UC Irvine machine learning data repository hosts a collection of data on breast cancer diagnostics, donated by Olvi Mangasarian, Nick Street, and William H. Wolberg. You can find this data at [http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)). For each record, there is an id number, 10 continuous variables, and a class (benign or malignant). There are 569 examples. Separate this dataset randomly into 100 validation, 100

test, and 369 training examples.

Write a program to train a support vector machine on this data using stochastic gradient descent. You should not use a package to train the classifier (you don't really need one), but your own code. You should ignore the id number, and use the continuous variables as a feature vector. You should search for an appropriate value of the regularization constant, trying at least the values $\lambda = [1e - 3, 1e - 2, 1e - 1, 1]$. Use the validation set for this search

You should use at least 50 epochs of at least 100 steps each. In each epoch, you should separate out 50 training examples at random for evaluation. You should compute the accuracy of the current classifier on the set held out for the epoch every 10 steps. You should produce:

- (a) A plot of the accuracy every 10 steps, for each value of the regularization constant.
- (b) Your estimate of the best value of the regularization constant, together with a brief description of why you believe that is a good value.
- (c) Your estimate of the accuracy of the best classifier on held out data

- 3.5.** The UC Irvine machine learning data repository hosts a collection of data on adult income, donated by Ronny Kohavi and Barry Becker. You can find this data at <https://archive.ics.uci.edu/ml/datasets/Adult> For each record, there is a set of continuous attributes, and a class ($\leq 50K$ or $> 50K$). There are 48842 examples. You should use only the continuous attributes (see the description on the web page) and drop examples where there are missing values of the continuous attributes. Separate the resulting dataset randomly into 10% validation, 10% test, and 80% training examples.

Write a program to train a support vector machine on this data using stochastic gradient descent. You should not use a package to train the classifier (you don't really need one), but your own code. You should ignore the id number, and use the continuous variables as a feature vector. You should search for an appropriate value of the regularization constant, trying at least the values $\lambda = [1e - 3, 1e - 2, 1e - 1, 1]$. Use the validation set for this search

You should use at least 50 epochs of at least 300 steps each. In each epoch, you should separate out 50 training examples at random for evaluation. You should compute the accuracy of the current classifier on the set held out for the epoch every 30 steps. You should produce:

- (a) A plot of the accuracy every 30 steps, for each value of the regularization constant.
- (b) Your estimate of the best value of the regularization constant, together with a brief description of why you believe that is a good value.
- (c) Your estimate of the accuracy of the best classifier on held out data

- 3.6.** The UC Irvine machine learning data repository hosts a collection of data on the whether p53 expression is active or inactive.

You can find out what this means, and more information about the dataset, by reading: Danziger, S.A., Baronio, R., Ho, L., Hall, L., Salmon, K., Hatfield, G.W., Kaiser, P., and Lathrop, R.H. (2009) Predicting Positive p53 Cancer Rescue Regions Using Most Informative Positive (MIP) Active Learning, *PLOS Computational Biology*, 5(9); Danziger, S.A., Zeng, J., Wang, Y., Brachmann, R.K. and Lathrop, R.H. (2007) Choosing where to look next in a mutation sequence space: Active Learning of informative p53 cancer rescue mutants, *Bioinformatics*, 23(13), 104-114; and Danziger, S.A., Swamidass, S.J., Zeng, J., Dearth, L.R., Lu, Q., Chen, J.H., Cheng, J., Hoang, V.P., Saigo, H., Luo, R., Baldi, P., Brachmann, R.K. and Lathrop, R.H. (2006) Functional

census of mutation sequence spaces: the example of p53 cancer rescue mutants, IEEE/ACM transactions on computational biology and bioinformatics / IEEE, ACM, 3, 114-125.

You can find this data at <https://archive.ics.uci.edu/ml/datasets/p53+Mutants>. There are a total of 16772 instances, with 5409 attributes per instance. Attribute 5409 is the class attribute, which is either active or inactive. There are several versions of this dataset. You should use the version `K8.data`.

- (a) Train an SVM to classify this data, using stochastic gradient descent. You will need to drop data items with missing values. You should estimate a regularization constant using cross-validation, trying at least 3 values. Your training method should touch at least 50% of the training set data. You should produce an estimate of the accuracy of this classifier on held out data consisting of 10% of the dataset, chosen at random.
 - (b) Now train a naive bayes classifier to classify this data. You should produce an estimate of the accuracy of this classifier on held out data consisting of 10% of the dataset, chosen at random.
 - (c) Compare your classifiers. Which one is better? why?
- 3.7.** The UC Irvine machine learning data repository hosts a collection of data on whether a mushroom is edible, donated by Jeff Schlimmer and to be found at <http://archive.ics.uci.edu/ml/datasets/Mushroom>. This data has a set of categorical attributes of the mushroom, together with two labels (poisonous or edible). Use the R random forest package (as in the example in the chapter) to build a random forest to classify a mushroom as edible or poisonous based on its attributes.
- (a) Produce a class-confusion matrix for this problem. If you eat a mushroom based on your classifier's prediction it is edible, what is the probability of being poisoned?

CODE SNIPPETS

Listing 3.2: R code used for the naive bayes example of worked example 3.1

```

setwd('/users/daf/Current/courses/Probcourse/Classification/RCode/BreastTissue')
wdat<-read.csv('cleanedbreast.csv')
library(klaR)
library(caret)
bigx<-wdat[, -c(1:2)]
bigy<-wdat[, 2]
wtd<-createDataPartition(y=bigy, p=.8, list=FALSE)
trax<-bigx[wtd,]
tray<-bigy[wtd]
model<-train(trax, tray, 'nb', trControl=trainControl(method='cv', number=10))
teclasses<-predict(model, newdata=bigx[-wtd,])
confusionMatrix(data=teclasses, bigy[-wtd])

```

Listing 3.3: R code used for the naive bayes example of worked example 3.2

```

setwd('/users/daf/Current/courses/Probcourse/Classification/RCode/MouseProtein')
wdat<-read.csv('Data_Cortex_Nuclear.csv')
#install.packages('klaR')
#install.packages('caret')
library(klaR)
library(caret)
cci<-complete.cases(wdat)
bigx<-wdat[cci, -c(82)]
bigy<-wdat[cci, 82]
wtd<-createDataPartition(y=bigy, p=.8, list=FALSE)
trax<-bigx[wtd,]
tray<-bigy[wtd]
model<-train(trax, tray, 'nb', trControl=trainControl(method='cv', number=10))
teclasses<-predict(model, newdata=bigx[-wtd,])
confusionMatrix(data=teclasses, bigy[-wtd])

```

CHAPTER 4

Extracting Important Relationships in High Dimensions

Chapter 13 described methods to explore the relationship between two elements in a dataset. We could extract a pair of elements and construct various plots. For vector data, we could also compute the correlation between different pairs of elements. But if each data item is d -dimensional, there could be a lot of pairs to deal with.

We will think of our dataset as a collection of d dimensional vectors. It turns out that there are easy generalizations of our summaries. However, is hard to plot d -dimensional vectors. We need to find some way to make them fit on a 2-dimensional plot. Some simple methods can offer insights, but to really get what is going on we need methods that can at all pairs of relationships in a dataset in one go.

These methods visualize the dataset as a “blob” in a d -dimensional space. Many such blobs are flattened in some directions, because components of the data are strongly correlated. Finding the directions in which the blobs are flat yields methods to compute lower dimensional representations of the dataset.

4.1 SOME PLOTS OF HIGH DIMENSIONAL DATA

4.1.1 Understanding Blobs with Scatterplot Matrices - CLEANUP

Plotting high dimensional data is tricky.

4.1.2 Parallel Plots

Parallel plots can sometimes reveal information, particularly when the dimension of the dataset is low. To construct a parallel plot, you compute a normalized representation of each component of each data item. The component is normalized by translating and scaling so that the minimum value over the dataset is zero, and the maximum value over the dataset is one. Now write the i 'th normalised data item as (n_1, n_2, \dots, n_d) . For this data item, you plot a broken line joining $(1, n_1)$ to $(2, n_2)$ to $(3, n_3)$, etc. These plots are superimposed on one another. In the case of the bodyfat dataset, this yields the plot of figure 4.1.

Some structures in the parallel plot are revealing. Outliers often stick out (in figure 4.1, it's pretty clear that there's a data point with a very low height value, and also one with a very large weight value). Outliers affect the scaling, and so make other structures difficult to spot. I have removed them for figure 4.2. In this figure, you can see that two negatively correlated components next to one another produce a butterfly like shape (bodyfat and density). In this plot, you can also see that there are still several data points that are very different from others (two data

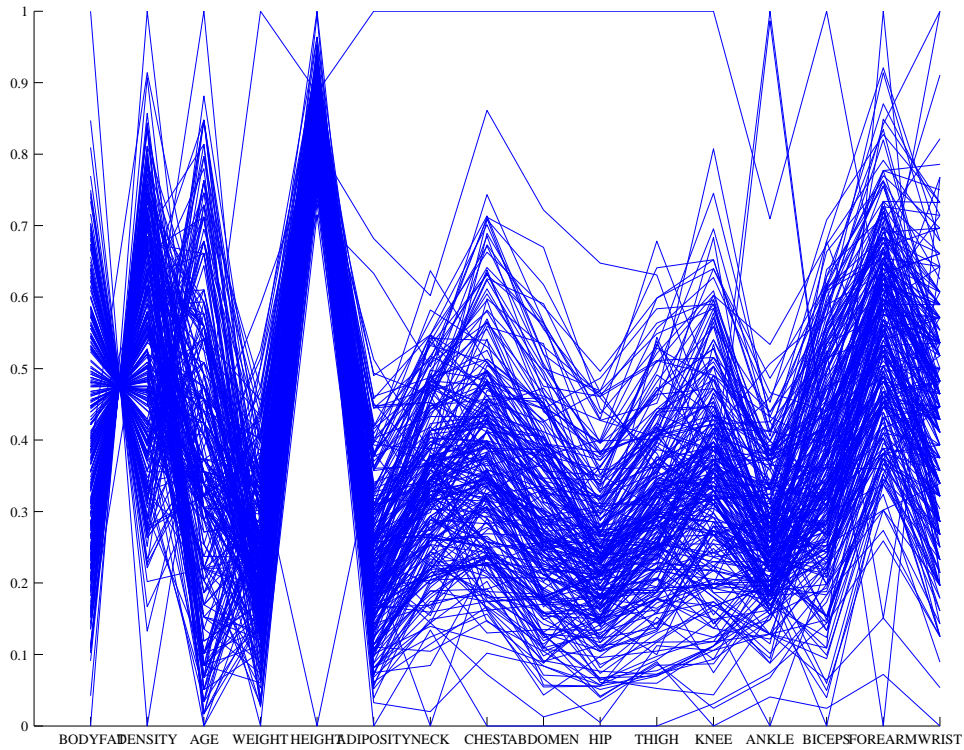


FIGURE 4.1: A parallel plot of the bodyfat dataset, including all data points. I have named the components on the horizontal axis. It is easy to see that large values of bodyfat correspond to small values of density, and vice versa. Notice that one datapoint has height very different from all others; similarly, one datapoint has weight very different from all others.

items have ankle values that are very different from the others, for example).

4.1.3 Scatterplot Matrices

One strategy that is very useful when there aren't too many dimensions is to use a scatterplot matrix. To build one, you lay out scatterplots for each pair of variables in a matrix. On the diagonal, you name the variable that is the vertical axis for each plot in the row, and the horizontal axis in the column. This sounds more complicated than it is; look at the example of figure 4.3, which shows a scatterplot matrix for four of the variables in the height weight dataset of <http://www2.stetson.edu/~jrasp/data.htm>; look for bodyfat.xls at that URL). This is originally a 16-dimensional dataset, but a 16 by 16 scatterplot matrix is squashed and hard to interpret.

What is nice about this kind of plot is that it's quite easy to spot correlations between pairs of variables, though you do need to take into account the coordinates have not been normalized. For figure 4.3, you can see that weight and adiposity

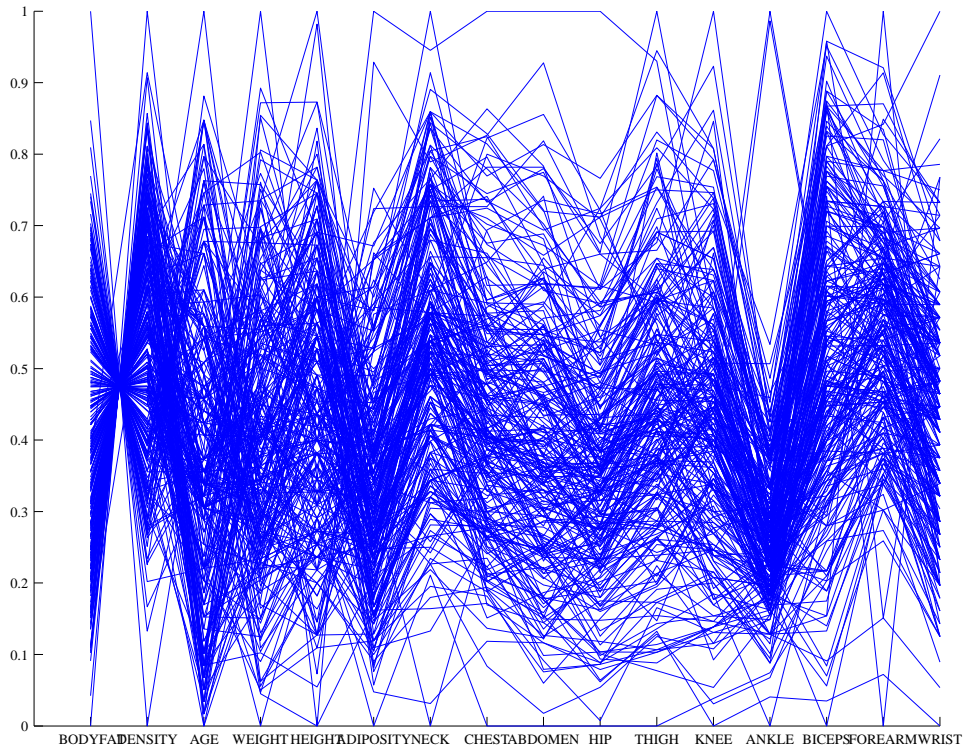


FIGURE 4.2: A plot with those data items removed, so that those components are renormalized. Two datapoints have rather distinct ankle measurements. Generally, you can see that large knees go with large ankles and large biceps (the v structure).

appear to show quite strong correlations, but weight and age are pretty weakly correlated. Height and age seem to have a low correlation. It is also easy to visualize unusual data points. Usually one has an interactive process to do so — you can move a “brush” over the plot to change the color of data points under the brush. To show what might happen, figure 4.4 shows a scatter plot matrix with some points shown as circles. Notice how they lie inside the “blob” of data in some views, and outside in others. This is an effect of projection.

UC Irvine keeps a large repository of datasets that are important in machine learning. You can find the repository at <http://archive.ics.uci.edu/ml/index.html>. Figures 4.5 and 4.6 show visualizations of a famous dataset to do with the botanical classification of irises.

Figures ??, ?? and 4.9 show visualizations of another dataset to do with forest fires in Portugal, also from the UC Irvine repository (look at <http://archive.ics.uci.edu/ml/datasets/Forest+Fires>). In this dataset, there are a variety of measurements of location, time, temperature, etc. together with the area burned by a wildfire. It would be nice to know what leads to large fires, and a visualization is the place to start. Many fires are tiny (or perhaps there was no area measurement?) and so

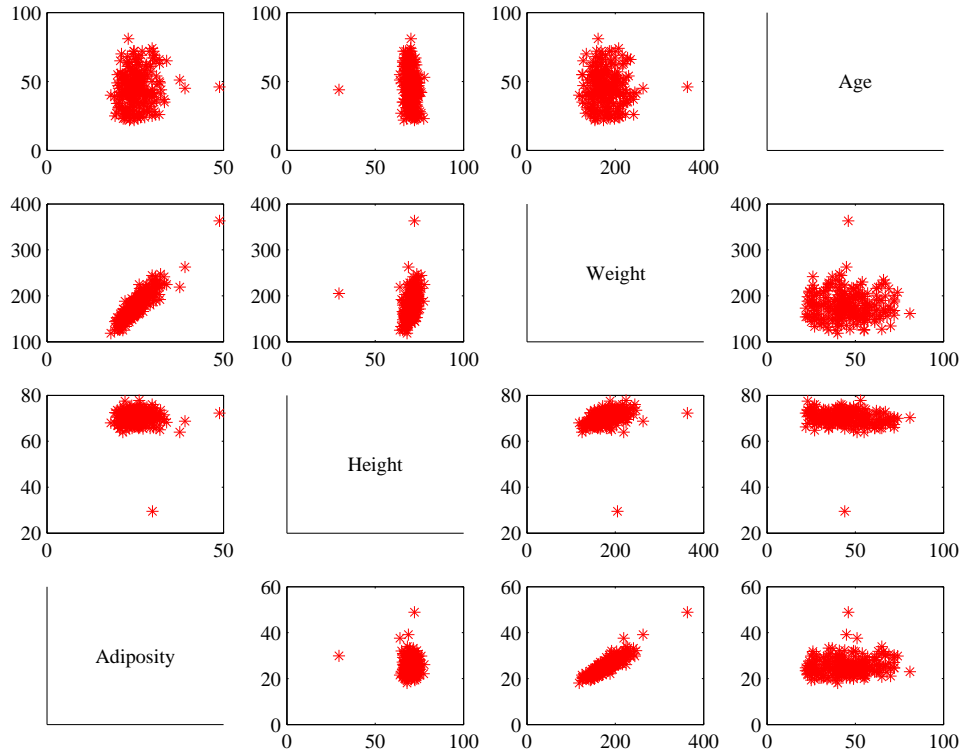


FIGURE 4.3: *This is a scatterplot matrix for four of the variables in the height weight dataset of <http://www2.stetson.edu/~jrasp/data.htm>. Each plot is a scatterplot of a pair of variables. The name of the variable for the horizontal axis is obtained by running your eye down the column; for the vertical axis, along the row. Although this plot is redundant (half of the plots are just flipped versions of the other half), that redundancy makes it easier to follow points by eye. You can look at a column, move down to a row, move across to a column, etc. Notice how you can spot correlations between variables and outliers (the arrows).*

many values of the area are zero. I found it helpful to take the log of area, and then to divide the values of the logarithm into seven categories. I ignored the first four variables, because I didn't think they'd be too important. **Exercise:** was I right? I made two scatterplot matrices, because an eight by eight matrix is too big to view. Generally, this visualization suggests that it would be hard to predict the size of a fire from these variables.

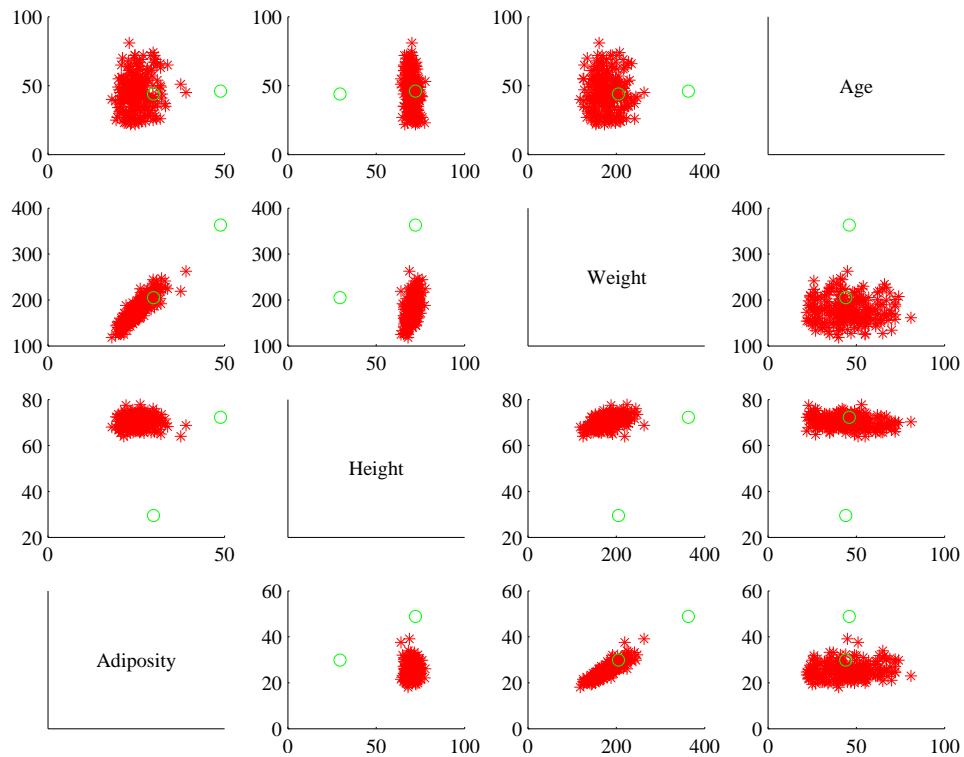


FIGURE 4.4: You should compare this figure with figure 4.3. I have marked two data points with circles in this figure; notice that in some panels these are far from the rest of the data, in others close by. A “brush” in an interactive application can be used to mark data like this to allow a user to explore a dataset.

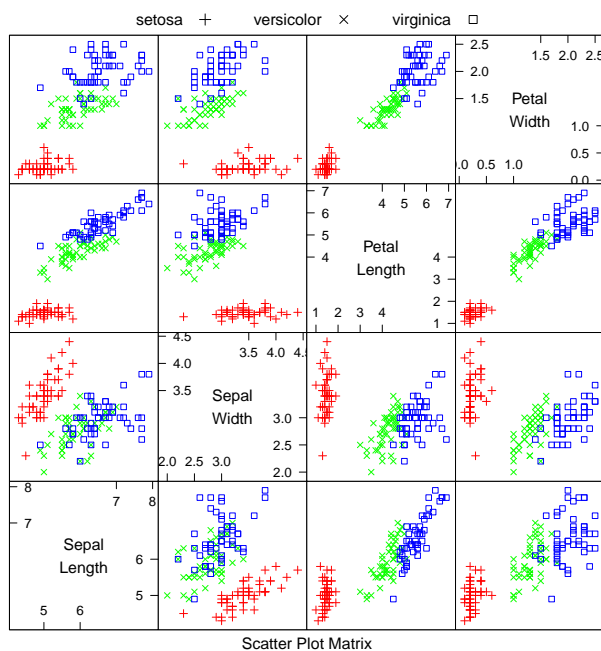


FIGURE 4.5: *This is a scatterplot matrix for the famous Iris data, originally due to ***. There are four variables, measured for each of three species of iris. I have plotted each species with a different marker. You can see from the plot that the species cluster quite tightly, and are different from one another. R code for this plot is on the website.*

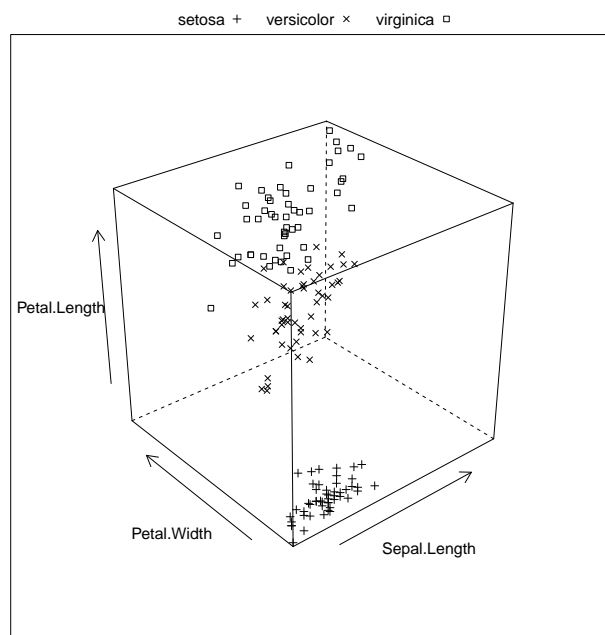


FIGURE 4.6: *This is a 3D scatterplot for the famous Iris data, originally due to ***. I have chosen three variables from the four, and have plotted each species with a different marker. You can see from the plot that the species cluster quite tightly, and are different from one another. R code for this plot is on the website.*

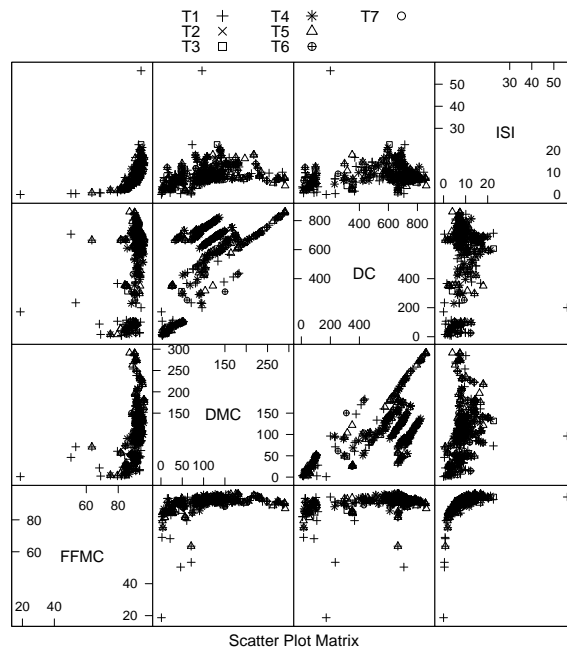


FIGURE 4.7: *This is a scatterplot matrix for the fire dataset from the UC Irvine repository. The smallest area fire is 'T1', and the largest is 'T7'; each is plotted with a different marker. These plots show severity of the fire, plotted against variables 5-8 of the dataset. You should notice that there isn't much separation between the markers. It might be very hard to predict the severity of a fire from these variables. R code for this plot is on the website.*

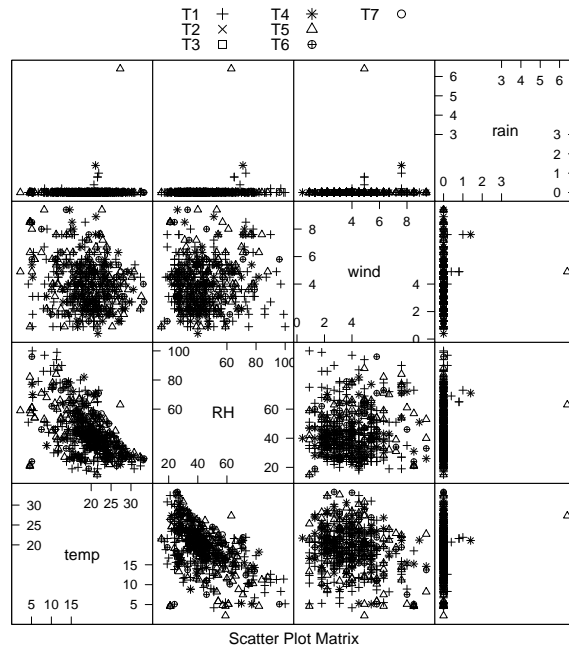


FIGURE 4.8: This is a scatterplot matrix for the fire dataset from the UC Irvine repository. The smallest area fire is 'T1', and the largest is 'T7'; each is plotted with a different marker. These plots show severity of the fire, plotted against variables 9-12 of the dataset. You should notice that there isn't much separation between the markers. It might be very hard to predict the severity of a fire from these variables. R code for this plot is on the website.

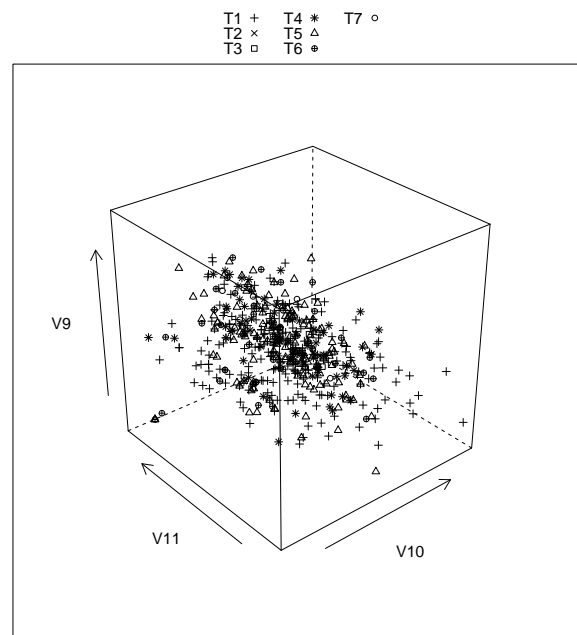


FIGURE 4.9: *This is a 3D scatterplot for the fire dataset from the UC Irvine repository. The smallest area fire is 'T1', and the largest is 'T7'; each is plotted with a different marker. These plots show severity of the fire, plotted against variables 9-11 of the dataset. You should notice that there isn't much separation between the markers. It might be very hard to predict the severity of a fire from these variables. R code for this plot is on the website.*

4.2 SUMMARIES OF HIGH DIMENSIONAL DATA

In this chapter, we assume that our data items are vectors. This means that we can add and subtract values and multiply values by a scalar without any distress. This is an important assumption, but it doesn't necessarily mean that data is continuous (for example, you can meaningfully add the number of children in one family to the number of children in another family). It does rule out a lot of discrete data. For example, you can't add "sports" to "grades" and expect a sensible answer.

Notation: Our data items are vectors, and we write a vector as \mathbf{x} . The data items are d -dimensional, and there are N of them. The entire data set is $\{\mathbf{x}\}$. When we need to refer to the i 'th data item, we write \mathbf{x}_i . We write $\{\mathbf{x}_i\}$ for a new dataset made up of N items, where the i 'th item is \mathbf{x}_i . If we need to refer to the j 'th component of a vector \mathbf{x}_i , we will write $x_i^{(j)}$ (notice this isn't in bold, because it is a component not a vector, and the j is in parentheses because it isn't a power). Vectors are always column vectors.

4.2.1 The Mean

For one-dimensional data, we wrote

$$\text{mean}(\{x\}) = \frac{\sum_i x_i}{N}.$$

This expression is meaningful for vectors, too, because we can add vectors and divide by scalars. We write

$$\text{mean}(\{\mathbf{x}\}) = \frac{\sum_i \mathbf{x}_i}{N}$$

and call this the mean of the data. Notice that each component of $\text{mean}(\{\mathbf{x}\})$ is the mean of that component of the data. There is not an easy analogue of the median, however (how do you order high dimensional data?) and this is a nuisance. Notice that, just as for the one-dimensional mean, we have

$$\text{mean}(\{\mathbf{x} - \text{mean}(\{\mathbf{x}\})\}) = 0$$

(i.e. if you subtract the mean from a data set, the resulting data set has zero mean).

4.2.2 Using Covariance to encode Variance and Correlation

Variance, standard deviation and correlation can each be seen as an instance of a more general operation on data. Assume that we have two one dimensional data sets $\{x\}$ and $\{y\}$. Then we can define the **covariance** of $\{x\}$ and $\{y\}$.

Definition: 4.1 *Covariance*

Assume we have two sets of N data items, $\{x\}$ and $\{y\}$. We compute the covariance by

$$\text{cov}(\{x\}, \{y\}) = \frac{\sum_i (x_i - \text{mean}(\{x\}))(y_i - \text{mean}(\{y\}))}{N}$$

Covariance measures the tendency of corresponding elements of $\{x\}$ and of $\{y\}$ to be larger than (resp. smaller than) the mean. Just like mean, standard deviation and variance, covariance can refer either to a property of a dataset (as in the definition here) or a particular expectation (as in chapter ??). The correspondence is defined by the order of elements in the data set, so that x_1 corresponds to y_1 , x_2 corresponds to y_2 , and so on. If $\{x\}$ tends to be larger (resp. smaller) than its mean for data points where $\{y\}$ is also larger (resp. smaller) than its mean, then the covariance should be positive. If $\{x\}$ tends to be larger (resp. smaller) than its mean for data points where $\{y\}$ is smaller (resp. larger) than its mean, then the covariance should be negative.

From this description, it should be clear we have seen examples of covariance already. Notice that

$$\text{std}(x)^2 = \text{var}(\{x\}) = \text{cov}(\{x\}, \{x\})$$

which you can prove by substituting the expressions. Recall that variance measures the tendency of a dataset to be different from the mean, so the covariance of a dataset with itself is a measure of its tendency not to be constant.

More important, notice that

$$\text{corr}(\{x, y\}) = \frac{\text{cov}(\{x\}, \{y\})}{\sqrt{\text{cov}(\{x\}, \{x\})} \sqrt{\text{cov}(\{y\}, \{y\})}}.$$

This is occasionally a useful way to think about correlation. It says that the correlation measures the tendency of $\{x\}$ and $\{y\}$ to be larger (resp. smaller) than their means for the same data points, *compared to* how much they change on their own.

Working with covariance (rather than correlation) allows us to unify some ideas. In particular, for data items which are d dimensional vectors, it is straightforward to compute a single matrix that captures all covariances between all pairs of components — this is the **covariance matrix**.

Definition: 4.2 *Covariance Matrix*

The covariance matrix is:

$$\text{Covmat}(\{\mathbf{x}\}) = \frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N}$$

Notice that it is quite usual to write a covariance matrix as Σ , and we will follow this convention.

Properties of the Covariance Matrix Covariance matrices are often written as Σ , whatever the dataset (you get to figure out precisely which dataset is intended, from context). Generally, when we want to refer to the j , k 'th entry of a matrix \mathcal{A} , we will write \mathcal{A}_{jk} , so Σ_{jk} is the covariance between the j 'th and k 'th components of the data.

- The j , k 'th entry of the covariance matrix is the covariance of the j 'th and the k 'th components of \mathbf{x} , which we write $\text{cov}(\{x^{(j)}\}, \{x^{(k)}\})$.
- The j , j 'th entry of the covariance matrix is the variance of the j 'th component of \mathbf{x} .
- The covariance matrix is symmetric.
- The covariance matrix is always positive semi-definite; it is positive definite, *unless* there is some vector \mathbf{a} such that $\mathbf{a}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}_i\})) = 0$ for all i .

Proposition:

$$\text{Covmat}(\{\mathbf{x}\})_{jk} = \text{cov}(\{x^{(j)}\}, \{x^{(k)}\})$$

Proof: Recall

$$\text{Covmat}(\{\mathbf{x}\}) = \frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N}$$

and the j , k 'th entry in this matrix will be

$$\frac{\sum_i (x_i^{(j)} - \text{mean}(\{x^{(j)}\}))(x_i^{(k)} - \text{mean}(\{x^{(k)}\}))^T}{N}$$

which is $\text{cov}(\{x^{(j)}\}, \{x^{(k)}\})$.

Proposition:

$$\text{Covmat}(\{\mathbf{x}_i\})_{jj} = \Sigma_{jj} = \text{var}\left(\{x^{(j)}\}\right)$$

Proof:

$$\begin{aligned} \text{Covmat}(\{\mathbf{x}\})_{jj} &= \text{cov}\left(\{x^{(j)}\}, \{x^{(j)}\}\right) \\ &= \text{var}\left(\{x^{(j)}\}\right) \end{aligned}$$

Proposition:

$$\text{Covmat}(\{\mathbf{x}\}) = \text{Covmat}(\{\mathbf{x}\})^T$$

Proof: We have

$$\begin{aligned} \text{Covmat}(\{\mathbf{x}\})_{jk} &= \text{cov}\left(\{x^{(j)}\}, \{x^{(k)}\}\right) \\ &= \text{cov}\left(\{x^{(k)}\}, \{x^{(j)}\}\right) \\ &= \text{Covmat}(\{\mathbf{x}\})_{kj} \end{aligned}$$

Proposition: Write $\Sigma = \text{Covmat}(\{\mathbf{x}\})$. If there is no vector \mathbf{a} such that $\mathbf{a}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})) = 0$ for all i , then for any vector \mathbf{u} , such that $\|\mathbf{u}\| > 0$,

$$\mathbf{u}^T \Sigma \mathbf{u} > 0.$$

If there is such a vector \mathbf{a} , then

$$\mathbf{u}^T \Sigma \mathbf{u} \geq 0.$$

Proof: We have

$$\begin{aligned} \mathbf{u}^T \Sigma \mathbf{u} &= \frac{1}{N} \sum_i [\mathbf{u}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))] [(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T \mathbf{u}] \\ &= \frac{1}{N} \sum_i [\mathbf{u}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))]^2. \end{aligned}$$

Now this is a sum of squares. If there is some \mathbf{a} such that $\mathbf{a}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})) = 0$ for every i , then the covariance matrix must be positive semidefinite (because the sum of squares could be zero in this case). Otherwise, it is positive definite, because the sum of squares will always be positive.

4.3 BLOB ANALYSIS OF HIGH-DIMENSIONAL DATA

When we plotted histograms, we saw that mean and variance were a very helpful description of data that had a unimodal histogram. If the histogram had more than one mode, one needed to be somewhat careful to interpret the mean and variance; in the pizza example, we plotted diameters for different manufacturers to try and see the data as a collection of unimodal histograms.

Generally, mean and covariance are a good description of data that lies in a “blob” (Figure 4.10). You might not believe that this is a technical term, but it’s quite widely used. This is because mean and covariance supply a natural coordinate system in which to interpret the blob. Mean and covariance are less useful as descriptions of data that forms multiple blobs (Figure 4.10). In chapter 16.20, we discuss methods to model data that forms multiple blobs, or other shapes that we will interpret as a set of blobs. But many datasets really are single blobs, and we concentrate on such data here. The way to understand a blob is to think about the coordinate transformations that place a blob into a particularly convenient form.

4.3.1 Transforming High Dimensional Data

Assume we apply an affine transformation to our data set $\{\mathbf{x}\}$, to obtain a new dataset $\{\mathbf{u}\}$, where $\mathbf{u}_i = \mathcal{A}\mathbf{x}_i + \mathbf{b}$. Here \mathcal{A} is any matrix (it doesn’t have to be

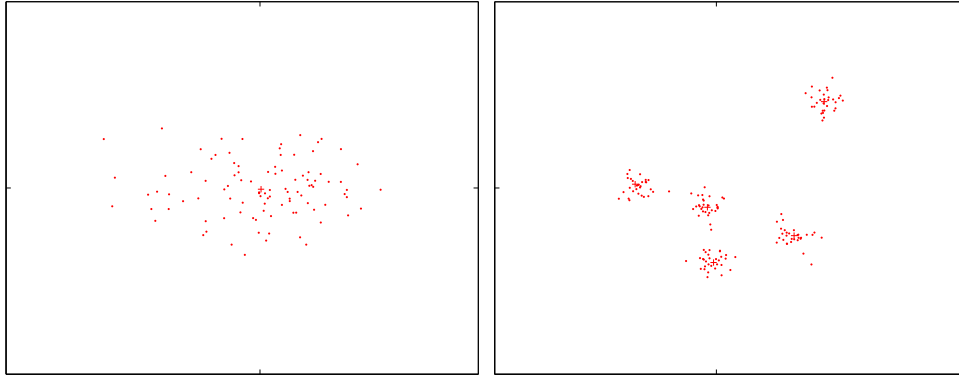


FIGURE 4.10: On the **left**, a “blob” in two dimensions. This is a set of data points that lie somewhat clustered around a single center, given by the mean. I have plotted the mean of these data points with a ‘+’. On the **right**, a data set that is best thought of as a collection of five blobs. I have plotted the mean of each with a ‘+’. We could compute the mean and covariance of this data, but it would be less revealing than the mean and covariance of a single blob. In chapter 16.20, I will describe automatic methods to describe this dataset as a series of blobs.

square, or symmetric, or anything else; it just has to have second dimension d). It is easy to compute the mean and covariance of $\{\mathbf{u}\}$. We have

$$\begin{aligned}\text{mean}(\{\mathbf{u}\}) &= \text{mean}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) \\ &= \mathcal{A}\text{mean}(\{\mathbf{x}\}) + \mathbf{b},\end{aligned}$$

so you get the new mean by multiplying the original mean by \mathcal{A} and adding \mathbf{b} .

The new covariance matrix is easy to compute as well. We have:

$$\begin{aligned}\text{Covmat}(\{\mathbf{u}\}) &= \text{Covmat}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) \\ &= \frac{\sum_i (\mathbf{u}_i - \text{mean}(\{\mathbf{u}\}))(\mathbf{u}_i - \text{mean}(\{\mathbf{u}\}))^T}{N} \\ &= \frac{\sum_i (\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})(\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})^T}{N} \\ &= \frac{\mathcal{A} \sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T \mathcal{A}^T}{N} \\ &= \mathcal{A} \text{Covmat}(\{\mathbf{x}\}) \mathcal{A}^T.\end{aligned}$$

4.3.2 Transforming Blobs

The trick to interpreting high dimensional data is to use the mean and covariance to understand the blob. Figure 4.11 shows a two-dimensional data set. Notice that there is obviously some correlation between the x and y coordinates (it’s a diagonal blob), and that neither x nor y has zero mean. We can easily compute the mean and subtract it from the data points, and this translates the blob so that the origin is at the center (Figure 4.11). In coordinates, this means we compute the new

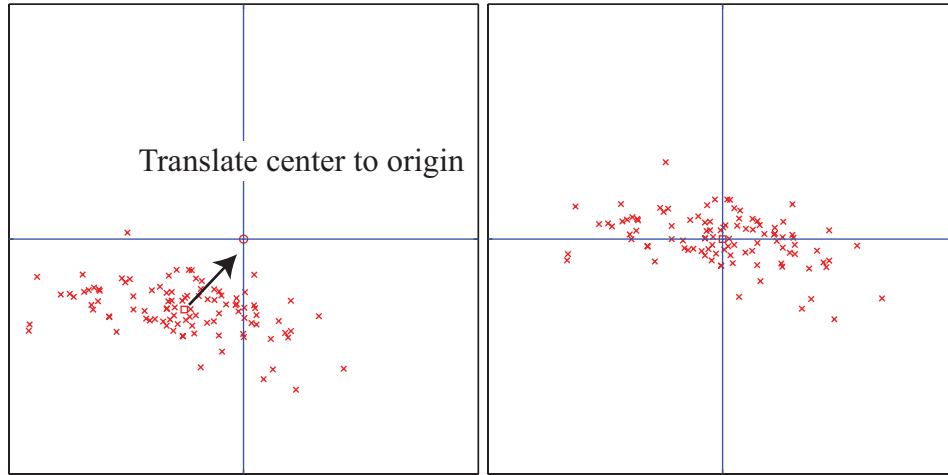


FIGURE 4.11: On the **left**, a “blob” in two dimensions. This is a set of data points that lie somewhat clustered around a single center, given by the mean. I have plotted the mean of these data points with a hollow square (it’s easier to see when there is a lot of data). To translate the blob to the origin, we just subtract the mean from each datapoint, yielding the blob on the **right**.

dataset $\{\mathbf{u}\}$ from the old dataset $\{\mathbf{x}\}$ by the rule $\mathbf{u}_i = \mathbf{x}_i - \text{mean}(\{\mathbf{x}\})$. This new dataset has been translated so that the mean is zero.

Once this blob is translated (Figure 4.12, left), we can rotate it as well. It is natural to try to rotate the blob so that there is no correlation between distinct pairs of dimensions. We can do so by diagonalizing the covariance matrix. In particular, let \mathcal{U} be the matrix formed by stacking the eigenvectors of $\text{Covmat}(\{\mathbf{x}\})$ into a matrix (i.e. $\mathcal{U} = [\mathbf{v}_1, \dots, \mathbf{v}_d]$, where \mathbf{v}_j are eigenvectors of the covariance matrix). We now form the dataset $\{\mathbf{n}\}$, using the rule

$$\mathbf{n}_i = \mathcal{U}^T \mathbf{u}_i = \mathcal{U}^T (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})).$$

The mean of this new dataset is clearly $\mathbf{0}$. The covariance of this dataset is

$$\begin{aligned} \text{Covmat}(\{\mathbf{n}\}) &= \text{Covmat}(\{\mathcal{U}^T \mathbf{x}\}) \\ &= \mathcal{U}^T \text{Covmat}(\{\mathbf{x}\}) \mathcal{U} \\ &= \Lambda, \end{aligned}$$

where Λ is a diagonal matrix of eigenvalues of $\text{Covmat}(\{\mathbf{x}\})$. Remember that, in describing diagonalization, we adopted the convention that the eigenvectors of the matrix being diagonalized were ordered so that the eigenvalues are sorted in descending order along the diagonal of Λ . We now have two very useful facts about $\{\mathbf{n}\}$: (a) every pair of distinct components has covariance zero, and so has correlation zero; (b) the first component has the highest variance, the second component has the second highest variance, and so on. We can rotate and translate any blob into a coordinate system that has these properties. *In this coordinate system, we*

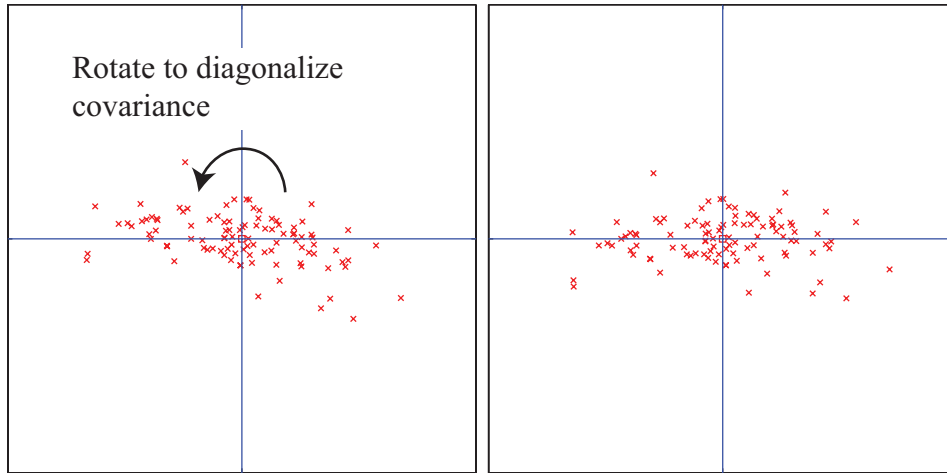


FIGURE 4.12: *On the left, the translated blob of figure 4.11. This blob lies somewhat diagonally, because the vertical and horizontal components are correlated. On the right, that blob of data rotated so that there is no correlation between these components. We can now describe the blob by the vertical and horizontal variances alone, as long as we do so in the new coordinate system. In this coordinate system, the vertical variance is significantly larger than the horizontal variance — the blob is short and wide.*

can describe the blob simply by giving the variances of each component — the covariances are zero.

Translating a blob of data doesn't change the scatterplot matrix in any interesting way (the axes change, but the picture doesn't). Rotating a blob produces really interesting results, however. Figure 4.14 shows the dataset of figure 4.3, translated to the origin and rotated to diagonalize it. Now we do not have names for each component of the data (they're linear combinations of the original components), but each pair is now not correlated. This blob has some interesting shape features. Figure 4.14 shows the gross shape of the blob best. Each panel of this figure has the same scale in each direction. You can see the blob extends about 80 units in direction 1, but only about 15 units in direction 2, and much less in the other two directions. You should think of this blob as being rather cigar-shaped; it's long in one direction, but there isn't much in the others. The cigar metaphor isn't perfect because there aren't any 4 dimensional cigars, but it's helpful. You can think of each panel of this figure as showing views down each of the four axes of the cigar.

Now look at figure ???. This shows the same rotation of the same blob of data, but now the scales on the axis have changed to get the best look at the detailed shape of the blob. First, you can see that blob is a little curved (look at the projection onto direction 2 and direction 4). There might be some effect here worth studying. Second, you can see that some points seem to lie away from the main blob. I have plotted each data point with a dot, and the interesting points

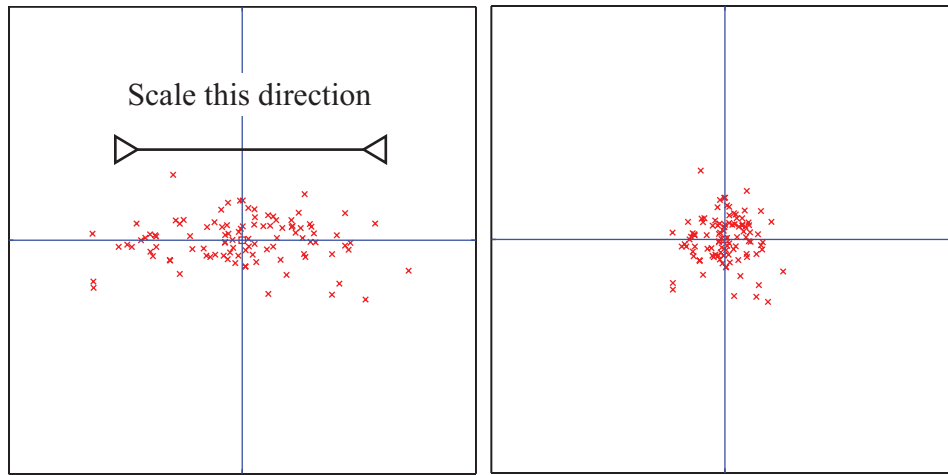


FIGURE 4.13: *On the left, the translated and rotated blob of figure 4.12. This blob is stretched — one direction has more variance than another. Because all covariances are zero, it is easy to scale the blob so that all variances are one (the blob on the right). You can think of this as a standard blob. All blobs can be reduced to a standard blob, by relatively straightforward linear algebra.*

with a number. These points are clearly special in some way.

We could now scale the data in this new coordinate system so that all the variances are either one (if there is any variation in that direction) or zero (directions where the data doesn't vary — these occur only if some directions are functions of others). Figure 4.13 shows the final scaling. The result is a standard blob. Our approach applies to any dimension — I gave 2D figures because they're much easier to understand. There is a crucial point here: we can reduce any blob of data, in any dimension, to a standard blob of that dimension. All blobs are the same, except for some stretching, some rotation, and some translation. This is why blobs are so well-liked.

4.3.3 Whitening Data

It is sometimes useful to actually reduce a dataset to a standard blob. Doing so is known as **whitening the data** (for reasons I find obscure). This can be a sensible thing to do when we don't have a clear sense of the relative scales of the components of each data vector. For example, if we have a dataset where one component ranges from $1e5$ to $2e5$, and the other component ranges from $-1e-5$ to $1e-5$, we are likely to face numerical problems in many computations (adding small numbers to big numbers is often unwise). Often, this kind of thing follows from a poor choice of units, rather than any kind of useful property of the data. In such a case, it could be quite helpful to whiten the data. Another reason to whiten the data might be that we know relatively little about the meaning of each component. In this case, the original choice of coordinate system was somewhat arbitrary anyhow, and

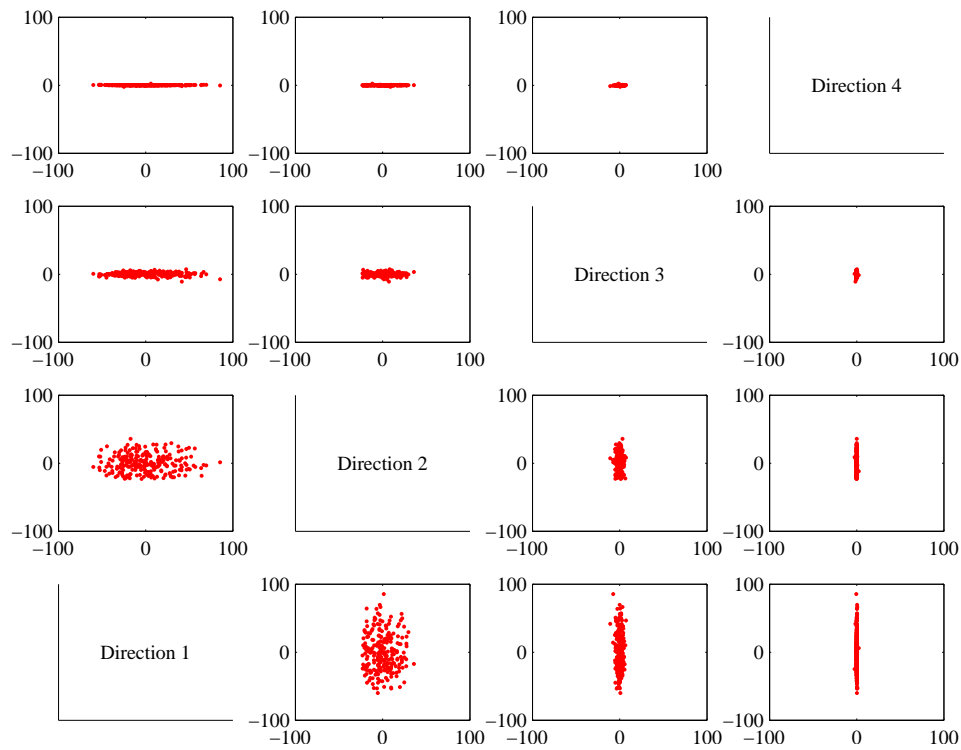


FIGURE 4.14: A panel plot of the bodyfat dataset of figure 4.3, now rotated so that the covariance between all pairs of distinct dimensions is zero. Now we do not know names for the directions — they’re linear combinations of the original variables. Each scatterplot is on the same set of axes, so you can see that the dataset extends more in some directions than in others.

transforming data to a uniform blob could be helpful.

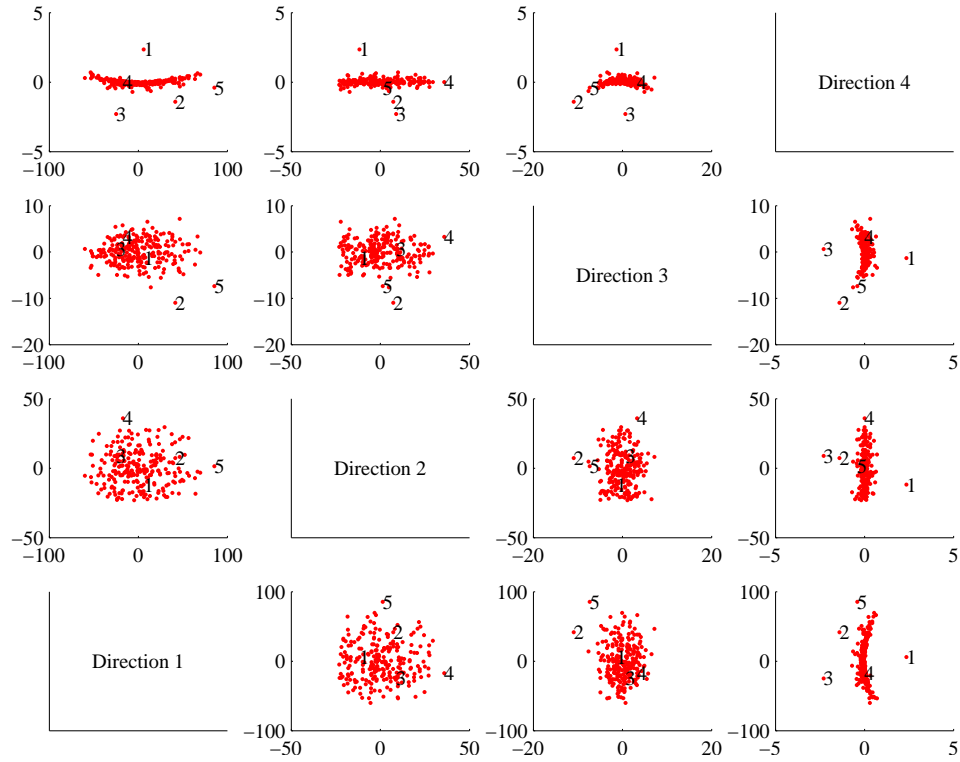


FIGURE 4.15: A panel plot of the bodyfat dataset of figure 4.3, now rotated so that the covariance between all pairs of distinct dimensions is zero. Now we do not know names for the directions — they’re linear combinations of the original variables. I have scaled the axes so you can see details; notice that the blob is a little curved, and there are several data points that seem to lie some way away from the blob, which I have numbered.

Useful Facts: 4.1 *Whitening a dataset*

For a dataset $\{\mathbf{x}\}$, compute:

- \mathcal{U} , the matrix of eigenvectors of $\text{Covmat}(\{\mathbf{x}\})$;
- and $\text{mean}(\{\mathbf{x}\})$.

Now compute $\{\mathbf{n}\}$ using the rule

$$\mathbf{n}_i = \mathcal{U}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})).$$

Then $\text{mean}(\{\mathbf{n}\}) = \mathbf{0}$ and $\text{Covmat}(\{\mathbf{n}\})$ is diagonal.

Now write Λ for the diagonal matrix of eigenvalues of $\text{Covmat}(\{\mathbf{x}\})$ (so that $\text{Covmat}(\{\mathbf{x}\})\mathcal{U} = \mathcal{U}\Lambda$). Assume that each of the diagonal entries of Λ is greater than zero (otherwise there is a redundant dimension in the data). Write λ_i for the i 'th diagonal entry of Λ , and write $\Lambda^{-(1/2)}$ for the diagonal matrix whose i 'th diagonal entry is $1/\sqrt{\lambda_i}$. Compute $\{\mathbf{z}\}$ using the rule

$$\mathbf{z}_i = \Lambda^{-(1/2)}\mathcal{U}(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})).$$

We have that $\text{mean}(\{\mathbf{z}\}) = \mathbf{0}$ and $\text{Covmat}(\{\mathbf{z}\}) = \mathcal{I}$. The dataset $\{\mathbf{z}\}$ is often known as **whitened data**.

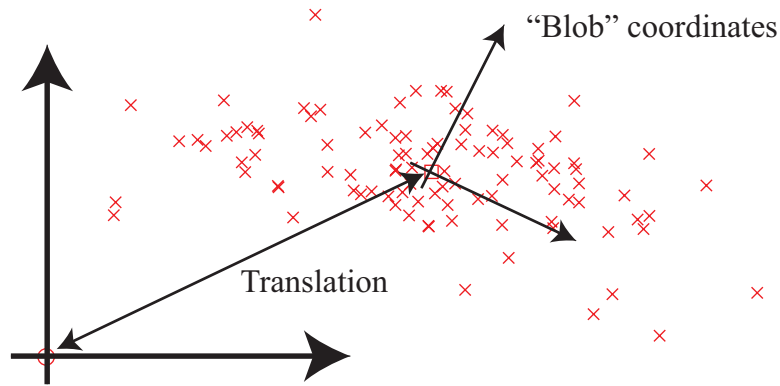


FIGURE 4.16: A 2D blob, with its natural blob coordinate system. The origin of this coordinate system is at the mean of the data. The coordinate axes are (a) at right angles to one another and (b) are directions that have no covariance.

It isn't always a good idea to whiten data. In some circumstances, each separate component is meaningful, and in a meaningful set of units. For example, one of the components might be a length using a natural scale and the other might be a time on a natural scale. When this happens, we might be reluctant to transform the data, either because we don't want to add lengths to times or because we want to preserve the scales.

4.4 PRINCIPAL COMPONENTS ANALYSIS

Mostly, when one deals with high dimensional data, it isn't clear which individual components are important. As we have seen with the height weight dataset (for example, in the case of density and weight) some components can be quite strongly correlated. Equivalently, as in Figure 16.20, the blob is not aligned with the coordinate axes.

4.4.1 The Blob Coordinate System and Smoothing

We can use the fact that we *could* rotate, translate and scale the blob to define a coordinate system within the blob. The origin of that coordinate system is the mean of the data, and the coordinate axes are given by the eigenvectors of the covariance matrix. These are orthonormal, so they form a set of unit basis vectors at right angles to one another (i.e. a coordinate system). You should think of these as blob coordinates; Figure 4.16 illustrates a set of blob coordinates.

The blob coordinate system is important because, *once we know the blob coordinate system*, we can identify important scales of variation in the data. For example, if you look at Figure 4.16, you can see that this blob is extended much further along one direction than along the other. We can use this information to identify the most significant forms of variation in very high dimensional data. In

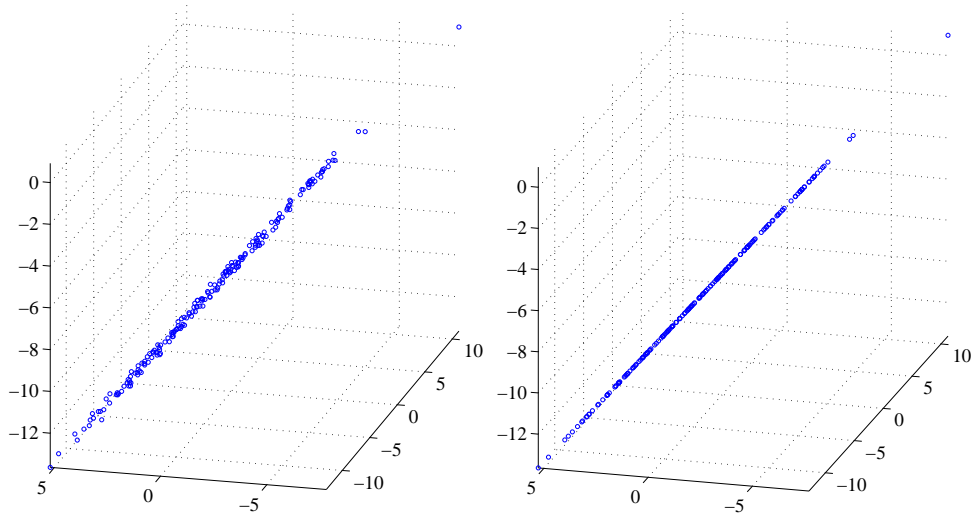


FIGURE 4.17: On the **left**, a blob of 3D data that has very low variance in two directions in the blob coordinates. As a result, all the data points are very close to a 1D blob. Experience shows that this is a common phenomenon. Although there might be many components in the data items, all data points are very close to a much lower dimensional object in the high dimensional space. When this is the case, we could obtain a lower dimensional representation of the data by working in blob coordinates, or we could smooth the data (as on the **right**), by projecting each data point onto the lower dimensional space.

some directions in the blob coordinate system, the blob will be spread out — ie have large variance — but in others, it might not be.

Equivalently, imagine we choose to represent each data item *in blob coordinates*. Then the mean over the dataset will be zero. Each pair of distinct coordinates will be uncorrelated. Some coordinates — corresponding to directions where the blob is spread out — will have a large range of values. Other coordinates — directions in which the blob is small — will have a small range of values. We could choose to replace these coordinates with zeros, with little significant loss in accuracy. The advantage of doing so is that we would have lower dimensional data to deal with.

However, it isn't particularly natural to work in blob coordinates. Each component of a data item may have a distinct meaning and scale (i.e. feet, pounds, and so on), but this is not preserved in any easy way in blob coordinates. Instead, we should like to (a) compute a lower dimensional representation in blob coordinates then (b) transform that representation into the original coordinate system of the data item. Doing so is a form of **smoothing** — suppressing small, irrelevant variations by exploiting multiple data items.

For example, look at Figure 4.17. Imagine we transform the blob on the left to blob coordinates. The covariance matrix in these coordinates is a 3×3 diagonal matrix. One of the values on the diagonal is large, because the blob is extended on

one direction; but the other two are small. This means that, in blob coordinates, the data varies significantly in one direction, but very little in the other two directions.

Now imagine we project the data points onto the high-variation direction; equivalently, we set the other two directions to zero for each data point. Each of the new data points is very close to the corresponding old data point, because by setting the small directions to zero we haven't moved the point very much. In blob coordinates, the covariance matrix of this new dataset has changed very little. It is again a 3×3 diagonal matrix, but now two of the diagonal values are zero, because there isn't any variance in those directions. The third value is large, because the blob is extended in that direction. We take the new dataset, and rotate and translate it into the original coordinate system. Each point must lie close to the corresponding point in the original dataset. However, the new dataset lies along a straight line (because it lay on a straight line in the blob coordinates). This process gets us the blob on the right in Figure 4.17. This blob is a smoothed version of the original blob.

Smoothing works because when two data items are strongly correlated, the value of one is a good guide to the value of the other. This principle works for more than two data items. Section 16.20 describes an example where the data items have dimension 101, but all values are extremely tightly correlated. In a case like this, there may be very few dimensions in blob coordinates that have any significant variation (3-6 for this case, depending on some details of what one believes is a small number, and so on). The components are so strongly correlated in this case that the 101-dimensional blob really looks like a slightly thickened 3 (or slightly more) dimensional blob that has been inserted into a 101-dimensional space (Figure 4.17). If we project the 101-dimensional data onto that structure in the original, 101-dimensional space, we may get much better estimates of the components of each data item than the original measuring device can supply. This occurs because each component is now estimated using correlations between all the measurements.

4.4.2 The Low-Dimensional Representation of a Blob

We wish to construct an r dimensional representation of a blob, where we have chosen r in advance. First, we compute $\{\mathbf{v}\}$ by translating the blob so its mean is at the origin, so that $\mathbf{v}_i = \mathbf{x}_i - \text{mean}(\{\mathbf{x}\})$. Now write $\mathcal{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N]$. The covariance matrix of $\{\mathbf{v}\}$ is then

$$\text{Covmat}(\{\mathbf{v}\}) = \frac{1}{N} \mathcal{V} \mathcal{V}^T = \text{Covmat}(\{\mathbf{x}\}).$$

Now write Λ for the diagonal matrix of eigenvalues of $\text{Covmat}(\{\mathbf{x}\})$ and \mathcal{U} for the matrix of eigenvectors, so that $\text{Covmat}(\{\mathbf{x}\})\mathcal{U} = \mathcal{U}\Lambda$. We assume that the elements of Λ are sorted in decreasing order along the diagonal. The covariance matrix for the dataset transformed into blob coordinates will be Λ . Notice that

$$\begin{aligned} \Lambda &= \mathcal{U}^T \text{Covmat}(\{\mathbf{x}\}) \mathcal{U} \\ &= \mathcal{U}^T \mathcal{V} \mathcal{V}^T \mathcal{U} \\ &= (\mathcal{U}^T \mathcal{V})(\mathcal{U}^T \mathcal{V})^T. \end{aligned}$$

This means we can interpret $(\mathcal{U}^T \mathcal{V})$ as a new dataset $\{\mathbf{b}\}$. This is our data, rotated into blob coordinates.

Now write Π_r for the $d \times d$ matrix

$$\begin{bmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{bmatrix}$$

which projects a d dimensional vector onto its first r components, and replaces the others with zeros. Then we have that

$$\Lambda_r = \Pi_r \Lambda \Pi_r^T$$

is the covariance matrix for the reduced dimensional data in blob coordinates. Notice that Λ_r keeps the r largest eigenvalues on the diagonal of Λ , and replaces all others with zero.

We have

$$\begin{aligned} \Lambda_r &= \Pi_r \Lambda \Pi_r^T \\ &= \Pi_r \mathcal{U}^T \text{Covmat}(\{\mathbf{x}\}) \mathcal{U} \Pi_r^T \\ &= (\Pi_r \mathcal{U}^T \mathcal{V})(\mathcal{V}^T \mathcal{U} \Pi_r^T) \\ &= \mathcal{P} \mathcal{P}^T \end{aligned}$$

where $\mathcal{P} = (\Pi_r \mathcal{U}^T \mathcal{V})$. This represents our data, rotated into blob coordinates, and then projected down to r dimensions, with remaining terms replaced by zeros. Write $\{\mathbf{b}_r\}$ for this new dataset.

Occasionally, we need to refer to this representation, and we give it a special name. Write

$$\text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\}) = \Pi_r \mathcal{U}^T (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))$$

where the notation seeks to explicitly keep track of the fact that the low dimensional representation of a particular data item *depends on the whole dataset* (because you have to be able to compute the mean, and the eigenvectors of the covariance). Notice that $\text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\})$ is a representation of the dataset with important properties:

- The representation is r -dimensional (i.e. the last $d - r$ components are zero).
- Each pair of distinct components of $\{\text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\})\}$ has zero covariance.
- The first component of $\{\text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\})\}$ has largest variance; the second component has second largest variance; and so on.

4.4.3 Smoothing Data with a Low-Dimensional Representation

We would now like to construct a low dimensional representation of the blob, in the original coordinates. We do so by rotating the low-dimensional representation back to the original coordinate system, then adding back the mean to translate the origin back to where it started. We can write this as

$$\begin{aligned} \text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\}) &= \mathcal{U} \Pi_r^T (\Pi_r \mathcal{U}^T (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))) + \text{mean}(\{\mathbf{x}\}) \\ &= \mathcal{U} \Pi_r^T \text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\}) + \text{mean}(\{\mathbf{x}\}) \end{aligned}$$

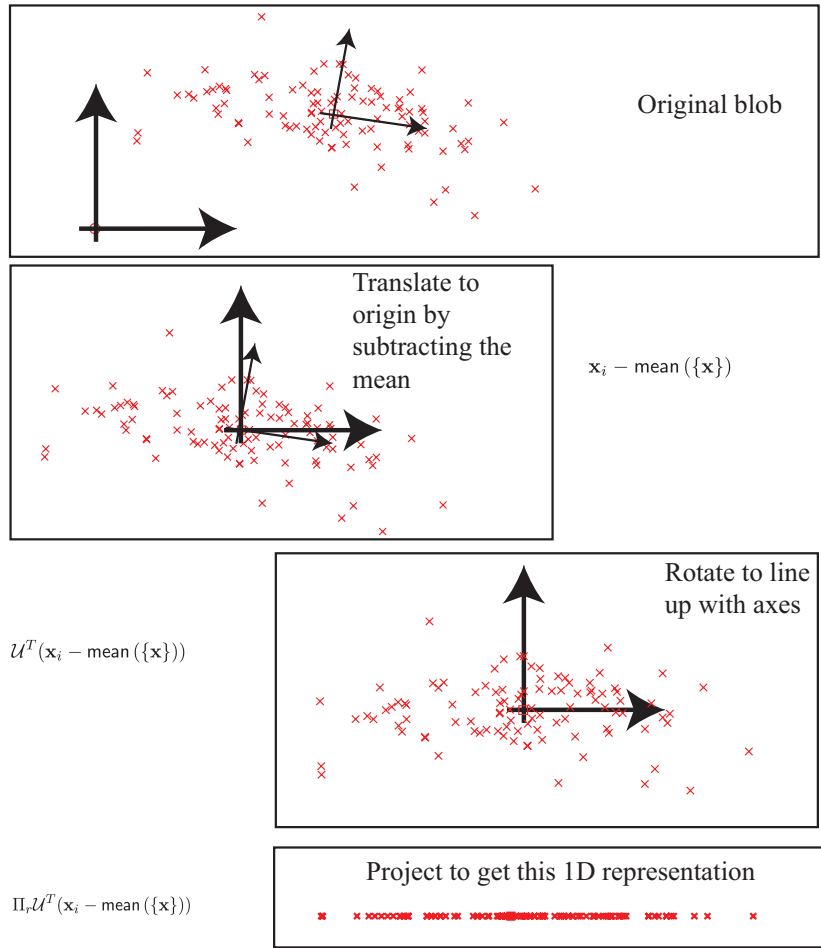


FIGURE 4.18: *Computing a low dimensional representation for principal components analysis.*

we have a new representation of the i 'th data item *in the original space* (Figure 16.20). Now consider the dataset obtained by smoothing each of our data items. We write this dataset as $\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\}$.

You should think of $\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\}$ as a smoothed version of the original data. One way to think of this process is that we have chosen a low-dimensional basis that represents the main variance in the data rather well. It is quite usual to think of a data item as being given by a the mean plus a weighted sum of these basis elements. In this view, the first weight has larger variance than the second, and so on. By construction, this dataset lies in an r dimensional affine subspace of the original space. We constructed this r -dimensional space to preserve the largest variance directions of the data. Each column of this matrix is known as a **principal component**. In particular, we have constructed this dataset so that

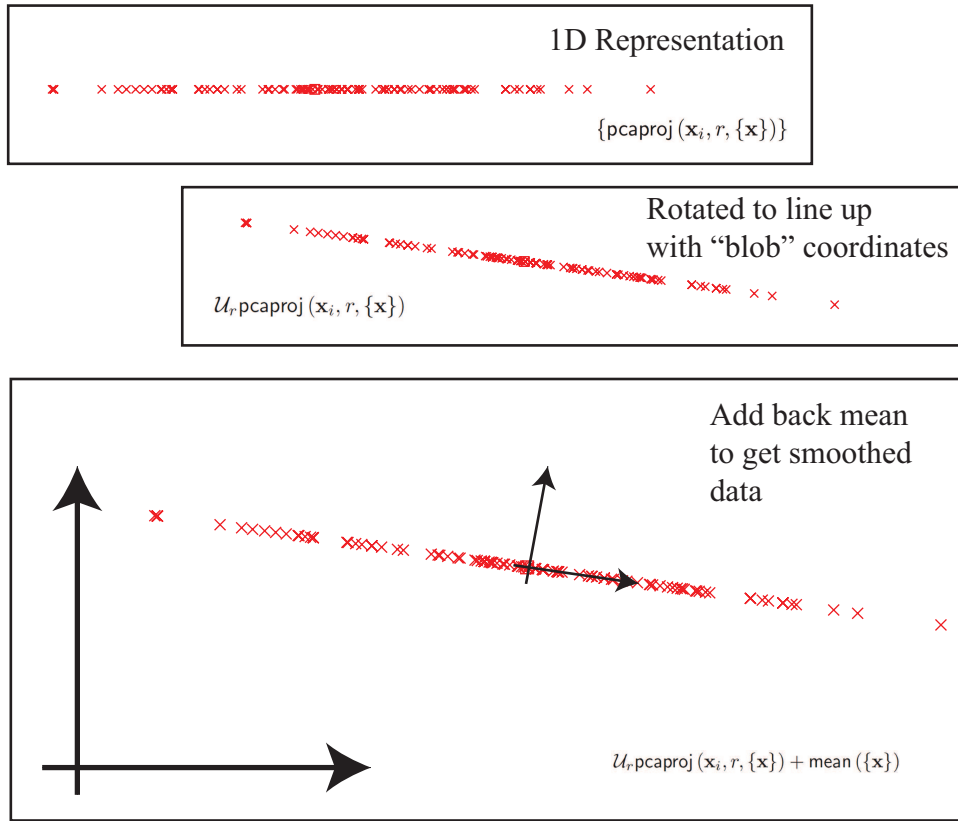


FIGURE 4.19: *Smoothing data with principal components analysis.*

- $\text{mean}(\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\}) = \text{mean}(\{\mathbf{x}\})$;
- $\text{Covmat}(\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\})$ has rank r ;
- $\text{Covmat}(\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\})$ is the best approximation of $\text{Covmat}(\{\mathbf{x}\})$ with rank r .

Figure 4.19 gives a visualization of the smoothing process. By comparing figures 4.14 and 4.20, you can see that a real dataset can lose several dimensions without much significant going wrong. As we shall see in the examples, some datasets can lose many dimensions without anything bad happening.

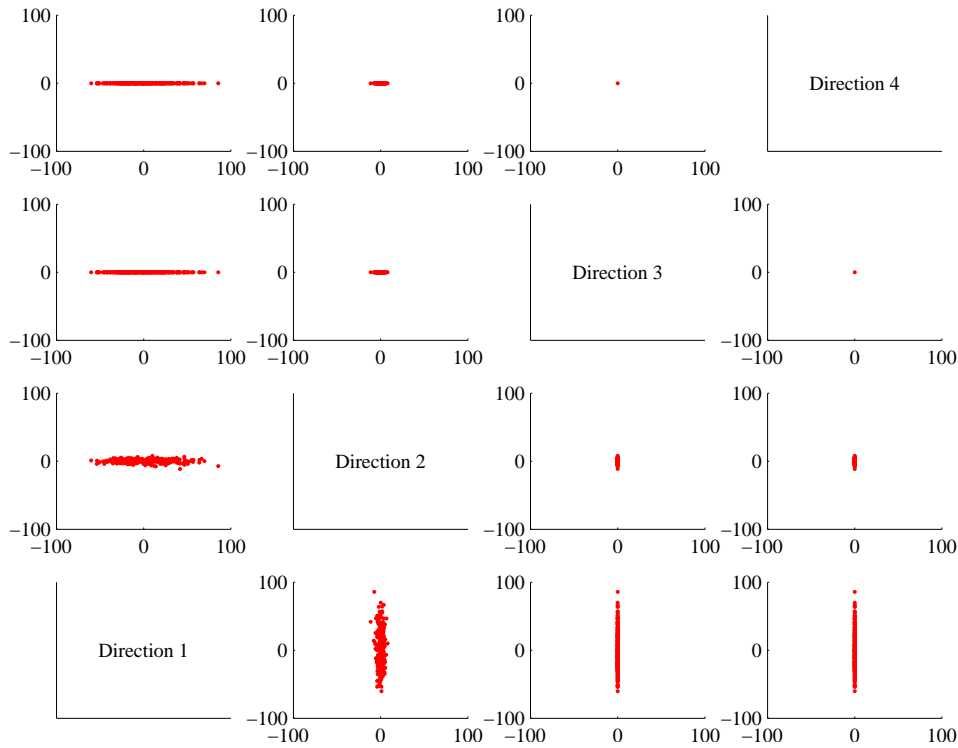


FIGURE 4.20: A panel plot of the bodyfat dataset of figure 4.3, with the dimension reduced to two using principal components analysis. Compare this figure to figure 4.14, which is on the same set of axes. You can see that the blob has been squashed in direction 3 and direction 4. But not much has really happened, because there wasn't very much variation in those directions in the first place.

Procedure: 4.1 Principal Components Analysis

Assume we have a general data set \mathbf{x}_i , consisting of N d -dimensional vectors. Now write $\Sigma = \text{Covmat}(\{\mathbf{x}\})$ for the covariance matrix.

Form \mathcal{U} , Λ , such that $\Sigma\mathcal{U} = \mathcal{U}\Lambda$ (these are the eigenvectors and eigenvalues of Σ). Ensure that the entries of Λ are sorted in decreasing order. Choose r , the number of dimensions you wish to represent. Typically, we do this by plotting the eigenvalues and looking for a “knee” (Figure ??). It is quite usual to do this by hand.

Constructing a low-dimensional representation: Form \mathcal{U}_r , a matrix consisting of the first r columns of \mathcal{U} . Now compute $\{\text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\})\} = \{(\Pi_r \mathcal{U}^T (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})))\}$. This is a set of data vectors which are r dimensional, and where each component is independent of each other component (i.e. the covariances of distinct components are zero).

Smoothing the data: Form $\{\text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\} = \{(\mathcal{U}_r \text{pcproj}(\mathbf{x}_i, r, \{\mathbf{x}\}) + \text{mean}(\{\mathbf{x}\}))\}$. These are d dimensional vectors that lie in a r -dimensional subspace of d -dimensional space. The “missing dimensions” have the lowest variance, and are independent.

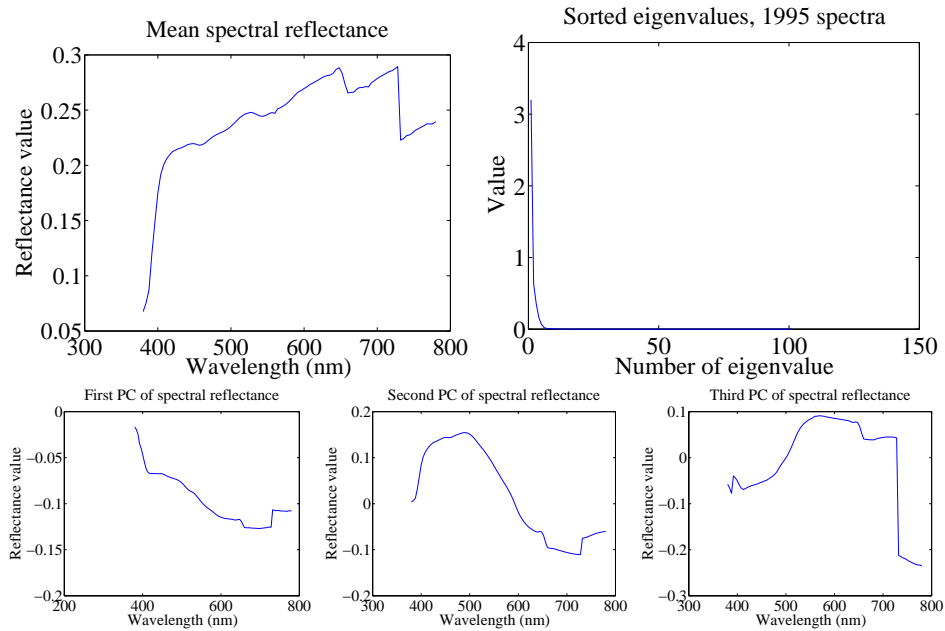


FIGURE 4.21: On the **top left**, the mean spectral reflectance of a dataset of 1995 spectral reflectances, collected by Kobus Barnard (at <http://www.cs.sfu.ca/~colour/data/>). On the **top right**, eigenvalues of the covariance matrix of spectral reflectance data, from a dataset of 1995 spectral reflectances, collected by Kobus Barnard (at <http://www.cs.sfu.ca/~colour/data/>). Notice how the first few eigenvalues are large, but most are very small; this suggests that a good representation using few principal components is available. The **bottom row** shows the first three principal components. A linear combination of these, with appropriate weights, added to the mean of figure ??, gives a good representation of the dataset.

4.4.4 The Error of the Low-Dimensional Representation

We took a dataset, $\{\mathbf{x}\}$, and constructed a d -dimensional dataset $\{\mathbf{b}\}$ in blob coordinates. We did so by translating, then rotating, the data, so no information was lost; we could reconstruct our original dataset by rotating, then translating $\{\mathbf{b}\}$. But in blob coordinates we projected each data item down to the first r components to get an r -dimensional dataset $\{\mathbf{b}_r\}$. We then reconstructed a smoothed dataset by rotating, then translating, $\{\mathbf{b}_r\}$. Information has been lost here, but how much?

The answer is easy to get if you recall that rotations and translations do not change lengths. This means that

$$\|\mathbf{x}_i - \text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\|^2 = \|\mathbf{b}_i - \mathbf{b}_{r,i}\|^2.$$

This expression is easy to evaluate, because \mathbf{b}_i and $\mathbf{b}_{r,i}$ agree in their first r com-

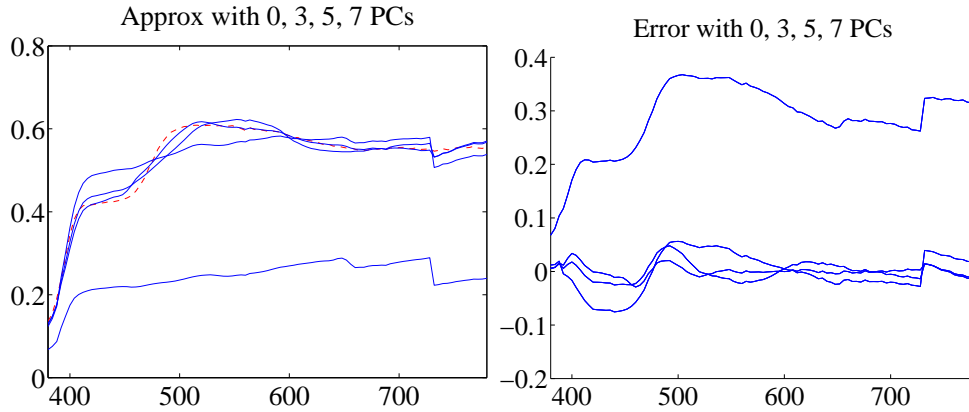


FIGURE 4.22: On the **left**, a spectral reflectance curve (dashed) and approximations using the mean, the mean and 3 principal components, the mean and 5 principal components, and the mean and 7 principal components. Notice the mean is a relatively poor approximation, but as the number of principal components goes up, the error falls rather quickly. On the **right** is the error for these approximations. Figure plotted from a dataset of 1995 spectral reflectances, collected by Kobus Barnard (at <http://www.cs.sfu.ca/~colour/data/>).

ponents. The remaining $d - r$ components of $\mathbf{b}_{r,i}$ are zero. So we can write

$$\|\mathbf{x}_i - \text{pcsmooth}(\mathbf{x}_i, r, \{\mathbf{x}\})\|^2 = \sum_{u=r+1}^d (\mathbf{b}_i^{(u)})^2.$$

Now a natural measure of error is the average over the dataset of this term. We have that

$$\frac{1}{N} \sum_{u=r+1}^d (\mathbf{b}_i^{(u)})^2 = \sum_{u=r+1}^d \text{var}(\{\mathbf{b}^{(u)}\})$$

which is easy to evaluate, because we know these variances — they are the values of the $d - r$ eigenvalues that we decided to ignore. So the mean error can be written as

$$\mathbf{1}^T (\Lambda - \Lambda_r) \mathbf{1}.$$

Now we could choose r by identifying how much error we can tolerate. More usual is to plot the eigenvalues of the covariance matrix, and look for a “knee”, like that in Figure 16.20. You can see that the sum of remaining eigenvalues is small.

4.4.5 Example: Representing Spectral Reflectances

Diffuse surfaces reflect light uniformly in all directions. Examples of diffuse surfaces include matte paint, many styles of cloth, many rough materials (bark, cement, stone, etc.). One way to tell a diffuse surface is that it does not look brighter (or darker) when you look at it along different directions. Diffuse surfaces can

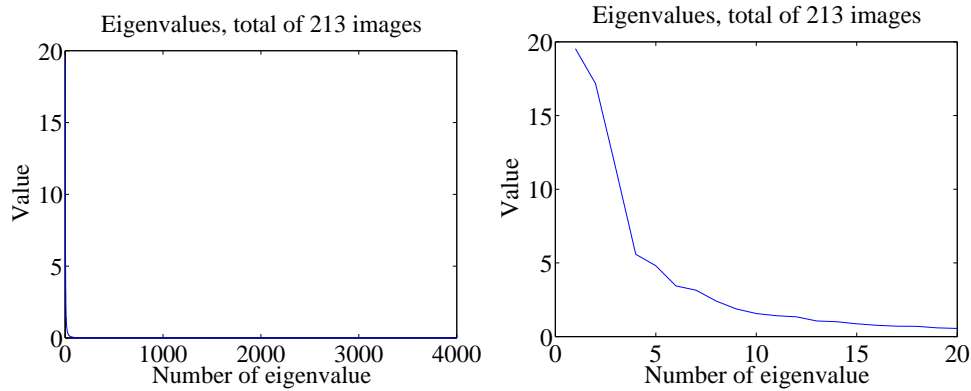


FIGURE 4.23: On the **left**, the eigenvalues of the covariance of the Japanese facial expression dataset; there are 4096, so it's hard to see the curve (which is packed to the left). On the **right**, a zoomed version of the curve, showing how quickly the values of the eigenvalues get small.

be colored, because the surface reflects different fractions of the light falling on it at different wavelengths. This effect can be represented by measuring the spectral reflectance of a surface, which is the fraction of light the surface reflects as a function of wavelength. This is usually measured in the visual range of wavelengths (about 380nm to about 770 nm). Typical measurements are every few nm, depending on the measurement device. I obtained data for 1995 different surfaces from <http://www.cs.sfu.ca/~colour/data/> (there are a variety of great datasets here, from Kobus Barnard).

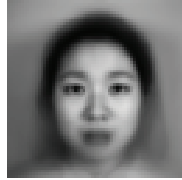
Each spectrum has 101 measurements, which are spaced 4nm apart. This represents surface properties to far greater precision than is really useful. Physical properties of surfaces suggest that the reflectance can't change too fast from wavelength to wavelength. It turns out that very few principal components are sufficient to describe almost any spectral reflectance function. Figure 4.21 shows the mean spectral reflectance of this dataset, and Figure 4.21 shows the eigenvalues of the covariance matrix.

This is tremendously useful in practice. One should think of a spectral reflectance as a function, usually written $\rho(\lambda)$. What the principal components analysis tells us is that we can represent this function rather accurately on a (really small) finite dimensional basis. This basis is shown in figure 4.21. This means that there is a mean function $r(\lambda)$ and k functions $\phi_m(\lambda)$ such that, for any $\rho(\lambda)$,

$$\rho(\lambda) = r(\lambda) + \sum_{i=1}^k c_i \phi_i(\lambda) + e(\lambda)$$

where $e(\lambda)$ is the error of the representation, which we know is small (because it consists of all the other principal components, which have tiny variance). In the case of spectral reflectances, using a value of k around 3-5 works fine for most applications (Figure 4.22). This is useful, because when we want to predict what

Mean image from Japanese Facial Expression dataset



First sixteen principal components of the Japanese Facial Expression dat



FIGURE 4.24: *The mean and first 16 principal components of the Japanese facial expression dataset.*

a particular object will look like under a particular light, we don't need to use a detailed spectral reflectance model; instead, it's enough to know the c_i for that object. This comes in useful in a variety of rendering applications in computer graphics. It is also the key step in an important computer vision problem, called **color constancy**. In this problem, we see a picture of a world of colored objects under unknown colored lights, and must determine what color the objects are. Modern color constancy systems are quite accurate, even though the problem sounds unconstrained. This is because they are able to exploit the fact that relatively few c_i are enough to accurately describe a surface reflectance.

4.4.6 Example: Representing Faces with Principal Components

An image is usually represented as an array of values. We will consider intensity images, so there is a single intensity value in each cell. You can turn the image into a vector by rearranging it, for example stacking the columns onto one another

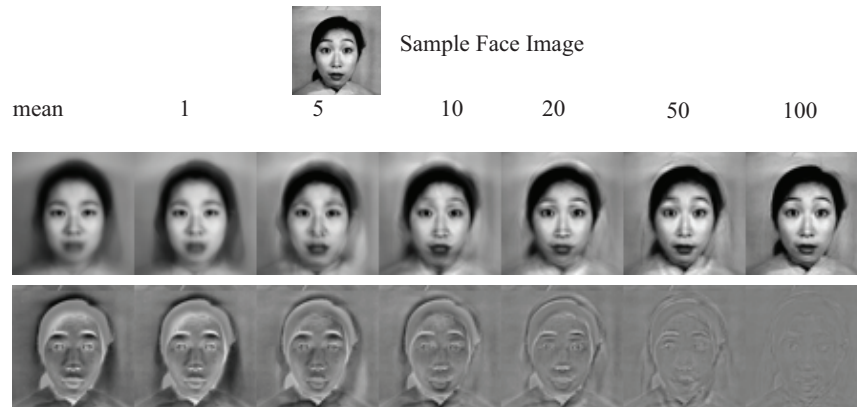


FIGURE 4.25: *Approximating a face image by the mean and some principal components; notice how good the approximation becomes with relatively few components.*

(use `reshape` in Matlab). This means you can take the principal components of a set of images. Doing so was something of a fashionable pastime in computer vision for a while, though there are some reasons that this is not a great representation of pictures. However, the representation yields pictures that can give great intuition into a dataset.

Figure ?? shows the mean of a set of face images encoding facial expressions of Japanese women (available at <http://www.kasrl.org/jaffe.html>; there are tons of face datasets at <http://www.face-rec.org/databases/>). I reduced the images to 64×64 , which gives a 4096 dimensional vector. The eigenvalues of the covariance of this dataset are shown in figure 4.23; there are 4096 of them, so it's hard to see a trend, but the zoomed figure suggests that the first couple of hundred contain most of the variance. Once we have constructed the principal components, they can be rearranged into images; these images are shown in figure 4.24. Principal components give quite good approximations to real images (figure 4.25).

The principal components sketch out the main kinds of variation in facial expression. Notice how the mean face in Figure 4.24 looks like a relaxed face, but with fuzzy boundaries. This is because the faces can't be precisely aligned, because each face has a slightly different shape. The way to interpret the components is to remember one adjusts the mean towards a data point by adding (or subtracting) some scale times the component. So the first few principal components have to do with the shape of the haircut; by the fourth, we are dealing with taller/shorter faces; then several components have to do with the height of the eyebrows, the shape of the chin, and the position of the mouth; and so on. These are all images of women who are not wearing spectacles. In face pictures taken from a wider set of models, moustaches, beards and spectacles all typically appear in the first couple of dozen principal components.

4.5 HIGH DIMENSIONS, SVD AND NIPALS

If you remember the curse of dimension, you should have noticed something of a problem in my account of PCA. When I described the curse, I said one consequence was that forming a covariance matrix for high dimensional data is hard or impossible. Then I described PCA as a method to understand the important dimensions in high dimensional datasets. But PCA appears to rely on covariance, so I should not be able to form the principal components in the first place. In fact, we can form principal components without computing a covariance matrix.

4.5.1 Principal Components by SVD

I will now assume the dataset has zero mean, to simplify notation. This is easily achieved. You subtract the mean from each data item at the start, and add the mean back once you've finished smoothing. As usual, we have N data items, each a d dimensional column vector. We will now arrange these into a matrix,

$$\mathcal{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \dots \mathbf{x}_N^T \end{pmatrix}$$

where each *row* of the matrix is a data vector. Now notice that the covariance matrix for this dataset can be formed by constructing $\mathcal{X}^T \mathcal{X}$, so that

$$\text{Covmat}(\{X\}) = \mathcal{X}^T \mathcal{X}$$

and if we form the SVD (see the math notes at the end if you don't remember this) of \mathcal{X} , we have $\mathcal{X} = \mathcal{U} \Sigma \mathcal{V}^T$. But we have $\mathcal{X}^T \mathcal{X} = \mathcal{V} \Sigma^T \Sigma \mathcal{V}^T$ so that

$$\mathcal{V}^T \mathcal{X}^T \mathcal{X} \mathcal{V} = \Sigma^T \Sigma$$

and $\Sigma^T \Sigma$ is diagonal. So we can recover the principal components of the dataset without actually forming the covariance matrix - we just form the SVD of \mathcal{X} .

4.5.2 Just a few Principal Components with NIPALS

For really big datasets, even taking the SVD is hard. Usually, we don't really want to recover all the principal components, because we want to recover a reasonably accurate low dimensional representation of the data. We continue to work with a data matrix \mathcal{X} , whose *rows* are data items. Now assume we wish to recover the first principal component. This means we are seeking a vector \mathbf{u} and a set of N numbers w_i such that $w_i \mathbf{u}$ is a good approximation to \mathbf{x}_i . In particular, we would like the dataset made of $w_i \mathbf{u}$ to encode as much of the variance of the original dataset as possible. Now we can stack the w_i into a column vector \mathbf{w} . The **Frobenius norm** is a term for the matrix norm obtained by summing squared entries of the matrix. We write

$$\|\mathcal{A}\|_F = \sum_{i,j} a_{ij}^2.$$

In the exercises, you will show that the choice of \mathbf{w} and \mathbf{u} minimizes the cost

$$\|\mathcal{X} - \mathbf{w} \mathbf{u}^T\|_F$$

which we can write as

$$C(\mathbf{w}, \mathbf{u}) = \sum_{ij} (x_{ij} - w_i u_j)^2.$$

Now we need to *find* the relevant \mathbf{w} and \mathbf{u} . Notice there is not a unique choice, because the pair $(s\mathbf{w}, (1/s)\mathbf{u})$ works as well as the pair (\mathbf{w}, \mathbf{u}) . We will choose \mathbf{u} such that $\|\mathbf{u}\| = 1$. There is still not a unique choice, because you can flip the signs in \mathbf{u} and \mathbf{w} , but this doesn't matter. The gradient of the cost function is a set of partial derivatives with respect to components of \mathbf{w} and \mathbf{u} . The partial with respect to w_k is

$$\frac{\partial C}{\partial w_k} = \sum_j (x_{kj} - w_k u_j) u_j$$

which can be written in matrix vector form as

$$\nabla_{\mathbf{w}} C = (\mathcal{X} - \mathbf{w}\mathbf{u}^T)\mathbf{u}.$$

Similarly, the partial with respect to u_l is

$$\frac{\partial C}{\partial u_l} = \sum_i (x_{il} - w_i u_l) w_i$$

which can be written in matrix vector form as

$$\nabla_{\mathbf{u}} C = (\mathcal{X}^T - \mathbf{u}\mathbf{w}^T)\mathbf{w}.$$

At the solution, these partial derivatives are zero. This suggests an algorithm. First, assume we have an estimate of \mathbf{u} , say $\mathbf{u}^{(n)}$. Then we could choose the \mathbf{w} that makes the partial wrt \mathbf{w} zero, so

$$\mathbf{w}^{(n+\frac{1}{2})} = \frac{\mathcal{X}\mathbf{u}^{(n)}}{(\mathbf{u}^{(n)})^T \mathbf{u}^{(n)}}.$$

Now we can update the estimate of \mathbf{u} by choosing a value that makes the partial wrt \mathbf{u} zero, using our estimate $\mathbf{w}^{(n+\frac{1}{2})}$, to get

$$\mathbf{u}^{(n+\frac{1}{2})} = \frac{\mathcal{X}^T \mathbf{w}^{(n+\frac{1}{2})}}{(\mathbf{w}^{(n+\frac{1}{2})})^T \mathbf{w}^{(n+\frac{1}{2})}}.$$

We need to rescale to ensure that our estimate of \mathbf{u} has unit length. Write $s = ((\mathbf{w}^{(n+\frac{1}{2})})^T \mathbf{w}^{(n+\frac{1}{2})})^{\frac{1}{2}}$. We get

$$\mathbf{u}^{(n+1)} = \frac{\mathbf{u}^{(n+\frac{1}{2})}}{s}$$

and

$$\mathbf{w}^{(n+1)} = s\mathbf{w}^{(n+\frac{1}{2})}.$$

This iteration can be started by choosing some row of \mathcal{X} as $\mathbf{u}^{(0)}$. You can test for convergence by checking $\|\mathbf{u}^{(n+1)} - \mathbf{u}^{(n)}\|$. If this is small enough, then the algorithm has converged.

To obtain a second principal component, you form $\mathcal{X}^{(1)} = \mathcal{X} - \mathbf{w}\mathbf{u}^T$ and apply the algorithm to that. You can get many principal components like this, but it's not a good way to get all of them (eventually numerical issues mean the estimates are poor). The algorithm is widely known as NIPALS (for Non-linear Iterative Partial Least Squares).

NIPALS is quite forgiving of missing values, though missing values make it hard to use matrix notation. Recall I wrote the cost function as $C(\mathbf{w}, \mathbf{u}) = \sum_{ij} (x_{ij} - w_i u_j)^2$. We change the sum so that it ranges over only the known values, to get

$$C(\mathbf{w}, \mathbf{u}) = \sum_{ij \in \text{known values}} (x_{ij} - w_i u_j)^2$$

then write

$$\frac{\partial C}{\partial w_k} = \sum_{j \in \text{known values for } k} (x_{kj} - w_k u_j) u_j$$

and

$$\frac{\partial C}{\partial u_l} = \sum_{i \in \text{known values for } l} (x_{il} - w_i u_l) w_i.$$

These partial derivatives must be zero at the solution, so we can estimate

$$w_k^{(n+\frac{1}{2})} = \frac{\sum_{j \in \text{known values for } k} x_{kj} u_j^{(n+\frac{1}{2})}}{\sum_{j \in \text{known values for } k} u_j^{(n+\frac{1}{2})}}$$

and

$$u_l^{(n+\frac{1}{2})} = \frac{\sum_{i \in \text{known values for } l} x_{il} w_l^{(n+\frac{1}{2})}}{\sum_{i \in \text{known values for } l} w_l^{(n+\frac{1}{2})}}$$

We then normalize as before.

Procedure: 4.2 *Obtaining some principal components with NIPALS*

4.5.3 Projection and Discriminative Problems

You should remind yourself how Lagrange multipliers work before reading this section.

Principal components analysis identifies the directions in high dimensional space that best describe the dataset. This may not always be what we are looking for. For a simple example, look at figure ???. In this case, we have two classes of data item. Each varies a lot in the x direction and little in the y direction. If we project onto the first principal component, the outcome represents variance in the data well, but suppresses the difference between the two classes. For many applications – for example, classification – we are interested in directions that emphasize the

differences between classes. One very important construction for such directions can be used for classification and for regression.

Assume we have a dataset of N items, each with two parts \mathbf{x}_i and \mathbf{y}_i . We assume that at least \mathbf{x}_i has high dimension. We also assume that $\text{mean}(\{\mathbf{x}\}) = 0$ and $\text{mean}(\{\mathbf{y}\}) = 0$. This simplifies notation, and is easy to achieve (subtract the mean). This is a mild generalization of what occurred in classification, where we had a feature vector \mathbf{x}_i and a label y_i for each data item; but now instead of having just a label, we have a vector. This situation arises in practice quite often. For example, \mathbf{x}_i might be a vector that describes an image and \mathbf{y}_i might be a vector that describes a caption for that image. If we have a classification problem with C classes, we might choose \mathbf{y}_i to be a vector with zero in all components except the one corresponding to the class (this is sometimes known as a **one-hot vector**). We wish to choose projections of \mathbf{x} and \mathbf{y} to a shared low dimensional space, so that the projections in this space are strongly correlated to one another.

For the moment, assume the low dimensional space is one dimensional. Then there is some unit vector \mathbf{a} so that projecting \mathbf{x}_i to that space is given by $\mathbf{a}^T \mathbf{x}_i$; similarly, there is some unit vector \mathbf{b} so that projecting \mathbf{y}_i to the space is given by $\mathbf{b}^T \mathbf{y}_i$. Now we stack the vectors into data matrices whose *rows* are data items, as above, so the i 'th row of \mathcal{X} is \mathbf{x}_i^T and the i 'th row of \mathcal{Y} is \mathbf{y}_i^T . Then $\mathbf{p}_x = \text{mat} \mathcal{X} \mathbf{a}$ is a vector containing all the projections of the x part, and $\mathbf{p}_y = \mathcal{Y} \mathbf{b}$ is a vector containing all the projections of the y part. We want these projections to be "like" one another to the extent possible.

One criterion we can use is to maximize $\mathbf{p}_x^T \mathbf{p}_y$ by choice of \mathbf{a} and \mathbf{b} . We must maximise

$$\mathbf{a}^T \max X^T \mathcal{Y} \mathbf{b} \text{ subject to } \mathbf{a}^T \mathbf{a} = 1 \text{ and } \mathbf{b}^T \mathbf{b} = 1.$$

Write λ_a and λ_b for Lagrange multipliers. To solve this problem, we must find \mathbf{a} and \mathbf{b} such that

$$\mathcal{X}^T \mathcal{Y} \mathbf{b} = \lambda_a \mathbf{a} \text{ and } \mathcal{Y}^T \mathcal{X} \mathbf{a} = \lambda_b \mathbf{b}.$$

We can substitute using the second equation to get

$$(\mathcal{X}^T \mathcal{Y} \mathcal{Y}^T \mathcal{X}) \mathbf{a} = \lambda_a \lambda_b \mathbf{a}$$

which means that \mathbf{a} is an eigenvector of $(\mathcal{X}^T \mathcal{Y} \mathcal{Y}^T \mathcal{X})$. It turns out we seek the eigenvector corresponding to the largest eigenvalue; equivalently, we must solve

$$\mathbf{a}^T (\mathcal{X}^T \mathcal{Y} \mathcal{Y}^T \mathcal{X}) \mathbf{a} \text{ subject to } \mathbf{a}^T \mathbf{a} = 1$$

There is a straightforward way to obtain a second, third, etc. dimension. We take the dataset, and subtract the \mathbf{a} component from each \mathbf{x}_i and the \mathbf{b} component from each \mathbf{y}_i ; we now have a new dataset, and seek another \mathbf{a} , \mathbf{b} using our procedure. These new directions must (a) maximise the covariance we are interested in and (b) are orthogonal to the original directions. As Figure ?? suggests, these projections of the data separate blobs of data with different labels.

****PROBLEM FIGURE

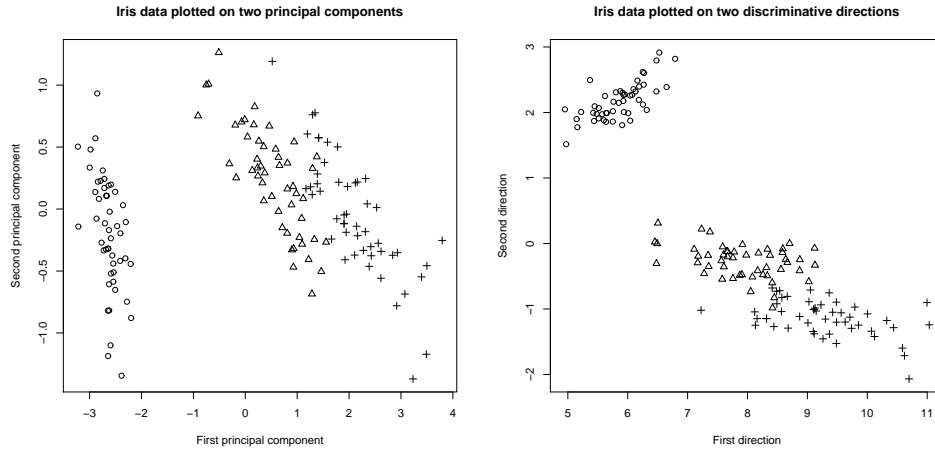


FIGURE 4.26:

4.5.4 Just a few Discriminative Directions with PLS1

Notice that we have, again, formed a (kind of) covariance matrix here, meaning that this method might not apply if the dimension of the data is big. But we could obtain \mathbf{a} from an SVD. An SVD yields $\mathcal{U}\Sigma\mathcal{V}^T = \mathcal{Y}^T\mathcal{X}$, where \mathcal{U} and \mathcal{V} are orthonormal and Σ is diagonal. In turn, we have

$$\mathbf{a}^T(\mathcal{X}^T\mathcal{Y}\mathcal{Y}^T\mathcal{X})\mathbf{a} = \mathbf{a}^T\mathcal{V}\Sigma^2\mathcal{V}^T\mathbf{a}$$

so that \mathbf{a} must be the column of \mathcal{V} corresponding to the largest singular value. Now the argument I used for extracting principal components works for this problem, too. Recall I was looking for the unit \mathbf{a} that maximized $\mathbf{a}^T\mathcal{X}^T\mathcal{X}\mathbf{a}$; I showed I could obtain this from an SVD of \mathcal{X} ; then I argued that I could recover \mathbf{a} , \mathbf{w} such that $\|\mathcal{X} - \mathbf{w}\mathbf{a}\|_F^2$ was minimized.

When we are looking for a discriminative direction, we would look for \mathbf{a} , \mathbf{w} such that $\|\mathcal{Y}^T\mathcal{X} - \mathbf{a}\mathbf{w}^T\|_F^2$ is minimized. The procedure above applies. This algorithm is usually called **PLS1** (for partial least squares one).

4.6 MULTI-DIMENSIONAL SCALING

One way to get insight into a dataset is to plot it. But choosing what to plot for a high dimensional dataset could be difficult. Assume we must plot the dataset in two dimensions (by far the most common choice). We wish to build a scatter plot in two dimensions — but where should we plot each data point? One natural requirement is that the points be laid out in two dimensions in a way that reflects how they sit in many dimensions. In particular, we would like points that are far apart in the high dimensional space to be far apart in the plot, and points that are close in the high dimensional space to be close in the plot.

4.6.1 Principal Coordinate Analysis

We will plot the high dimensional point \mathbf{x}_i at \mathbf{v}_i , which is a two-dimensional vector. Now the squared distance between points i and j in the high dimensional space is

$$D_{ij}^{(2)}(\mathbf{x}) = (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j)$$

(where the superscript is to remind you that this is a squared distance). We could build an $N \times N$ matrix of squared distances, which we write $\mathcal{D}^{(2)}(\mathbf{x})$. The i, j 'th entry in this matrix is $D_{ij}^{(2)}(\mathbf{x})$, and the \mathbf{x} argument means that the distances are between points in the high-dimensional space. Now we could choose the \mathbf{v}_i to make

$$\sum_{ij} \left(D_{ij}^{(2)}(\mathbf{x}) - D_{ij}^{(2)}(\mathbf{v}) \right)^2$$

as small as possible. Doing so should mean that points that are far apart in the high dimensional space are far apart in the plot, and that points that are close in the high dimensional space are close in the plot.

In its current form, the expression is difficult to deal with, but we can refine it. Because translation does not change the distances between points, it cannot change either of the $\mathcal{D}^{(2)}$ matrices. So it is enough to solve the case when the mean of the points \mathbf{x}_i is zero. We can assume that $\frac{1}{N} \sum_i \mathbf{x}_i = \mathbf{0}$. Now write $\mathbf{1}$ for the n -dimensional vector containing all ones, and \mathcal{I} for the identity matrix. Notice that

$$D_{ij}^{(2)} = (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_i - 2\mathbf{x}_i \cdot \mathbf{x}_j + \mathbf{x}_j \cdot \mathbf{x}_j.$$

Now write

$$\mathcal{A} = \left[\mathcal{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^T \right].$$

Using this expression, you can show that the matrix \mathcal{M} , defined below,

$$\mathcal{M}(\mathbf{x}) = -\frac{1}{2} \mathcal{A} \mathcal{D}^{(2)}(\mathbf{x}) \mathcal{A}^T$$

has i, j th entry $\mathbf{x}_i \cdot \mathbf{x}_j$ (exercises). I now argue that, to make $\mathcal{D}^{(2)}(\mathbf{v})$ is close to $\mathcal{D}^{(2)}(\mathbf{x})$, it is enough to make $\mathcal{M}(\mathbf{v})$ close to $\mathcal{M}(\mathbf{x})$. Proving this will take us out of our way unnecessarily, so I omit a proof.

We can choose a set of \mathbf{v}_i that makes $\mathcal{D}^{(2)}(\mathbf{v})$ close to $\mathcal{D}^{(2)}(\mathbf{x})$ quite easily, using the method of the previous section. Take the dataset of N d -dimensional column vectors \mathbf{x}_i , and form a matrix \mathcal{X} by stacking the vectors, so

$$\mathcal{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N].$$

In this notation, we have

$$\mathcal{M}(\mathbf{x}) = \mathcal{X}^T \mathcal{X}.$$

This matrix is symmetric, and it is positive semidefinite. It can't be positive definite, because the data is zero mean, so $\mathcal{M}(\mathbf{x})\mathbf{1} = 0$. The $\mathcal{M}(\mathbf{v})$ we seek must (a) be as close as possible to $\mathcal{M}(\mathbf{x})$ and (b) have rank 2. It must have rank 2 because

there must be some \mathcal{V} which is $2 \times N$ so that $\mathcal{M}(\mathbf{v}) = \mathcal{V}^T \mathcal{V}$. The columns of this \mathcal{V} are our \mathbf{v}_i .

We can use the method of section 17.1.2 to construct $\mathcal{M}(\mathbf{v})$ and \mathcal{V} . As usual, we write \mathcal{U} for the matrix of eigenvectors of $\mathcal{M}(\mathbf{x})$, Λ for the diagonal matrix of eigenvalues sorted in descending order, Λ_2 for the 2×2 upper left hand block of Λ , and $\Lambda_2^{(1/2)}$ for the matrix of positive square roots of the eigenvalues. Then our methods yield

$$\mathcal{M}(\mathbf{v}) = \mathcal{U}_2 \Lambda_2^{(1/2)} \Lambda_2^{(1/2)} \mathcal{U}_2^T$$

and

$$\mathcal{V} = \Lambda_2^{(1/2)} \mathcal{U}_2^T$$

and we can plot these \mathbf{v}_i (example in section 16.20). This method for constructing a plot is known as **principal coordinate analysis**.

This plot might not be perfect, because reducing the dimension of the data points should cause some distortions. In many cases, the distortions are tolerable. In other cases, we might need to use a more sophisticated scoring system that penalizes some kinds of distortion more strongly than others. There are many ways to do this; the general problem is known as **multidimensional scaling**.

Procedure: 4.3 *Principal Coordinate Analysis*

Assume we have a matrix $D^{(2)}$ consisting of the squared differences between each pair of N points. We do not need to know the points. We wish to compute a set of points in r dimensions, such that the distances between these points are as similar as possible to the distances in $D^{(2)}$. Form $\mathcal{A} = [\mathcal{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^T]$. Form $\mathcal{W} = \frac{1}{2} \mathcal{A} D^{(2)} \mathcal{A}^T$. Form \mathcal{U} , Λ , such that $\mathcal{W}\mathcal{U} = \mathcal{U}\Lambda$ (these are the eigenvectors and eigenvalues of \mathcal{W}). Ensure that the entries of Λ are sorted in decreasing order. Choose r , the number of dimensions you wish to represent. Form Λ_r , the top left $r \times r$ block of Λ . Form $\Lambda_r^{(1/2)}$, whose entries are the positive square roots of Λ_r . Form \mathcal{U}_r , the matrix consisting of the first r columns of \mathcal{U} . Then $\mathcal{V} = \Lambda_r^{(1/2)} \mathcal{U}_r^T = [\mathbf{v}_1, \dots, \mathbf{v}_N]$ is the set of points to plot.

4.6.2 Example: Mapping with Multidimensional Scaling

Multidimensional scaling gets positions (the \mathcal{V} of section 4.6.1) from distances (the $\mathcal{D}^{(2)}(\mathbf{x})$ of section 4.6.1). This means we can use the method to build maps from distances alone. I collected distance information from the web (I used <http://www.distancefromto.net>, but a google search on “city distances” yields a wide range of possible sources), then apply multidimensional scaling. Table ?? shows distances between the South African provincial capitals, in kilometers, rounded to the nearest kilometer. I then used principal coordinate analysis to find positions for each capital, and rotated, translated and scaled the resulting plot to check it against a real map (Figure 4.27).

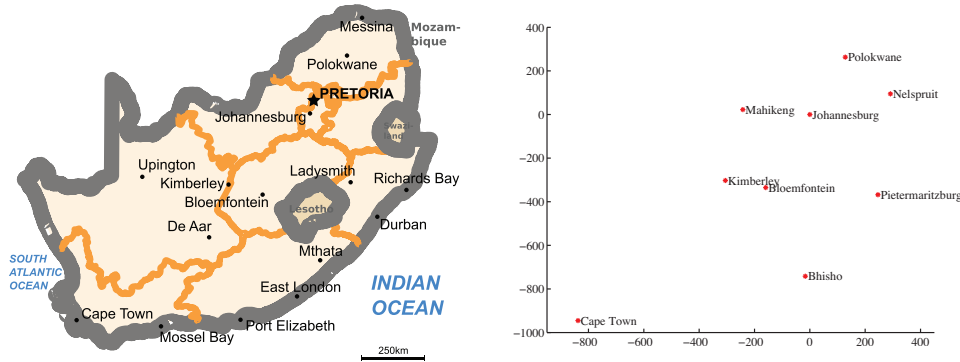


FIGURE 4.27: On the left, a public domain map of South Africa, obtained from http://commons.wikimedia.org/wiki/File:Map_of_South_Africa.svg, and edited to remove surrounding countries. On the right, the locations of the cities inferred by multidimensional scaling, rotated, translated and scaled to allow a comparison to the map by eye. The map doesn't have all the provincial capitals on it, but it's easy to see that MDS has placed the ones that are there in the right places (use a piece of ruled tracing paper to check).

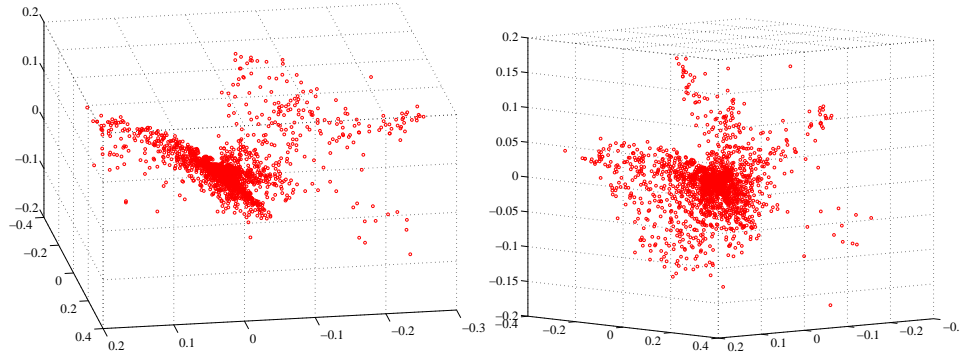


FIGURE 4.28: Two views of the spectral data of section 4.4.5, plotted as a scatter plot by applying principal coordinate analysis to obtain a 3D set of points. Notice that the data spreads out in 3D, but seems to lie on some structure; it certainly isn't a single blob. This suggests that further investigation would be fruitful.

One natural use of principal coordinate analysis is to see if one can spot any structure in a dataset. Does the dataset form a blob, or is it clumpy? This isn't a perfect test, but it's a good way to look and see if anything interesting is happening. In figure 4.28, I show a 3D plot of the spectral data, reduced to three dimensions using principal coordinate analysis. The plot is quite interesting. You should notice that the data points are spread out in 3D, but actually seem to lie on a complicated curved surface — they very clearly don't form a uniform blob. To me, the structure looks somewhat like a butterfly. I don't know why this occurs, but it certainly suggests that something worth investigating is going on. Perhaps the

choice of samples that were measured is funny; perhaps the measuring instrument doesn't make certain kinds of measurement; or perhaps there are physical processes that prevent the data from spreading out over the space.

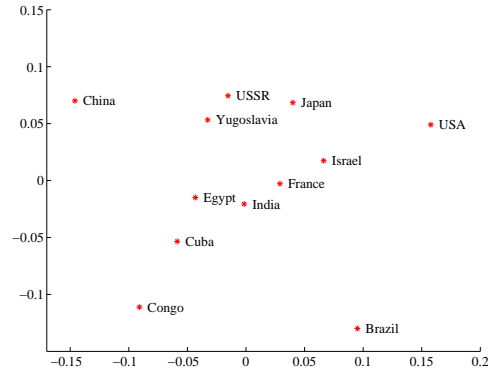


FIGURE 4.29: A map of country similarity, prepared from the data of figure ???. The map is often interpreted as showing a variation in development or wealth (poorest at bottom left to richest at top right); and freedom (most repressed at top left and freest at bottom right). I haven't plotted these axes, because the interpretation wouldn't be consistent with current intuition (the similarity data is forty years old, and quite a lot has happened in that time).

Our algorithm has one really interesting property. In some cases, we do not actually know the datapoints as vectors. Instead, we *just* know distances between the datapoints. This happens often in the social sciences, but there are important cases in computer science as well. As a rather contrived example, one could survey people about breakfast foods (say, eggs, bacon, cereal, oatmeal, pancakes, toast, muffins, kippers and sausages for a total of 9 items). We ask each person to rate the similarity of each pair of distinct items on some scale. We advise people that similar items are ones where, if they were offered both, they would have no particular preference; but, for dissimilar items, they would have a strong preference for one over the other. The scale might be “very similar”, “quite similar”, “similar”, “quite dissimilar”, and “very dissimilar” (scales like this are often called **Likert scales**). We collect these similarities from many people for each pair of distinct items, and then average the similarity over all respondents. We compute distances from the similarities in a way that makes very similar items close and very dissimilar items distant. Now we have a table of distances between items, and can compute a \mathbf{V} and produce a scatter plot. This plot is quite revealing, because items that most people think are easily substituted appear close together, and items that are hard to substitute are far apart. The neat trick here is that we did not start with a \mathcal{X} , but with just a set of distances; but we were able to associate a vector with “eggs”, and produce a meaningful plot.

Table ?? shows data from one such example. Students were interviewed (in 1971! things may have changed since then) about their perceptions of the similarity of countries. The averaged perceived similarity is shown in table ??. Large numbers

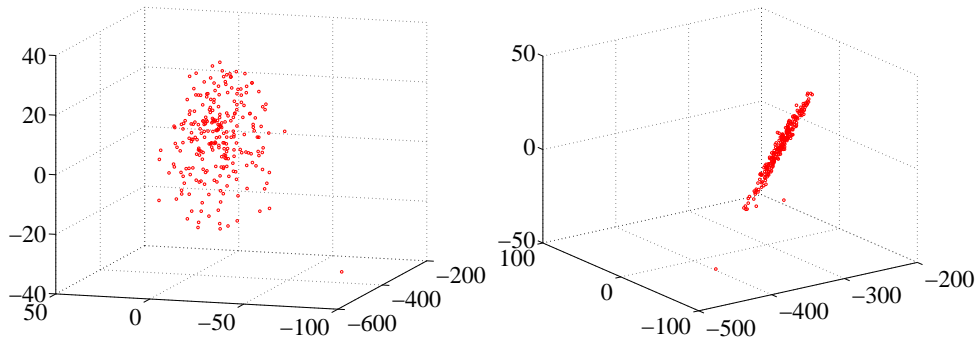


FIGURE 4.30: *Two views of a multidimensional scaling to three dimensions of the height-weight dataset. Notice how the data seems to lie in a flat structure in 3D, with one outlying data point. This means that the distances between data points can be (largely) explained by a 2D representation.*

reflect high similarity, so we can't use these numbers directly. It is reasonable to turn these numbers into distances by (a) using 0 as the distance between a country and itself and (b) using $e^{-s_{ij}}$ as the distance between countries i and j (where s_{ij} is the similarity between them). Once we have distances, we can apply the procedure of section 4.6.1 to get points, then plot a scatter plot (Figure 4.29).

4.7 EXAMPLE: UNDERSTANDING HEIGHT AND WEIGHT

Recall the height-weight data set of section ?? (from <http://www2.stetson.edu/~jrasp/data.htm>; look for bodyfat.xls at that URL). This is, in fact, a 16-dimensional dataset. The entries are (in this order): *bodyfat; density; age; weight; height; adiposity; neck; chest; abdomen; hip; thigh; knee; ankle; biceps; forearm; wrist*. We know already that many of these entries are correlated, but it's hard to grasp a 16 dimensional dataset in one go. The first step is to investigate with a multidimensional scaling.

Figure ?? shows a multidimensional scaling of this dataset down to three dimensions. The dataset seems to lie on a (fairly) flat structure in 3D, meaning that inter-point distances are relatively well explained by a 2D representation. Two points seem to be special, and lie far away from the flat structure. The structure isn't perfectly flat, so there will be small errors in a 2D representation; but it's clear that a lot of dimensions are redundant. Figure 4.31 shows a 2D representation of these points. They form a blob that is stretched along one axis, and there is no sign of multiple blobs. There's still at least one special point, which we shall ignore but might be worth investigating further. The distortions involved in squashing this dataset down to 2D seem to have made the second special point less obvious than it was in figure ??.

The next step is to try a principal component analysis. Figure 4.32 shows the mean of the dataset. The components of the dataset have different units, and shouldn't really be compared. But it is difficult to interpret a table of 16 numbers, so I have plotted the mean by showing a vertical bar for each component. Figure 4.33

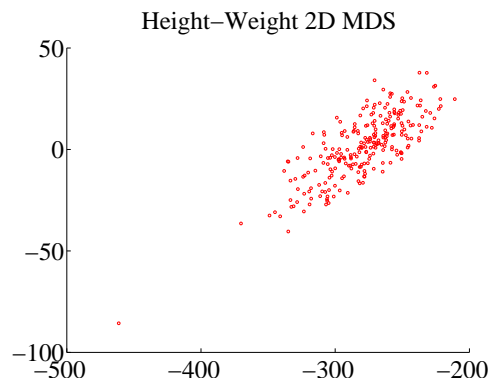


FIGURE 4.31: A multidimensional scaling to two dimensions of the height-weight dataset. One data point is clearly special, and another looks pretty special. The data seems to form a blob, with one axis quite a lot more important than another.

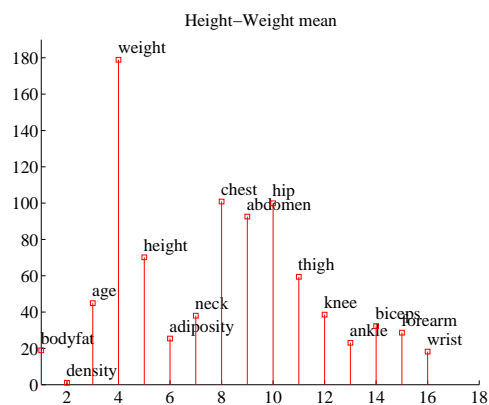


FIGURE 4.32: The mean of the *bodyfat.xls* dataset. Each component is likely in a different unit (though I don't know the units), making it difficult to plot the data without being misleading. I've adopted one solution here, by plotting each component as a vertical bar, and labelling the bar. You shouldn't try to compare the values to one another. Instead, think of this plot as a compact version of a table.

shows the eigenvalues of the covariance for this dataset. Notice how one dimension is very important, and after the third principal component, the contributions become small. Of course, I could have said “fourth”, or “fifth”, or whatever — the precise choice depends on how small a number you think is “small”.

Figure 4.33 also shows the first principal component. The eigenvalues justify thinking of each data item as (roughly) the mean plus some weight times this principal component. From this plot you can see that data items with a larger value of *weight* will also have larger values of most other measurements, except *age* and *density*. You can also see how much larger; if the weight goes up by 8.5 units, then the abdomen will go up by 3 units, and so on. This explains the main variation

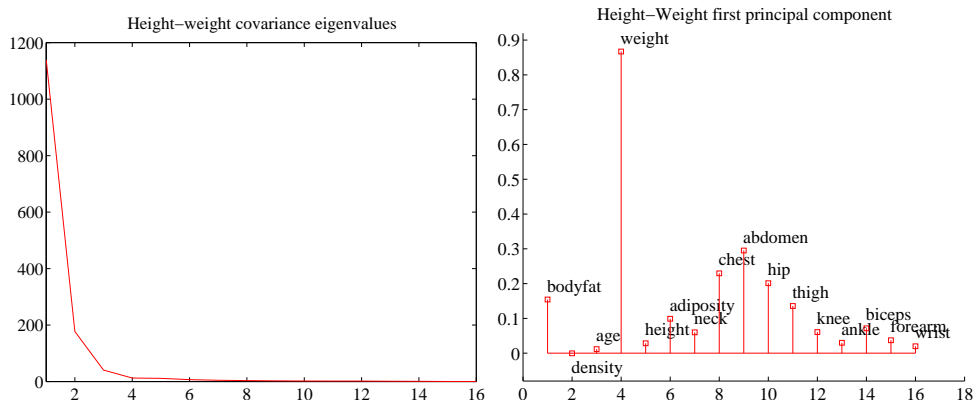


FIGURE 4.33: On the **left**, the eigenvalues of the covariance matrix for the bodyfat data set. Notice how fast the eigenvalues fall off; this means that most principal components have very small variance, so that data can be represented well with a small number of principal components. On the **right**, the first principal component for this dataset, plotted using the same convention as for figure 4.32.

in the dataset.

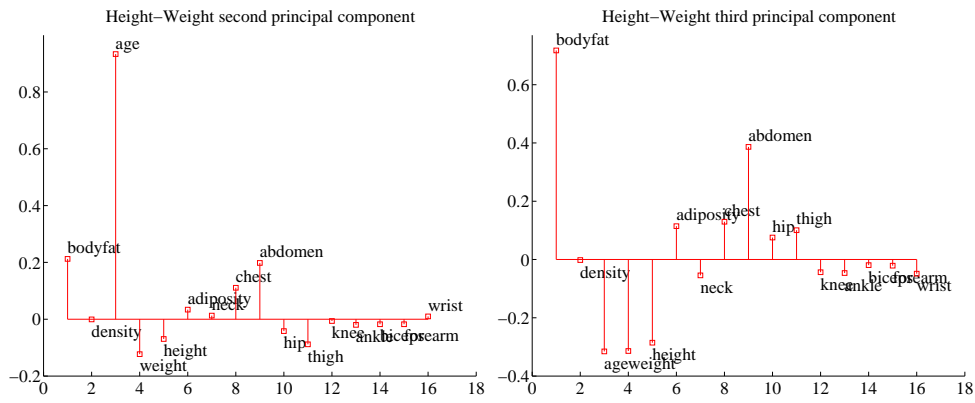


FIGURE 4.34: On the **left**, the second principal component, and on the **right** the third principal component of the height-weight dataset.

In the rotated coordinate system, the components are not correlated, and they have different variances (which are the eigenvalues of the covariance matrix). You can get some sense of the data by adding these variances; in this case, we get 1404. This means that, in the translated and rotated coordinate system, the average data point is about $37 = \sqrt{1404}$ units away from the center (the origin). Translations and rotations do not change distances, so the average data point is about 37 units from the center in the original dataset, too. If we represent a datapoint by using the mean and the first three principal components, there will be some error. We

can estimate the average error from the component variances. In this case, the sum of the first three eigenvalues is 1357, so the mean square error in representing a datapoint by the first three principal components is $\sqrt{(1404 - 1357)}$, or 6.8. The relative error is $6.8/37 = 0.18$. Another way to represent this information, which is more widely used, is to say that the first three principal components explain all but $(1404 - 1357)/1404 = 0.034$, or 3.4% of the variance; notice that this is the square of the relative error, which will be a much smaller number.

All this means that explaining a data point as the mean and the first three principal components produces relatively small errors. Figure 4.35 shows the second and third principal component of the data. These two principal components suggest some further conclusions. As *age* gets larger, *height* and *weight* get slightly smaller, but the weight is redistributed; *abdomen* gets larger, whereas *thigh* gets smaller. A smaller effect (the third principal component) links *bodyfat* and *abdomen*. As *bodyfat* goes up, so does *abdomen*.

4.8 WHAT YOU SHOULD REMEMBER - NEED SOMETHING

PROBLEMS

Summaries

- 4.1. You have a dataset $\{\mathbf{x}\}$ of N vectors, \mathbf{x}_i , each of which is d -dimensional. We will consider a linear function of this dataset. Write \mathbf{a} for a constant vector; then the value of this linear function evaluated on the i 'th data item is $\mathbf{a}^T \mathbf{x}_i$. Write $f_i = \mathbf{a}^T \mathbf{x}_i$. We can make a new dataset $\{f\}$ out of the values of this linear function.
- Show that $\text{mean}(\{f\}) = \mathbf{a}^T \text{mean}(\{\mathbf{x}\})$ (easy).
 - Show that $\text{var}(\{f\}) = \mathbf{a}^T \text{Covmat}(\{\mathbf{x}\}) \mathbf{a}$ (harder, but just push it through the definition).
 - Assume the dataset has the special property that there exists some \mathbf{a} so that $\mathbf{a}^T \text{Covmat}(\{\mathbf{x}\}) \mathbf{a} = 0$. Show that this means that the dataset lies on a hyperplane.

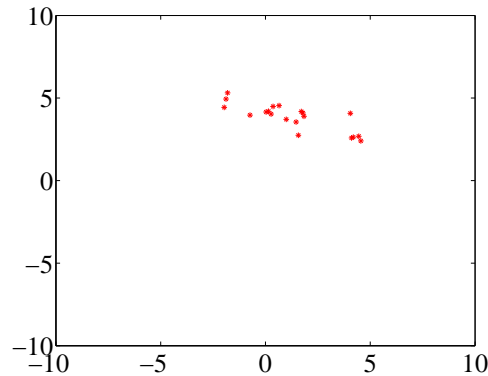


FIGURE 4.35: Figure for the question

- 4.2. On Figure 4.35, mark the mean of the dataset, the first principal component,

Listing 4.1: R code for iris example.

```

# r code for scatterplot of iris data
irisdat<-read.csv('iris.dat', header=FALSE);
library('lattice')
numiris=irisdat[,c(1, 2, 3, 4)]
postscript("irisscatterplot.eps")
# so that I get a postscript file
speciesnames<-c('setosa', 'versicolor', 'virginica')
pchr<-c(1, 2, 3)
colr<-c('red', 'green', 'blue', 'yellow', 'orange')
ss<-expand.grid(species=1:3)
parset<-with(ss, simpleTheme(pch=pchr[species],
                             col=colr[species]))
splom(irisdat[, c(1:4)], groups=irisdat$V5,
      par.settings=parset,
      varnames=c('Sepal\nLength', 'Sepal\nWidth',
                 'Petal\nLength', 'Petal\nWidth'),
      key=list(text=list(speciesnames),
               points=list(pch=pchr, columns=3))
dev.off()

```

and the second principal component.

- 4.3. You have a dataset $\{\mathbf{x}\}$ of N vectors, \mathbf{x}_i , each of which is d -dimensional. Assume that $\text{Covmat}(\{\mathbf{x}\})$ has one non-zero eigenvalue. Assume that \mathbf{x}_1 and \mathbf{x}_2 do not have the same value.

- (a) Show that you can choose a set of t_i so that you can represent *every* data item \mathbf{x}_i *exactly*

$$\mathbf{x}_i = \mathbf{x}_1 + t_i(\mathbf{x}_2 - \mathbf{x}_1).$$

- (b) Now consider the dataset of these t values. What is the relationship between (a) $\text{std}(t)$ and (b) the non-zero eigenvalue of $\text{Covmat}(\{\mathbf{x}\})$? Why?

PROGRAMMING EXERCISES

- 4.4. Obtain the iris dataset from the UC Irvine machine learning data repository at <http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data>.
- (a) Plot a scatterplot matrix of this dataset, showing each species with a different marker. The fragment of R code in listing 4.1 should take you most of the way.
- (b) Now obtain the first two principal components of the data. Plot the data on those two principal components alone, again showing each species with a different marker. Has this plot introduced significant distortions? Explain
- (c) Now use PLS1 to obtain two discriminative directions, and project the data on to those directions. Does the plot look better? Explain *Keep in mind that the most common error here is to forget that the \mathcal{X} and the \mathcal{Y} in PLS1 are centered - i.e. we subtract the mean.*
- 4.5. Take the wine dataset from the UC Irvine machine learning data repository at <https://archive.ics.uci.edu/ml/datasets/Wine>.
- (a) Plot the eigenvalues of the covariance matrix in sorted order. How many principal components should be used to represent this dataset? Why?

- (b) Construct a stem plot of each of the first 3 principal components (i.e. the eigenvectors of the covariance matrix with largest eigenvalues). What do you see?
 - (c) Compute the first two principal components of this dataset, and project it onto those components. Now produce a scatter plot of this two dimensional dataset, where data items of class 1 are plotted as a '1', class 2 as a '2', and so on.
- 4.6. Take the wheat kernel dataset from the UC Irvine machine learning data repository at <http://archive.ics.uci.edu/ml/datasets/seeds>. Compute the first two principal components of this dataset, and project it onto those components.
- (a) Produce a scatterplot of this projection. Do you see any interesting phenomena?
 - (b) Plot the eigenvalues of the covariance matrix in sorted order. How many principal components should be used to represent this dataset? why?
- 4.7. The UC Irvine machine learning data repository hosts a collection of data on breast cancer diagnostics, donated by Olvi Mangasarian, Nick Street, and William H. Wolberg. You can find this data at [http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)). For each record, there is an id number, 10 continuous variables, and a class (benign or malignant). There are 569 examples. Separate this dataset randomly into 100 validation, 100 test, and 369 training examples.
- (a) Plot this dataset on the first three principal components, using different markers for benign and malignant cases. What do you see?
 - (b) Now use PLS1 to obtain three discriminative directions, and project the data on to those directions. Does the plot look better? Explain *Keep in mind that the most common error here is to forget that the \mathcal{X} and the \mathcal{Y} in PLS1 are centered - i.e. we subtract the mean. Once you have computed the first (etc. direction, you should subtract it from \mathcal{X} , but leave \mathcal{Y} alone.*

CHAPTER 5

Clustering: Models of High Dimensional Data

High-dimensional data comes with problems. Data points tend not to be where you think; they can scattered quite far apart, and can be quite far from the mean. Except in special cases, the only really reliable probability model is the Gaussian (or Gaussian blob, or blob).

There is an important rule of thumb for coping with high dimensional data: **Use simple models.** A blob is a good simple model. Modelling data as a blob involves computing its mean and its covariance. Sometimes, as we shall see, the covariance can be hard to compute. Even so, a blob model is really useful. It is natural to try and extend this model to cover datasets that don't obviously consist of a single blob.

One very good, very simple, model for high dimensional data is to assume that it consists of multiple blobs. To build models like this, we must determine (a) what the blob parameters are and (b) which datapoints belong to which blob. Generally, we will collect together data points that are close and form blobs out of them. This process is known as **clustering**.

Clustering is a somewhat puzzling activity. It is extremely useful to cluster data, and it seems to be quite important to do it reasonably well. But it surprisingly hard to give crisp criteria for a good (resp. bad) clustering of a dataset. Usually, clustering is part of building a model, and the main way to know that the clustering algorithm is bad is that the model is bad.

5.1 AGGLOMERATIVE AND DIVISIVE CLUSTERING

There are two natural algorithms for clustering. In **divisive clustering**, the entire data set is regarded as a cluster, and then clusters are recursively split to yield a good clustering (Algorithm 5.2). In **agglomerative clustering**, each data item is regarded as a cluster, and clusters are recursively merged to yield a good clustering (Algorithm 5.1).

```
Make each point a separate cluster
Until the clustering is satisfactory
    Merge the two clusters with the
        smallest inter-cluster distance
end
```

Algorithm 5.1: *Agglomerative Clustering or Clustering by Merging.*

```

Construct a single cluster containing all points
Until the clustering is satisfactory
    Split the cluster that yields the two
        components with the largest inter-cluster distance
end

```

Algorithm 5.2: *Divisive Clustering, or Clustering by Splitting.*

There are two major issues in thinking about clustering:

- *What is a good inter-cluster distance?* Agglomerative clustering uses an inter-cluster distance to fuse nearby clusters; divisive clustering uses it to split insufficiently coherent clusters. Even if a natural distance between data points is available (which might not be the case for vision problems), there is no canonical inter-cluster distance. Generally, one chooses a distance that seems appropriate for the data set. For example, one might choose the distance between the closest elements as the inter-cluster distance, which tends to yield extended clusters (statisticians call this method **single-link clustering**). Another natural choice is the maximum distance between an element of the first cluster and one of the second, which tends to yield rounded clusters (statisticians call this method **complete-link clustering**). Finally, one could use an average of distances between elements in the cluster, which also tends to yield “rounded” clusters (statisticians call this method **group average clustering**).
- *How many clusters are there?* This is an intrinsically difficult task if there is no model for the process that generated the clusters. The algorithms we have described generate a hierarchy of clusters. Usually, this hierarchy is displayed to a user in the form of a **dendrogram**—a representation of the structure of the hierarchy of clusters that displays inter-cluster distances—and an appropriate choice of clusters is made from the dendrogram (see the example in Figure 5.1).

The main difficulty in using a divisive model is knowing where to split. This is sometimes made easier for particular kinds of data. For example, we could segment an image by clustering pixel values. In this case, you can sometimes find good splits by constructing a histogram of intensities, or of color values.

Another important thing to notice about clustering from the example of figure 5.1 is that there is no right answer. There are a variety of different clusterings of the same data. For example, depending on what scales in that figure mean, it might be right to zoom out and regard all of the data as a single cluster, or to zoom in and regard each data point as a cluster. Each of these representations may be useful.

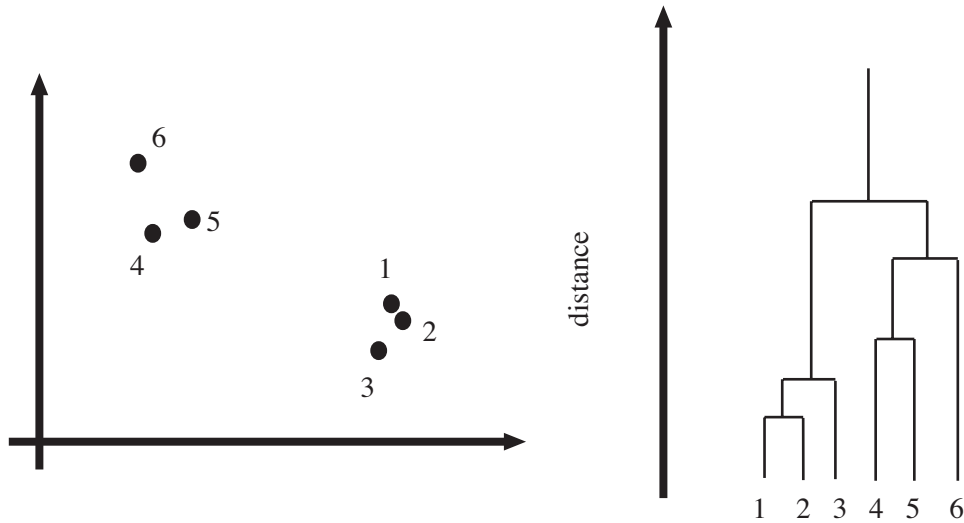


FIGURE 5.1: **Left**, a data set; **right**, a dendrogram obtained by agglomerative clustering using single-link clustering. If one selects a particular value of distance, then a horizontal line at that distance splits the dendrogram into clusters. This representation makes it possible to guess how many clusters there are and to get some insight into how good the clusters are.

5.1.1 Clustering and Distance

In the algorithms above, and in what follows, we assume that the features are scaled so that distances (measured in the usual way) between data points are a good representation of their similarity. This is quite an important point. For example, imagine we are clustering data representing brick walls. The features might contain several distances: the spacing between the bricks, the length of the wall, the height of the wall, and so on. If these distances are given in the same set of units, we could have real trouble. For example, assume that the units are centimeters. Then the spacing between bricks is of the order of one or two centimeters, but the heights of the walls will be in the hundreds of centimeters. In turn, this means that the distance between two datapoints is likely to be completely dominated by the height and length data. This could be what we want, but it might also not be a good thing.

There are some ways to manage this issue. One is to know what the features measure, and know how they should be scaled. Usually, this happens because you have a deep understanding of your data. If you don't (which happens!), then it is often a good idea to try and normalize the scale of the data set. There are two good strategies. The simplest is to translate the data so that it has zero mean (this is just for neatness - translation doesn't change distances), then scale each direction so that it has unit variance. More sophisticated is to translate the data so that it has zero mean, then transform it so that each direction is independent and has unit variance. Doing so is sometimes referred to as **decorrelation** or **whitening**

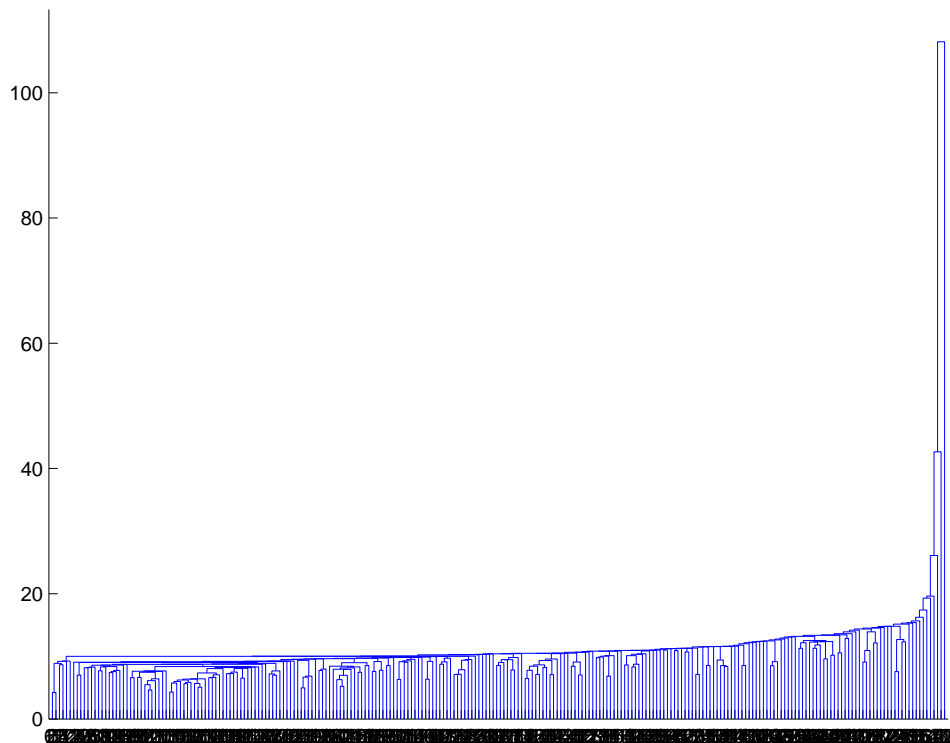


FIGURE 5.2: A dendrogram obtained from the body fat dataset, using single link clustering. Recall that the data points are on the horizontal axis, and that the vertical axis is distance; there is a horizontal line linking two clusters that get merged, established at the height at which they're merged. I have plotted the entire dendrogram, despite the fact it's a bit crowded at the bottom, because it shows that most data points are relatively close (i.e. there are lots of horizontal branches at about the same height).

(because you make the data more like white noise); I described how to do this in section ??.

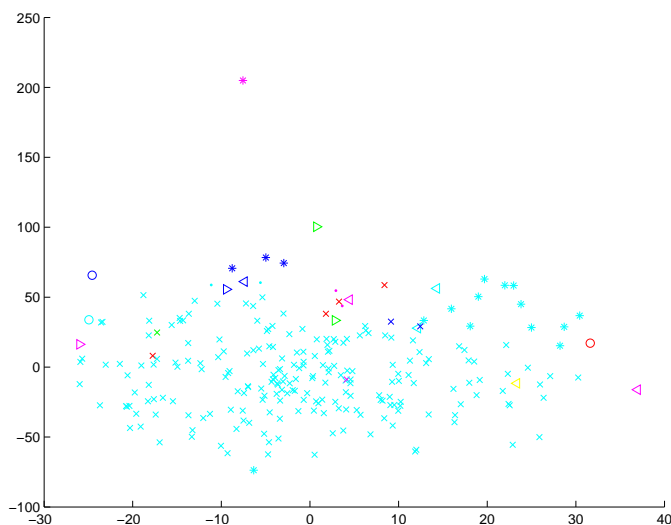


FIGURE 5.3: A clustering of the body fat dataset, using agglomerative clustering, single link distance, and requiring a maximum of 30 clusters. I have plotted each cluster with a distinct marker (though some markers differ only by color; you might need to look at the PDF version to see this figure at its best). Notice that one cluster contains much of the data, and that there are a set of small isolated clusters. The original data is 16 dimensional, which presents plotting problems; I show a scatter plot on the first two principal components (though I computed distances for clustering in the original 16 dimensional space).

Worked example 5.1 *Agglomerative clustering in Matlab*

Cluster the height-weight dataset of <http://www2.stetson.edu/~jrjsp/data.htm> (look for bodyfat.xls) using an agglomerative clusterer, and describe the results.

Solution: Matlab provides some tools that are useful for agglomerative clustering. These functions use a scheme where one first builds the whole tree of merges, then analyzes that tree to decide which clustering to report. `linkage` will determine which pairs of clusters should be merged at which step (there are arguments that allow you to choose what type of inter-cluster distance it should use); `dendrogram` will plot you a dendrogram; and `cluster` will extract the clusters from the linkage, using a variety of options for choosing the clusters. I used these functions to prepare the dendrogram of figure 5.2 for the height-weight dataset of section ?? (from <http://www2.stetson.edu/~jrjsp/data.htm>; look for bodyfat.xls). I deliberately forced Matlab to plot the whole dendrogram, which accounts for the crowded look of the figure (you can allow it to merge small leaves, etc.). I used a single-link strategy. In particular, notice that many data points are about the same distance from one another, which suggests a single big cluster with a smaller set of nearby clusters. The clustering of figure 5.3 supports this view. I plotted the data points on the first two principal components, using different colors and shapes of marker to indicate different clusters. There are a total of 30 clusters here, though most are small.

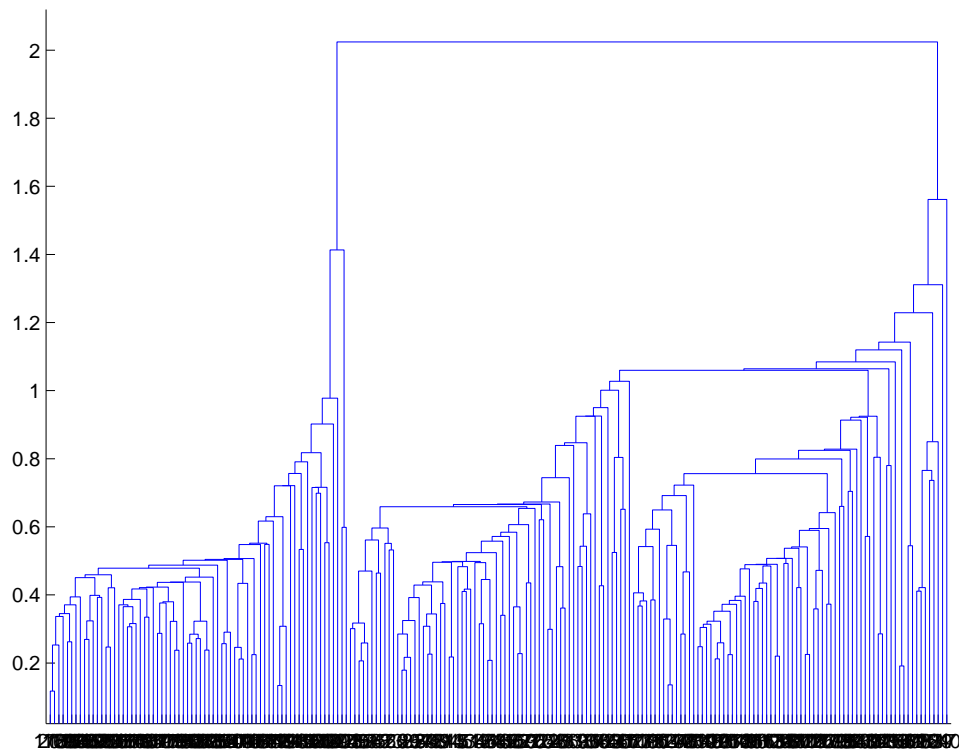


FIGURE 5.4: A dendrogram obtained from the seed dataset, using single link clustering. Recall that the data points are on the horizontal axis, and that the vertical axis is distance; there is a horizontal line linking two clusters that get merged, established at the height at which they're merged. I have plotted the entire dendrogram, despite the fact it's a bit crowded at the bottom, because you can now see how clearly the data set clusters into a small set of clusters — there are a small number of vertical “runs”.

Worked example 5.2 *Agglomerative clustering in Matlab, 2*

Cluster the seed dataset from the UC Irvine Machine Learning Dataset Repository (you can find it at <http://archive.ics.uci.edu/ml/datasets/seeds>).

Solution: Each item consists of seven measurements of a wheat kernel; there are three types of wheat represented in this dataset. As you can see in figures 5.4 and 5.5, this data clusters rather well.

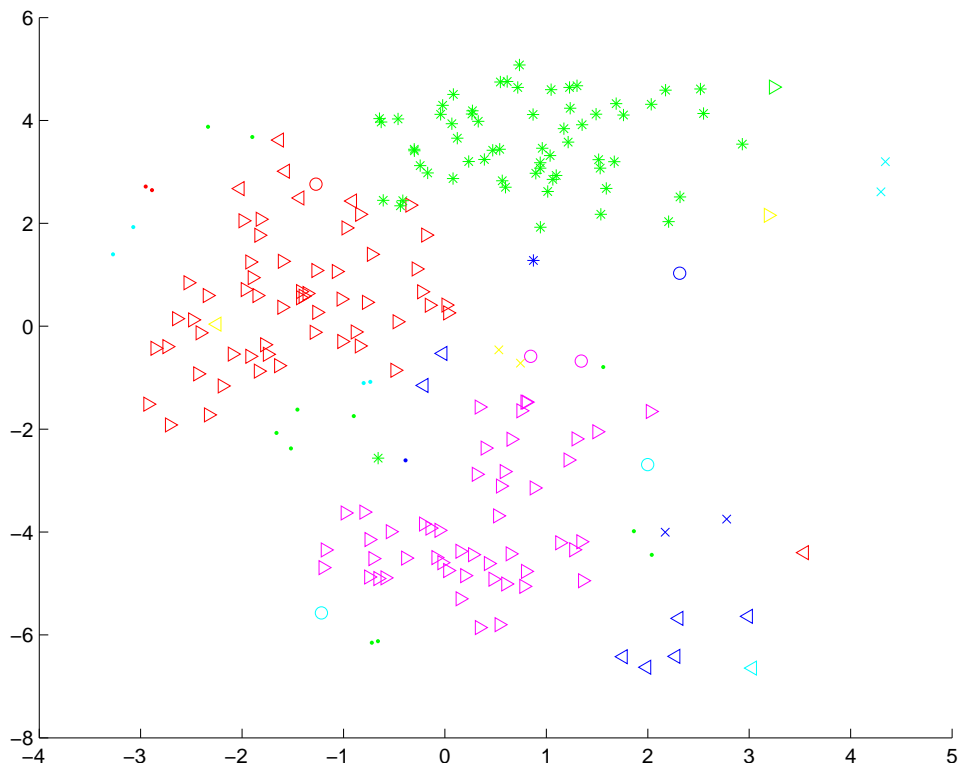


FIGURE 5.5: A clustering of the seed dataset, using agglomerative clustering, single link distance, and requiring a maximum of 30 clusters. I have plotted each cluster with a distinct marker (though some markers differ only by color; you might need to look at the PDF version to see this figure at its best). Notice that there are a set of fairly natural isolated clusters. The original data is 8 dimensional, which presents plotting problems; I show a scatter plot on the first two principal components (though I computed distances for clustering in the original 8 dimensional space).

Worked example 5.3 Agglomerative clustering in R

Cluster the questions in the student evaluation dataset from the UC Irvine Machine Learning Dataset Repository (you can find it at <https://archive.ics.uci.edu/ml/datasets/Turkiye+Student+Evaluation>; this dataset was donated by G. Gunduz and E. Fokoue), and display it on the first two principal components. You should use 5 clusters, but investigate the results of choosing others.

Solution: R provides tools for agglomerative clustering, too. I used the block of code shown in listing 5.1 to produce figure 5.6. The clustering shown in that figure is to 5 clusters. To get some idea of what the clusters are like, you can compute the per-cluster means, using a line in that listing. I found means where: all students in the cluster gave moderately low, low, moderately high and high answers to all questions, respectively; and one where all students answered 1 to all questions. There are 5820 students in this collection. The clusters suggest that answers tend to be quite strongly correlated — a student who gives a low answer to a question will likely give low answers to others, too. Choosing other numbers of clusters wasn't particularly revealing, though there were more levels to the answers.

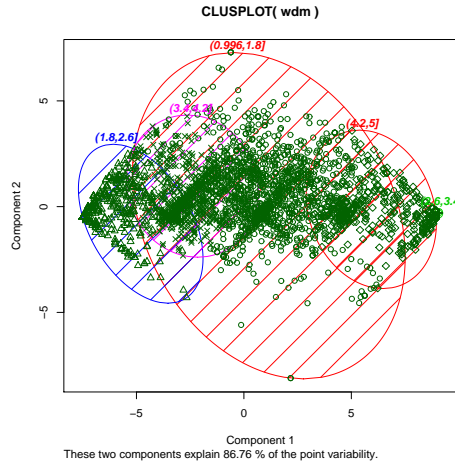


FIGURE 5.6: A clustering of the student evaluation dataset, using agglomerative clustering and 5 clusters. I produced this clustering with the code fragment in listing 5.1. One difficulty with this dataset is there are 28 questions (Q1-Q28), meaning the data is quite difficult to plot. You should be able to see the clusters on the figure, but notice that they are smeared on top of one another; this is because I had to project off 26 dimensions to produce the plot. The plotting function I used (`clusplot`) plots a covariance ellipse for each cluster.

5.2 THE K-MEANS ALGORITHM AND VARIANTS

Assume we have a dataset that, we believe, forms many clusters that look like blobs. If we knew where the center of each of the clusters was, it would be easy to tell which cluster each data item belonged to — it would belong to the cluster with the closest center. Similarly, if we knew which cluster each data item belonged to, it would be easy to tell where the cluster centers were — they'd be the mean of the data items in the cluster. This is the point closest to every point in the cluster.

We can turn these observations into an algorithm. Assume that we know how many clusters there are in the data, and write k for this number. The j th data item to be clustered is described by a feature vector \mathbf{x}_j . We write \mathbf{c}_i for the center of the i th cluster. We assume that most data items are close to the center of their cluster. This suggests that we cluster the data by minimizing the the cost function

$$\Phi(\text{clusters}, \text{data}) = \sum_{i \in \text{clusters}} \left\{ \sum_{j \in i\text{th cluster}} (\mathbf{x}_j - \mathbf{c}_i)^T (\mathbf{x}_j - \mathbf{c}_i) \right\}.$$

Notice that if we know the center for each cluster, it is easy to determine which cluster is the best choice for each point. Similarly, if the allocation of points to clusters is known, it is easy to compute the best center for each cluster. However, there are far too many possible allocations of points to clusters to search this space for a minimum. Instead, we define an algorithm that iterates through two activities:

Listing 5.1: R code for student example.

```

setwd('/users/daf/Current/courses/Probcourse/Clustering/RCode')
wdat<-read.csv('turkiye-student-evaluation_R_Specific.csv')
wdm<-wdat[,c(6:33)] # choose the questions
d <- dist(wdm, method = "euclidean") # distance matrix
fit <- hclust(d, method="ward") # the clustering
cm<-cov(wdm)
ev<-eigen(cm)
princ<-ev$vectors
proj<-t(rbind(t(princ[, 1]), t(princ[, 2])))
wdzm<-scale(wdm, center=TRUE, scale=FALSE) #subtract the mean
wdzm<-data.matrix(wdzm)
wdpcs<-wdzm%*%proj
kv<-c(5)
groups<-cutree(fit, k=kv)
wdbig<-wdm
wdbig$group<-groups
cg<-cut(groups, kv)
p<-ggplot(as.data.frame(wdpcs))+geom_point(aes(y=V1, x=V2, shape=cg))
setEPS()
postscript("clusters.eps")
p
dev.off()
cm1<-colMeans(subset(wdbig, group==1)) # mean of the 1st cluster

```

- Assume the cluster centers are known and, allocate each point to the closest cluster center.
- Assume the allocation is known, and choose a new set of cluster centers. Each center is the mean of the points allocated to that cluster.

We then choose a start point by randomly choosing cluster centers, and then iterate these stages alternately. This process eventually converges to a local minimum of the objective function (the value either goes down or is fixed at each step, and it is bounded below). It is not guaranteed to converge to the global minimum of the objective function, however. It is also not guaranteed to produce k clusters, unless we modify the allocation phase to ensure that each cluster has some nonzero number of points. This algorithm is usually referred to as **k-means** (summarized in Algorithm 5.3).

Usually, we are clustering high dimensional data, so that visualizing clusters can present a challenge. If the dimension isn't too high, then we can use panel plots. An alternative is to project the data onto two principal components, and plot the clusters there. A natural dataset to use to explore k-means is the iris data, where we know that the data should form three clusters (because there are three species). Recall this dataset from section ???. I reproduce figure 4.5 from that section as figure 5.11, for comparison. Figures 5.8, ??? and ??? show different k-means clusterings of that data.

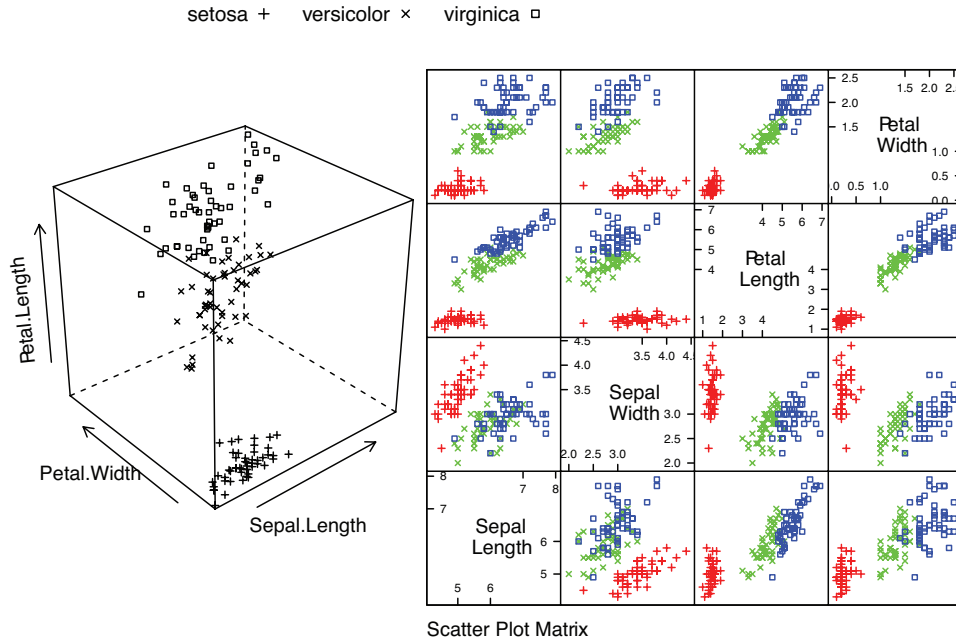


FIGURE 5.7: **Left:** a 3D scatterplot for the famous Iris data, originally due to *****. I have chosen three variables from the four, and have plotted each species with a different marker. You can see from the plot that the species cluster quite tightly, and are different from one another. **Right:** a scatterplot matrix for the famous Iris data, originally due to *****. There are four variables, measured for each of three species of iris. I have plotted each species with a different marker. You can see from the plot that the species cluster quite tightly, and are different from one another.

Worked example 5.4 *K*-means clustering in *R*

Cluster the iris dataset into two clusters using k-means, then plot the results on the first two principal components

Solution: I used the code fragment in listing 5.2, which produced figure ??

5.2.1 How to choose K

The iris data is just a simple example. We know that the data forms clean clusters, and we know there should be three of them. Usually, we don't know how many clusters there should be, and we need to choose this by experiment. One strategy is to cluster for a variety of different values of *k*, then look at the value of the cost function for each. If there are more centers, each data point can find a center that is close to it, so we expect the value to go down as *k* goes up. This means that

```

Choose  $k$  data points to act as cluster centers
Until the cluster centers change very little
  Allocate each data point to cluster whose center is nearest.
  Now ensure that every cluster has at least
  one data point; one way to do this is by
  supplying empty clusters with a point chosen at random from
  points far from their cluster center.
  Replace the cluster centers with the mean of the elements
  in their clusters.
end

```

Algorithm 5.3: *Clustering by K-Means.*

Listing 5.2: R code for iris example.

```

setwd('/users/daf/Current/courses/Probcourse/Clustering/RCode')
#library('lattice')
# work with iris dataset this is famous, and included in R
# there are three species
head(iris)
#
library('cluster')
numiris=iris[,c(1, 2, 3, 4)] #the numerical values
#scalediris<-scale(numiris) # scale to unit variance
nclus<-2
sfit<-kmeans(numiris, nclus)
colr<-c('red', 'green', 'blue', 'yellow', 'orange')
clusplot(numiris, sfit$cluster, color=TRUE, shade=TRUE,
         labels=0, lines=0)

```

looking for the k that gives the smallest value of the cost function is not helpful, because that k is always the same as the number of data points (and the value is then zero). However, it can be very helpful to plot the value as a function of k , then look at the “knee” of the curve. Figure 5.11 shows this plot for the iris data. Notice that $k = 3$ — the “true” answer — doesn’t look particularly special, but $k = 2$, $k = 3$, or $k = 4$ all seem like reasonable choices. It is possible to come up with a procedure that makes a more precise recommendation by penalizing clusterings that use a large k , because they may represent inefficient encodings of the data. However, this is often not worth the bother.

In some special cases (like the iris example), we might know the right answer to check our clustering against. In such cases, one can evaluate the clustering by looking at the number of different labels in a cluster (sometimes called the purity), and the number of clusters. A good solution will have few clusters, all of which have high purity.

Mostly, we don’t have a right answer to check against. An alternative strategy, which might seem crude to you, for choosing k is extremely important in practice. Usually, one clusters data to use the clusters in an application (one of the most

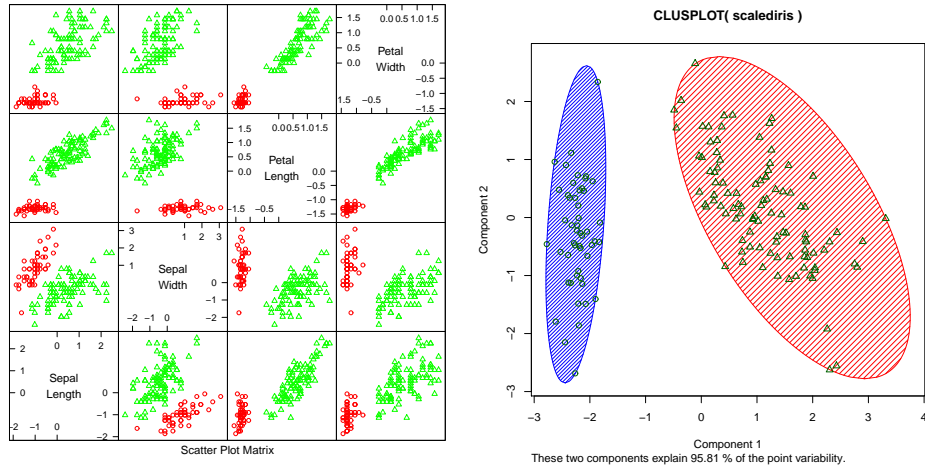


FIGURE 5.8: On the **left**, a panel plot of the iris data clustered using k -means with $k = 2$. By comparison with figure 5.11, notice how the versicolor and virginica clusters appear to have been merged. On the **right**, this data set projected onto the first two principal components, with one blob drawn over each cluster.

important, vector quantization, is described in section 5.3). There are usually natural ways to evaluate this application. For example, vector quantization is often used as an early step in texture recognition or in image matching; here one can evaluate the error rate of the recognizer, or the accuracy of the image matcher. One then chooses the k that gets the best evaluation score on validation data. In this view, the issue is not how good the clustering is; it's how well the system that uses the clustering works.

5.2.2 Soft Assignment

One difficulty with k -means is that each point must belong to exactly one cluster. But, given we don't know how many clusters there are, this seems wrong. If a point is close to more than one cluster, why should it be forced to choose? This reasoning suggests we assign points to cluster centers with weights.

We allow each point to carry a total weight of 1. In the conventional k -means algorithm, it must choose a single cluster, and assign its weight to that cluster alone. In soft-assignment k -means, the point can allocate some weight to each cluster center, as long as (a) the weights are all non-negative and (b) the weights sum to one. Write $w_{i,j}$ for the weight connecting point i to cluster center j . We interpret these weights as the degree to which the point participates in a particular cluster. We require $w_{i,j} \geq 0$ and $\sum_j w_{i,j} = 1$.

We would like $w_{i,j}$ to be large when \mathbf{x}_i is close to \mathbf{c}_j , and small otherwise. Write $d_{i,j}$ for the distance $\|\mathbf{x}_i - \mathbf{c}_j\|$ between these two. Write

$$s_{i,j} = e^{\frac{-d_{i,j}^2}{2\sigma^2}}$$

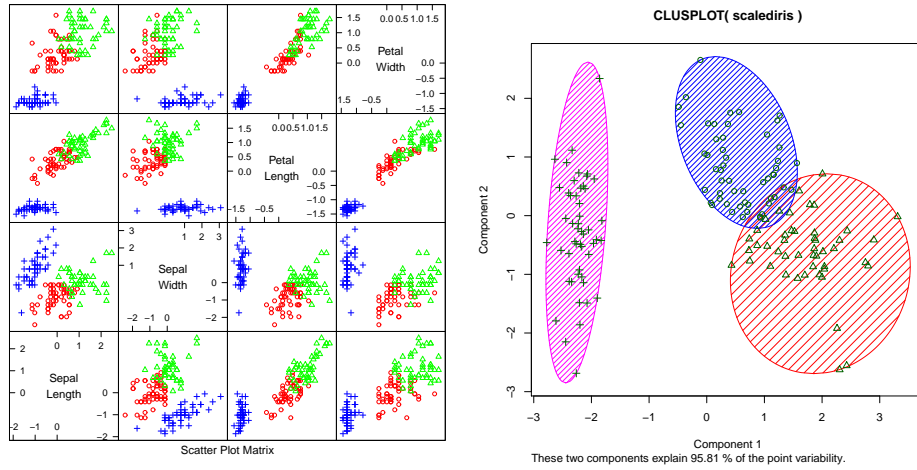


FIGURE 5.9: On the **left**, a panel plot of the iris data clustered using k -means with $k = 3$. By comparison with figure 5.11, notice how the clusters appear to follow the species labels. On the **right**, this data set projected onto the first two principal components, with one blob drawn over each cluster.

where $\sigma > 0$ is a choice of scaling parameter. This is often called the affinity between the point i and the center j . Now a natural choice of weights is

$$w_{i,j} = \frac{s_{i,j}}{\sum_{l=1}^k s_{i,l}}.$$

All these weights are non-negative, they sum to one, and the weight is large if the point is much closer to one center than to any other. The scaling parameter σ sets the meaning of “much closer” — we measure distance in units of σ .

Once we have weights, re-estimating the cluster centers is easy. We use a weights to compute a weighted average of the points. In particular, we re-estimate the j 'th cluster center by

$$\frac{\sum_i w_{i,j} \mathbf{x}_i}{\sum_i w_{i,j}}.$$

Notice that k -means is a special case of this algorithm where σ limits to zero. In this case, each point has a weight of one for some cluster, and zero for all others, and the weighted mean becomes an ordinary mean. I have collected the description into Algorithm 5.4 for convenience.

5.2.3 General Comments on K-Means

If you experiment with k -means, you will notice one irritating habit of the algorithm. It almost always produces either some rather spread out clusters, or some single element clusters. Most clusters are usually rather tight and blobby clusters, but there is usually one or more bad cluster. This is fairly easily explained. Because every data point must belong to some cluster, data points that are far from all

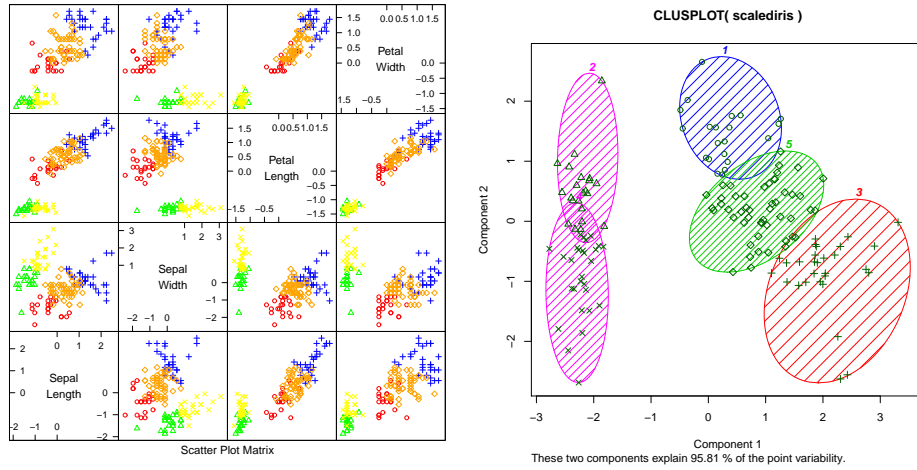


FIGURE 5.10: *On the left*, a panel plot of the iris data clustered using k -means with $k = 5$. *By comparison with figure 5.11*, notice how setosa seems to have been broken in two groups, and versicolor and virginica into a total of three. *On the right*, this data set projected onto the first two principal components, with one blob drawn over each cluster.

others (a) belong to some cluster and (b) very likely “drag” the cluster center into a poor location. This applies even if you use soft assignment, because every point must have total weight one. If the point is far from all others, then it will be assigned to the closest with a weight very close to one, and so may drag it into a poor location, or it will be in a cluster on its own.

There are ways to deal with this. If k is very big, the problem is often not significant, because then you simply have many single element clusters that you can ignore. It isn’t always a good idea to have too large a k , because then some larger clusters might break up. An alternative is to have a junk cluster. Any point that is too far from the closest true cluster center is assigned to the junk cluster, and the center of the junk cluster is not estimated. Notice that points should not be assigned to the junk cluster permanently; they should be able to move in and out of the junk cluster as the cluster centers move.

In some cases, we want to cluster objects that can’t be averaged. For example, you can compute distances between two trees but you can’t meaningfully average them. In some cases, you might have a table of distances between objects, but not know vectors representing the objects. For example, one could collect data on the similarities between countries (as in Section 4.6.2, particularly Figure 4.29), then try and cluster using this data (similarities can be turned into distances by, for example, taking the negative logarithm). A variant of k -means, known as k -medoids, applies to this case.

In k -medoids, the cluster centers are data items rather than averages, but the rest of the algorithm has a familiar form. We assume the number of medoids is known, and initialize these randomly. We then iterate two procedures. In the first,

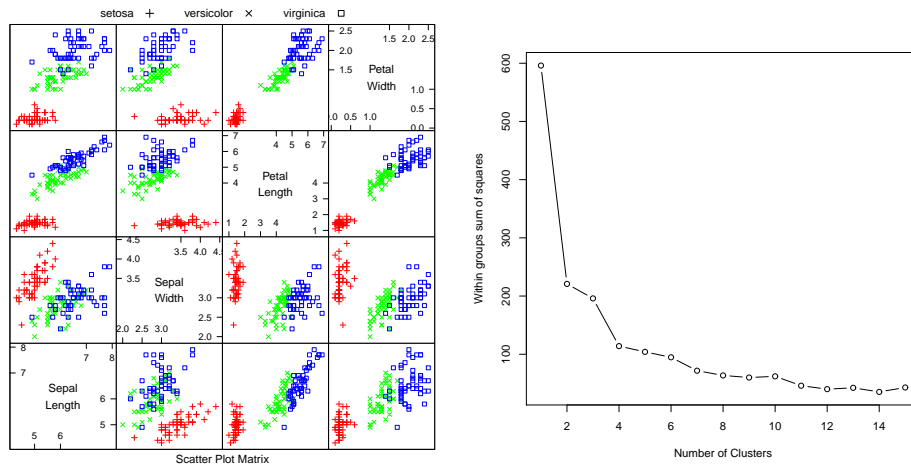


FIGURE 5.11: On the left, the scatterplot matrix for the Iris data, for reference. On the right, a plot of the value of the cost function for each of several different values of k . Notice how there is a sharp drop in cost going from $k = 1$ to $k = 2$, and again at $k = 4$; after that, the cost falls off slowly. This suggests using $k = 2$, $k = 3$, or $k = 4$, depending on the precise application.

we allocate data points to medoids. In the second, we choose the best medoid for each cluster by finding the medoid that minimizes the sum of distances of points in the cluster to that medoid (blank search is fine).

5.3 DESCRIBING REPETITION WITH VECTOR QUANTIZATION

Repetition is an important feature of many interesting signals. For example, images contain *textures*, which are orderly patterns that look like large numbers of small structures that are repeated. Examples include the spots of animals such as leopards or cheetahs; the stripes of animals such as tigers or zebras; the patterns on bark, wood, and skin. Similarly, speech signals contain *phonemes* — characteristic, stylised sounds that people assemble together to produce speech (for example, the “ka” sound followed by the “tuh” sound leading to “cat”). Another example comes from accelerometers. If a subject wears an accelerometer while moving around, the signals record the accelerations during their movements. So, for example, brushing one’s teeth involves a lot of repeated twisting movements at the wrist, and walking involves swinging the hand back and forth.

Repetition occurs in subtle forms. The essence is that a small number of local patterns can be used to represent a large number of examples. You see this effect in pictures of scenes. If you collect many pictures of, say, a beach scene, you will expect most to contain some waves, some sky, and some sand. The individual patches of wave, sky or sand can be surprisingly similar, and different images are made by selecting some patches from a vocabulary of patches, then placing them down to form an image. Similarly, pictures of living rooms contain chair patches, TV patches, and carpet patches. Many different living rooms can be made from

Choose k data points to act as cluster centers

Until the cluster centers change very little

First, we estimate the weights

For i indexing data points

For j indexing cluster centers

$$\text{Compute } s_{i,j} = e^{\frac{-\|\mathbf{x}_i - \mathbf{c}_j\|}{2\sigma^2}}$$

For i indexing data points

For j indexing cluster centers

$$\text{Compute } w_{i,j} = s_{i,j} / \sum_{l=1}^k s_{i,l}$$

Now, we re-estimate the centers

For j indexing cluster centers

$$\text{Compute } \mathbf{c}_j = \frac{\sum_i w_{i,j} \mathbf{x}_i}{\sum_i w_{i,j}}$$

end

Algorithm 5.4: *Soft Clustering by K-Means.*

small vocabularies of patches; but you won't often see wave patches in living rooms, or carpet patches in beach scenes.

An important part of representing signals that repeat is building a vocabulary of patterns that repeat, then describing the signal in terms of those patterns. For many problems, problems, knowing what vocabulary elements appear and how often is much more important than knowing where they appear. For example, if you want to tell the difference between zebra's and leopards, you need to know whether stripes or spots are more common, but you don't particularly need to know where they appear. As another example, if you want to tell the difference between brushing teeth and walking using accelerometer signals, knowing that there are lots of (or few) twisting movements is important, but knowing how the movements are linked together in time may not be.

5.3.1 Vector Quantization

It is natural to try and find patterns by looking for small pieces of signal of fixed size that appear often. In an image, a piece of signal might be a 10x10 patch; in a sound file, which is likely represented as a vector, it might be a subvector of fixed size. But finding patterns that appear often is hard to do, because the signal is continuous — each pattern will be slightly different, so we cannot simply count how many times a particular pattern occurs.

Here is a strategy. First, we take a training set of signals, and cut each signal into vectors of fixed dimension (say d). It doesn't seem to matter too much if these

vectors overlap or not. We then build a set of clusters out of these vectors; this set of clusters is often thought of as a dictionary. We can now describe any new vector with the cluster center closest to that vector. This means that a vector in a continuous space is described with a number in the range $[1, \dots, k]$ (where you get to choose k), and two vectors that are close should be described by the same number. This strategy is known as **vector quantization**.

We can now build features that represent important repeated structure in signals. We take a signal, and cut it up into vectors of length d . These might overlap, or be disjoint; we follow whatever strategy we used in building the dictionary. We then take each vector, and compute the number that describes it (i.e. the number of the closest cluster center, as above). We then compute a histogram of the numbers we obtained for all the vectors in the signal. This histogram describes the signal.

Notice several nice features to this construction. First, it can be applied to anything that can be thought of in terms of vectors, so it will work for speech signals, sound signals, accelerometer signals, images, and so on. You might need to adjust some indices. For example, you cut the image into patches, then rearrange the patch to form a vector. As another example, accelerometer signals are three dimensional vectors that depend on time, so you cut out windows of a fixed number of time samples (say t), then rearrange to get a $3t$ dimensional vector.

Another nice feature is the construction can accept signals of different length, and produce a description of fixed length. One accelerometer signal might cover 100 time intervals; another might cover 200; but the description is always a histogram with k buckets, so it's always a vector of length k .

Yet another nice feature is that we don't need to be all that careful how we cut the signal into fixed length vectors. This is because it is hard to hide repetition. This point is easier to make with a figure than in text, so look at figure ??.

5.3.2 Example: Groceries in Portugal

At <http://archive.ics.uci.edu/ml/datasets/Wholesale+customers>, you will find a dataset giving sums of money spent annually on different commodities by customers in Portugal. The commodities are divided into a set of categories (fresh; milk; grocery; frozen; detergents and paper; and delicatessen) relevant for the study. These customers are divided by channel (two channels) and by region (three regions). You can think of the data as being divided into six groups (one for each channel-region pair). There are 440 records, and so there are many customers per group. Figure 5.12 shows a panel plot of the customer data; the data has been clustered, and I gave each of 20 clusters its own marker. Relatively little structure is apparent in this scatter plot. You can't, for example, see evidence of six groups that are cleanly separated.

It's unlikely that all the customers in a group are the same. Instead, we expect that there might be different "types" of customer. For example, customers who prepare food at home might spend more money on fresh or on grocery, and those who mainly buy prepared food might spend more money on delicatessen; similarly, coffee drinkers with cats or children might spend more on milk than the lactose-intolerant, and so on. Because some of these effects are driven by things like wealth and the tendency of people to like to have neighbors who are similar to

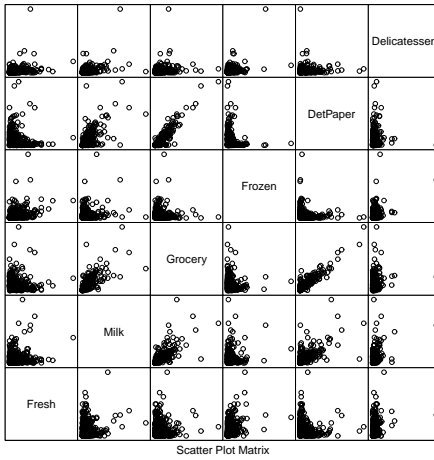


FIGURE 5.12: A panel plot of the wholesale customer data of <http://archive.ics.uci.edu/ml/datasets/Wholesale+customers>, which records sums of money spent annually on different commodities by customers in Portugal. This data is recorded for six different groups (two channels each within three regions). I have plotted each group with a different marker, but you can't really see much structure here, for reasons explained in the text.

them, you could expect that different groups contain different numbers of each type of customer. There might be more deli-spenders in wealthier regions; more milk-spenders and detergent-spenders in regions where it is customary to have many children; and so on.

An effect like this is hard to see on a panel plot (Figure 5.12). The plot for this dataset is hard to read, because the dimension is fairly high for a panel plot and the data is squashed together in the bottom left corner. There is another effect. If customers are grouped in the way I suggested above, then each group might look the same in a panel plot. A group of some milk-spenders and more detergent-spenders will have many data points with high milk expenditure values (and low other values) and also many data points with high detergent expenditure values (and low other values). In a panel plot, this will look like two blobs; but another group with more milk-spenders and some detergent-spenders will also look like two blobs, in about the same place. It will be hard to spot the difference. A histogram of the types within each group will make this difference obvious.

I used k-means clustering to cluster the customer data to 20 different clusters (Figure 5.14). I chose 20 rather arbitrarily, but with the plot of error against k in mind. Then I described the each group of data by the histogram of customer types that appeared in that group (Figure ??). Notice how the distinction between the groups is now apparent — the groups do appear to contain quite different distributions of customer type. It looks as though the channels (rows in this figure) are more different than the regions (columns in this figure). To be more confident in this analysis, we would need to be sure that different types of customer really

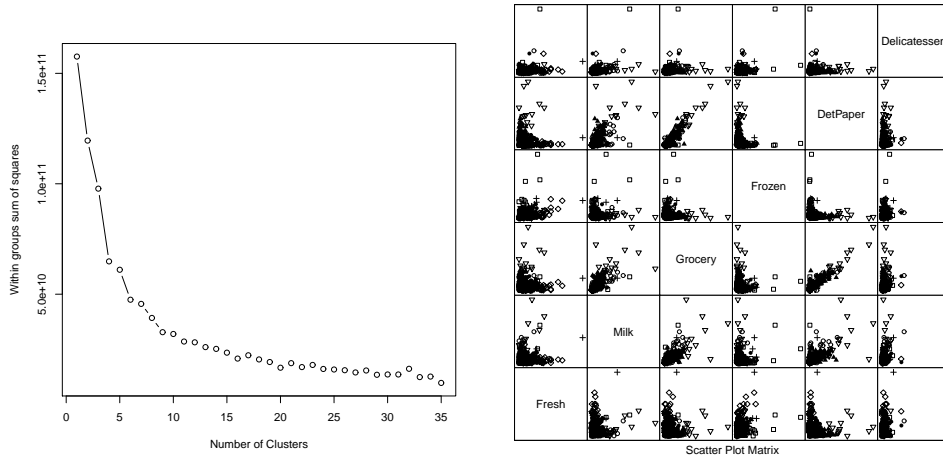


FIGURE 5.13: *On the left, the sum of squared error for clusterings of the customer data with k-means for k running from 2 to 35. This suggests using a k somewhere in the range 10-30; I chose 20. On the right, I have clustered this data to 20 cluster centers with k-means. The clusters do seem to be squashed together, but the plot on the left suggests that clusters do capture some important information. Using too few clusters will clearly lead to problems. Notice that I did not scale the data, because each of the measurements is in a comparable unit. For example, it wouldn't make sense to scale expenditures on fresh and expenditures on grocery with a different scale.*

are different. We could do this by repeating the analysis for fewer clusters, or by looking at the similarity of customer types.

5.3.3 Efficient Clustering and Hierarchical K Means

One important difficulty occurs in applications. We might need to have an enormous dataset (millions of image patches are a real possibility), and so a very large k . In this case, k means clustering becomes difficult because identifying which cluster center is closest to a particular data point scales linearly with k (and we have to do this for every data point at every iteration). There are two useful strategies for dealing with this problem.

The first is to notice that, if we can be reasonably confident that each cluster contains many data points, some of the data is redundant. We could randomly subsample the data, cluster that, then keep the cluster centers. This works, but doesn't scale particularly well.

A more effective strategy is to build a hierarchy of k-means clusters. We randomly subsample the data (typically, quite aggressively), then cluster this with a small value of k . Each data item is then allocated to the closest cluster center, and the data in each cluster is clustered again with k-means. We now have something that looks like a two-level tree of clusters. Of course, this process can be repeated to produce a multi-level tree of clusters. It is easy to use this tree to vector quantize a

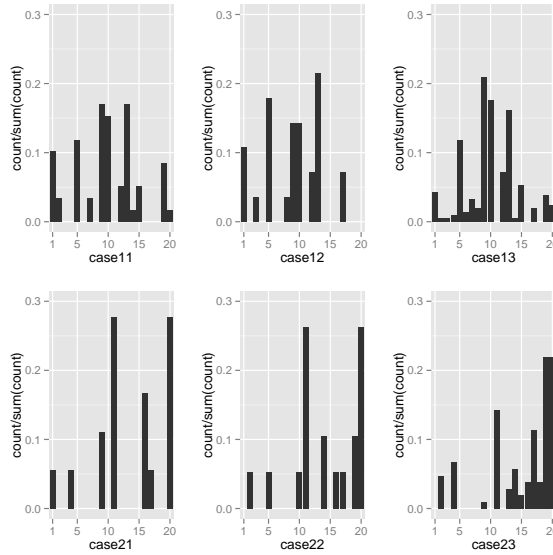


FIGURE 5.14: *The histogram of different types of customer, by group, for the customer data. Notice how the distinction between the groups is now apparent — the groups do appear to contain quite different distributions of customer type. It looks as though the channels (rows in this figure) are more different than the regions (columns in this figure).*

query data item. We vector quantize at the first level. Doing so chooses a branch of the tree, and we pass the data item to this branch. It is either a leaf, in which case we report the number of the leaf, or it is a set of clusters, in which case we vector quantize, and pass the data item down. This procedure is efficient both when one clusters and at run time.

5.3.4 Example: Activity from Accelerometer Data

A complex example dataset appears at <https://archive.ics.uci.edu/ml/datasets/Dataset+for+ADL+Recognition+>. This dataset consists of examples of the signal from a wrist mounted accelerometer, produced as different subjects engaged in different activities of daily life. Activities include: brushing teeth, climbing stairs, combing hair, descending stairs, and so on. Each is performed by sixteen volunteers. The accelerometer samples the data at 32Hz (i.e. this data samples and reports the acceleration 32 times per second). The accelerations are in the x, y and z-directions. Figure 5.15 shows the x-component of various examples of toothbrushing.

There is an important problem with using data like this. Different subjects take quite different amounts of time to perform these activities. For example, some subjects might be more thorough tooth-brushers than other subjects. As another example, people with longer legs walk at somewhat different frequencies than people with shorter legs. This means that the same activity performed by different subjects will produce data vectors *that are of different lengths*. It's not a good idea to deal

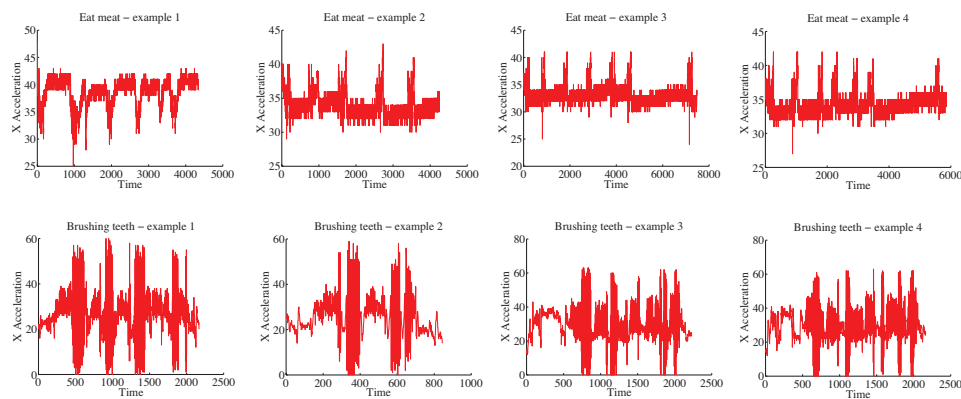


FIGURE 5.15: *Some examples from the accelerometer dataset at <https://archive.ics.uci.edu/ml/datasets/Dataset+for+ADL+Recognition+with+Wrist-worn+Accelerometer>. I have labelled each signal by the activity. These show acceleration in the X direction (Y and Z are in the dataset, too). There are four examples for **brushing teeth** and four for **eat meat**. You should notice that the examples don't have the same length in time (some are slower and some faster eaters, etc.), but that there seem to be characteristic features that are shared within a category (brushing teeth seems to involve faster movements than eating meat).*

with this by warping time and resampling the signal. For example, doing so will make a thorough toothbrusher look as though they are moving their hands very fast (or a careless toothbrusher look ludicrously slow: think speeding up or slowing down a movie). So we need a representation that can cope with signals that are a bit longer or shorter than other signals.

Another important property of these signals is that all examples of a particular activity should contain repeated patterns. For example, brushing teeth should show fast accelerations up and down; walking should show a strong signal at somewhere around 2 Hz; and so on. These two points should suggest vector quantization to you. Representing the signal in terms of stylized, repeated structures is probably a good idea because the signals probably contain these structures. And if we represent the signal in terms of the relative frequency with which these structures occur, the representation will have a fixed length, even if the signal doesn't. To do so, we need to consider (a) over what time scale we will see these repeated structures and (b) how to ensure we segment the signal into pieces so that we see these structures.

Generally, repetition in activity signals is so obvious that we don't need to be smart about segment boundaries. I broke these signals into 32 sample segments, one following the other. Each segment represents one second of activity. This is long enough for the body to do something interesting, but not so long that our representation will suffer if we put the segment boundaries in the wrong place. This resulted in about 40,000 segments. I then used hierarchical k-means to cluster these segments. I used two levels, with 40 cluster centers at the first level, and 12 at the second. Figure 5.16 shows some cluster centers at the second level.

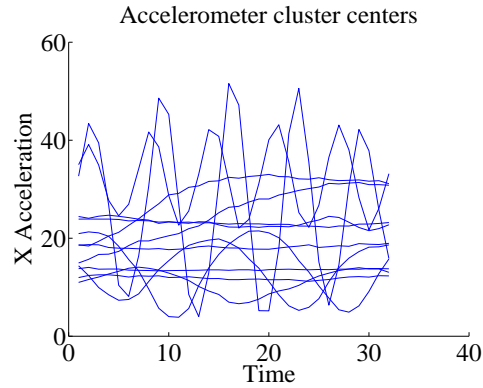


FIGURE 5.16: Some cluster centers from the accelerometer dataset. Each cluster center represents a one-second burst of activity. There are a total of 480 in my model, which I built using hierarchical k -means. Notice there are a couple of centers that appear to represent movement at about 5Hz; another few that represent movement at about 2Hz; some that look like 0.5Hz movement; and some that seem to represent much lower frequency movement. These cluster centers are samples (rather than chosen to have this property).

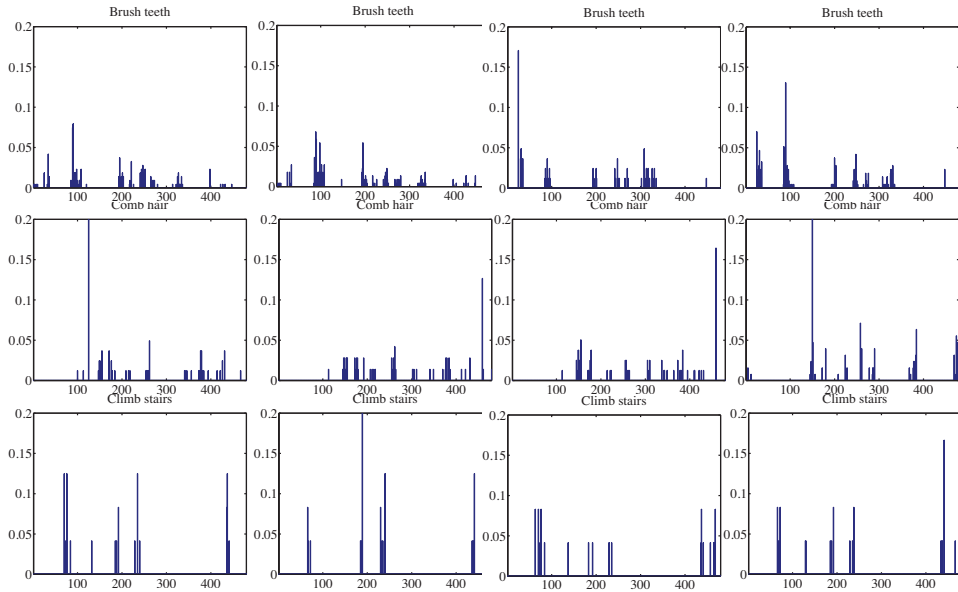


FIGURE 5.17: Histograms of cluster centers for the accelerometer dataset, for different activities. You should notice that (a) these histograms look somewhat similar for different actors performing the same activity and (b) these histograms look somewhat different for different activities.

I then computed histogram representations for different example signals (Figure 5.17). You should notice that when the activity label is different, the histogram looks different, too.

Another useful way to check this representation is to compare the average within class chi-squared distance with the average between class chi-squared distance. I computed the histogram for each example. Then, for each pair of examples, I computed the chi-squared distance between the pair. Finally, for each pair of *activity labels*, I computed the average distance between pairs of examples where one example has one of the activity labels and the other example has the other activity label. In the ideal case, all the examples with the same label would be very close to one another, and all examples with different labels would be rather different. Table 5.1 shows what happens with the real data. You should notice that for some pairs of activity label, the mean distance between examples is smaller than one would hope for (perhaps some pairs of examples are quite close?). But generally, examples of activities with different labels tend to be further apart than examples of activities with the same label.

0.9	2.0	1.9	2.0	2.0	2.0	1.9	2.0	1.9	1.9	2.0	2.0	2.0	2.0
	1.6	2.0	1.8	2.0	2.0	2.0	1.9	1.9	2.0	1.9	1.9	2.0	1.7
		1.5	2.0	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	2.0
			1.4	2.0	2.0	2.0	2.0	2.0	2.0	2.0	2.0	2.0	1.8
				1.5	1.8	1.7	1.9	1.9	1.8	1.9	1.9	1.8	2.0
					0.9	1.7	1.9	1.9	1.8	1.9	1.9	1.9	2.0
						0.3	1.9	1.9	1.5	1.9	1.9	1.9	2.0
							1.8	1.8	1.9	1.9	1.9	1.9	1.9
								1.7	1.9	1.9	1.9	1.9	1.9
									1.6	1.9	1.9	1.9	2.0
										1.8	1.9	1.9	1.9
											1.8	2.0	1.9
												1.5	2.0
													1.5

TABLE 5.1: Each column of the table represents an activity for the activity dataset <https://archive.ics.uci.edu/ml/datasets/Dataset+for+ADL+Recognition+with+Wrist-worn+Accelerometer>, as does each row. In each of the upper diagonal cells, I have placed the average chi-squared distance between histograms of examples from that pair of classes (I dropped the lower diagonal for clarity). Notice that in general the diagonal terms (average within class distance) are rather smaller than the off diagonal terms. This quite strongly suggests we can use these histograms to classify examples successfully.

Yet another way to check the representation is to try classification with nearest neighbors, using the chi-squared distance to compute distances. I split the dataset into 80 test pairs and 360 training pairs; using 1-nearest neighbors, I was able to get a held-out error rate of 0.79. This suggests that the representation is fairly good at exposing what is important.

5.4 YOU SHOULD

5.4.1 remember:

New term: clustering	97
New term: decorrelation	99
New term: whitening	99
New term: k-means	105
New term: vector quantization	113

PROGRAMMING EXERCISES

- 5.1.** Obtain the activities of daily life dataset from the UC Irvine machine learning website (<https://archive.ics.uci.edu/ml/datasets/Dataset+for+ADL+Recognition+with+Wrist-worn+Accelerometer>) data provided by Barbara Bruno, Fulvio Mastrogiiovanni and Antonio Sgorbissa).
- (a) Build a classifier that classifies sequences into one of the 14 activities provided. To make features, you should vector quantize, then use a histogram of cluster centers (as described in the subsection; this gives a pretty explicit set of steps to follow). You will find it helpful to use hierarchical k-means to vector quantize. You may use whatever multi-class classifier you wish, though I'd start with R's decision forest, because it's easy to use and effective. You should report (a) the total error rate and (b) the class confusion matrix of your classifier.
 - (b) Now see if you can improve your classifier by (a) modifying the number of cluster centers in your hierarchical k-means and (b) modifying the size of the fixed length samples that you use.

CHAPTER 6

Clustering using Probability Models

6.1 THE MULTIVARIATE NORMAL DISTRIBUTION

All the nasty facts about high dimensional data, above, suggest that we need to use quite simple probability models. By far the most important model is the multivariate normal distribution, which is quite often known as the multivariate gaussian distribution. There are two sets of parameters in this model, the mean μ and the covariance Σ . For a d -dimensional model, the mean is a d -dimensional column vector and the covariance is a $d \times d$ dimensional matrix. The covariance is a symmetric matrix. For our definitions to be meaningful, the covariance matrix must be positive definite.

The form of the distribution $p(\mathbf{x}|\mu, \Sigma)$ is

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right),$$

where Σ is a positive definite matrix. Notice that if Σ is not positive definite, then we cannot have a probability distribution, because there are some directions \mathbf{d} such that $\exp(-\frac{1}{2}(\mathbf{d} - \mu)^T \Sigma^{-1}(\mathbf{d} - \mu))$ does not fall off to zero as t limits to infinity. In turn, this means we can't compute the integral, and so can't normalize.

The following facts explain the names of the parameters:

Useful Facts: 6.1 *Parameters of a Multivariate Normal Distribution*

Assuming a multivariate normal distribution, we have

- $\mathbb{E}[\mathbf{x}] = \mu$, meaning that the mean of the distribution is μ .
- $\mathbb{E}[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] = \Sigma$, meaning that the entries in Σ represent covariances.

Assume I know have a dataset of items \mathbf{x}_i , where i runs from 1 to N , and we wish to model this data with a multivariate normal distribution. The maximum likelihood estimate of the mean, $\hat{\mu}$, is

$$\hat{\mu} = \frac{\sum_i \mathbf{x}_i}{N}$$

(which is quite easy to show). The maximum likelihood estimate of the covariance, $\hat{\Sigma}$, is

$$\hat{\Sigma} = \frac{\sum_i (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T}{N}$$

(which is rather a nuisance to show, because you need to know how to differentiate a determinant). You should be aware that this estimate is not guaranteed to be positive definite, even though the covariance matrix of a gaussian must be positive definite. We deal with this problem by checking the estimate. If its smallest eigenvalue is too close to zero, then we add some small positive constant times the identity to get a positive definite matrix.

6.1.1 Affine Transformations and Gaussians

Gaussians behave very well under affine transformations. In fact, we've already worked out all the math. Assume I have a dataset $\{\mathbf{x}\}$ with N data points \mathbf{x}_i , each a d -dimensional vector. The mean of the maximum likelihood gaussian model is $\text{mean}(\{\mathbf{x}\})$, and the covariance is $\text{Covmat}(\{\mathbf{x}\})$. We assume that this is positive definite, or adjust it as above.

I can now transform the data with an affine transformation, to get $\mathbf{y}_i = \mathcal{A}\mathbf{x}_i + \mathbf{b}$. We assume that \mathcal{A} is a square matrix with full rank, so that this transformation is 1-1. The mean of the maximum likelihood gaussian model for the transformed dataset is $\text{mean}(\{\mathbf{y}\}) = \mathcal{A}\text{mean}(\{\mathbf{x}\}) + \mathbf{b}$. Similarly, the covariance is $\text{Covmat}(\{\mathbf{y}\}) = \mathcal{A}\text{Covmat}(\{\mathbf{x}\})\mathcal{A}^T$.

A very important point follows in an obvious way. I can apply an affine transformation to any multivariate gaussian to obtain one with (a) zero mean and (b) independent components. In turn, this means that, *in the right coordinate system*, any gaussian is a product of zero mean, unit standard deviation, one-dimensional normal distributions. This fact is quite useful. For example, it means that simulating multivariate normal distributions is quite straightforward — you could simulate a standard normal distribution for each component, then apply an affine transformation.

6.1.2 Plotting a 2D Gaussian: Covariance Ellipses

There are some useful tricks for plotting a 2D Gaussian, which are worth knowing both because they're useful, and they help to understand Gaussians. Assume we are working in 2D; we have a Gaussian with mean $\boldsymbol{\mu}$ (which is a 2D vector), and covariance Σ (which is a 2x2 matrix). We could plot the collection of points \mathbf{x} that has some fixed value of $p(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$. This set of points is given by:

$$\frac{1}{2} ((\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})) = c^2$$

where c is some constant. I will choose $c^2 = \frac{1}{2}$, because the choice doesn't matter, and this choice simplifies some algebra. You might recall that a set of points \mathbf{x} that satisfies a quadratic like this is a conic section. Because Σ (and so Σ^{-1}) is positive definite, the curve is an ellipse. There is a useful relationship between the geometry of this ellipse and the Gaussian.

This ellipse — like all ellipses — has a major axis and a minor axis. These are at right angles, and meet at the center of the ellipse. We can determine the properties of the ellipse in terms of the Gaussian quite easily. The geometry of the ellipse isn't affected by rotation or translation, so we will translate the ellipse so that $\mu = \mathbf{0}$ (i.e. the mean is at the origin) and rotate it so that Σ^{-1} is diagonal. Writing $\mathbf{x} = [x, y]$ we get that the set of points on the ellipse satisfies

$$\frac{1}{2} \left(\frac{1}{k_1^2} x^2 + \frac{1}{k_2^2} y^2 \right) = \frac{1}{2}$$

where $\frac{1}{k_1^2}$ and $\frac{1}{k_2^2}$ are the diagonal elements of Σ^{-1} . We will assume that the ellipse has been rotated so that $k_1 < k_2$. The points $(k_1, 0)$ and $(-k_1, 0)$ lie on the ellipse, as do the points $(0, k_2)$ and $(0, -k_2)$. The major axis of the ellipse, in this coordinate system, is the x-axis, and the minor axis is the y-axis. In this coordinate system, x and y are independent. If you do a little algebra, you will see that the standard deviation of x is $\text{abs}(k_1)$ and the standard deviation of y is $\text{abs}(k_2)$. So the ellipse is longer in the direction of largest standard deviation and shorter in the direction of smallest standard deviation.

Now rotating the ellipse means we will pre- and post-multiply the covariance matrix with some rotation matrix. Translating it will move the origin to the mean. As a result, the ellipse has its center at the mean, its major axis is in the direction of the eigenvector of the covariance with largest eigenvalue, and its minor axis is in the direction of the eigenvector with smallest eigenvalue. A plot of this ellipse, which can be coaxed out of most programming environments with relatively little effort, gives us a great deal of information about the underlying Gaussian. These ellipses are known as **covariance ellipses**.

6.2 MIXTURE MODELS AND CLUSTERING

It is natural to think of clustering in the following way. The data was created by a collection of distinct models (one per cluster). For each data item, something (nature?) chose which model was to produce a point, and then the model produced a point. We see the results: crucially, we'd like to know what the models were, but we don't know which model produced which point. If we knew the models, it would be easy to decide which model produced which point. Similarly, if we knew which point went to which model, we could determine what the models were.

One encounters this situation — or problems that can be mapped to this situation — again and again. It is very deeply embedded in clustering problems. It is pretty clear that a natural algorithm is to iterate between estimating which model gets which point, and the model parameters. We have seen this approach before, in the case of k-means.

A particularly interesting case occurs when the models are probabilistic. There is a standard, and very important, algorithm for estimation here, called **EM** (or **expectation maximization**, if you want the long version). I will develop this algorithm in two simple cases, and we will see it in a more general form later.

Notation: This topic lends itself to a glorious festival of indices, limits of sums and products, etc. I will do one example in quite gory detail; the other follows the same form, and for that we'll proceed more expeditiously. Writing the

limits of sums or products explicitly is usually even more confusing than adopting a compact notation. When I write \sum_i or \prod_i , I mean a sum (or product) over all values of i . When I write $\sum_{i,\hat{j}}$ or $\prod_{i,\hat{j}}$, I mean a sum (or product) over all values of i *except* for the j 'th item. I will write vectors, as usual, as \mathbf{x} ; the i 'th such vector in a collection is \mathbf{x}_i , and the k 'th component of the i 'th vector in a collection is x_{ik} . In what follows, I will construct a vector δ_i corresponding to the i 'th data item \mathbf{x}_i (it will tell us what cluster that item belongs to). I will write δ to mean all the δ_i (one for each data item). The j 'th component of this vector is δ_{ij} . When I write \sum_{δ_u} , I mean a sum over all values that δ_u can take. When I write \sum_{δ} , I mean a sum over all values that each δ can take. When I write \sum_{δ,δ_v} , I mean a sum over all values that all δ can take, *omitting* all cases for the v 'th vector δ_v .

6.2.1 A Finite Mixture of Blobs

A blob of data points is quite easily modelled with a single normal distribution. Obtaining the parameters is straightforward (estimate the mean and covariance matrix with the usual expressions). Now imagine I have t blobs of data, and I know t . A normal distribution is likely a poor model, but I could think of the data as being produced by t normal distributions. I will assume that each normal distribution has a fixed, *known* covariance matrix Σ , but the mean of each is unknown. Because the covariance matrix is fixed, and *known*, we can compute a factorization $\Sigma = \mathcal{A}\mathcal{A}^T$. The factors must have full rank, because the covariance matrix must be positive definite. This means that we can apply \mathcal{A}^{-1} to all the data, so that each blob covariance matrix (and so each normal distribution) is the identity.

Write μ_j for the mean of the j 'th normal distribution. We can model a distribution that consists of t distinct blobs by forming a weighted sum of the blobs, where the j 'th blob gets weight π_j . We ensure that $\sum_j \pi_j = 1$, so that we can think of the overall model as a probability distribution. We can then model the data as samples from the probability distribution

$$p(\mathbf{x}|\mu_1, \dots, \mu_k, \pi_1, \dots, \pi_k) = \sum_j \pi_j \left[\frac{1}{\sqrt{(2\pi)^d}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_j)^T(\mathbf{x} - \mu_j)\right) \right].$$

The way to think about this probability distribution is that a point is generated by first choosing one of the normal distributions (the j 'th is chosen with probability π_j), then generating a point from that distribution. This is a pretty natural model of clustered data. Each mean is the center of a blob. Blobs with many points in them have a high value of π_j , and blobs with few points have a low value of π_j . We must now use the data points to estimate the values of π_j and μ_j (again, I am assuming that the blobs – and the normal distribution modelling each – have the identity as a covariance matrix). A distribution of this form is known as a **mixture of normal distributions**.

Writing out the likelihood will reveal a problem: we have a product of many sums. The usual trick of taking the log will not work, because then you have a sum of logs of sums, which is hard to differentiate and hard to work with. A much more productive approach is to think about a set of hidden variables which tell us which blob each data item comes from. For the i 'th data item, we construct a vector δ_i . The j 'th component of this vector is δ_{ij} , where $\delta_{ij} = 1$ if \mathbf{x}_i comes from blob

(equivalently, normal distribution) j and zero otherwise. Notice there is exactly one 1 in δ_i , because each data item comes from one blob. I will write δ to mean all the δ_i (one for each data item). Assume we know the values of these terms. I will write $\theta = (\mu_1, \dots, \mu_k, \pi_1, \dots, \pi_k)$ for the unknown parameters. Then we can write

$$p(\mathbf{x}_i | \delta_i, \theta) = \prod_j \left[\frac{1}{\sqrt{(2\pi)^d}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) \right]^{\delta_{ij}}$$

(because $\delta_{ij} = 1$ means that \mathbf{x}_i comes from blob j , so the terms in the product are a collection of 1's and the probability we want). We also have

$$p(\delta_{ij} = 1 | \theta) = \pi_j$$

allowing us to write

$$p(\delta_i | \theta) = \prod_j [\pi_j]^{\delta_{ij}}$$

(because this is the probability that we select blob j to produce a data item; again, the terms in the product are a collection of 1's and the probability we want). This means that

$$p(\mathbf{x}_i, \delta_i | \theta) = \prod_j \left\{ \left[\frac{1}{\sqrt{(2\pi)^d}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) \right] \pi_j \right\}^{\delta_{ij}}$$

and we can write a log-likelihood. The data are the observed values of \mathbf{x} and δ (remember, we pretend we know these; I'll fix this in a moment), and the parameters are the unknown values of μ_1, \dots, μ_k and π_1, \dots, π_k . We have

$$\begin{aligned} \mathcal{L}(\mu_1, \dots, \mu_k, \pi_1, \dots, \pi_k; \mathbf{x}, \delta) &= \mathcal{L}(\theta; \mathbf{x}, \delta) \\ &= \sum_{ij} \left\{ \left[\left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) \right] + \log \pi_j \right\} \delta_{ij} + K \end{aligned}$$

where K is a constant that absorbs the normalizing constants for the normal distributions. You should check this expression gives the right answer. I have used the δ_{ij} as a “switch” – for one term, $\delta_{ij} = 1$ and the term in curly brackets is “on”, and for all others that term is multiplied by zero. The problem with all this is that we don't know δ . I will deal with this when we have another example.

6.2.2 Topics and Topic Models

A real attraction of probabilistic clustering methods is that we can cluster data where there isn't a clear distance function. One example occurs in document processing. For many kinds of document, we obtain a good representation by (a) choosing a list of different words then (b) representing the document by a vector of word counts, where we simply ignore every word outside the list. This is a viable representation for many applications because quite often, most of the words people actually use come from a relatively short list (typically 100s to 1000s, depending on the particular application). The vector has one component for each word in the

list, and that component contains the number of times that particular word is used. The problem is to cluster the documents.

It isn't a particularly good idea to cluster on the distance between word vectors. This is because quite small changes in word use might lead to large differences between count vectors. For example, some authors might write "car" when others write "auto". In turn, two documents might have a large (resp. small) count for "car" and a small (resp. large) count for "auto". Just looking at the counts would significantly overstate the difference between the vectors. However, the counts are informative: a document that uses the word "car" often, and the word "lipstick" seldom, is likely quite different from a document that uses "lipstick" often and "car" seldom.

We get a useful notion of the differences between documents by pretending that the count vector for each document comes from one of a small set of underlying topics. Each topic generates words as independent, identically distributed samples from a multinomial distribution, with one probability per word in the vocabulary. You should think of each topic as being like a cluster center. If two documents come from the same topic, they should have "similar" word distributions. Topics are one way to deal with changes in word use. For example, one topic might have quite high probability of generating the word "car" *and* a high probability of generating the word "auto"; another might have low probability of generating those words, but a high probability of generating "lipstick".

We cluster documents together if they come from the same topic. Imagine we know which document comes from which topic. Then we could estimate the word probabilities using the documents in each topic. Now imagine we know the word probabilities for each topic. Then we could tell (at least in principle) which topic a document comes from by looking at the probability each topic generates the document, and choosing the topic with the highest probability. This should strike you as being a circular argument. It has a form you should recognize from k-means, though the details of the distance have changed.

To construct a probabilistic model, we will assume that a document is generated in two steps. We will have t topics. First, we choose a topic, choosing the j 'th topic with probability π_j . Then we will obtain a set of words by repeatedly drawing IID samples from that topic, and record the count of each word in a count vector. Assume we have N vectors of word counts, and write \mathbf{x}_i for the i 'th such vector. Each topic is a multinomial probability distribution. Write \mathbf{p}_j for the vector of word probabilities for the j 'th topic. We assume that words are generated independently, conditioned on the topic. Write x_{ik} for the k 'th component of \mathbf{x}_i , and so on. Notice that $\mathbf{x}_i^T \mathbf{1}$ is the sum of entries in \mathbf{x}_i , and so the number of words in document i . Then the probability of observing the counts in \mathbf{x}_i when the document was generated by topic j is

$$p(\mathbf{x}_i | \mathbf{p}_j) = \left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right) \prod_u p_{ju}^{x_{iu}}.$$

We can now write the probability of observing a document. Again, we write $\theta =$

$(\mathbf{p}_1, \dots, \mathbf{p}_t, \pi_1, \dots, \pi_t)$ for the vector of unknown parameters. We have

$$\begin{aligned} p(\mathbf{x}_i|\theta) &= \sum_l p(\mathbf{x}_i|\text{topic is } l)p(\text{topic is } l|\theta) \\ &= \sum_l \left[\left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right) \prod_u p_{lu}^{x_{iu}} \right] \pi_l. \end{aligned}$$

This model is widely called a **topic model**; be aware that there are many kinds of topic model, and this is a simple one. The expression should look unpromising, in a familiar way. If you write out a likelihood, you will see a product of sums; and if you write out a log-likelihood, you will see a sum of logs of sums. Neither is enticing. We could use the same trick we used for a mixture of normals. Write $\delta_{ij} = 1$ if \mathbf{x}_i comes from topic j , and $\delta_{ij} = 0$ otherwise. Then we have

$$p(\mathbf{x}_i|\delta_{ij} = 1, \theta) = \left[\left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right) \prod_u p_{ju}^{x_{iu}} \right]$$

(because $\delta_{ij} = 1$ means that \mathbf{x}_i comes from topic j). This means we can write

$$p(\mathbf{x}_i|\delta_i, \theta) = \prod_j \left\{ \left[\left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right) \prod_u p_{ju}^{x_{iu}} \right] \right\}^{\delta_{ij}}$$

(because $\delta_{ij} = 1$ means that \mathbf{x}_i comes from topic j , so the terms in the product are a collection of 1's and the probability we want). We also have

$$p(\delta_{ij} = 1|\theta) = \pi_j$$

(because this is the probability that we select topic j to produce a data item), allowing us to write

$$p(\delta_i|\theta) = \prod_j [\pi_j]^{\delta_{ij}}$$

(again, the terms in the product are a collection of 1's and the probability we want). This means that

$$p(\mathbf{x}_i, \delta_i|\theta) = \prod_j \left[\left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right) \prod_u (p_{ju}^{x_{iu}}) \pi_j \right]^{\delta_{ij}}$$

and we can write a log-likelihood. The data are the observed values of \mathbf{x} and δ (remember, we pretend we know these for the moment), and the parameters are the unknown values collected in θ . We have

$$\mathcal{L}(\theta; \mathbf{x}, \delta) = \sum_i \left\{ \sum_j \left[\sum_u x_{iu} \log p_{ju} + \log \pi_j \right] \delta_{ij} \right\} + K$$

where K is a term that contains all the

$$\log \left(\frac{(\mathbf{x}_i^T \mathbf{1})!}{\prod_v x_{iv}!} \right)$$

terms. This is of no interest to us, because it doesn't depend on any of our parameters. It takes a fixed value for each dataset. You should check this expression, noticing that, again, I have used the δ_{ij} as a "switch" – for one term, $\delta_{ij} = 1$ and the term in curly brackets is "on", and for all others that term is multiplied by zero. The problem with all this, as before, is that we don't know δ_{ij} . But there is a recipe.

6.3 THE EM ALGORITHM

There is a straightforward, natural, and very powerful recipe. In essence, we will average out the things we don't know. But this average will depend on our estimate of the parameters, so we will average, then re-estimate parameters, then re-average, and so on. If you lose track of what's going on here, think of the example of k-means with soft weights (section 16.20; this is what the equations for the case of a mixture of normals will boil down to). In this analogy, the δ tell us which cluster center a data item came from. Because we don't know the values of the δ , we assume we have a set of cluster centers; these allow us to make a soft estimate of the δ ; then we use this estimate to re-estimate the centers; and so on.

This is an instance of a general recipe. Recall we wrote θ for a vector of parameters. In the mixture of normals case, θ contained the means and the mixing weights; in the topic model case, it contained the topic distributions and the mixing weights. Assume we have an estimate of the value of this vector, say $\theta^{(n)}$. We could then compute $p(\delta|\theta^{(n)}, \mathbf{x})$. In the mixture of normals case, this is a guide to which example goes to which cluster. In the topic case, it is a guide to which example goes to which topic.

We could use this to compute the expected value of the likelihood with respect to δ . We compute

$$Q(\theta; \theta^{(n)}) = \sum_{\delta} \mathcal{L}(\theta; \mathbf{x}, \delta) p(\delta|\theta^{(n)}, \mathbf{x})$$

(where the sum is over all values of δ). Notice that $Q(\theta; \theta^{(n)})$ is a *function* of θ (because \mathcal{L} was), but now does not have any unknown δ terms in it. This $Q(\theta; \theta^{(n)})$ encodes what we know about δ .

For example, assume that $p(\delta|\theta^{(n)}, \mathbf{x})$ has a single, narrow peak in it, at (say) $\delta = \delta^0$. In the mixture of normals case, this would mean that there is one allocation of points to clusters that is significantly better than all others, given $\theta^{(n)}$. For this example, $Q(\theta; \theta^{(n)})$ will be approximately $\mathcal{L}(\theta; \mathbf{x}, \delta^0)$.

Now assume that $p(\delta|\theta^{(n)}, \mathbf{x})$ is about uniform. In the mixture of normals case, this would mean that any particular allocation of points to clusters is about as good as any other. For this example, $Q(\theta; \theta^{(n)})$ will average \mathcal{L} over all possible δ values with about the same weight for each.

We obtain the next estimate of θ by computing

$$\theta^{(n+1)} = \operatorname{argmax}_{\theta} Q(\theta; \theta^{(n)})$$

and iterate this procedure until it converges (which it does, though I shall not prove that). The algorithm I have described is extremely general and powerful, and is

known as **expectation maximization** or (more usually) **EM**. The step where we compute $Q(\theta; \theta^{(n)})$ is called the **E step**; the step where we compute the new estimate of θ is known as the **M step**.

One trick to be aware of: it is quite usual to ignore additive constants in the log-likelihood, because they have no effect. When you do the E-step, taking the expectation of a constant gets you a constant; in the M-step, the constant can't change the outcome. As a result, I will tend to be careless about it. In the mixture of normals example, below, I've tried to keep track of it; for the mixture of multinomials, I've ignored it.

6.3.1 Example: Mixture of Normals: The E-step

Now let us do the actual calculations for a mixture of normal distributions. The E step requires a little work. We have

$$Q(\theta; \theta^{(n)}) = \sum_{\delta} \mathcal{L}(\theta; \mathbf{x}, \delta) p(\delta | \theta^{(n)}, \mathbf{x})$$

If you look at this expression, it should strike you as deeply worrying. There are a very large number of different possible values of δ . In this case, there are $N \times t$ cases (there is one δ_i for each data item, and each of these can have a one in each of t locations). It isn't obvious how we could compute this average.

But notice

$$p(\delta | \theta^{(n)}, \mathbf{x}) = \frac{p(\delta, \mathbf{x} | \theta^{(n)})}{p(\mathbf{x} | \theta^{(n)})}$$

and let us deal with numerator and denominator separately. For the numerator, notice that the \mathbf{x}_i and the δ_i are independent, identically distributed samples, so that

$$p(\delta, \mathbf{x} | \theta^{(n)}) = \prod_i p(\delta_i, \mathbf{x}_i | \theta^{(n)}).$$

The denominator is slightly more work. We have

$$\begin{aligned} p(\mathbf{x} | \theta^{(n)}) &= \sum_{\delta} p(\delta, \mathbf{x} | \theta^{(n)}) \\ &= \sum_{\delta} \left[\prod_i p(\delta_i, \mathbf{x}_i | \theta^{(n)}) \right] \\ &= \prod_i \left[\sum_{\delta_i} p(\delta_i, \mathbf{x}_i | \theta^{(n)}) \right]. \end{aligned}$$

You should check the last step; one natural thing to do is check with $N = 2$ and

$t = 2$. This means that we can write

$$\begin{aligned} p(\delta|\theta^{(n)}, \mathbf{x}) &= \frac{p(\delta, \mathbf{x}|\theta^{(n)})}{p(\mathbf{x}|\theta^{(n)})} \\ &= \frac{\prod_i p(\delta_i, \mathbf{x}_i|\theta^{(n)})}{\prod_i \left[\sum_{\delta_i} p(\delta_i, \mathbf{x}_i|\theta^{(n)}) \right]} \\ &= \prod_i \frac{p(\delta_i, \mathbf{x}_i|\theta^{(n)})}{\sum_{\delta_i} p(\delta_i, \mathbf{x}_i|\theta^{(n)})} \\ &= \prod_i p(\delta_i|\mathbf{x}_i, \theta^{(n)}) \end{aligned}$$

Now we need to look at the log-likelihood. We have

$$\mathcal{L}(\theta; \mathbf{x}, \delta) = \sum_{ij} \left\{ \left[\left(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^T(\mathbf{x}_i - \mu_j) \right) \right] + \log \pi_j \right\} \delta_{ij} + K.$$

The K term is of no interest – it will result in a constant – but we will try to keep track of it. To simplify the equations we need to write, I will construct a t dimensional vector \mathbf{c}_i for the i 'th data point. The j 'th component of this vector will be

$$\left\{ \left[\left(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^T(\mathbf{x}_i - \mu_j) \right) \right] + \log \pi_j \right\}$$

so we can write

$$\mathcal{L}(\theta; \mathbf{x}, \delta) = \sum_i \mathbf{c}_i^T \delta_i + K.$$

Now all this means that

$$\begin{aligned} \mathcal{Q}(\theta; \theta^{(n)}) &= \sum_{\delta} \mathcal{L}(\theta; \mathbf{x}, \delta) p(\delta|\theta^{(n)}, \mathbf{x}) \\ &= \sum_{\delta} \left(\sum_i \mathbf{c}_i^T \delta_i + K \right) p(\delta|\theta^{(n)}, \mathbf{x}) \\ &= \sum_{\delta} \left(\sum_i \mathbf{c}_i^T \delta_i + K \right) \prod_u p(\delta_u|\theta^{(n)}, \mathbf{x}) \\ &= \sum_{\delta} \left(\mathbf{c}_1^T \delta_1 \prod_u p(\delta_u|\theta^{(n)}, \mathbf{x}) + \dots + \mathbf{c}_N^T \delta_N \prod_u p(\delta_u|\theta^{(n)}, \mathbf{x}) \right). \end{aligned}$$

We can simplify further. We have that $\sum_{\delta_i} p(\delta_i|\mathbf{x}_i, \theta^{(n)}) = 1$, because this is a probability distribution. Notice that, for any index v ,

$$\begin{aligned} \sum_{\delta} \left(\mathbf{c}_v^T \delta_v \prod_u p(\delta_u|\theta^{(n)}, \mathbf{x}) \right) &= \sum_{\delta_v} \left(\mathbf{c}_v^T \delta_v p(\delta_v|\theta^{(n)}, \mathbf{x}) \right) \left[\sum_{\delta, \hat{\delta}_v} \prod_{u, \hat{v}} p(\delta_u|\theta^{(n)}, \mathbf{x}) \right] \\ &= \sum_{\delta_v} \left(\mathbf{c}_v^T \delta_v p(\delta_v|\theta^{(n)}, \mathbf{x}) \right) \end{aligned}$$

So we can write

$$\begin{aligned}
 \mathcal{Q}(\theta; \theta^{(n)}) &= \mathcal{L}(\theta; \mathbf{x}, \delta) p(\delta | \theta^{(n)}, \mathbf{x}) \\
 &= \sum_i \left[\sum_{\delta_i} \mathbf{c}_i^T \delta_i p(\delta_i | \theta^{(n)}, \mathbf{x}) \right] + K \\
 &= \sum_i \left[\left(\sum_j \left\{ \left[\left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) + \log \pi_j \right] w_{ij} \right\} \right) \right] + K
 \end{aligned}$$

where

$$\begin{aligned}
 w_{ij} &= 1p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) + 0p(\delta_{ij} = 0 | \theta^{(n)}, \mathbf{x}) \\
 &= p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}).
 \end{aligned}$$

Now

$$\begin{aligned}
 p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) &= \frac{p(\mathbf{x}, \delta_{ij} = 1 | \theta^{(n)})}{p(\mathbf{x} | \theta^{(n)})} \\
 &= \frac{p(\mathbf{x}, \delta_{ij} = 1 | \theta^{(n)})}{\sum_l p(\mathbf{x}, \delta_{il} = 1 | \theta^{(n)})} \\
 &= \frac{p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)}) \prod_{u, \hat{i}} p(\mathbf{x}_u, \delta_u | \theta)}{(\sum_l p(\mathbf{x}, \delta_{il} = 1 | \theta^{(n)})) \prod_{u, \hat{i}} p(\mathbf{x}_u, \delta_u | \theta)} \\
 &= \frac{p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)})}{\sum_l p(\mathbf{x}, \delta_{il} = 1 | \theta^{(n)})}
 \end{aligned}$$

If the last couple of steps puzzle you, remember we obtained $p(\mathbf{x}, \delta | \theta) = \prod_i p(\mathbf{x}_i, \delta_i | \theta)$. Also, look closely at the denominator; it expresses the fact that the data must have come from somewhere. So the main question is to obtain $p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)})$. But

$$\begin{aligned}
 p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)}) &= p(\mathbf{x}_i | \delta_{ij} = 1, \theta^{(n)}) p(\delta_{ij} = 1 | \theta^{(n)}) \\
 &= \left[\frac{1}{\sqrt{(2\pi)^d}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) \right] \pi_j.
 \end{aligned}$$

Substituting yields

$$p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) = \frac{[\exp(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j))] \pi_j}{\sum_k [\exp(-\frac{1}{2}(\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k))] \pi_k} = w_{ij}.$$

6.3.2 Example: Mixture of Normals: The M-step

The M-step is more straightforward. Recall

$$\mathcal{Q}(\theta; \theta^{(n)}) = \left(\sum_{ij} \left\{ \left[\left(-\frac{1}{2} (\mathbf{x}_i - \mu_j)^T (\mathbf{x}_i - \mu_j) \right) + \log \pi_j \right] w_{ij} + K \right\} \right)$$

and we have to maximise this with respect to μ and π , and the terms w_{ij} are known. This maximization is easy. We compute

$$\mu_j^{(n+1)} = \frac{\sum_i x_i w_{ij}}{\sum_i w_{ij}}$$

and

$$\pi_j^{(n+1)} = \frac{\sum_i w_{ij}}{N}.$$

You should check these expressions by differentiating and setting to zero. When you do so, remember that, because π is a probability distribution, $\sum_j \pi_j = 1$ (otherwise you'll get the wrong answer).

6.3.3 Example: Topic Model: The E-Step

We need to work out two steps. The E step requires a little calculation. We have

$$\begin{aligned} \mathcal{Q}(\theta; \theta^{(n)}) &= \sum_{\delta} \mathcal{L}(\theta; \mathbf{x}, \delta) p(\delta | \theta^{(n)}, \mathbf{x}) \\ &= \sum_{\delta} \left(\sum_{ij} \left\{ \left[\sum_u x_{iu} \log p_{ju} \right] + \log \pi_j \right\} \delta_{ij} \right) p(\delta | \theta^{(n)}, \mathbf{x}) \\ &= \left(\sum_{ij} \left\{ \left[\sum_k x_{i,k} \log p_{j,k} \right] + \log \pi_j \right\} w_{ij} \right) \end{aligned}$$

Here the last two steps follow from the same considerations as in the mixture of normals. The \mathbf{x}_i and δ_i are IID samples, and so the expectation simplifies as in that case. If you're uncertain, rewrite the steps of section 6.3.1. The form of this Q function is the same as that (a sum of $\mathbf{c}_i^T \delta_i$ terms, but using a different expression for \mathbf{c}_i). In this case, as above,

$$\begin{aligned} w_{ij} &= 1p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) + 0p(\delta_{ij} = 0 | \theta^{(n)}, \mathbf{x}) \\ &= p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}). \end{aligned}$$

Again, we have

$$\begin{aligned} p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) &= \frac{p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)})}{p(\mathbf{x}_i | \theta^{(n)})} \\ &= \frac{p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)})}{\sum_l p(\mathbf{x}_i, \delta_{il} = 1 | \theta^{(n)})} \end{aligned}$$

and so the main question is to obtain $p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)})$. But

$$\begin{aligned} p(\mathbf{x}_i, \delta_{ij} = 1 | \theta^{(n)}) &= p(\mathbf{x}_i | \delta_{ij} = 1, \theta^{(n)}) p(\delta_{ij} = 1 | \theta^{(n)}) \\ &= \left[\prod_k p_{j,k}^{x_{i,k}} \right] \pi_j. \end{aligned}$$

Substituting yields

$$p(\delta_{ij} = 1 | \theta^{(n)}, \mathbf{x}) = \frac{\left[\prod_k p_{j,k}^{x_k} \right] \pi_j}{\sum_l \left[\prod_k p_{l,k}^{x_k} \right] \pi_l}$$

6.3.4 Example: Topic Model: The M-step

The M-step is more straightforward. Recall

$$\mathcal{Q}(\theta; \theta^{(n)}) = \left(\sum_{ij} \left\{ \left[\sum_k x_{i,k} \log p_{j,k} \right] + \log \pi_j \right\} w_{ij} \right)$$

and we have to maximise this with respect to μ and π , and the terms w_{ij} are known. This maximization is easy, but remember that the probabilities sum to one, so you need either to use a Lagrange multiplier or to set one probability to $(1 - \text{all others})$. You should get

$$\mathbf{p}_j^{(n+1)} = \frac{\sum_i \mathbf{x}_i w_{ij}}{\sum_i \mathbf{x}_i^T \mathbf{1} w_{ij}}$$

and

$$\pi_j^{(n+1)} = \frac{\sum_i w_{ij}}{N}.$$

You should check these expressions by differentiating and setting to zero.

6.3.5 EM in Practice

The algorithm we have seen is amazingly powerful; I will use it again, ideally with less notation. One could reasonably ask whether it produces a “good” answer. Slightly surprisingly, the answer is yes. The algorithm produces a local minimum of $p(\mathbf{x}|\theta)$, the likelihood of the data conditioned on parameters. This is rather surprising because we engaged in all the activity with δ to avoid directly dealing with this likelihood (which in our cases was an unattractive product of sums). I did not prove this, but it’s true anyway.

There are some practical issues. First, how many cluster centers should there be? Mostly, the answer is a practical one. We are usually clustering data for a reason (vector quantization is a really good reason), and then we search for a k that yields the best results.

Second, how should one start the iteration? This depends on the problem you want to solve, but for the two cases I have described, a rough clustering using k-means usually provides an excellent start. In the mixture of normals problem, you can take the cluster centers as initial values for the means, and the fraction of points in each cluster as initial values for the mixture weights. In the topic model problem, you can cluster the count vectors with k-means, use the overall counts within a cluster to get an initial estimate of the multinomial model probabilities, and use the fraction of documents within a cluster to get mixture weights. You need to be careful here, though. You really don’t want to initialize a topic probability with a zero value for any word (otherwise no document containing that word can ever go into the cluster, which is a bit extreme). For our purposes, it will be enough

to allocate a small value to each zero count, then adjust all the word probabilities to be sure they sum to one. More complicated approaches are possible.

Third, we need to avoid numerical problems in the implementation. Notice that you will be evaluating terms that look like

$$\frac{\pi_k e^{-(\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k)/2}}{\sum_u \pi_u e^{-(\mathbf{x}_i - \mu_u)^T (\mathbf{x}_i - \mu_u)/2}}.$$

Imagine you have a point that is far from all cluster means. If you just blithely exponentiate the negative distances, you could find yourself dividing zero by zero, or a tiny number by a tiny number. This can lead to trouble. There's an easy alternative. Find the center the point is closest to. Now subtract the square of this distance (d_{\min}^2 for concreteness) from all the distances. Then evaluate

$$\frac{\pi_k e^{-\left[(\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) - d_{\min}^2\right]/2}}{\sum_u \pi_u e^{-\left[(\mathbf{x}_i - \mu_u)^T (\mathbf{x}_i - \mu_u) - d_{\min}^2\right]/2}}$$

which is a better way of estimating the same number (notice the $e^{-d_{\min}^2/2}$ terms cancel top and bottom).

The last problem is more substantial. EM will get me to a local minimum of $p(\mathbf{x}|\theta)$, but there might be more than one local minimum. For clustering problems, the usual case is there are lots of them. One doesn't really expect a clustering problem to have a single best solution, as opposed to a lot of quite good solutions. Points that are far from all clusters are a particular source of local minima; placing these points in different clusters yields somewhat different sets of cluster centers, each about as good as the other. It's not usual to worry much about this point. A natural strategy is to start the method in a variety of different places (use k means with different start points), and choose the one that has the best value of Q when it has converged.

However, EM isn't magic. There are problems where computing the expectation is hard, typically because you have to sum over a large number of cases which don't have the nice independence structure that helped in the examples I showed. There are strategies for dealing with this problem — essentially, you can get away with an approximate expectation — but they're beyond our reach at present.

6.4 YOU SHOULD

6.4.1 remember:

Useful facts: Parameters of a Multivariate Normal Distribution . . .	122
New term: covariance ellipses	124
New term: EM	124
New term: expectation maximization	124
New term: mixture of normal distributions	125
New term: topic model	128
New term: expectation maximization	130
New term: EM	130
New term: E step	130
New term: M step	130

PROGRAMMING EXERCISES

- 6.1.** Obtain the activities of daily life dataset from the UC Irvine machine learning website (<https://archive.ics.uci.edu/ml/datasets/Dataset+for+ADL+Recognition+with+Wrist-worn+Accelerometer>) data provided by Barbara Bruno, Fulvio Mastrogiovanni and Antonio Sgorbissa).
- (a) Build a classifier that classifies sequences into one of the 14 activities provided. To make features, you should vector quantize, then use a histogram of cluster centers (as described in the subsection; this gives a pretty explicit set of steps to follow). You will find it helpful to use hierarchical k-means to vector quantize. You may use whatever multi-class classifier you wish, though I'd start with R's decision forest, because it's easy to use and effective. You should report (a) the total error rate and (b) the class confusion matrix of your classifier.
 - (b) Now see if you can improve your classifier by (a) modifying the number of cluster centers in your hierarchical k-means and (b) modifying the size of the fixed length samples that you use.

CHAPTER 7

Regression

Classification tries to predict a class from a data item. **Regression** tries to predict a value. For example, we know the zip code of a house, the square footage of its lot, the number of rooms and the square footage of the house, and we wish to predict its likely sale price. As another example, we know the cost and condition of a trading card for sale, and we wish to predict a likely profit in buying it and then reselling it. As yet another example, we have a picture with some missing pixels – perhaps there was text covering them, and we want to replace it – and we want to fill in the missing values. As a final example, you can think of classification as a special case of regression, where we want to predict either $+1$ or -1 ; this isn't usually the best way to proceed, however. Predicting values is very useful, and so there are many examples like this.

7.1 OVERVIEW

Some formalities are helpful here. In the simplest case, we have a dataset consisting of a set of N pairs (\mathbf{x}_i, y_i) . We think of y_i as the value of some function evaluated at \mathbf{x}_i , but with some random component. This means there might be two data items where the \mathbf{x}_i are the same, and the y_i are different. We refer to the \mathbf{x}_i as **explanatory variables** and the y_i is a **dependent variable**. We regularly say that we are regressing the dependent variable against the explanatory variables. We want to use the examples we have — the **training examples** — to build a model of the dependence between y and \mathbf{x} . This model will be used to predict values of y for new values of \mathbf{x} , which are usually called **test examples**. By far the most important model has the form $y = \mathbf{x}^T \beta + \xi$, where β are some set of parameters we need to choose and ξ are random effects. Now imagine that we have one independent variable. An appropriate choice of \mathbf{x} (details below) will mean that the predictions made by this model will lie on a straight line. Figure 7.1 shows two regressions. The data are plotted with a scatter plot, and the line gives the prediction of the model for each value on the x axis.

We do not guarantee that different values of \mathbf{x} produce different values of y . Data just isn't like this (see the crickets example Figure 7.1). Traditionally, regression produces some representation of a probability distribution for y conditioned on \mathbf{x} , so that we would get (say) some representation of a distribution on the houses likely sale value. The best prediction would then be the expected value of that distribution.

It should be clear that none of this will work if there is not some relationship between the training examples and the test examples. If I collect training data on the height and weight of children, I'm unlikely to get good predictions of the weight of adults from their height. We can be more precise with a probabilistic framework. We think of \mathbf{x}_i as IID samples from some (usually unknown) probability distribution $P(X)$. Then the test examples should also be IID samples from $P(X)$,

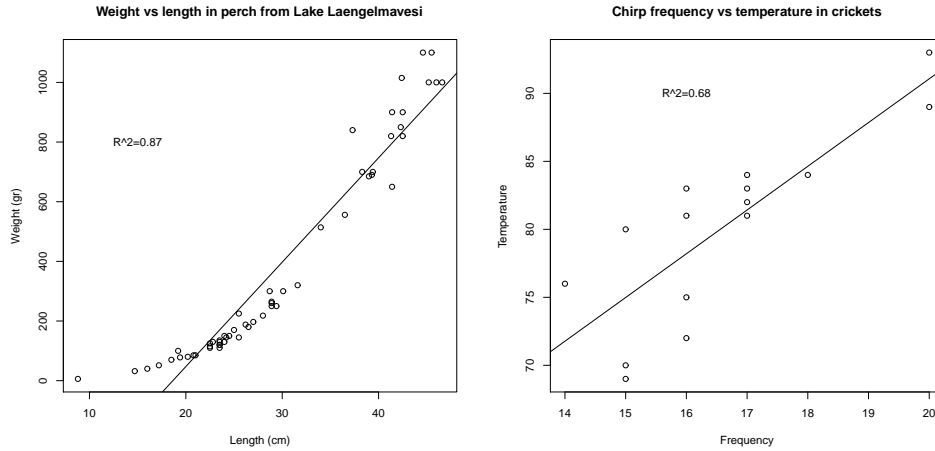


FIGURE 7.1: *On the left*, a regression of weight against length for perch from a Finnish lake (you can find this dataset, and the back story at http://www.amstat.org/publications/jse/jse_data_archive.htm; look for “fishcatch” on that page). Notice that the linear regression fits the data fairly well, meaning that you should be able to predict the weight of a perch from its length fairly well. *On the right*, a regression of air temperature against chirp frequency for crickets. The data is fairly close to the line, meaning that you should be able to tell the temperature from the pitch of cricket’s chirp fairly well. This data is from <http://mste.illinois.edu/patel/amar430/keyprob1.html>. The R^2 you see on each figure is a measure of the goodness of fit of the regression (section 7.2.4).

or, at least, rather like them – you usually can’t check this point with any certainty. A probabilistic formalism can help be precise about the y_i , too. Assume another random variable Y has joint distribution with X given by $P(Y, X)$. We think of each y_i as a sample from $P(Y | \{X = \mathbf{x}_i\})$. Then our modelling problem would be: given the training data, build a model that takes a test example \mathbf{x} and yields a model of $P(Y | \{X = \mathbf{x}_i\})$.

Thinking about the problem this way should make it clear that we’re not relying on any exact, physical, or causal relationship between Y and X . It’s enough that their joint probability makes useful predictions possible, something we will test by experiment. This means that you can build regressions that work in somewhat surprising circumstances. For example, regressing childrens’ reading ability against their foot size can be quite successful. This isn’t because having big feet somehow helps you read; it’s because on the whole, older children read better, and also have bigger feet.

To do anything useful with this formalism requires some aggressive simplifying assumptions. There are very few circumstances that require a comprehensive representation of $P(Y | \{X = \mathbf{x}_i\})$. Usually, we are interested in $\mathbb{E}[Y | \{X = \mathbf{x}_i\}]$ (the mean of $P(Y | \{X = \mathbf{x}_i\})$) and in $\text{var}(\{P(Y | \{X = \mathbf{x}_i\})\})$. To recover this representation, we assume that, for any pair of examples (\mathbf{x}, y) , the value of y is obtained

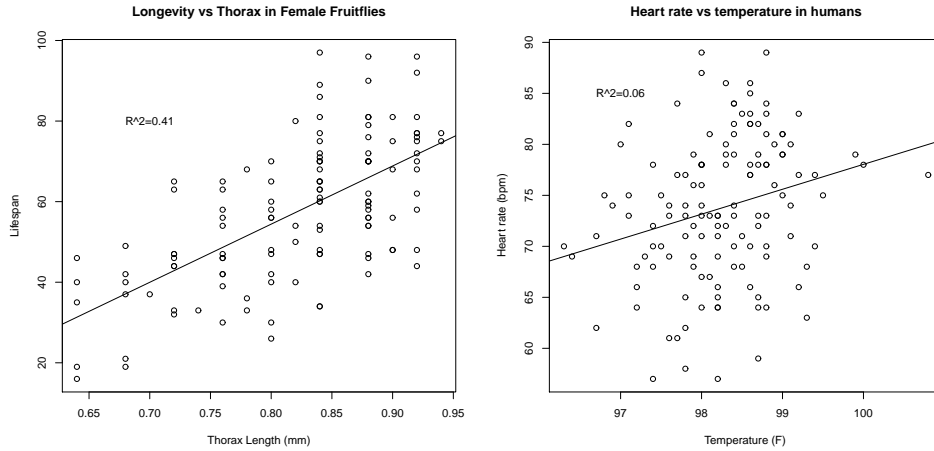


FIGURE 7.2: Regressions do not necessarily yield good predictions or good model fits. On the **left**, a regression of the lifespan of female fruitflies against the length of their torso as adults (apparently, this doesn't change as a fruitfly ages; you can find this dataset, and the back story at http://www.amstat.org/publications/jse/jse_data_archive.htm; look for “fruitfly” on that page). The figure suggests you can make some prediction of how long your fruitfly will last by measuring its torso, but not a particularly accurate one. On the **right**, a regression of heart rate against body temperature for adults. You can find the data at http://www.amstat.org/publications/jse/jse_data_archive.htm as well; look for “temperature” on that page. Notice that predicting heart rate from body temperature isn't going to work that well, either.

by applying some (unknown) function f to \mathbf{x} , then adding some random variable ξ with zero mean. We can write $y(\mathbf{x}) = f(\mathbf{x}) + \xi$, though it's worth remembering that there can be many different values of y associated with a single \mathbf{x} . Now we must make some estimate of f — which yields $\mathbb{E}[Y | \{X = \mathbf{x}_i\}]$ — and estimate the variance of ξ . The variance of ξ might be constant, or might vary with \mathbf{x} .

7.1.1 Regression to Spot Trends

Regression isn't only used to predict values. Another reason to build a regression model is to compare trends in data. Doing so can make it clear what is really happening. Here is an example from Efron (“Computer-Intensive methods in statistical regression”, B. Efron, SIAM Review, 1988). The table in the appendix shows some data from medical devices, which sit in the body and release a hormone. The data shows the amount of hormone currently in a device after it has spent some time in service, and the time the device spent in service. The data describes devices from three production lots (A, B, and C). Each device, from each lot, is supposed to have the same behavior. The important question is: Are the lots the same? The amount of hormone changes over time, so we can't just compare the amounts currently in each device. Instead, we need to determine the relationship between time in service

and hormone, and see if this relationship is different between batches. We can do so by regressing hormone against time.

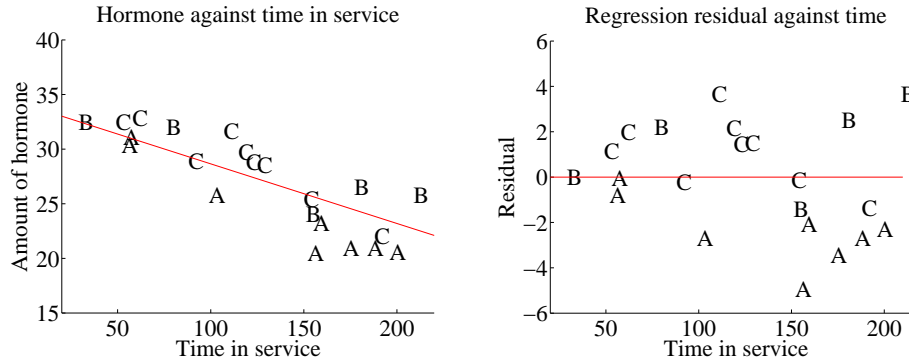


FIGURE 7.3: On the **left**, a scatter plot of hormone against time for devices from tables 7.1 and 7.1. Notice that there is a pretty clear relationship between time and amount of hormone (the longer the device has been in service the less hormone there is). The issue now is to understand that relationship so that we can tell whether lots A, B and C are the same or different. The best fit line to all the data is shown as well, fitted using the methods of section 7.2. On the **right**, a scatter plot of residual — the distance between each data point and the best fit line — against time for the devices from tables 7.1 and 7.1. Now you should notice a clear difference; some devices from lots B and C have positive and some negative residuals, but all lot A devices have negative residuals. This means that, when we account for loss of hormone over time, lot A devices still have less hormone in them. This is pretty good evidence that there is a problem with this lot.

Figure 7.3 shows how a regression can help. In this case, we have modelled the amount of hormone in the device as

$$a \times (\text{time in service}) + b$$

for a , b chosen to get the best fit (much more on this point later!). This means we can plot each data point on a scatter plot, together with the best fitting line. This plot allows us to ask whether any particular batch behaves differently from the overall model in any interesting way.

However, it is hard to evaluate the distances between data points and the best fitting line by eye. A sensible alternative is to subtract the amount of hormone predicted by the model from the amount that was measured. Doing so yields a **residual** — the difference between a measurement and a prediction. We can then plot those residuals (Figure 7.3). In this case, the plot suggests that lot A is special — all devices from this lot contain less hormone than our model predicts.

Definition: 7.1 *Regression*

Regression accepts a feature vector and produces a prediction, which is usually a number, but can sometimes have other forms. You can use these predictions as predictions, or to study trends in data. It is possible, but not usually particularly helpful, to see classification as a form of regression.

7.2 LINEAR REGRESSION AND LEAST SQUARES

Assume we have a dataset consisting of a set of N pairs (\mathbf{x}_i, y_i) . We think of y_i as the value of some function evaluated at \mathbf{x}_i , with some random component added. This means there might be two data items where the \mathbf{x}_i are the same, and the y_i are different. We refer to the \mathbf{x}_i as **explanatory variables** and the y_i is a **dependent variable**. We want to use the examples we have — the **training examples** — to build a model of the dependence between y and \mathbf{x} . This model will be used to predict values of y for new values of \mathbf{x} , which are usually called **test examples**. It can also be used to understand the relationships between the \mathbf{x} . The model needs to have some probabilistic component; we do not expect that y is a function of \mathbf{x} , and there is likely some error in evaluating y anyhow.

7.2.1 Linear Regression

We cannot expect that our model makes perfect predictions. Furthermore, y may not be a function of \mathbf{x} — it is quite possible that the same value of \mathbf{x} could lead to different y 's. One way that this could occur is that y is a measurement (and so subject to some measurement noise). Another is that there is some randomness in y . For example, we expect that two houses with the same set of features (the \mathbf{x}) might still sell for different prices (the y 's).

A good, simple model is to assume that the dependent variable (i.e. y) is obtained by evaluating a linear function of the explanatory variables (i.e. \mathbf{x}), then adding a zero-mean normal random variable. We can write this model as

$$y = \mathbf{x}^T \beta + \xi$$

where ξ represents random (or at least, unmodelled) effects. We will always assume that ξ has zero mean. In this expression, β is a vector of weights, which we must estimate. When we use this model to predict a value of y for a particular set of explanatory variables \mathbf{x}^* , we cannot predict the value that ξ will take. Our best available prediction is the mean value (which is zero). Notice that if $\mathbf{x} = \mathbf{0}$, the model predicts $y = 0$. This may seem like a problem to you — you might be concerned that we can fit only lines through the origin — but remember that \mathbf{x} contains explanatory variables, and we can choose what appears in \mathbf{x} . The two

examples show how a sensible choice of \mathbf{x} allows us to fit a line with an arbitrary y -intercept.

Definition: 7.2 *Linear regression*

A linear regression takes the feature vector \mathbf{x} and predicts $\mathbf{x}^T\beta$, for some vector of coefficients β . The coefficients are adjusted, using data, to produce the best predictions.

Example: 7.1 *A linear model fitted to a single explanatory variable*

Assume we fit a linear model to a single explanatory variable. Then the model has the form $y = x\beta + \xi$, where ξ is a zero mean random variable. For any value x^* of the explanatory variable, our best estimate of y is βx^* . In particular, if $x^* = 0$, the model predicts $y = 0$, which is unfortunate. We can draw the model by drawing a line through the origin with slope β in the x, y plane. The y -intercept of this line must be zero.

Example: 7.2 *A linear model with a non-zero y -intercept*

Assume we have a single explanatory variable, which we write u . We can then create a *vector* $\mathbf{x} = [u, 1]^T$ from the explanatory variable. We now fit a linear model to this vector. Then the model has the form $y = \mathbf{x}^T\beta + \xi$, where ξ is a zero mean random variable. For any value $\mathbf{x}^* = [u^*, 1]^T$ of the explanatory variable, our best estimate of y is $(\mathbf{x}^*)^T\beta$, which can be written as $y = \beta_1 u^* + \beta_2$. If $x^* = 0$, the model predicts $y = \beta_2$. We can draw the model by drawing a line through the origin with slope β_1 and y -intercept β_2 in the x, y plane.

7.2.2 Choosing β

We must determine β . We can proceed in two ways. I show both because different people find different lines of reasoning more compelling. Each will get us to the same solution. One is probabilistic, the other isn't. Generally, I'll proceed as if they're interchangeable, although at least in principle they're different.

Probabilistic approach: we could assume that ξ is a zero mean normal random variable with unknown variance. Then $P(y|x, \beta)$ is normal, with mean $\mathbf{x}^T\beta$, and so we can write out the log-likelihood of the data. Write σ^2 for the variance of ξ , which we don't know, but will not worry about right now. We have that

$$\begin{aligned}\log \mathcal{L}(\beta) &= -\sum_i \log P(y_i|\mathbf{x}_i, \beta) \\ &= \frac{1}{2\sigma^2} \sum_i (y_i - \mathbf{x}_i^T\beta)^2 + \text{term not depending on } \beta\end{aligned}$$

Maximizing the log-likelihood of the data is equivalent to minimizing the negative log-likelihood of the data. Furthermore, the term $\frac{1}{2\sigma^2}$ does not affect the location of the minimum, so we must have that β minimizes $\sum_i (y_i - \mathbf{x}_i^T\beta)^2$, or anything proportional to it. It is helpful to minimize an expression that is an average of squared errors, because (hopefully) this doesn't grow much when we add data. We therefore minimize

$$\left(\frac{1}{N}\right) \left(\sum_i (y_i - \mathbf{x}_i^T\beta)^2\right).$$

Direct approach: notice that, if we have an estimate of β , we have an estimate of the values of the unmodelled effects ξ_i for each example. We just take $\xi_i = y_i - \mathbf{x}_i^T\beta$. It is quite natural to make the unmodelled effects “small”. A good measure of size is the mean of the squared values, which means we want to minimize

$$\left(\frac{1}{N}\right) \left(\sum_i (y_i - \mathbf{x}_i^T\beta)^2\right).$$

We can write all this more conveniently using vectors and matrices. Write \mathbf{y} for the vector

$$\begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}$$

and \mathcal{X} for the matrix

$$\begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \dots \mathbf{x}_n^T \end{pmatrix}.$$

Then we want to minimize

$$\left(\frac{1}{N}\right) (\mathbf{y} - \mathcal{X}\beta)^T (\mathbf{y} - \mathcal{X}\beta)$$

which means that we must have

$$\mathcal{X}^T \mathcal{X}\beta - \mathcal{X}^T \mathbf{y} = 0.$$

For reasonable choices of features, we could expect that $\mathcal{X}^T \mathcal{X}$ — which should strike you as being a lot like a covariance matrix — has full rank. If it does, which

Listing 7.1: R code used for the linear regression example of worked example 7.1

```

efd<-read.table('efrontable.txt',header=TRUE)
# the table has the form
#N1 Ah Bh Ch N2 At Bt Ct
# now we need to construct a new dataset
hor<-stack(efd, select=2:4)
tim<-stack(efd, select=6:8)
foo<-data.frame(time=tim[, c("values")],
                hormone=hor[, c("values")])
foo.lm<-lm(hormone~time, data=foo)
plot(foo)
abline(foo.lm)

```

is the usual case, this equation is easy to solve. If it does not, there is more to do, which we will do in section 7.4.2.

Remember this: *The vector of coefficients β for a linear regression is usually estimated using a least-squares procedure.*

Worked example 7.1 *Simple Linear Regression with R*

Regress the hormone data against time for all the devices in the Efron example.

Solution: This example is mainly used to demonstrate how to regress in R. There is sample code in listing 7.1. The summary in the listing produces a great deal of information (try it). Most of it won't mean anything to you yet. You can get a figure by doing `plot(foo.lm)`, but these figures will not mean anything yet, either. In the code, I've shown how to plot the data and a line on top of it.

7.2.3 Residuals

Assume we have produced a regression by solving

$$\mathcal{X}^T \mathcal{X} \hat{\beta} - \mathcal{X}^T \mathbf{y} = 0$$

for the value of $\hat{\beta}$. I write $\hat{\beta}$ because this is an *estimate*; we likely don't have the true value of the β that generated the data (the model might be wrong; etc.). We cannot expect that $\mathcal{X} \hat{\beta}$ is the same as \mathbf{y} . Instead, there is likely to be some error. The **residual** is the vector

$$\mathbf{e} = \mathbf{y} - \mathcal{X} \hat{\beta}$$

which gives the difference between the true value and the model's prediction at each point. Each component of the residual is an estimate of the unmodelled effects for that data point. The **mean square error** is

$$m = \frac{\mathbf{e}^T \mathbf{e}}{N}$$

and this gives the average of the squared error of prediction on the training examples.

Notice that the mean squared error is not a great measure of how good the regression is. This is because the value depends on the units in which the dependent variable is measured. So, for example, if you measure y in meters you will get a different mean squared error than if you measure y in kilometers.

7.2.4 R-squared

There is an important quantitative measure of how good a regression is which doesn't depend on units. Unless the dependent variable is a constant (which would make prediction easy), it has some variance. If our model is of any use, it should explain some aspects of the value of the dependent variable. This means that the variance of the residual should be smaller than the variance of the dependent variable. If the model made perfect predictions, then the variance of the residual should be zero.

We can formalize all this in a relatively straightforward way. We will ensure that \mathcal{X} always has a column of ones in it, so that the regression can have a non-zero y -intercept. We now fit a model

$$\mathbf{y} = \mathcal{X}\beta + \mathbf{e}$$

(where \mathbf{e} is the vector of residual values) by choosing β such that $\mathbf{e}^T \mathbf{e}$ is minimized. Then we get some useful technical results.

Useful Facts: 7.1 *Regression*

We write $\mathbf{y} = \mathcal{X}\hat{\beta} + \mathbf{e}$, where \mathbf{e} is the residual. Assume \mathcal{X} has a column of ones, and $\hat{\beta}$ is chosen to minimize $\mathbf{e}^T \mathbf{e}$. Then we have

1. $\mathbf{e}^T \mathcal{X} = \mathbf{0}$, i.e. that \mathbf{e} is orthogonal to any column of \mathcal{X} . This is because, if \mathbf{e} is not orthogonal to some column of \mathcal{X} , we can increase or decrease the $\hat{\beta}$ term corresponding to that column to make the error smaller. Another way to see this is to notice that $\hat{\beta}$ is chosen to minimize $\frac{1}{N} \mathbf{e}^T \mathbf{e}$, which is $\frac{1}{N} (\mathbf{y} - \mathcal{X}\hat{\beta})^T (\mathbf{y} - \mathcal{X}\hat{\beta})$. Now because this is a minimum, the gradient with respect to $\hat{\beta}$ is zero, so $(\mathbf{y} - \mathcal{X}\hat{\beta})^T (-\mathcal{X}) = -\mathbf{e}^T \mathcal{X} = 0$.
2. $\mathbf{e}^T \mathbf{1} = 0$ (recall that \mathcal{X} has a column of all ones, and apply the previous result).
3. $\mathbf{1}^T (\mathbf{y} - \mathcal{X}\hat{\beta}) = 0$ (same as previous result).
4. $\mathbf{e}^T \mathcal{X}\hat{\beta} = 0$ (first result means that this is true).

Now \mathbf{y} is a one dimensional dataset arranged into a vector, so we can compute $\text{mean}(\{y\})$ and $\text{var}[y]$. Similarly, $\mathcal{X}\hat{\beta}$ is a one dimensional dataset arranged into a vector (its elements are $\mathbf{x}_i^T \hat{\beta}$), as is \mathbf{e} , so we know the meaning of mean and variance for each. We have a particularly important result:

$$\text{var}[y] = \text{var}[\mathcal{X}\hat{\beta}] + \text{var}[e].$$

This is quite easy to show, with a little more notation. Write $\bar{\mathbf{y}} = (1/N)(\mathbf{1}^T \mathbf{y})\mathbf{1}$ for the vector whose entries are all $\text{mean}(\{y\})$; similarly for $\bar{\mathbf{e}}$ and for $\overline{\mathcal{X}\hat{\beta}}$. We have

$$\text{var}[y] = (1/N)(\mathbf{y} - \bar{\mathbf{y}})^T (\mathbf{y} - \bar{\mathbf{y}})$$

and so on for $\text{var}[e_i]$, etc. Notice from the facts that $\bar{\mathbf{y}} = \overline{\mathcal{X}\hat{\beta}}$. Now

$$\begin{aligned} \text{var}[y] &= (1/N) \left(\left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right] + [\mathbf{e} - \bar{\mathbf{e}}] \right)^T \left(\left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right] + [\mathbf{e} - \bar{\mathbf{e}}] \right) \\ &= (1/N) \left(\left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right]^T \left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right] + 2[\mathbf{e} - \bar{\mathbf{e}}]^T \left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right] + [\mathbf{e} - \bar{\mathbf{e}}]^T [\mathbf{e} - \bar{\mathbf{e}}] \right) \\ &= (1/N) \left(\left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right]^T \left[\mathcal{X}\hat{\beta} - \overline{\mathcal{X}\hat{\beta}} \right] + [\mathbf{e} - \bar{\mathbf{e}}]^T [\mathbf{e} - \bar{\mathbf{e}}] \right) \\ &\quad \text{because } \bar{\mathbf{e}} = \mathbf{0} \text{ and } \mathbf{e}^T \mathcal{X}\hat{\beta} = 0 \text{ and } \mathbf{e}^T \mathbf{1} = 0 \\ &= \text{var}[\mathcal{X}\hat{\beta}] + \text{var}[e]. \end{aligned}$$

This is extremely important, because it allows us to think about a regression as explaining variance in \mathbf{y} . As we are better at explaining \mathbf{y} , $\text{var}[e]$ goes down. In

turn, a natural measure of the goodness of a regression is what percentage of the variance of \mathbf{y} it explains. This is known as R^2 (the r-squared measure). We have

$$R^2 = \frac{\text{var}[\mathbf{x}_i^T \hat{\boldsymbol{\beta}}]}{\text{var}[y_i]}$$

which gives some sense of how well the regression explains the training data. Notice that the value of R^2 is not affected by the units of \mathbf{y} (exercises)

Good predictions result in high values of R^2 , and a perfect model will have $R^2 = 1$ (which doesn't usually happen). For example, the regression of figure 7.3 has an R^2 value of 0.87. Figures 7.1 and 7.2 show the R^2 values for the regressions plotted there; notice how better models yield larger values of R^2 . Notice that if you look at the summary that R provides for a linear regression, it will offer you *two* estimates of the value for R^2 . These estimates are obtained in ways that try to account for (a) the amount of data in the regression, and (b) the number of variables in the regression. For our purposes, the differences between these numbers and the R^2 I defined are not significant. For the figures, I computed R^2 as I described in the text above, but if you substitute one of R's numbers nothing terrible will happen.

Remember this: *The quality of predictions made by a regression can be evaluated by looking at the fraction of the variance in the dependent variable that is explained by the regression. This number is called R^2 , and lies between zero and one; regressions with larger values make better predictions.*

7.2.5 Transforming Variables

Sometimes the data isn't in a form that leads to a good linear regression. In this case, transforming explanatory variables, the dependent variable, or both can lead to big improvements. Figure 7.4 shows one example, based on the idea of word frequencies. Some words are used very often in text; most are used seldom. The dataset for this figure consists of counts of the number of times a word occurred for the 100 most common words in Shakespeare's printed works. It was originally collected from a concordance, and has been used to attack a variety of interesting questions, including an attempt to assess how many words Shakespeare knew. This is hard, because he likely knew many words that he didn't use in his works, so one can't just count. If you look at the plot of Figure 7.4, you can see that a linear regression of count (the number of times a word is used) against rank (how common a word is, 1-100) is not really useful. The most common words are used very often, and the number of times a word is used falls off very sharply as one looks at less common words. You can see this effect in the scatter plot of residual against dependent variable in Figure 7.4 — the residual depends rather strongly on the dependent variable. This is an extreme example that illustrates how poor linear regressions can be.

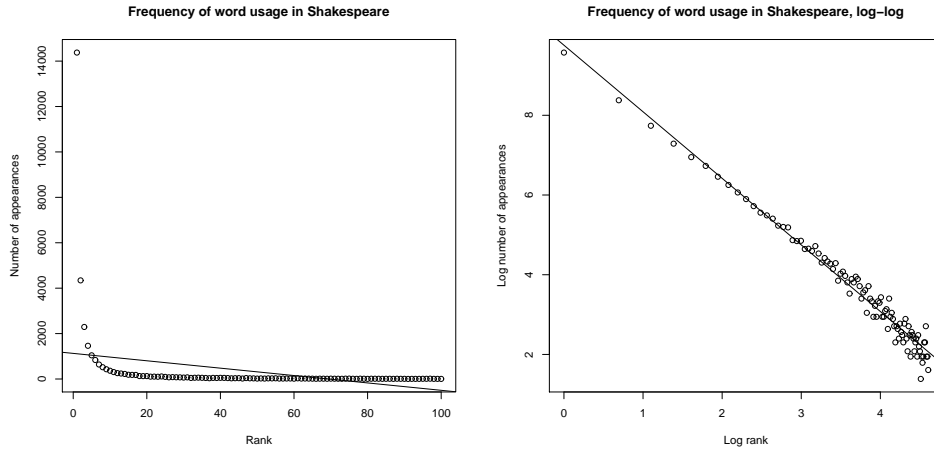


FIGURE 7.4: *On the left*, word count plotted against rank for the 100 most common words in Shakespeare, using a dataset that comes with R (called “bard”, and quite likely originating in an unpublished report by J. Gani and I. Saunders). I show a regression line too. This is a poor fit by eye, and the R^2 is poor, too ($R^2 = 0.1$). *On the right*, log word count plotted against log rank for the 100 most common words in Shakespeare, using a dataset that comes with R (called “bard”, and quite likely originating in an unpublished report by J. Gani and I. Saunders). The regression line is very close to the data.

However, if we regress log-count against log-rank, we get a very good fit indeed. This suggests that Shakespeare’s word usage (at least for the 100 most common words) is consistent with **Zipf’s law**. This gives the relation between frequency f and rank r for a word as

$$f \propto \frac{1}{r^s}$$

where s is a constant characterizing the distribution. Our linear regression suggests that s is approximately 1.67 for this data.

In some cases, the natural logic of the problem will suggest variable transformations that improve regression performance. For example, one could argue that humans have approximately the same density, and so that weight should scale as the cube of height; in turn, this suggests that one regress weight against the cube root of height. Generally, shorter people tend not to be scaled versions of taller people, so the cube root might be too aggressive, and so one thinks of the square root.

Remember this: *The performance of a regression can be improved by transforming variables. Transformations can follow from looking at plots, or thinking about the logic of the problem*

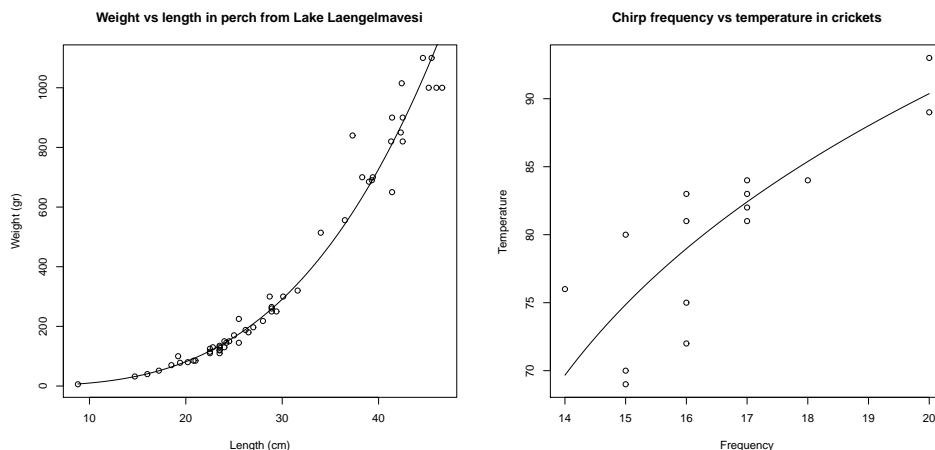


FIGURE 7.5: The Box-Cox transformation suggests a value of $\lambda = 0.303$ for the regression of weight against height for the perch data of Figure 7.1. You can find this dataset, and the back story at http://www.amstat.org/publications/jse/jse_data_archive.htm; look for “fishcatch” on that page). On the **left**, a plot of the resulting curve overlaid on the data. For the cricket temperature data of that figure (from <http://mste.illinois.edu/patel/amar430/keyprob1.html>), the transformation suggests a value of $\lambda = 4.75$. On the **right**, a plot of the resulting curve overlaid on the data.

The **Box-Cox transformation** is a method that can search for a transformation of the dependent variable that improves the regression. The method uses a one-parameter family of transformations, with parameter λ , then searches for the best value of this parameter using maximum likelihood. A clever choice of transformation means that this search is relatively straightforward. We define the Box-Cox transformation of the dependent variable to be

$$y_i^{(bc)} = \begin{cases} \frac{y_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log y_i & \text{if } \lambda = 0 \end{cases} .$$

It turns out to be straightforward to estimate a good value of λ using maximum likelihood. One searches for a value of λ that makes residuals look most like a normal distribution. Statistical software will do it for you; the exercises sketch out the method. This transformation can produce significant improvements in a

regression. For example, the transformation suggests a value of $\lambda = 0.303$ for the fish example of Figure 7.1. It isn't natural to plot $\text{weight}^{0.303}$ against height, because we don't really want to predict $\text{weight}^{0.303}$. Instead, we plot the predictions of weight that come from this model, which will lie on a curve with the form $(ax + b)^{\frac{1}{0.303}}$, rather than on a straight line. Similarly, the transformation suggests a value of $\lambda = 0.475$ for the cricket data. Figure 7.5 shows the result of these transforms.

7.2.6 Can you Trust Your Regression?

Linear regression is useful, but it isn't magic. Some regressions make poor predictions (recall the regressions of figure 7.2). As another example, regressing the first digit of your telephone number against the length of your foot won't work.

We have some straightforward tests to tell whether a regression is working. You can **look at a plot** for a dataset with one explanatory variable and one dependent variable. You plot the data on a scatter plot, then plot the model as a line on that scatterplot. Just looking at the picture can be informative (compare Figure 7.1 and Figure 7.2).

You can check if the regression **predicts a constant**. This is usually a bad sign. You can check this by looking at the predictions for each of the training data items. If the variance of these predictions is small compared to the variance of the independent variable, the regression isn't working well. If you have only one explanatory variable, then you can plot the regression line. If the line is horizontal, or close, then the value of the explanatory variable makes very little contribution to the prediction. This suggests that there is no particular relationship between the explanatory variable and the independent variable.

You can also check, by eye, if **the residual isn't random**. If $y - \mathbf{x}^T\beta$ is a zero mean normal random variable, then the value of the residual vector should not depend on the corresponding y -value. Similarly, if $y - \mathbf{x}^T\beta$ is just a zero mean collection of unmodelled effects, we want the value of the residual vector to not depend on the corresponding y -value either. If it does, that means there is some phenomenon we are not modelling. Looking at a scatter plot of \mathbf{e} against \mathbf{y} will often reveal trouble in a regression (Figure 7.7). In the case of Figure 7.7, the trouble is caused by a few data points that are very different from the others severely affecting the regression. We will discuss how to identify and deal with such points in Section ???. Once they have been removed, the regression improves markedly (Figure 7.8).

Remember this: *Linear regressions can make bad predictions. You can check for trouble by: evaluating R^2 ; looking at a plot; looking to see if the regression makes a constant prediction; or checking whether the residual is random. Other strategies exist, but are beyond the scope of this book.*

Procedure: 7.1 *Linear Regression using Least Squares*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each x_i is a d -dimensional explanatory vector, and each y_i is a single dependent variable. We assume that each data point conforms to the model

$$y_i = \mathbf{x}_i^T \beta + \xi_i$$

where ξ_i represents unmodelled effects. We assume that ξ_i are samples of a random variable with 0 mean and unknown variance. Sometimes, we assume the random variable is normal. Write

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}$$

and

$$\mathcal{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \dots \\ \mathbf{x}_n^T \end{pmatrix}.$$

We estimate $\hat{\beta}$ (the value of β) by solving the linear system

$$\mathcal{X}^T \mathcal{X} \hat{\beta} - \mathcal{X}^T \mathbf{y} = 0.$$

For a data point \mathbf{x} , our model predicts $\mathbf{x}^T \hat{\beta}$. The residuals are

$$\mathbf{e} = \mathbf{y} - \mathcal{X} \hat{\beta}.$$

We have that $\mathbf{e}^T \mathbf{1} = 0$. The mean square error is given by

$$m = \frac{\mathbf{e}^T \mathbf{e}}{N}.$$

The R^2 is given by

$$\frac{\text{var}(\{\mathbf{x}_i^T \hat{\beta}\})}{\text{var}(\{\mathbf{y}\})}.$$

Values of R^2 range from 0 to 1; a larger value means the regression is better at explaining the data.

7.3 PROBLEM DATA POINTS

I have described regressions on a single explanatory variable, because it is easy to plot the line in this case. You can find most problems by looking at the line and the data points. But a single explanatory variable isn't the most common or useful

case. If we have many explanatory variables, it can be hard to plot the regression in a way that exposes problems. This section mainly describes methods to identify and solve difficulties that don't involve looking at the line.

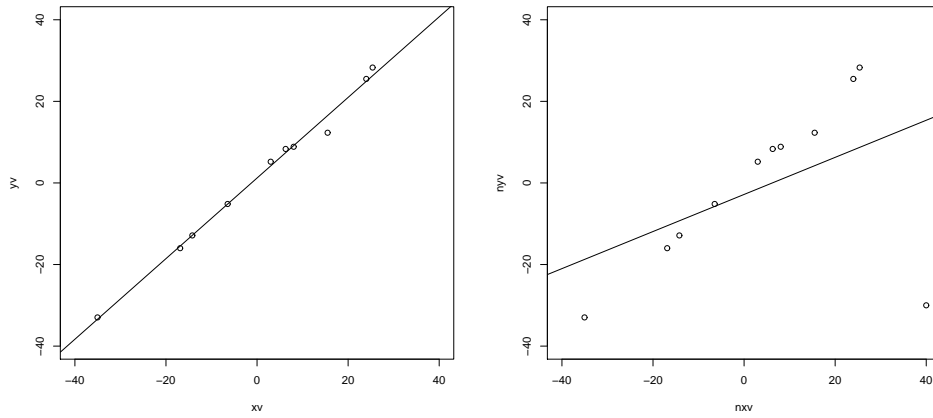


FIGURE 7.6: On the **left**, a synthetic dataset with one independent and one explanatory variable, with the regression line plotted. Notice the line is close to the data points, and its predictions seem likely to be reliable. On the **right**, the result of adding a single outlying datapoint to that dataset. The regression line has changed significantly, because the regression line tries to minimize the sum of squared vertical distances between the data points and the line. Because the outlying datapoint is far from the line, the squared vertical distance to this point is enormous. The line has moved to reduce this distance, at the cost of making the other points further from the line.

7.3.1 Problem Data Points have Significant Impact

Outlying data points can significantly weaken the usefulness of a regression. For some regression problems, we can identify data points that might be a problem, and then resolve how to deal with them. One possibility is that they are true outliers — someone recorded a data item wrong, or they represent an effect that just doesn't occur all that often. Another is that they are important data, and our linear model may not be good enough. If the data points really are outliers, we can drop them from the data set. If they aren't, we may be able to improve the regression by transforming features or by finding a new explanatory variable.

When we construct a regression, we are solving for the β that minimizes $\sum_i (y_i - \mathbf{x}_i^T \beta)^2$, equivalently for the β that produces the smallest value of $\sum_i e_i^2$. This means that residuals with large value can have a very strong influence on the outcome — we are squaring that large value, resulting in an enormous value. Generally, many residuals of medium size will have a smaller cost than one large residual and the rest tiny. As figure 7.6 illustrates, this means that a data point that lies far from the others can swing the regression line significantly.

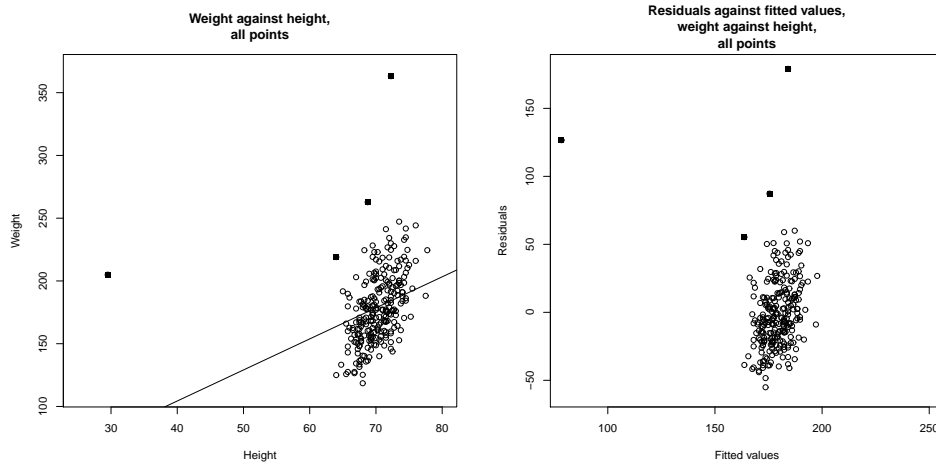


FIGURE 7.7: On the **left**, weight regressed against height for the bodyfat dataset. The line doesn't describe the data particularly well, because it has been strongly affected by a few data points (filled-in markers). On the **right**, a scatter plot of the residual against the value predicted by the regression. This doesn't look like noise, which is a sign of trouble.

This creates a problem, because data points that are very different from most others (sometimes called **outliers**) can also have the highest influence on the outcome of the regression. Figure 7.8 shows this effect for a simple case. When we have only one explanatory variable, there's an easy method to spot problem data points. We produce a scatter plot and a regression line, and the difficulty is usually obvious. In particularly tricky cases, printing the plot and using a see-through ruler to draw a line by eye can help (if you use an opaque ruler, you may not see some errors).

These data points can come from many sources. They may simply be errors. Failures of equipment, transcription errors, someone guessing a value to replace lost data, and so on are some methods that might produce outliers. Another possibility is your understanding of the problem is wrong. If there are some rare effects that are very different than the most common case, you might see outliers. Major scientific discoveries have resulted from investigators taking outliers seriously, and trying to find out what caused them (though you shouldn't see a Nobel prize lurking behind every outlier).

What to do about outliers is even more fraught. The simplest strategy is to find them, then remove them from the data. I will describe some methods that can identify outliers, but you should be aware that this strategy can get dangerous fairly quickly. First, you might find that each time you remove a few problematic data points, some more data points look strange to you. This process is unlikely to end well. Second, you should be aware that throwing out outliers can *increase* your future prediction error, particularly if they're caused by real effects. An alternative strategy is to build methods that can either discount the effects of outliers, or

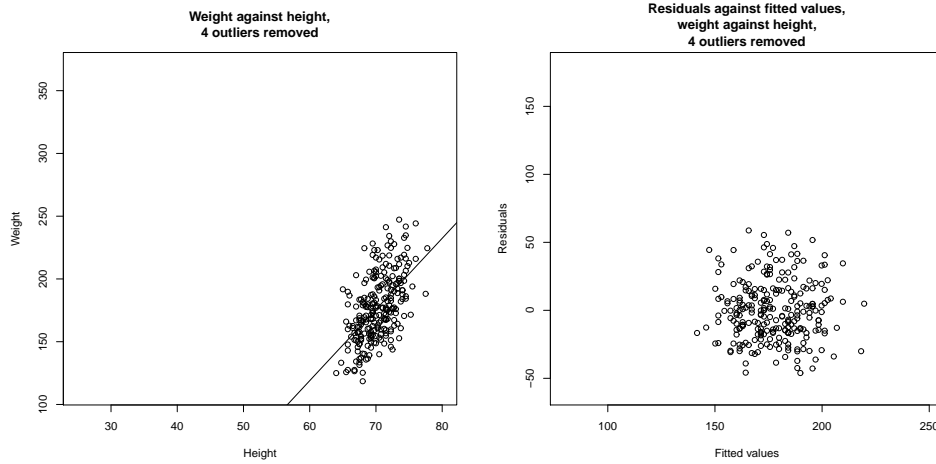


FIGURE 7.8: On the **left**, weight regressed against height for the bodyfat dataset. I have now removed the four suspicious looking data points, identified in Figure 7.7 with filled-in markers; these seemed the most likely to be outliers. On the **right**, a scatter plot of the residual against the value predicted by the regression. Notice that the residual looks like noise. The residual seems to be uncorrelated to the predicted value; the mean of the residual seems to be zero; and the variance of the residual doesn't depend on the predicted value. All these are good signs, consistent with our model, and suggest the regression will yield good predictions.

model them; I describe some such methods, which can be technically complex, in the following chapter.

Remember this: *Outliers can affect linear regressions significantly. Usually, if you can plot the regression, you can look for outliers by eyeballing the plot. Other methods exist, but are beyond the scope of this text.*

7.3.2 The Hat Matrix and Leverage

Write $\hat{\beta}$ for the estimated value of β , and $\mathbf{y}^{(p)} = \mathcal{X}\hat{\beta}$ for the predicted y values. Then we have

$$\hat{\beta} = (\mathcal{X}^T \mathcal{X})^{-1} (\mathcal{X}^T \mathbf{y})$$

so that

$$\mathbf{y}^{(p)} = (\mathcal{X} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T) \mathbf{y}.$$

What this means is that the values the model predicts at training points are a linear function of the true values at the training points. The matrix $(\mathcal{X} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T)$ is

sometimes called the **hat matrix**. The hat matrix is written \mathcal{H} , and I shall write the i, j 'th component of the hat matrix h_{ij} .

Remember this: *The predictions of a linear regression at training points are a linear function of the y -values at the training points. The linear function is given by the hat matrix.*

The hat matrix has a variety of important properties. I won't prove any here, but the proofs are in the exercises. It is a symmetric matrix. The eigenvalues can be only 1 or 0. And the row sums have the important property that

$$\sum_j h_{ij}^2 \leq 1.$$

This is important, because it can be used to find data points that have values that are hard to predict. The **leverage** of the i 'th training point is the i 'th diagonal element, h_{ii} , of the hat matrix \mathcal{H} . Now we can write the prediction at the i 'th training point $y_{p,i} = h_{ii}y_i + \sum_{j \neq i} h_{ij}y_j$. But if h_{ii} has large absolute value, then all the other entries in that row of the hat matrix must have small absolute value. This means that, if a data point has high leverage, the model's value at that point is predicted almost entirely by the observed value at that point. Alternatively, it's hard to use the other training data to predict a value at that point.

Here is another way to see this importance of h_{ii} . Imagine we change the value of y_i by adding Δ ; then $y_i^{(p)}$ becomes $y_i^{(p)} + h_{ii}\Delta$. In turn, a large value of h_{ii} means that the predictions at the i 'th point are very sensitive to the value of y_i .

Remember this: *Ideally, the value predicted for a particular data point depends on many other data points. Leverage measures the importance of a data point in producing a prediction at that data point. If the leverage of a point is high, other points are not contributing much to the prediction for that point, and it may well be an outlier.*

7.3.3 Cook's Distance

Another way to find points that may be creating problems is to look at the effect of omitting the point from the regression. We could compute $\mathbf{y}^{(p)}$ using the whole data set. We then omit the i 'th point from the dataset, compute the regression coefficients from the remaining data (which I will write $\hat{\beta}_i$), then compare $\mathbf{y}^{(p)}$ to $\mathcal{X}\hat{\beta}_i$. If there is a large difference, the point is suspect, because omitting it

strongly changes the predictions. The score for the comparison is called **Cook's distance**. If a point has a large value of Cook's distance, then it has a strong influence on the regression and might well be an outlier. Typically, one computes Cook's distance for each point, and takes a closer look at any point with a large value. This procedure is described in more detail in procedure 96

Notice the rough similarity to cross-validation (omit some data and recompute). But in this case, we are using the procedure to identify points we might not trust, rather than to get an unbiased estimate of the error.

Procedure: 7.2 *Computing Cook's distance*

We have a dataset containing N pairs (\mathbf{x}_i, y_i) . Each x_i is a d -dimensional explanatory vector, and each y_i is a single dependent variable. Write $\hat{\beta}$ for the coefficients of a linear regression (see procedure 7.1), and $\hat{\beta}_i$ for the coefficients of the linear regression computed by omitting the i 'th data point, $\mathbf{y}^{(p)}$ for $\mathcal{X}\hat{\beta}$, and m for the mean square error. The Cook's distance of the i 'th data point is

$$\frac{(\mathbf{y}^{(p)} - \mathcal{X}\hat{\beta}_i)^T (\mathbf{y}^{(p)} - \mathcal{X}\hat{\beta}_i)}{dm}.$$

Large values of this distance suggest a point may present problems. Statistical software will compute and plot this distance for you.

Remember this: *The Cook's distance of a training data point measures the effect on predictions of leaving that point out of the regression. A large value of Cook's distance suggests other points are poor at predicting the value at a given point, so a point with a large value of Cook's distance may be an outlier.*

7.3.4 Standardized Residuals

The hat matrix has another use. It can be used to tell how “large” a residual is. The residuals that we measure depend on the units in which y was expressed, meaning we have no idea what a “large” residual is. For example, if we were to express y in kilograms, then we might want to think of 0.1 as a small residual. Using exactly the same dataset, but now with y expressed in grams, that residual value becomes 100 — is it really “large” because we changed units?

Now recall that we assumed, in section 7.2.1, that $y - \mathbf{x}^T\beta$ was a zero mean normal random variable, but we didn't know its variance. It can be shown that,

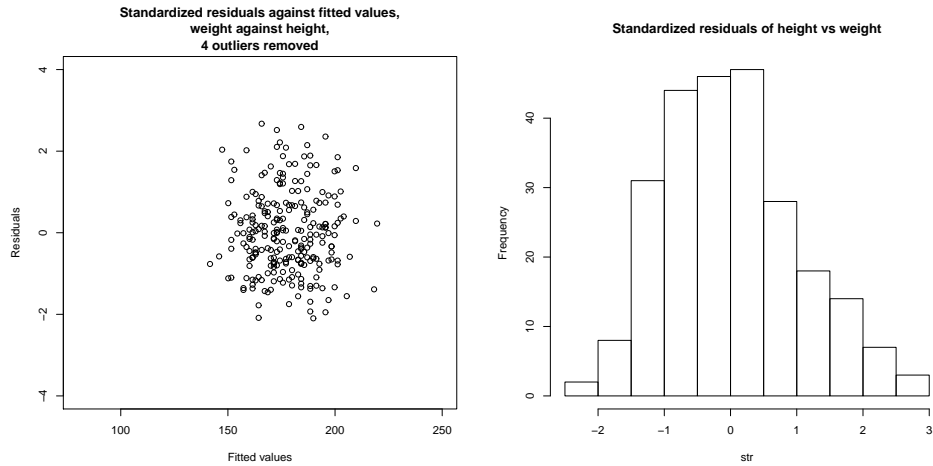


FIGURE 7.9: On the **left**, standardized residuals plotted against predicted value for weight regressed against height for the bodyfat dataset. I removed the four suspicious looking data points, identified in Figure 7.7 with filled-in markers ; these seemed the most likely to be outliers. You should compare this plot with the residuals in figure 7.8, which are not standardized. On the **right**, a histogram of the residual values. Notice this looks rather like a histogram of a standard normal random variable, though there are slightly more large positive residuals than one would like. This suggests the regression is working tolerably.

under our assumption, the i 'th residual value, e_i , is a sample of a normal random variable whose variance is

$$\left(\frac{\mathbf{e}^T \mathbf{e}}{N}\right) (1 - h_{ii}).$$

This means we can tell whether a residual is large by **standardizing** it – that is, dividing by its standard deviation. Write s_i for the standard residual at the i 'th training point. Then we have that

$$s_i = \frac{e_i}{\sqrt{\left(\frac{\mathbf{e}^T \mathbf{e}}{N}\right) (1 - h_{ii})}}.$$

When the regression is behaving, this standard residual should look like a sample of a standard normal random variable. In turn, this means that if all is going well, about 66% of the residuals should have values in the range $[-1, 1]$, and so on. Large values of the standard residuals are a sign of trouble.

R produces a nice diagnostic plot that can be used to look for problem data points (code and details in the appendix). The plot is a scatter plot of the standardized residuals against leverage, with level curves of Cook's distance superimposed. Figure 7.10 shows an example. Some bad points that are likely to present problems are identified with a number (you can control how many, and the number, with arguments to `plot`; appendix). Problem points will have high leverage and/or high

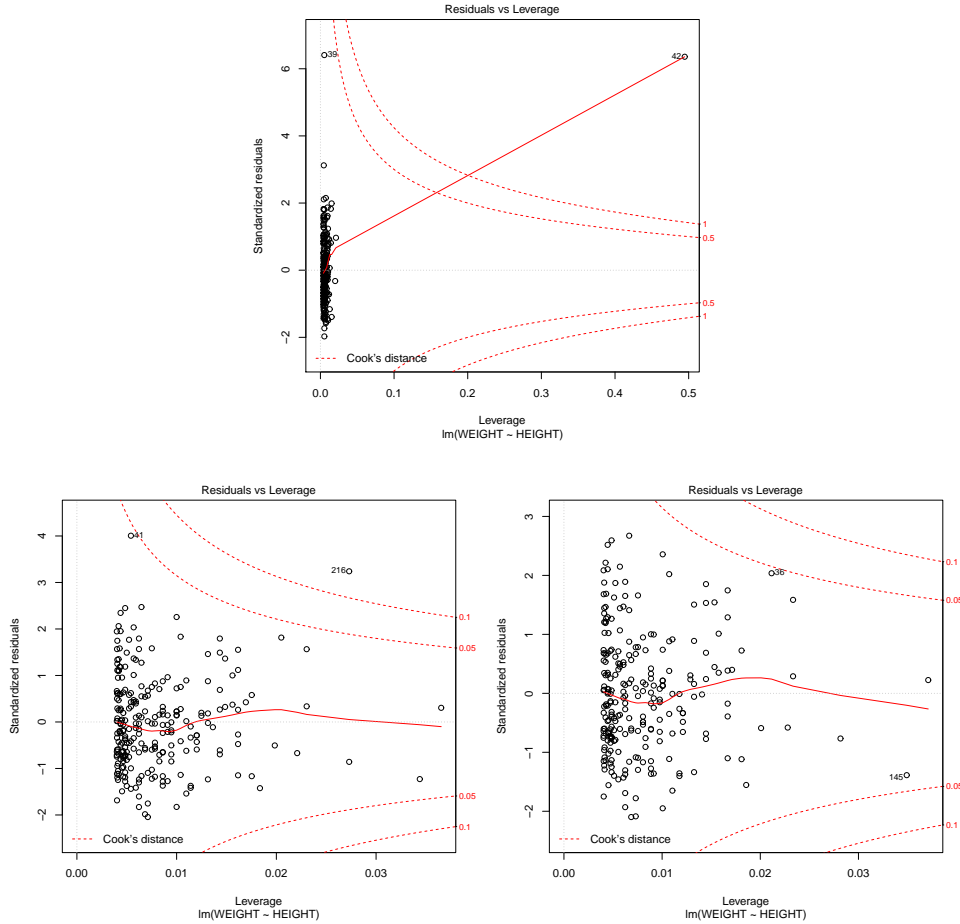


FIGURE 7.10: A diagnostic plot, produced by R, of a linear regression of weight against height for the bodyfat dataset. **Top:** the whole dataset; **bottom left:** with the two most extreme points in the top figure removed; **bottom right:** with two further points (highest residual) removed. Details in text.

Cook's distance and/or high residual. The figure shows this plot for three different versions of the dataset (original; two problem points removed; and two further problem points removed).

7.4 MANY EXPLANATORY VARIABLES

In earlier sections, I implied you could put anything into the explanatory variables. This is correct, and makes it easy to do the math for the general case. However, I have plotted only cases where there was one explanatory variable (together with a constant, which hardly counts). In some cases (section 7.4.1), we can add explana-

tory variables and still have an easy plot. Adding explanatory variables can cause the matrix $\mathcal{X}^T \mathcal{X}$ to have poor condition number; there's an easy strategy to deal with this (section 7.4.2).

Most cases are hard to plot successfully, and one needs better ways to visualize the regression than just plotting. The value of R^2 is still a useful guide to the goodness of the regression, but the way to get more insight is to use the tools of the previous section.

7.4.1 Functions of One Explanatory Variable

Imagine we have only one measurement to form explanatory variables. For example, in the perch data of Figure 7.1, we have only the length of the fish. If we evaluate functions of that measurement, and insert them into the vector of explanatory variables, the resulting regression is still easy to plot. It may also offer better predictions. The fitted line of Figure 7.1 looks quite good, but the data points look as though they might be willing to follow a curve. We can get a curve quite easily. Our current model gives the weight as a linear function of the length with a noise term (which we wrote $y_i = \beta_1 x_i + \beta_0 + \xi_i$). But we could expand this model to incorporate other functions of the length. In fact, it's quite surprising that the weight of a fish should be predicted by its length. If the fish doubled in each direction, say, its weight should go up by a factor of eight. The success of our regression suggests that fish do not just scale in each direction as they grow. But we might try the model $y_i = \beta_2 x_i^2 + \beta_1 x_i + \beta_0 + \xi_i$. This is easy to do. The i 'th row of the matrix \mathcal{X} currently looks like $[x_i, 1]$. We build a new matrix $\mathcal{X}^{(b)}$, where the i 'th row is $[x_i^2, x_i, 1]$, and proceed as before. This gets us a new model. The nice thing about this model is that it is easy to plot – our predicted weight is still a function of the length, it's just not a linear function of the length. Several such models are plotted in Figure 7.11.

You should notice that it can be quite easy to add a lot of functions like this (in the case of the fish, I tried x_i^3 as well). However, it's hard to decide whether the regression has actually gotten better. The least-squares error *on the training data* will *never* go up when you add new explanatory variables, so the R^2 will *never* get worse. This is easy to see, because you could always use a coefficient of zero with the new variables and get back the previous regression. However, the models that you choose are likely to produce worse and worse predictions as you add explanatory variables. Knowing when to stop can be tough (Section 8.1), though it's sometimes obvious that the model is untrustworthy (Figure 7.11).

Remember this: *If you have only one measurement, you can construct a high dimensional \mathbf{x} by using functions of that measurement. This produces a regression that has many explanatory variables, but is still easy to plot. Knowing when to stop is hard. An understanding of the problem is helpful.*

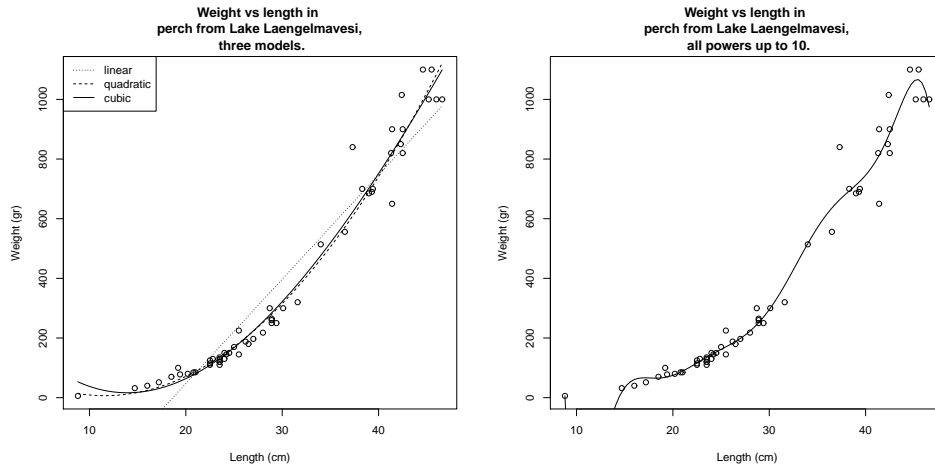


FIGURE 7.11: On the **left**, several different models predicting fish weight from length. The line uses the explanatory variables 1 and x_i ; and the curves use other monomials in x_i as well, as shown by the legend. This allows the models to predict curves that lie closer to the data. It is important to understand that, while you can make a curve go closer to the data by inserting monomials, that doesn't mean you necessarily have a better model. On the **right**, I have used monomials up to x_i^{10} . This curve lies very much closer to the data points than any on the other side, at the cost of some very odd looking wiggles inbetween data points (look at small lengths; the model goes quite strongly negative there, but I can't bring myself to change the axes and show predictions that are obvious nonsense). I can't think of any reason that these structures would come from true properties of fish, and it would be hard to trust predictions from this model.

7.4.2 Regularizing Linear Regressions

When we have many explanatory variables, some might be significantly correlated. This means that we can predict, quite accurately, the value of one explanatory variable using the values of the other variables. This means there must be a vector \mathbf{w} so that $\mathcal{X}\mathbf{w}$ is small (exercises). In turn, that $\mathbf{w}^T \mathcal{X}^T \mathcal{X} \mathbf{w}$ must be small, so that $\mathcal{X}^T \mathcal{X}$ has some small eigenvalues. These small eigenvalues lead to bad predictions, as follows. The vector \mathbf{w} has the property that $\mathcal{X}^T \mathcal{X} \mathbf{w}$ is small. This means that $\mathcal{X}^T \mathcal{X}(\hat{\beta} + \mathbf{w})$ is not much different from $\mathcal{X}^T \mathcal{X} \hat{\beta}$ (equivalently, the matrix can turn large vectors into small ones). All this means that $(\mathcal{X}^T \mathcal{X})^{-1}$ will turn some small vectors into big ones. A small change in $\mathcal{X}^T \mathbf{Y}$ can lead to a large change in the estimate of $\hat{\beta}$.

This is a problem, because we can expect that different samples from the same data will have somewhat different values of $\mathcal{X}^T \mathbf{Y}$. For example, imagine the person recording fish measurements in Lake Laengelmavesi recorded a different set of fish; we expect changes in \mathcal{X} and \mathbf{Y} . But, if $\mathcal{X}^T \mathcal{X}$ has small eigenvalues, these changes could produce large changes in our model.

The problem is relatively easy to control. When there are small eigenvalues

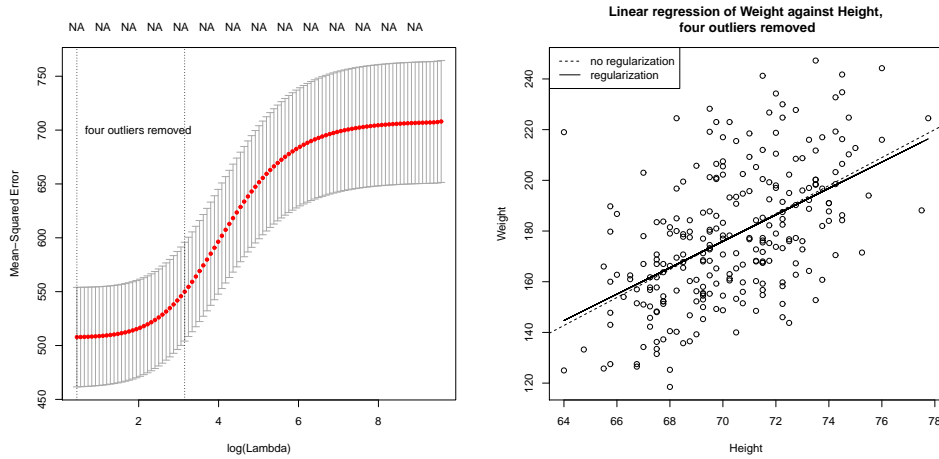


FIGURE 7.12: On the **left**, cross-validated error estimated for different choices of regularization constant for a linear regression of weight against height for the bodyfat dataset, with four outliers removed. The horizontal axis is log regression constant; the vertical is cross-validated error. The mean of the error is shown as a spot, with vertical error bars. The vertical lines show a range of reasonable choices of regularization constant (**left** yields the lowest observed error, **right** the error whose mean is within one standard error of the minimum). On the **right**, two regression lines on a scatter plot of this dataset; one is the line computed without regularization, the other is obtained using the regularization parameter that yields the lowest observed error. In this case, the regularizer doesn't change the line much, but may produce improved values on new data (notice how the cross-validated error is fairly flat with low values of the regularization constant).

in $\mathcal{X}^T \mathcal{X}$, we expect that $\hat{\beta}$ will be large (because we can add components in the direction of \mathbf{w} without changing all that much), and the largest components in $\hat{\beta}$ might be very inaccurately estimated. If we are trying to predict new y values, we expect that large components in $\hat{\beta}$ turn into large errors in prediction (exercises).

An important and useful way to suppress these errors is to try to find a $\hat{\beta}$ that isn't large, and also gives a low error. We can do this by regularizing, using the same trick we saw in the case of classification. Instead of choosing the value of β that minimizes

$$\left(\frac{1}{N}\right) (\mathbf{y} - \mathcal{X}\beta)^T (\mathbf{y} - \mathcal{X}\beta)$$

we minimize

$$\left(\frac{1}{N}\right) (\mathbf{y} - \mathcal{X}\beta)^T (\mathbf{y} - \mathcal{X}\beta) + \lambda \beta^T \beta$$

Error + Regularizer

Here $\lambda > 0$ is a constant that weights the two requirements (small error; small $\hat{\beta}$) relative to one another. Notice also that dividing the total error by the number of

data points means that our choice of λ shouldn't be affected by changes in the size of the data set.

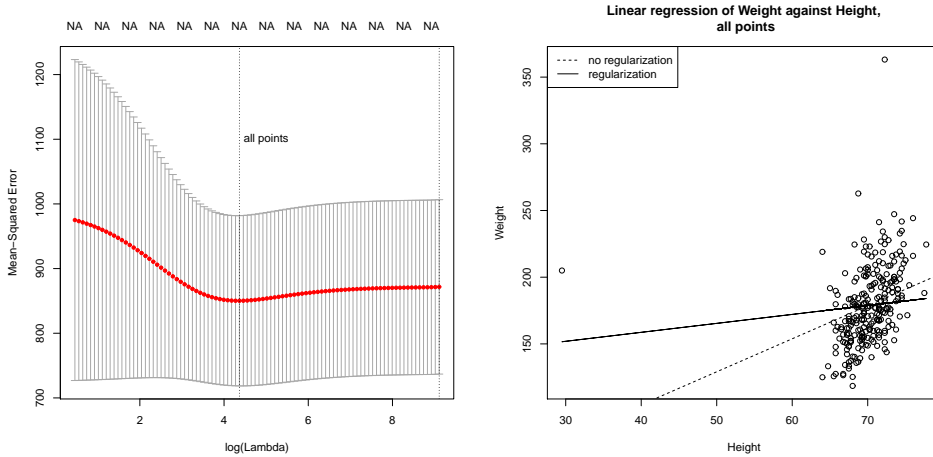


FIGURE 7.13: *Regularization doesn't make outliers go away. On the left, cross-validated error estimated for different choices of regularization constant for a linear regression of weight against height for the bodyfat dataset, with all points. The horizontal axis is log regression constant; the vertical is cross-validated error. The mean of the error is shown as a spot, with vertical error bars. The vertical lines show a range of reasonable choices of regularization constant (left yields the lowest observed error, right the error whose mean is within one standard error of the minimum). On the right, two regression lines on a scatter plot of this dataset; one is the line computed without regularization, the other is obtained using the regularization parameter that yields the lowest observed error. In this case, the regularizer doesn't change the line much, but may produce improved values on new data (notice how the cross-validated error is fairly flat with low values of the regularization constant).*

Regularization helps deal with the small eigenvalue, because to solve for β we must solve the equation

$$\left[\left(\frac{1}{N} \right) \mathcal{X}^T \mathcal{X} + \lambda \mathcal{I} \right] \hat{\beta} = \left(\frac{1}{N} \right) \mathcal{X}^T \mathbf{y}$$

(obtained by differentiating with respect to β and setting to zero) and the smallest eigenvalue of the matrix $\left(\left(\frac{1}{N} \right) (\mathcal{X}^T \mathcal{X} + \lambda \mathcal{I}) \right)$ will be at least λ (exercises). Penalizing a regression with the size of β in this way is sometimes known as **ridge regression**.

We choose λ in the same way we used for classification; split the training set into a training piece and a validation piece, train for different values of λ , and test the resulting regressions on the validation piece. The error is a random variable, random because of the random split. It is a fair model of the error that would occur on a randomly chosen test example (assuming that the training set is “like” the test set, in a way that I do not wish to make precise yet). We could use multiple

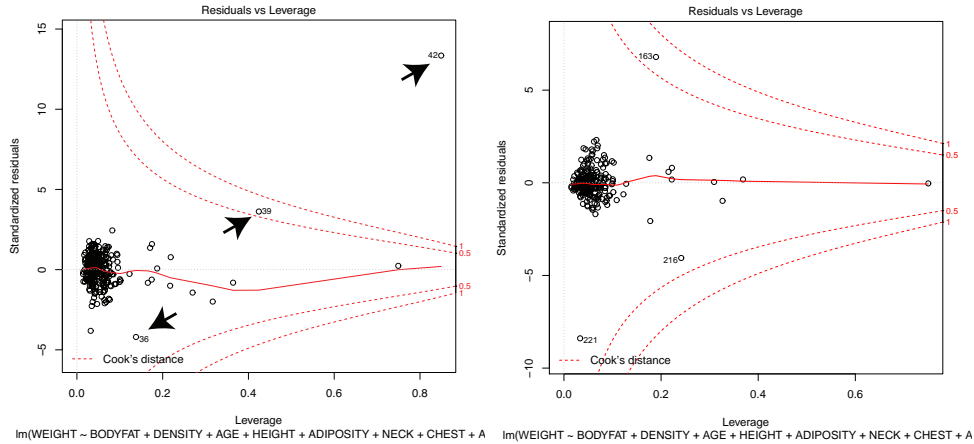


FIGURE 7.14: On the **left**, residuals plotted against leverage for a regression of weight against all other measurements for the bodyfat dataset. I did not remove the outliers. The contours on the plot are contours of Cook's distance; I have overlaid arrows showing points with suspiciously large Cook's distance. Notice also that several points have high leverage, without having a large residual value. These points may or may not present problems. On the **right**, the same plot for this dataset with points 36, 39, 41 and 42 removed (these are the points I have been removing for each such plot). Notice that another point now has high Cook's distance, but mostly the residual is much smaller.

splits, and average over the splits. Doing so yields both an average error for a value of λ and an estimate of the standard deviation of error.

Statistical software will do all the work for you. I used the `glmnet` package in R (see exercises for details). Figure 7.12 shows an example, for weight regressed against height. Notice the regularization doesn't change the model (plotted in the figure) all that much. For each value of λ (horizontal axis), the method has computed the mean error and standard deviation of error using cross-validation splits, and displays these with error bars. Notice that $\lambda = 0$ yields poorer predictions than a larger value; large $\hat{\beta}$ really are unreliable. Notice that now there is now no λ that yields the smallest validation error, because the value of error depends on the random splits used in cross-validation. A reasonable choice of λ lies between the one that yields the smallest error encountered (one vertical line in the plot) and the largest value whose mean error is within one standard deviation of the minimum (the other vertical line in the plot).

All this is quite similar to regularizing a classification problem. We started with a cost function that evaluated the errors caused by a choice of β , then added a term that penalized β for being "large". This term is the squared length of β , as a vector. It is sometimes known as the L_2 **norm** of the vector. In section 16.20, I describe the consequences of using other norms.

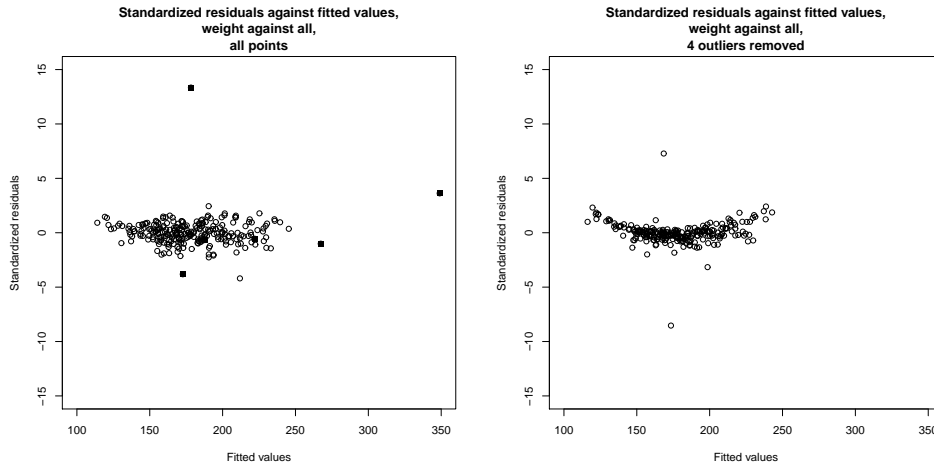


FIGURE 7.15: *On the left, standardized residuals plotted against predicted value for weight regressed against all variables for the bodyfat dataset. Four data points appear suspicious, and I have marked these with a filled in marker. On the right, standardized residuals plotted against predicted value for weight regressed against all variables for the bodyfat dataset, but with the four suspicious looking data points removed. Notice two other points stick out markedly.*

Remember this: *The performance of a regression can be improved by regularizing, particularly if some explanatory variables are correlated. The procedure is similar to that used for classification.*

7.4.3 Example: Weight against Body Measurements

We can now look at regressing weight against all body measurements for the bodyfat dataset. We can't plot this regression (too many independent variables), but we can approach the problem in a series of steps.

Finding suspect points: Figure 7.14 shows the R diagnostic plots for a regression of weight against all body measurements for the bodyfat dataset. We've already seen there are outliers, so the odd structure of this plot should be no particular surprise. There are several really worrying points here. As the figure shows, removing the four points identified in the caption, based on their very high standardized residuals, high leverage, and high Cook's distance, yields improvements. We can get some insight by plotting standardized residuals against predicted value (Figure 7.9). There is clearly a problem here; the residual seems to depend quite strongly on the predicted value. Removing the four outliers we have already identified leads to a much improved plot, also shown in Figure 7.15. This is banana-shaped, which is suspicious. There are two points that seem to come from some

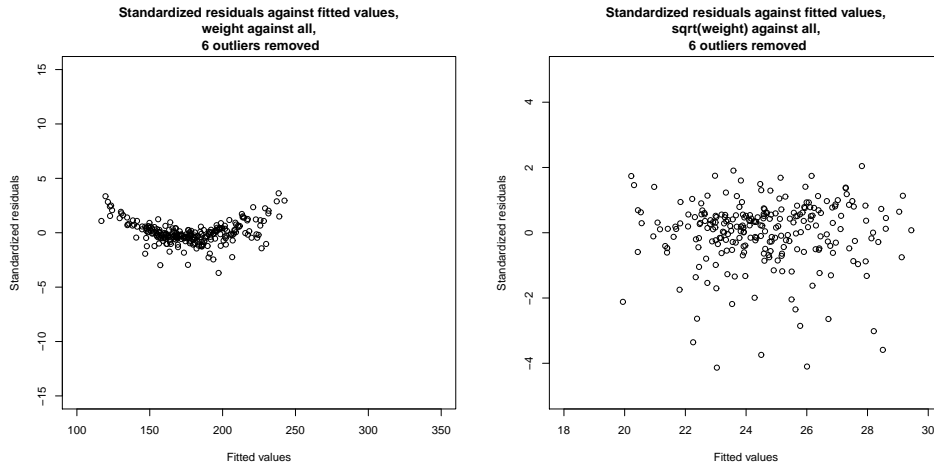


FIGURE 7.16: *On the left*, standardized residuals plotted against predicted value for weight regressed against all variables for the bodyfat dataset. I removed the four suspicious data points of Figure 7.15, and the two others identified in that figure. Notice a suspicious “banana” shape – the residuals are distinctly larger for small and for large predicted values. This suggests a non-linear transformation of something might be helpful. I used a Box-Cox transformation, which suggested a value of 0.5 (i.e. regress $2(\sqrt{\text{weight}} - 1)$) against all variables. *On the right*, the standardized residuals for this regression. Notice that the “banana” has gone, though there is a suspicious tendency for the residuals to be smaller rather than larger. Notice also the plots are on different axes. It’s fair to compare these plots by eye; but it’s not fair to compare details, because the residual of a predicted square root means something different than the residual of a predicted value.

other model (one above the center of the banana, one below). Removing these points gives the residual plot shown in Figure 7.16.

Transforming variables: The banana shape of the plot of standardized residuals against value is a suggestion that some non-linearity somewhere would improve the regression. One option is a non-linear transformation of the independent variables. Finding the right one might require some work, so it’s natural to try a Box-Cox transformation first. This gives the best value of the parameter as 0.5 (i.e. the dependent variable should be $\sqrt{\text{weight}}$, which makes the residuals look much better (Figure 7.16).

Choosing a regularizing value: Figure 7.17 shows the `glmnet` plot of cross-validated error as a function of regularizer weight. A sensible choice of value here seems to be a bit smaller than -2 (between the value that yields the smallest error encountered – one vertical line in the plot – and the largest value whose mean error is within one standard deviation of the minimum – the other vertical line in the plot). I chose -2.2

How good are the resulting predictions likely to be: the standardized residuals don’t seem to depend on the predicted values, but how good are the

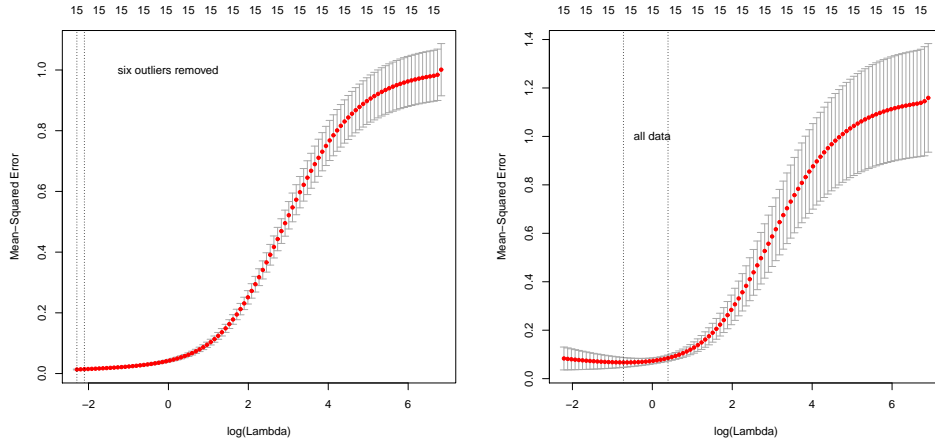


FIGURE 7.17: Plots of mean-squared error as a function of log regularization parameter (i.e. $\log \lambda$) for a regression of $\text{weight}^{1/2}$ against all variables for the bodyfat dataset. These plots show mean-squared error averaged over cross-validation folds with a vertical one standard deviation bar. On the **left**, the plot for the dataset with the six outliers identified in Figure 16.20 removed. On the **right**, the plot for the whole dataset. Notice how the outliers increase the variability of the error, and the best error.

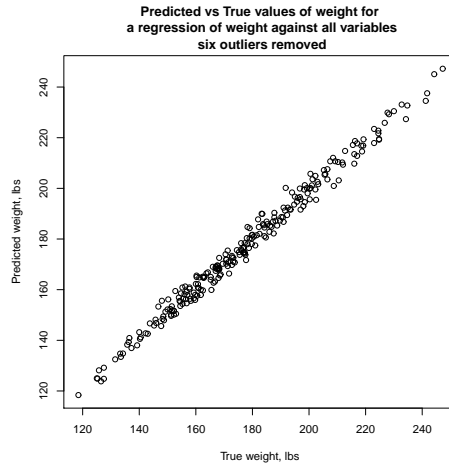


FIGURE 7.18: A scatter plot of the predicted weight against the true weight for the bodyfat dataset. The prediction is made with all variables, but the six outliers identified above are omitted. I used a Box-Cox transformation with parameter $1/2$, and the regularization parameter that yielded the smallest mean square error in Figure 7.17.

predictions? We already have some information on this point. Figure 7.17 shows cross-validation errors for regressions of $\text{weight}^{1/2}$ against height for different regularization weights, but some will find this slightly indirect. We want to predict weight, not $\text{weight}^{1/2}$. I chose the regularization weight that yielded the lowest mean-square-error for the model of Figure 7.17, omitting the six outliers previously mentioned. I then computed the predicted weight for each data point using that model (which predicts $\text{weight}^{1/2}$, remember; but squaring takes care of that). Figure 7.18 shows the predicted values plotted against the true values. You should not regard this plot as a safe way to estimate generalization (the points were used in training the model; Figure 7.17 is better for that), but it helps to visualize the errors. This regression looks as though it is quite good at predicting bodyweight from other measurements.

7.5 YOU SHOULD

7.5.1 remember:

New term: Regression	138
New term: explanatory variables	138
New term: dependent variable	138
New term: training examples	138
New term: test examples	138
New term: residual	141
Definition: Regression	142
New term: explanatory variables	142
New term: dependent variable	142
New term: training examples	142
New term: test examples	142
Definition: Linear regression	143
Estimating β	145
New term: residual	145
New term: mean square error	146
Useful facts: Regression	147
R^2 evaluates the quality of predictions made by a regression	148
New term: Zipf's law	149
Transforming variables is useful	150
New term: Box-Cox transformation	150
Linear regressions can fail.	151
New term: outliers	154
Outliers can affect linear regressions significantly.	155
New term: hat matrix	156
The hat matrix mixes training y -values to produce predictions.	156
New term: leverage	156
Be suspicious of points with high leverage.	156
New term: Cook's distance	157
Be suspicious of points with high Cook's distance.	157
New term: standardizing	158
Appending functions of a measurement to \mathbf{x} is useful.	160
New term: ridge regression	163
New term: L_2 norm	164
You can regularize a regression	165
New term: condition number	171
New term: condition number	171

APPENDIX: DATA

Batch A		Batch B		Batch C	
Amount of Hormone	Time in Service	Amount of Hormone	Time in Service	Amount of Hormone	Time in Service
25.8	99	16.3	376	28.8	119
20.5	152	11.6	385	22.0	188
14.3	293	11.8	402	29.7	115
23.2	155	32.5	29	28.9	88
20.6	196	32.0	76	32.8	58
31.1	53	18.0	296	32.5	49
20.9	184	24.1	151	25.4	150
20.9	171	26.5	177	31.7	107
30.4	52	25.8	209	28.5	125

TABLE 7.1: A table showing the amount of hormone remaining and the time in service for devices from lot A, lot B and lot C. The numbering is arbitrary (i.e. there's no relationship between device 3 in lot A and device 3 in lot B). We expect that the amount of hormone goes down as the device spends more time in service, so cannot compare batches just by comparing numbers.

PROBLEMS

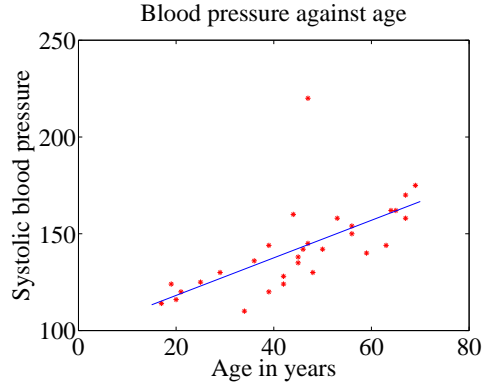


FIGURE 7.19: A regression of blood pressure against age, for 30 data points.

7.1. Figure 7.19 shows a linear regression of systolic blood pressure against age. There are 30 data points.

- (a) Write $e_i = y_i - \mathbf{x}_i^T \beta$ for the residual. What is the $\text{mean}(\{e\})$ for this regression?
- (b) For this regression, $\text{var}(\{y\}) = 509$ and the R^2 is 0.4324. What is $\text{var}(\{e\})$ for this regression?
- (c) How well does the regression explain the data?
- (d) What could you do to produce better predictions of blood pressure (without actually measuring blood pressure)?

7.2. This exercise investigates the effect of correlation on a regression. Assume we have N data items, (\mathbf{x}_i, y_i) . Write x_{i1} for the first component of \mathbf{x}_i , and $\mathbf{x}_{i,\hat{1}}$ for the vector obtained by deleting the first component of \mathbf{x}_i . Assume our data has the property that $x_{i1} = \mathbf{x}_{i,\hat{1}}^T \mathbf{u} + r_i$. Write \mathbf{r} for the vector of residuals (i.e. the i 'th component of \mathbf{r} is r_i). Now assume $\mathbf{r}^T \mathbf{1} = 0$ (i.e. the average of the r_i 's is zero) and $\mathbf{r}^T \mathbf{r} \leq \epsilon$. All this means that the data have the property that the first component is relatively accurately predicted by the other components.

- (a) Write $\mathbf{w} = [-1, \mathbf{u}]^T$. Show that

$$\mathbf{w}^T \mathcal{X}^T \mathcal{X} \mathbf{w} \leq \epsilon.$$

- (b) Now show that the smallest eigenvalue of $\mathcal{X}^T \mathcal{X}$ is less than or equal to ϵ .
- (c) Write $s_k = \sum_u x_{uk}^2$, and s_{\max} for $\max(s_1, \dots, s_d)$. Show that the largest eigenvalue of $\mathcal{X}^T \mathcal{X}$ is greater than or equal to s_{\max} .
- (d) The **condition number** of a matrix is the ratio of largest to smallest eigenvalue of a matrix. Use the information above to bound the **condition number** of $\mathcal{X}^T \mathcal{X}$.
- (e) Assume that $\hat{\beta}$ is the solution to $\mathcal{X}^T \mathcal{X} \hat{\beta} = \mathcal{X}^T \mathbf{Y}$. Show that the

$$(\mathcal{X}^T \mathbf{Y} - \mathcal{X}^T \mathcal{X}(\hat{\beta} + \mathbf{w}))^T (\mathcal{X}^T \mathbf{Y} - \mathcal{X}^T \mathcal{X}(\hat{\beta} + \mathbf{w}))$$

is bounded above by

$$\epsilon^2 (1 + \mathbf{u}^T \mathbf{u})$$

- (f) Use the last sub exercises to explain why correlated data will lead to a poor estimate of $\hat{\beta}$.
- 7.3.** In this exercise, I will show that the prediction process of chapter 13 (see page 316) is a linear regression with two independent variables. Assume we have N data items which are 2-vectors $(x_1, y_1), \dots, (x_N, y_N)$, where $N > 1$. These could be obtained, for example, by extracting components from larger vectors. As usual, we will write \hat{x}_i for x_i in normalized coordinates, and so on. The correlation coefficient is r (this is an important, traditional notation).
- (a) Show that $r = \frac{\text{mean}(\{(x - \text{mean}(\{x})) (y - \text{mean}(\{y}))\})}{(\text{std}(x)\text{std}(y))}$.
- (b) Now write $s = \frac{\text{std}(y)}{\text{std}(x)}$. Now assume that we have an x_0 , for which we wish to predict a y value. Show that the value of the prediction obtained using the method of page 317 is

$$sr(x_0 - \text{mean}(\{x\})) + \text{mean}(\{y\}).$$

- (c) Show that $sr = \text{mean}(\{(xy)\}) - \text{mean}(\{x\})\text{mean}(\{y\})$.
- (d) Now write

$$\mathcal{X} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \dots & \dots \\ x_n & 1 \end{pmatrix} \text{ and } \mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}.$$

The coefficients of the linear regression will be $\hat{\beta}$, where $\mathcal{X}^T \mathcal{X} \hat{\beta} = \mathcal{X}^T \mathbf{Y}$. Show that

$$\mathcal{X}^T \mathcal{X} = N \begin{pmatrix} \text{mean}(\{x^2\}) & \text{mean}(\{x\}) \\ \text{mean}(\{x\}) & 1 \end{pmatrix}$$

- (e) Now show that $\text{var}(\{x\}) = \text{mean}(\{(x - \text{mean}(\{x}))^2\}) = \text{mean}(\{x^2\}) - \text{mean}(\{x\})^2$.
- (f) Now show that $\text{std}(x)\text{std}(y)\text{corr}(\{(x, y)\}) = \text{mean}(\{(x - \text{mean}(\{x})) (y - \text{mean}(\{y}))\})$.

Regression: Choosing and Managing Models

8.1 MODEL SELECTION: WHICH MODEL IS BEST?

It is usually quite easy to have many explanatory variables in a regression problem. Even if you have only one measurement, you could always compute a variety of non-linear functions of that measurement. As we have seen, inserting variables into a model will reduce the fitting cost, but that doesn't mean that better predictions will result (section 7.4.1). We need to choose which explanatory variables we will use. A linear model with few explanatory variables may make poor predictions because the model itself is incapable of representing the independent variable accurately (an effect known as bias). A linear model with many explanatory variables may make poor predictions because we can't estimate the coefficients well (an effect known as variance). Choosing which explanatory variables we will use (and so which model we will use) requires that we balance these effects, described in greater detail in section 8.1.1. In the following sections, we describe straightforward methods of doing so.

8.1.1 Bias and Variance

We now look at the process of finding a model in a fairly abstract way. Doing so makes plain three distinct and important effects that cause models to make predictions that are wrong. One is **irreducible error**. Even a perfect choice of model can make mistake predictions, because more than one prediction could be correct for the same \mathbf{x} . Another way to think about this is that there could be many future data items, all of which have the same \mathbf{x} , but each of which has a different y . In this case some of our predictions must be wrong, and the effect is unavoidable.

A second effect is **bias**. We must use some collection of models. Even the best model in the collection may not be capable of predicting all the effects that occur in the data. Errors that are caused by the best model still not being able to predict the data accurately are attributed to bias.

The third effect is **variance**. We must choose our model from the collection of models. The model we choose is unlikely to be the best model. This might occur, for example, because our estimates of the parameters aren't exact because we have a limited amount of data. Errors that are caused by our choosing a model that is not the best in the family are attributed to variance.

All this can be written out in symbols. We have a vector of predictors \mathbf{x} , and a random variable Y . At any given point \mathbf{x} , we have

$$Y = f(\mathbf{x}) + \xi$$

where ξ is noise and f is an unknown function. We have $\mathbb{E}[\xi] = 0$, and $\mathbb{E}[\xi^2] = \text{var}(\{\xi\}) = \sigma_\xi^2$; furthermore, ξ is independent of X . We have some procedure that takes a selection of training data, consisting of pairs (\mathbf{x}_i, y_i) , and selects a model \hat{f} . We will use this model to predict values for future \mathbf{x} . It is highly unlikely that \hat{f} is the same as f ; assuming that it is involves assuming that we can perfectly estimate the best model with a finite dataset, which doesn't happen.

We need to understand the error that will occur when we use \hat{f} to predict for some data item that isn't in the training set. This is the error that we will encounter in practice. The error at any point \mathbf{x} is

$$\mathbb{E}[(Y - \hat{f}(\mathbf{X}))^2]$$

where the expectation is taken over $P(Y|\mathbf{x})$. This expectation can be written in an extremely useful form. Recall $\text{var}(\{U\}) = \mathbb{E}[U^2] - \mathbb{E}[U]^2$. This means we have

$$\begin{aligned} \mathbb{E}[(Y - \hat{f}(\mathbf{X}))^2] &= \mathbb{E}[Y^2] - 2\mathbb{E}[Y\hat{f}] + \mathbb{E}[\hat{f}^2] \\ &= \text{var}(\{Y\}) + \mathbb{E}[Y]^2 + \text{var}(\{\hat{f}\}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[Y\hat{f}]. \end{aligned}$$

Now $Y = f(X) + \xi$, $\mathbb{E}[\xi] = 0$, and ξ is independent of X so we have $\mathbb{E}[Y] = \mathbb{E}[f]$ and $\text{var}(\{Y\}) = \text{var}(\{\xi\}) = \sigma_\xi^2$. This yields

$$\begin{aligned} \mathbb{E}[(Y - \hat{f}(\mathbf{X}))^2] &= \text{var}(\{Y\}) + \mathbb{E}[f]^2 + \text{var}(\{\hat{f}\}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[f\hat{f}] \\ &= \sigma_\xi^2 + \mathbb{E}[(f - \hat{f})^2] + \text{var}(\{\hat{f}\}) \\ &= \sigma_\xi^2 + (f - \mathbb{E}[\hat{f}])^2 + \text{var}(\{\hat{f}\}) \quad (f \text{ isn't random}). \end{aligned}$$

The expected error on all future data is the sum of three terms. The irreducible error is σ_ξ^2 ; even the true model must produce this error, on average. The best model to choose would be $\mathbb{E}[\hat{f}]$ (remember, the expectation is over choices of training data; this model would be the one that best represented all possible attempts to train). But we don't have $\mathbb{E}[\hat{f}]$. Instead, we have \hat{f} . The variance is $\text{var}(\{\hat{f}\}) = \mathbb{E}[(\hat{f} - \mathbb{E}[\hat{f}])^2]$. This term represents the fact that the model we chose (\hat{f}) is different from the mean model ($\mathbb{E}[\hat{f}]$). The difference arises because our training data is a subset of all data, and our model is chosen to be good on the training data, rather than on every possible training set. The bias is $(f - \mathbb{E}[\hat{f}])^2$. This term reflects the fact that even the best choice of model ($\mathbb{E}[\hat{f}]$) may not be the same as the true source of data ($\mathbb{E}[f]$ which is the same as f , because f is deterministic).

There is usually a tradeoff between bias and variance. Generally, when a model comes from a "small" or "simple" family, we expect that (a) we can estimate the best model in the family reasonably accurately (so the variance will be low) but (b) the model may have real difficulty reproducing the data (meaning the bias

is large). Similarly, if the model comes from a “large” or “complex” family, the variance is likely to be high (because it will be hard to estimate the best model in the family accurately) but the bias will be low (because the model can more accurately reproduce the data). All modelling involves managing this tradeoff between bias and variance. I am avoiding being precise about the complexity of a model because it can be tricky to do. One reasonable proxy is the number of parameters we have to estimate to determine the model.

You can see a crude version this tradeoff in the perch example of section 7.4.1 and Figure 7.11. Recall that, as I added monomials to the regression of weight against length, the fitting error went down; but the model that uses length¹⁰ as an explanatory variable makes very odd predictions away from the training data. When I use low degree monomials, the dominant source of error is bias; and when I use high degree monomials, the dominant source of error is variance. A common mistake is to feel that the major difficulty is bias, and so to use extremely complex models. Usually the result is poor estimates of model parameters, and huge variance. Experienced modellers fear variance far more than they fear bias.

The bias-variance discussion suggests it isn’t a good idea simply to use all the explanatory variables that you can obtain (or think of). Doing so might lead to a model with serious variance problems. Instead, we must choose a model that uses a subset of the explanatory variables that is small enough to control variance, and large enough that the bias isn’t a problem. We need some strategy to choose explanatory variables. The simplest (but by no means the best; we’ll see better in this chapter) approach is to search sets of explanatory variables for a good set. The main difficulty is knowing when you have a good set.

8.1.2 Choosing a Model using Penalties: AIC and BIC

We would like to choose one of a set of models. We cannot do so using just the training error, because more complex models will tend to have lower training error, and so the model with the lowest training error will tend to be the most complex model. Training error is a poor guide to test error, because lower training error is evidence of lower bias on the models part; but with lower bias, we expect to see greater variance, and the training error doesn’t take that into account.

One strategy is to penalize the model for complexity. We add some penalty, reflecting the complexity of the model, to the training error. We then expect to see the general behavior of figure 8.1. The training error goes down, and the penalty goes up as the model gets more complex, so we expect to see a point where the sum is at a minimum.

There are a variety of ways of constructing penalties. **AIC** (short for An Information Criterion) is a method due originally to Akaike, in ****. Rather than using the training error, AIC uses the maximum value of the log-likelihood of the model. Write \mathcal{L} for this value. Write k for the number of parameters estimated to fit the model. Then the AIC is

$$2k - 2\mathcal{L}$$

and a better model has a smaller value of AIC (remember this by remembering that a larger log-likelihood corresponds to a better model). Estimating AIC is straightforward for regression models if you assume that the noise is a zero mean

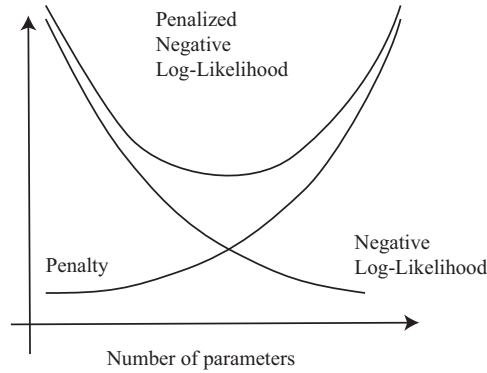


FIGURE 8.1: When we add explanatory variables (and so parameters) to a model, the value of the negative log-likelihood of the best model can't go up, and usually goes down. This means that we cannot use the value as a guide to how many explanatory variables there should be. Instead, we add a penalty that increases as a function of the number of parameters, and search for the model that minimizes the sum of negative log-likelihood and penalty. AIC and BIC grow linearly with the number of parameters, but I am following the usual convention of plotting the penalty as a curve rather than a straight line.

normal random variable. You estimate the mean-squared error, which gives the variance of the noise, and so the log-likelihood of the model. You do have to keep track of two points. First, k is the total number of parameters estimated to fit the model. For example, in a linear regression model, where you model y as $\mathbf{x}^T \beta + \xi$, you need to estimate d parameters to estimate $\hat{\beta}$ and the variance of ξ (to get the log-likelihood). So in this case $k = d + 1$. Second, log-likelihood is usually only known up to a constant, so that different software implementations often use different constants. This is wildly confusing when you don't know about it (why would AIC and `extractAIC` produce different numbers on the same model?) but of no real significance – you're looking for the smallest value of the number, and the actual value doesn't mean anything. Just be careful to compare only numbers computed with the same routine.

An alternative is **BIC** (Bayes' Information Criterion), given by

$$2k \log N - 2\mathcal{L}$$

(where N is the size of the training data set). You will often see this written as $2\mathcal{L} - 2k \log N$; I have given the form above so that one always wants the smaller value as with AIC. There is a considerable literature comparing AIC and BIC. AIC has a mild reputation for overestimating the number of parameters required, but is often argued to have firmer theoretical foundations.

Worked example 8.1 *AIC and BIC*

Write M_d for the model that predicts weight from length for the perch dataset as $\sum_{j=0}^{j=d} \beta_j \text{length}^j$. Choose an appropriate value of $d \in [1, 10]$ using AIC and BIC.

Solution: I used the R functions `AIC` and `BIC`, and got the table below.

	1	2	3	4	5	6	7	8	9	10
AIC	677	617	617	613	615	617	617	612	613	614
BIC	683	625	627	625	629	633	635	633	635	638

The best model by AIC has (rather startlingly!) $d = 8$. One should not take small differences in AIC too seriously, so models with $d = 4$ and $d = 9$ are fairly plausible, too. BIC suggests $d = 2$.

8.1.3 Choosing a Model using Cross-Validation

AIC and BIC are estimates of error on future data. An alternative is to measure this error on held out data, using a cross-validation strategy (as in section 3.1.4). One splits the training data into F **olds**, where each data item lies in exactly one fold. The case $F = N$ is sometimes called “leave-one-out” cross-validation. One then sets aside one fold in turn, fitting the model to the remaining data, and evaluating the model error on the left-out fold. The model error is then averaged. This process gives us an estimate of the performance of a model on held-out data. Numerous variants are available, particularly when lots of computation and lots of data are available. For example: one might not average over all folds; one might use fewer or more folds; and so on.

Worked example 8.2 *Cross-validation*

Write M_d for the model that predicts weight from length for the perch dataset as $\sum_{j=0}^{j=d} \beta_j \text{length}^j$. Choose an appropriate value of $d \in [1, 10]$ using leave-one-out cross validation.

Solution: I used the R functions `CV1m`, which takes a bit of getting used to. There is sample code in the exercises section. I found:

1	2	3	4	5	6	7	8	9	10
1.9e4	4.0e3	7.2e3	4.5e3	6.0e3	5.6e4	1.2e6	4.0e6	3.9e6	1.9e8

where the best model is $d = 2$.

8.1.4 A Search Process: Forward and Backward Stagewise Regression

Assume we have a set of explanatory variables and we wish to build a model, choosing some of those variables for our model. Our explanatory variables could be many distinct measurements, or they could be different non-linear functions of

the same measurement, or a combination of both. We can evaluate models relative to one another fairly easily (AIC, BIC or cross-validation, your choice). However, choosing which set of explanatory variables to use can be quite difficult, because there are so many sets. The problem is that you cannot predict easily what adding or removing an explanatory variable will do. Instead, when you add (or remove) an explanatory variable, the errors that the model makes change, and so the usefulness of all other variables changes too. This means that (at least in principle) you have to look at every subset of the explanatory variables. Imagine you start with a set of F possible explanatory variables (including the original measurement, and a constant). You don't know how many to use, so you might have to try every different group, of each size, and there are far too many groups to try. There are two useful alternatives.

In **forward stagewise regression**, you start with an empty working set of explanatory variables. You then iterate the following process. For each of the explanatory variables not in the working set, you construct a new model using the working set and that explanatory variable, and compute the model evaluation score. If the best of these models has a better score than the model based on the working set, you insert the appropriate variable into the working set and iterate. If no variable improves the working set, you decide you have the best model and stop. This is fairly obviously a greedy algorithm.

Backward stagewise regression is pretty similar, but you start with a working set containing all the variables, and remove variables one-by-one and greedily. As usual, greedy algorithms are very helpful but not capable of exact optimization. Each of these strategies can produce rather good models, but neither is guaranteed to produce the best model.

8.1.5 Significance: What Variables are Important?

Imagine you regress some measure of risk of death against blood pressure, whether someone smokes or not, and the length of their thumb. Because high blood pressure and smoking tend to increase risk of death, you would expect to see “large” coefficients for these explanatory variables. Since changes in the thumb length have no effect, you would expect to see “small” coefficients for these explanatory variables. This suggests a regression can be used to determine what effects are important in building a model.

One difficulty is the result of sampling variance. Imagine that we have an explanatory variable that has absolutely no relationship to the dependent variable. If we had an arbitrarily large amount of data, and could exactly identify the correct model, we'd find that, in the correct model, the coefficient of that variable was zero. But we don't have an arbitrarily large amount of data. Instead, we have a sample of data. Hopefully, our sample is random, so that the reasoning of section 16.20 can be applied. Using that reasoning, our estimate of the coefficient is the value of a random variable whose expected value is zero, but whose variance isn't. As a result, we are very unlikely to see a zero. This reasoning applies to each coefficient of the model. To be able to tell which ones are small, we would need to know the standard deviation of each, so we can tell whether the value we observe is a small number of standard deviations away from zero. This line of reasoning is very like hypothesis

testing. It turns out that the sampling variance of regression coefficients can be estimated in a straightforward way. In turn, we have an estimate of the extent to which their difference from zero could be a result of random sampling. R will produce this information routinely; use `summary` on the output of `lm`.

A second difficulty has to do with practical significance, and is rather harder. We could have explanatory variables that are genuinely linked to the independent variable, but might not matter very much. This is a common phenomenon, particularly in medical statistics. It requires considerable care to disentangle some of these issues. Here is an example. Bowel cancer is an unpleasant disease, which could kill you. Being screened for bowel cancer is at best embarrassing and unpleasant, and involves some startling risks. There is considerable doubt, from reasonable sources, about whether screening has value and if so, how much (as a start point, you could look at Ransohoff DF. How Much Does Colonoscopy Reduce Colon Cancer Mortality?. *Ann Intern Med.* 2009). There is some evidence linking eating red or processed meat to incidence of bowel cancer. A good practical question is: should one abstain from eating red or processed meat based on increased bowel cancer risk?

Coming to an answer is tough; the coefficient in any regression is clearly not zero, but it's pretty small as these numbers indicate. The UK population in 2012 was 63.7 million (this is a summary figure from Google, using World Bank data; there's no reason to believe that it's significantly wrong). I obtained the following figures from the UK cancer research institute website, at <http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/bowel-cancer>. There were 41,900 new cases of bowel cancer in the UK in 2012. Of these cases, 43% occurred in people aged 75 or over. 57% of people diagnosed with bowel cancer survive for ten years or more after diagnosis. Of diagnosed cases, an estimated 21% are linked to eating red or processed meat, and the best current estimate is that the risk of incidence is between 17% and 30% higher per 100g of red meat eaten per day (i.e. if you eat 100g of red meat per day, your risk increases by some number between 17% and 30%; 200g a day gets you twice that number; and – rather roughly – so on). These numbers are enough to confirm that there is a non-zero coefficient linking the amount of red or processed meat in your diet with your risk of bowel cancer (though you'd have a tough time estimating the exact value of that coefficient from the information here). If you eat more red meat, your risk of dying of bowel cancer really will go up. But the numbers I gave above suggest that (a) it won't go up much and (b) you might well die rather late in life, where the chances of dying of something are quite strong. The coefficient linking eating red meat and bowel cancer is clearly pretty small, because the incidence of the disease is about 1 in 1500 per year. Does it matter? you get to choose, and your choice has consequences.

8.2 ROBUST REGRESSION

We have seen that outlying data points can result in a poor model. This is caused by the squared error cost function: squaring a large error yields an enormous number. One way to resolve this problem is to identify and remove outliers before fitting a model. This can be difficult, because it can be hard to specify precisely when

a point is an outlier. Worse, in high dimensions most points will look somewhat like outliers, and we may end up removing all most all the data. The alternative solution I offer here is to come up with a cost function that is less susceptible to problems with outliers. The general term for a regression that can ignore some outliers is a **robust regression**.

8.2.1 M-Estimators and Iteratively Reweighted Least Squares

One way to reduce the effect of outliers on a least-squares solution would be to weight each point in the cost function. We need some method to estimate an appropriate set of weights. This would use a large weight for errors at points that are “trustworthy”, and a low weight for errors at “suspicious” points.

We can obtain such weights using an **M-estimator**, which estimates parameters by replacing the negative log-likelihood with a term that is better behaved. In our examples, the negative log-likelihood has always been squared error. Write β for the parameters of the model being fitted, and $r_i(\mathbf{x}_i, \beta)$ for the residual error of the model on the i th data point. For us, r_i will always be $y_i - \mathbf{x}_i^T \beta$. So rather than minimizing

$$\sum_i (r_i(\mathbf{x}_i, \beta))^2$$

as a function of β , we will minimize an expression of the form

$$\sum_i \rho(r_i(\mathbf{x}_i, \beta); \sigma),$$

for some appropriately chosen function ρ . Clearly, our negative log-likelihood is one such estimator (use $\rho(u; \sigma) = u^2$). The trick to M-estimators is to make $\rho(u; \sigma)$ look like u^2 for smaller values of u , but ensure that it grows more slowly than u^2 for larger values of u .

The **Huber loss** is one important M-estimator. We use

$$\rho(u; \sigma) = \begin{cases} \frac{u^2}{2} & |u| < \sigma \\ \sigma|u| - \frac{\sigma^2}{2} & |u| \geq \sigma \end{cases}$$

which is the same as u^2 for $-\sigma \leq u \leq \sigma$, and then switches to $|u|$ for larger (or smaller) σ (Figure ??). The Huber loss is convex (meaning that there will be a unique minimum for our models) and differentiable, but its derivative is not continuous. The choice of the parameter σ (which is known as **scale**) has an effect on the estimate. You should interpret this parameter as the distance that a point can lie from the fitted function while still being seen as an **inlier** (anything that isn't even partially an outlier).

Generally, M-estimators are discussed in terms of their **influence function**. This is

$$\frac{\partial \rho}{\partial u}.$$

Its importance becomes evidence when we consider algorithms to fit $\hat{\beta}$ using an

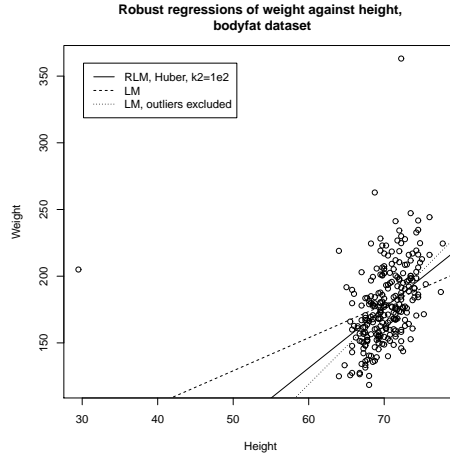


FIGURE 8.2: Comparing three different linear regression strategies on the bodyfat data, regressing weight against height. Notice that using an M-estimator gives an answer very like that obtained by rejecting outliers by hand. The answer may well be “better” because it isn’t certain that each of the four points rejected is an outlier, and the robust method may benefit from some of the information in these points. I tried a range of scales for the Huber loss (the ‘k2’ parameter), but found no difference in the line resulting over scales varying by a factor of $1e4$, which is why I plot only one scale.

M-estimator. Our minimization criterion is

$$\begin{aligned} \nabla_{\beta} \left(\sum_i \rho(y_i - \mathbf{x}_i^T \beta; \sigma) \right) &= \sum_i \left[\frac{\partial \rho}{\partial \mathbf{u}}(y_i - \mathbf{x}_i^T \beta; \sigma) \right] (-\mathbf{x}_i) \\ &= 0. \end{aligned}$$

Now write $w_i(\beta)$ for

$$\frac{\frac{\partial \rho}{\partial \mathbf{u}}(y_i - \mathbf{x}_i^T \beta; \sigma)}{y_i - \mathbf{x}_i^T \beta}.$$

We can write the minimization criterion as

$$\sum_i [w_i(\beta)] (y_i - \mathbf{x}_i^T \beta)(-\mathbf{x}_i) = 0.$$

Now write $\mathcal{W}(\beta)$ for the diagonal matrix whose i ’th diagonal entry is $w_i(\beta)$. Then our fitting criterion is equivalent to

$$\mathcal{X}^T [\mathcal{W}(\beta)] \mathbf{Y} = \mathcal{X}^T [\mathcal{W}(\beta)] \mathcal{X} \beta.$$

The difficulty in solving this is that $w_i(\beta)$ depend on β , so we can’t just solve a linear system in β . We could use the following strategy. Use \mathcal{W} tries to downweight points that are suspiciously inconsistent with our current estimate of β , then update

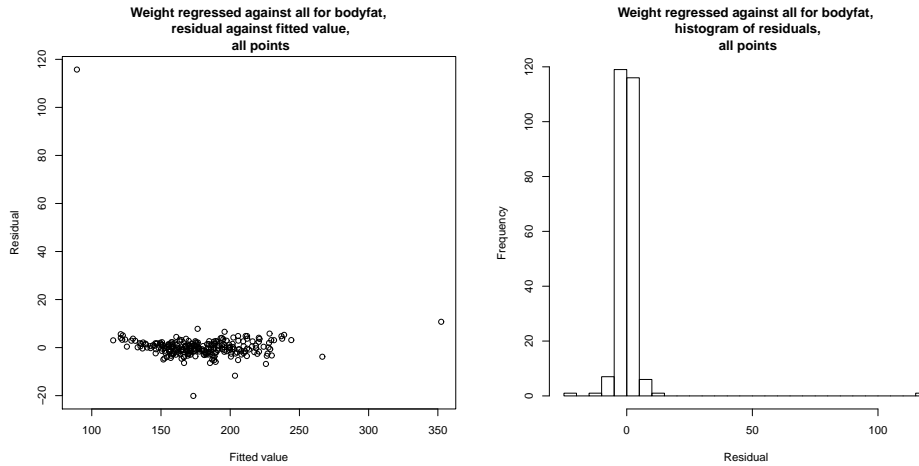


FIGURE 8.3: A robust linear regression of weight against all variables for the bodyfat dataset, using the Huber loss and all data points. On the **left**, residual plotted against fitted value (the residual is not standardized). Notice that there are some points with very large residual, but most have much smaller residual; this wouldn't happen with a squared error. On the **right**, a histogram of the residual. If one ignores the extreme residual values, this looks normal. The robust process has been able to discount the effect of the outliers, without us needing to identify and reject outliers by hand.

β using those weights. The strategy is known as **iteratively reweighted least squares**, and is very effective.

We assume we have an estimate of the correct parameters $\hat{\beta}^{(n)}$, and consider updating it to $\hat{\beta}^{(n+1)}$. We compute

$$w_i^{(n)} = w_i(\hat{\beta}^{(n)}) = \frac{\frac{\partial \rho}{\partial u}(y_i - \mathbf{x}_i^T \hat{\beta}^{(n)}; \sigma)}{y_i - \mathbf{x}_i^T \hat{\beta}^{(n)}}.$$

We then estimate $\hat{\beta}^{(n+1)}$ by solving

$$\mathcal{X}^T \mathcal{W}^{(n)} \mathbf{Y} = \mathcal{X}^T \mathcal{W}^{(n)} \mathcal{X} \hat{\beta}^{(n+1)}.$$

The key to this algorithm is finding good start points for the iteration. One strategy is randomized search. We select a small subset of points uniformly at random, and fit some $\hat{\beta}$ to these points, then use the result as a start point. If we do this often enough, one of the start points will be an estimate that is not contaminated by outliers.

8.2.2 Scale for M-Estimators

The estimators require a sensible estimate of σ , which is often referred to as **scale**. Typically, the scale estimate is supplied at each iteration of the solution method.

One reasonable estimate is the **MAD** or **median absolute deviation**, given by

$$\sigma^{(n)} = 1.4826 \operatorname{median}_i |r_i^{(n)}(x_i; \hat{\beta}^{(n-1)})|.$$

Another a popular estimate of scale is obtained with **Huber's proposal 2** (that is what everyone calls it!). Choose some constant $k_1 > 0$, and define $\Xi(u) = \min(|u|, k_1)^2$. Now solve the following equation for σ :

$$\sum_i \Xi\left(\frac{r_i^{(n)}(x_i; \hat{\beta}^{(n-1)})}{\sigma}\right) = Nk_2$$

where k_2 is another constant, usually chosen so that the estimator gives the right answer for a normal distribution (exercises). This equation needs to be solved with an iterative method; the MAD estimate is the usual start point. R provides `hubers`, which will compute this estimate of scale (and figures out k_2 for itself). The choice of k_1 depends somewhat on how contaminated you expect your data to be. As $k_1 \rightarrow \infty$, this estimate becomes more like the standard deviation of the data.

8.3 GENERALIZED LINEAR MODELS

We have used a linear regression to predict a value from a feature vector, but implicitly have assumed that this value is a real number. Other cases are important, and some of them can be dealt with using quite simple generalizations of linear regression. When we derived linear regression, I said one way to think about the model was

$$y = \mathbf{x}^T \beta + \xi$$

where ξ was a normal random variable with zero mean and variance σ_ξ^2 . Another way to write this is to think of y as the value of a random variable Y . In this case, Y has mean $\mathbf{x}^T \beta$ and variance σ_ξ^2 . This can be written as

$$Y \sim N(\mathbf{x}^T \beta, \sigma_\xi^2).$$

This offers a fruitful way to generalize: we replace the normal distribution with some other parametric distribution, and predict the parameter using $\mathbf{x}^T \beta$. Two examples are particularly important.

8.3.1 Logistic Regression

Assume the y values can be either 0 or 1. You could think of this as a two class classification problem, and deal with it using an SVM. There are sometimes advantages to seeing it as a regression problem. One is that we get to see a new classification method that explicitly models class posteriors, which an SVM doesn't do.

We build the model by asserting that the y values represent a draw from a Bernoulli random variable (definition below, for those who have forgotten). The parameter of this random variable is θ , the probability of getting a one. But $0 \leq \theta \leq 1$, so we can't just model θ as $\mathbf{x}^T \beta$. We will choose some **link function** g so that we can model $g(\theta)$ as $\mathbf{x}^T \beta$. This means that, in this case, g must map the interval between 0 and 1 to the whole line, and must be 1-1. The link function maps θ to $\mathbf{x}^T \beta$; the direction of the map is chosen by convention. We build our model by asserting that $g(\theta) = \mathbf{x}^T \beta$.

Definition: 8.1 *Bernoulli random variable*

A Bernoulli random variable with parameter θ takes the value 1 with probability θ and 0 with probability $1 - \theta$. This is a model for a coin toss, among other things.

Notice that, for a Bernoulli random variable, we have that

$$\log \left[\frac{P(y = 1|\theta)}{P(y = 0|\theta)} \right] = \log \left[\frac{\theta}{1 - \theta} \right]$$

and the **logit function** $g(u) = \log \left[\frac{u}{1-u} \right]$ meets our needs for a link function (it maps the interval between 0 and 1 to the whole line, and is 1-1). This means we can build our model by asserting that

$$\log \left[\frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} \right] = \mathbf{x}^T \beta$$

then solving for the β that maximizes the log-likelihood of the data. Simple manipulation yields

$$P(y = 1|\mathbf{x}) = \frac{e^{\mathbf{x}^T \beta}}{1 + e^{\mathbf{x}^T \beta}} \text{ and } P(y = 0|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{x}^T \beta}}.$$

In turn, this means the log-likelihood of a dataset will be

$$\mathcal{L}(\beta) = \sum_i \left[\mathbb{I}_{[y=1]}(y_i) \mathbf{x}_i^T \beta - \log \left(1 + e^{\mathbf{x}_i^T \beta} \right) \right].$$

You can obtain β from this log-likelihood by gradient ascent (or rather a lot faster by Newton's method, if you know that).

A regression of this form is known as a **logistic regression**. It has the attractive property that it produces estimates of posterior probabilities. Another interesting property is that a logistic regression is a lot like an SVM. To see this, we replace the labels with new ones. Write $\hat{y}_i = 2y_i - 1$; this means that \hat{y}_i takes the values -1 and 1 , rather than 0 and 1 . Now $\mathbb{I}_{[y=1]}(y_i) = \frac{\hat{y}_i + 1}{2}$, so we can write

$$\begin{aligned} -\mathcal{L}(\beta) &= -\sum_i \left[\frac{\hat{y}_i + 1}{2} \mathbf{x}_i^T \beta - \log \left(1 + e^{\mathbf{x}_i^T \beta} \right) \right] \\ &= \sum_i \left[\frac{\hat{y}_i + 1}{2} \mathbf{x}_i^T \beta - \log \left(1 + e^{\mathbf{x}_i^T \beta} \right) \right] \\ &= \sum_i \left[\log \left(\frac{1 + e^{\mathbf{x}_i^T \beta}}{e^{\frac{\hat{y}_i + 1}{2} \mathbf{x}_i^T \beta}} \right) \right] \\ &= \sum_i \left[\log \left(e^{-\frac{(\hat{y}_i + 1)}{2} \mathbf{x}_i^T \beta} + e^{\frac{1 - \hat{y}_i}{2} \mathbf{x}_i^T \beta} \right) \right] \end{aligned}$$

and we can interpret the term in square brackets as a loss function. If you plot it, you will notice that it behaves rather like the hinge loss. When $\hat{y}_i = 1$, if $\mathbf{x}^T \beta$ is positive the loss is very small, but if $\mathbf{x}^T \beta$ is strongly negative, the loss grows linearly in $\mathbf{x}^T \beta$. There is similar behavior when $\hat{y}_i = -1$. The transition is smooth, unlike the hinge loss. Logistic regression should (and does) behave well for the same reasons the SVM behaves well.

Be aware that logistic regression has one annoying quirk. When the data are linearly separable (i.e. there exists some β such that $y_i \mathbf{x}_i^T \beta > 0$ for all data items), logistic regression will behave badly. To see the problem, choose the β that separates the data. Now it is easy to show that increasing the magnitude of β will increase the log likelihood of the data; there isn't any limit. These situations arise fairly seldom in practical data.

8.3.2 Multiclass Logistic Regression

Imagine $y \in [0, 1, \dots, C-1]$. Then it is natural to model $p(y|\mathbf{x})$ with a discrete probability distribution on these values. This can be specified by choosing $(\theta_0, \theta_1, \dots, \theta_{C-1})$ where each term is between 0 and 1 and $\sum_i \theta_i = 1$. Our link function will need to map this constrained vector of θ values to a \Re^{C-1} . We can do this with a fairly straightforward variant of the logit function, too. Notice that there are $C-1$ probabilities we need to model (the C 'th comes from the constraint $\sum_i \theta_i = 1$). We choose one vector β for each probability, and write β_i for the vector used to model θ_i . Then we can write

$$\mathbf{x}^T \beta_i = \log \left(\frac{\theta_i}{1 - \sum_u \theta_u} \right)$$

and this yields the model

$$\begin{aligned} P(y = 0|\mathbf{x}, \beta) &= \frac{e^{\mathbf{x}^T \beta_0}}{1 + \sum_i e^{\mathbf{x}^T \beta_i}} \\ P(y = 1|\mathbf{x}, \beta) &= \frac{e^{\mathbf{x}^T \beta_1}}{1 + \sum_i e^{\mathbf{x}^T \beta_i}} \\ &\dots \\ P(y = C-1|\mathbf{x}, \beta) &= \frac{1}{1 + \sum_i e^{\mathbf{x}^T \beta_i}} \end{aligned}$$

and we would fit this model using maximum likelihood. The likelihood is easy to write out, and gradient descent is a good strategy for actually fitting models.

8.3.3 Regressing Count Data

Now imagine that the y_i values are counts. For example, y_i might have the count of the number of animals caught in a small square centered on \mathbf{x}_i in a study region. As another example, \mathbf{x}_i might be a set of features that represent a customer, and y_i might be the number of times that customer bought a particular product. The natural model for count data is a Poisson model, with parameter θ representing the intensity (reminder below).

Definition: 8.2 *Poisson distribution*

A non-negative, integer valued random variable X has a Poisson distribution when its probability distribution takes the form

$$P(\{X = k\}) = \frac{\theta^k e^{-\theta}}{k!},$$

where $\theta > 0$ is a parameter often known as the **intensity** of the distribution.

Now we need $\theta > 0$. A natural link function is to use

$$\mathbf{x}^T \beta = \log \theta$$

yielding a model

$$P(\{X = k\}) = \frac{e^{k\mathbf{x}^T \beta} e^{-e^{k\mathbf{x}^T \beta}}}{k!}.$$

Now assume we have a dataset. The negative log-likelihood can be written as

$$\begin{aligned} -\mathcal{L}(\beta) &= -\sum_i \log \left(\frac{e^{y_i \mathbf{x}_i^T \beta} e^{-e^{y_i \mathbf{x}_i^T \beta}}}{y_i!} \right) \\ &= -\sum_i \left(y_i \mathbf{x}_i^T \beta - e^{y_i \mathbf{x}_i^T \beta} - \log(y_i!) \right). \end{aligned}$$

There isn't a closed form minimum available, but the log-likelihood is convex, and gradient descent (or Newton's method) are enough to find a minimum. Notice that the $\log(y_i!)$ term isn't relevant to the minimization, and is usually dropped.

8.3.4 Deviance

Cross-validating a model is done by repeatedly splitting a data set into two pieces, training on one, evaluating some score on the other, and averaging the score. But we need to keep track of *what* to score. For earlier linear regression models (eg section 16.20), we have used the squared error of predictions. This doesn't really make sense for a generalized linear model, because predictions are of quite different form. It is usual to use the **deviance** of the model. Write y_t for the true prediction at a point, \mathbf{x}_p for the independent variables we want to obtain a prediction for, $\hat{\beta}$ for our estimated parameters; a generalized linear model yields $P(y|\mathbf{x}_p, \hat{\beta})$. For our purposes, you should think of the deviance as

$$-2 \log P(y_t|\mathbf{x}_p, \hat{\beta})$$

(this expression is sometimes adjusted in software to deal with extreme cases, etc.). Notice that this is quite like the least squares error for the linear regression case, because there

$$-2 \log P(y|\mathbf{x}_p, \hat{\beta}) = (\mathbf{x}_p^T \hat{\beta} - y_t)^2 / \sigma^2 + K$$

for K some constant.

8.4 L1 REGULARIZATION AND SPARSE MODELS

Forward and backward stagewise regression were strategies for adding independent variables to, or removing independent variables from, a model. An alternative, and very powerful, strategy is to construct a model with a method that forces some coefficients to be zero. The resulting model ignores the corresponding independent variables. Models built this way are often called **sparse models**, because (one hopes) that many independent variables will have zero coefficients, and so the model is using a sparse subset of the possible predictors.

In some situations, we are forced to use a sparse model. For example, imagine there are more independent variables than there are examples. In this case, the matrix $\mathcal{X}^T \mathcal{X}$ will be rank deficient. We could use a ridge regression (Section 7.4.2) and the rank deficiency problem will go away, but it would be hard to trust the resulting model, because it will likely use all the predictors (more detail below). We really want a model that uses a small subset of the predictors. Then, because the model ignores the other predictors, there will be more examples than there are predictors *that we use*.

There is now quite a strong belief amongst practitioners that using sparse models is the best way to deal with high dimensional problems (although there are lively debates about *which* sparse model to use, etc.). This is sometimes called the “bet on sparsity” principle: use a sparse model for high dimensional data, because dense models don’t work well for such problems.

8.4.1 Dropping Variables with L1 Regularization

We have a large set of explanatory variables, and we would like to choose a small set that explains most of the variance in the independent variable. We could do this by encouraging β to have many zero entries. In section 7.4.2, we saw we could regularize a regression by adding a term to the cost function that discouraged large values of β . Instead of solving for the value of β that minimized $\sum_i (y_i - \mathbf{x}_i^T \beta)^2 = (\mathbf{y} - \mathcal{X}\beta)^T (\mathbf{y} - \mathcal{X}\beta)$ (which I shall call the **error cost**), we minimized

$$\sum_i (y_i - \mathbf{x}_i^T \beta)^2 + \frac{\lambda}{2} \beta^T \beta = (\mathbf{y} - \mathcal{X}\beta)^T (\mathbf{y} - \mathcal{X}\beta) + \frac{\lambda}{2} \beta^T \beta$$

(which I shall call the **L2 regularized error**). Here $\lambda > 0$ was a constant chosen by cross-validation. Larger values of λ encourage entries of β to be small, but do not force them to be zero. The reason is worth understanding.

Write β_k for the k ’th component of β , and write β_{-k} for all the other components. Now we can write the L2 regularized error as a function of β_k :

$$(a + \lambda)\beta_k^2 - 2b(\beta_{-k})\beta_k + c(\beta_{-k})$$

where a is a function of the data and b and c are functions of the data and of β_{-k} . Now notice that the best value of β_k will be

$$\beta_k = \frac{b(\beta_{-k})}{(a + \lambda)}.$$

Notice that λ doesn't appear in the numerator. This means that, to force β_k to zero by increasing λ , we may have to make λ arbitrarily large. This is because the improvement in the penalty obtained by going from a small β_k to $\beta_k = 0$ is tiny – the penalty is proportional to β_k^2 .

To force some components of β to zero, we need a penalty that grows linearly around zero rather than quadratically. This means we should use the **L₁ norm** of β , given by

$$\|\beta\|_1 = \sum_k |\beta_k|.$$

To choose β , we must now solve

$$(\mathbf{y} - \mathcal{X}\beta)^T(\mathbf{y} - \mathcal{X}\beta) + \lambda\|\beta\|_1$$

for an appropriate choice of λ . An equivalent problem is to solve a constrained minimization problem, where one minimizes

$$(\mathbf{y} - \mathcal{X}\beta)^T(\mathbf{y} - \mathcal{X}\beta) \text{ subject to } \|\beta\|_1 \leq t$$

where t is some value chosen to get a good result, typically by cross-validation. There is a relationship between the choice of t and the choice of λ (with some thought, a smaller t will correspond to a bigger λ) but it isn't worth investigating in any detail.

Actually solving this system is quite involved, because the cost function is not differentiable. You should *not* attempt to use stochastic gradient descent, because this will not compel zeros to appear in $\hat{\beta}$ (exercises). There are several methods, which are beyond our scope. As the value of λ increases, the number of zeros in $\hat{\beta}$ will increase too. We can choose λ in the same way we used for classification; split the training set into a training piece and a validation piece, train for different values of λ , and test the resulting regressions on the validation piece. However, one consequence of modern methods is that we can generate a very good approximation to the path $\hat{\beta}(\lambda)$ for all values of $\lambda \geq 0$ about as easily as we can choose $\hat{\beta}$ for a particular value of λ .

One way to understand the models that result is to look at the behavior of cross-validated error as λ changes. The error is a random variable, random because of the random split. It is a fair model of the error that would occur on a randomly chosen test example (assuming that the training set is “like” the test set, in a way that I do not wish to make precise yet). We could use multiple splits, and average over the splits. Doing so yields both an average error for each value of λ and an estimate of the standard deviation of error. Figure 8.4 shows the result of doing so for two datasets. Again, there is no λ that yields the smallest validation error, because the value of error depends on the random split cross-validation. A reasonable choice of λ lies between the one that yields the smallest

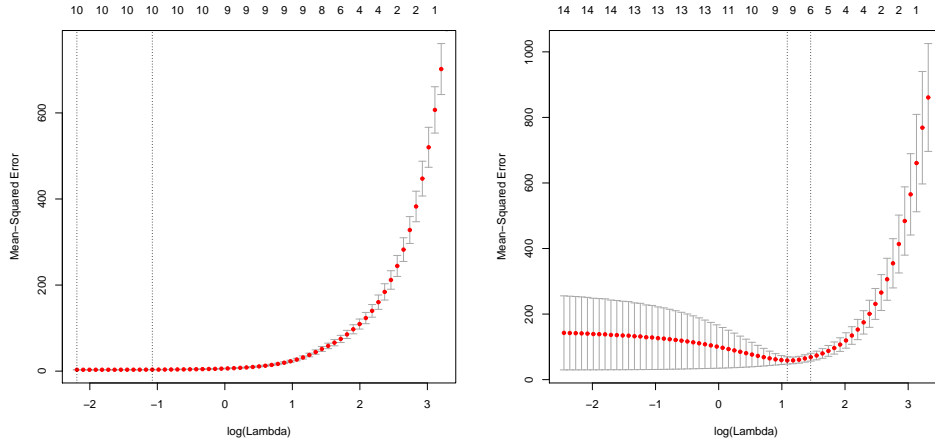


FIGURE 8.4: Plots of mean-squared error as a function of log regularization parameter (i.e. $\log \lambda$) for a regression of weight against all variables for the bodyfat dataset using an L1 regularizer (i.e. a lasso). These plots show mean-squared error averaged over cross-validation folds with a vertical one standard deviation bar. On the **left**, the plot for the dataset with the six outliers identified in Figure 16.20 removed. On the **right**, the plot for the whole dataset. Notice how the outliers increase the variability of the error, and the best error. The top row of numbers gives the number of non-zero components in $\hat{\beta}$. Notice how as λ increases, this number falls (there are 15 explanatory variables, so the largest model would have 15 variables). The penalty ensures that explanatory variables with small coefficients are dropped as λ gets bigger.

error encountered (one vertical line in the plot) and the largest value whose mean error is within one standard deviation of the minimum (the other vertical line in the plot). It is informative to keep track of the number of zeros in $\hat{\beta}$ as a function of λ , and this is shown in Figure 8.4.

Listing 8.1: R code used for the lasso regression example of worked example 8.4

```

setwd('/users/daf/Current/courses/Probcourse/Regression/RCode/HeightWeight');
library(gdata)
bfd<-read.xls('BodyFat.xls')
library(glmnet)
xmat<-as.matrix(bfd[, -c(1, 5, 18)])
ymat<-as.matrix(bfd[, 5])
# keeping in the outliers
dmodel<-cv.glmnet(xmat, ymat, alpha=0)
plot(dmodel)
# without outliers
cbfd<-bfd[-c(216, 39, 41, 42, 221, 163),]
xmat<-as.matrix(cbfd[, -c(1, 5, 18)])
ymat<-as.matrix(cbfd[, 5])
model<-cv.glmnet(xmat, ymat, alpha=0)
plot(model)

```

Worked example 8.3 *Building an L1 regularized regression*

Fit a linear regression to the bodyfat dataset, predicting weight as a function of all variables, and using the lasso to regularize. How good are the predictions? Do outliers affect the predictions?

Solution: I used the `glmnet` package, and I benefited a lot from example code by Trevor Hastie and Junyang Qian and published at https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html. You can see from Figure 8.7 that (a) for the case of outliers removed, the predictions are very good and (b) the outliers create problems. Note the magnitude of the error, and the low variance, for good cross validated choices. The main point of this example is to give you a start on producing R code, and I have put a code snippet in example 8.1.

Another way to understand the models is to look at how $\hat{\beta}$ changes as λ changes. We expect that, as λ gets smaller, more and more coefficients become non-zero. Figure 8.5 shows plots of coefficient values as a function of $\log \lambda$ for a regression of weight against all variables for the bodyfat dataset, penalised using the L_1 norm. For different values of λ , one gets different solutions for $\hat{\beta}$. When λ is very large, the penalty dominates, and so the norm of $\hat{\beta}$ must be small. In turn, most components of $\hat{\beta}$ are zero. As λ gets smaller, the norm of $\hat{\beta}$ falls and some components of become non-zero. At first glance, the variable whose coefficient grows very large seems important. Look more carefully; this is the last component introduced into the model. But Figure 8.4 implies that the right model has 7 components. This means that the right model has $\log \lambda \approx 1.3$, the vertical line shown in the detailed figure. In the best model, that coefficient is in fact zero.

The L_1 norm can sometimes produce an impressively small model from a large number of variables. In the UC Irvine Machine Learning repository, there is a dataset to do with the geographical origin of music (<https://archive.ics.uci.edu/>

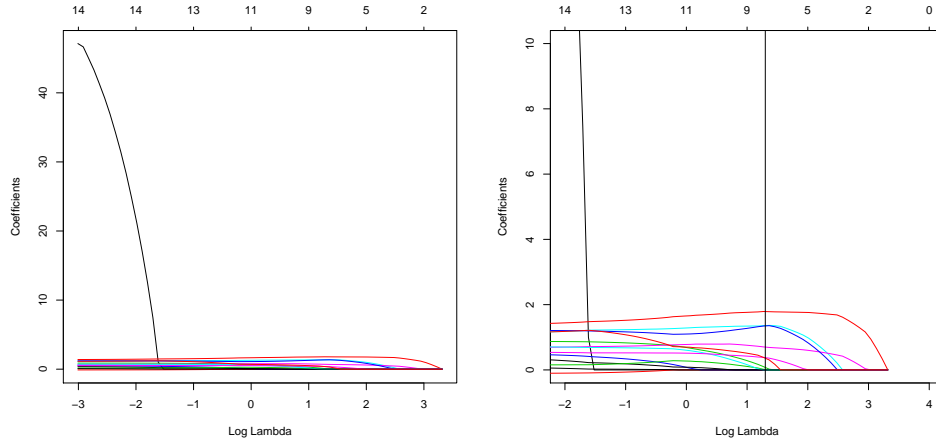


FIGURE 8.5: Plots of coefficient values as a function of $\log \lambda$ for a regression of weight against all variables for the bodyfat dataset, penalised using the L_1 norm. In each case, the six outliers identified in Figure 16.20 were removed. On the left, the plot of the whole path for each coefficient (each curve is one coefficient). On the right, a detailed version of the plot. The vertical line shows the value of $\log \lambda$ the produces the model with smallest cross-validated error (look at Figure 8.4). Notice that the variable that appears to be important, because it would have a large weight with $\lambda = 0$, does not appear in this model.

ml/datasets/Geographical+Original+of+Music). The dataset was prepared by Fang Zhou, and donors were Fang Zhou, Claire Q, and Ross D. King. Further details appear on that webpage, and in the paper: “Predicting the Geographical Origin of Music” by Fang Zhou, Claire Q and Ross. D. King, which appeared at ICDM in 2014. There are two versions of the dataset. One has 116 explanatory variables (which are various features representing music), and 2 independent variables (the latitude and longitude of the location where the music was collected). Figure 8.6 shows the results of a regression of latitude against the independent variables using L_1 regularization. Notice that the model that achieves the lowest cross-validated prediction error uses only 38 of the 116 variables.

Regularizing a regression with the L_1 norm is sometimes known as a **lasso**. A nuisance feature of the lasso is that, if several explanatory variables are correlated, it will tend to choose one for the model and omit the others (example in exercises). This can lead to models that have worse predictive error than models chosen using the L_2 penalty. One nice feature of good minimization algorithms for the lasso is that it is easy to use both an L_1 penalty and an L_2 penalty together. One can form

$$\left(\frac{1}{N}\right) \left(\sum_i (y_i - \mathbf{x}_i^T \beta)^2\right) + \lambda \left(\frac{(1-\alpha)}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1\right)$$

Error + Regularizer

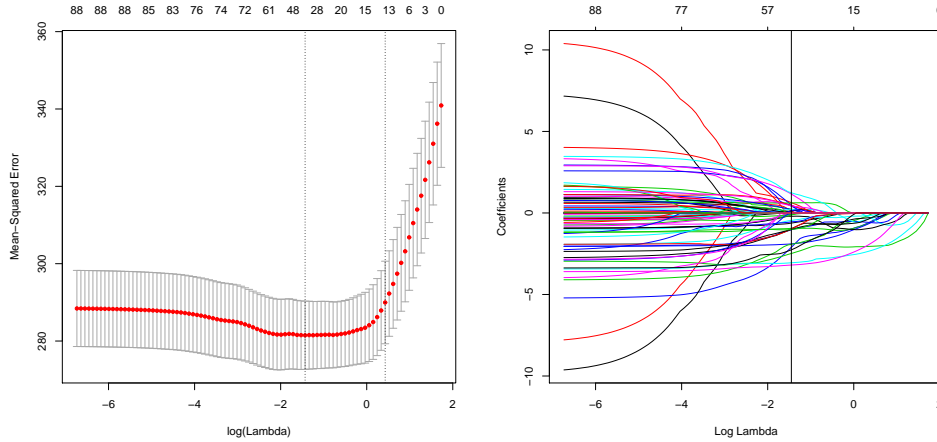


FIGURE 8.6: Mean-squared error as a function of log regularization parameter (i.e. $\log \lambda$) for a regression of latitude against features describing music (details in text), using the dataset at <https://archive.ics.uci.edu/ml/datasets/Geographical+Original+of+Music> and penalized with the L_1 norm. The plot on the left shows mean-squared error averaged over cross-validation folds with a vertical one standard deviation bar. The top row of numbers gives the number of non-zero components in $\hat{\beta}$. Notice how as λ increases, this number falls. The penalty ensures that explanatory variables with small coefficients are dropped as λ gets bigger. On the right, a plot of the coefficient values as a function of $\log \lambda$ for the same regression. The vertical line shows the value of $\log \lambda$ the produces the model with smallest cross-validated error. Only 38 of 116 explanatory variables are used by this model.

where one usually chooses $0 \leq \alpha \leq 1$ by hand. Doing so can both discourage large values in β and encourage zeros. Penalizing a regression with a mixed norm like this is sometimes known as **elastic net**. It can be shown that regressions penalized with elastic net tend to produce models with many zero coefficients, while not omitting correlated explanatory variables. All the computation can be done by the `glmnet` package in R (see exercises for details).

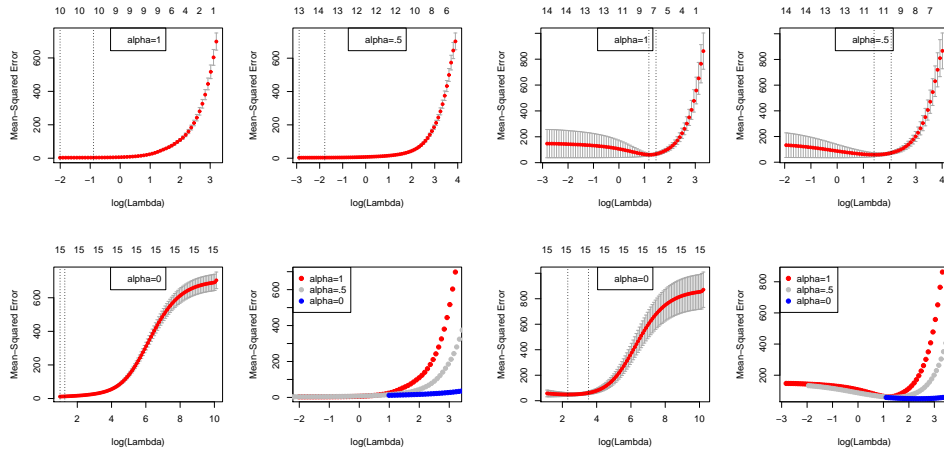


FIGURE 8.7: Plots of mean-squared error as a function of log regularization parameter (i.e. $\log \lambda$) for a regression of weight against all variables for the bodyfat dataset using an elastic net regularizer for various choices of α . The case $\alpha = 1$ corresponds to a lasso; $\alpha = 0$ corresponds to a ridge; and $\alpha = 0.5$ is one possible choice yielding an elastic net. These plots show mean-squared error averaged over cross-validation folds with a vertical one standard deviation bar. On the left, the plot for the dataset with the six outliers identified in Figure 16.20 removed. On the right, the plot for the whole dataset. Notice how the outliers increase the variability of the error, and the best error. The top row of numbers gives the number of non-zero components in $\hat{\beta}$. Notice how as λ increases, this number falls (there are 15 explanatory variables, so the largest model would have 15 variables). The penalty ensures that explanatory variables with small coefficients are dropped as λ gets bigger.

Worked example 8.4 *Building an elastic net regression*

Fit a linear regression to the bodyfat dataset, predicting weight as a function of all variables, and using the elastic net to regularize. How good are the predictions? Do outliers affect the predictions?

Solution: I used the `glmnet` package, and I benefited a lot from example code by Trevor Hastie and Junyang Qian and published at https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html. The package will do ridge, lasso and elastic net regressions. One adjusts a parameter in the function call, α , that balances the terms; $\alpha = 0$ is ridge and $\alpha = 1$ is lasso. You can see from Figure 8.7 that (a) for the case of outliers removed, the predictions are very good and (b) the outliers create problems. Note the magnitude of the error, and the low variance, for good cross validated choices. The main point of this example is to give you a start on producing R code, and I have put a code snippet in example 8.2.

Listing 8.2: R code used for the elastic net regression example of worked example 8.4

```

setwd( '/users/daf/Current/courses/Probcourse/Regression/RCode/HeightWeight' );
library( gdata )
bfd<-read.xls( 'BodyFat.xls' )
library( glmnet )
library( pls )
x<-as.matrix( bfd[, -c( 1, 5, 18 ) ] )
y<-as.matrix( bfd[, 5 ] )
foldid=sample( 1:10, size=length( y ), replace=TRUE )
cv1=cv.glmnet( x, y, foldid=foldid, alpha=1 )
cv.5=cv.glmnet( x, y, foldid=foldid, alpha=.5 )
cv0=cv.glmnet( x, y, foldid=foldid, alpha=0 )
par( mfrow=c( 2, 2 ) )
plot( cv1 );
legend( "top", legend="alpha=1" )
plot( cv.5 );
legend( "top", legend="alpha=.5" );
plot( cv0 );
legend( "top", legend="alpha=0" )
plot( log( cv1$lambda ), cv1$cvm, pch=19, col="red",
      xlab="log( Lambda )", ylab=cv1$name )
points( log( cv.5$lambda ), cv.5$cvm, pch=19, col="grey" )
points( log( cv0$lambda ), cv0$cvm, pch=19, col="blue" )
legend( "topleft",
       legend=c( "alpha=1", "alpha=.5", "alpha=0" ),
       pch=19, col=c( "red", "grey", "blue" ) )

```

8.4.2 Wide Datasets

Now imagine we have more independent variables than examples (this is sometimes referred to as a “wide” dataset). This occurs quite often for a wide range of datasets; it’s particularly common for biological datasets and natural language datasets. Unregularized linear regression must fail, because $\mathcal{X}^T \mathcal{X}$ must be rank deficient. Using an L2 (ridge) regularizer will produce an answer that should seem untrustworthy. The estimate of β is constrained by the data in some directions, but in other directions it is constrained only by the regularizer.

An estimate produced by L1 (lasso) regularization should look more reliable to you. Zeros in the estimate of β mean that the corresponding independent variables are ignored. Now if there are many zeros in the estimate of β , the model is being fit with a small subset of the independent variables. If this subset is small enough, then the number of independent variables that are actually being used is smaller than the number of examples. If the model gives low enough error, it should seem trustworthy in this case. There are some hard questions to face here (eg does the model choose the “right” set of variables?) that we can’t deal with.

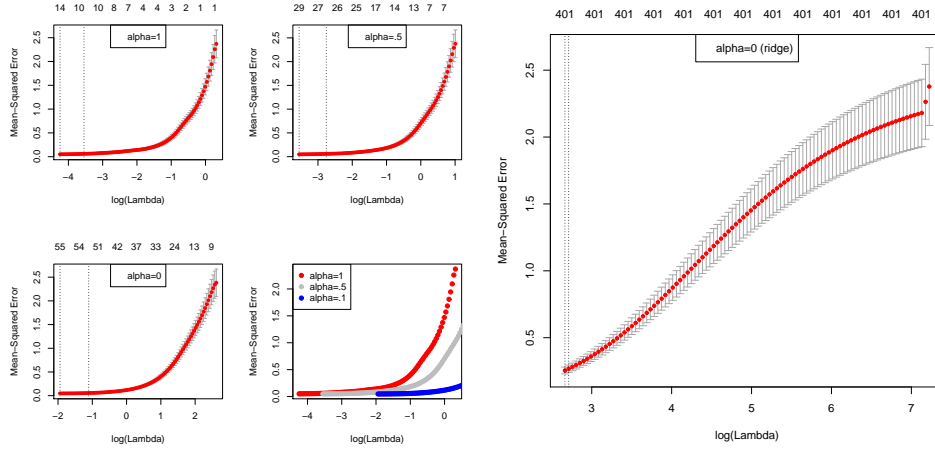


FIGURE 8.8: On the left, a comparison between three values of α in a `glmnet` regression predicting octane from NIR spectra (see Example 8.5). The plots show cross-validated error against log regularization coefficient for $\alpha = 1$ (lasso) and two elastic-net cases, $\alpha = 0.5$ and $\alpha = 0.1$. I have plotted these curves separately, with error bars, and on top of each other but without error bars. The values at the top of each separate plot show the number of independent variables with non-zero coefficients in the best model with that regularization parameter. On the right, a ridge regression for comparison. Notice that the error is considerably larger, even at the best value of the regularization parameter.

Worked example 8.5 *L1 regularized regression for a “wide” dataset*

The gasoline dataset has 60 examples of near infrared spectra for gasoline of different octane ratings. The dataset is due to John H. Kalivas, and was originally described in the article “Two Data Sets of Near Infrared Spectra”, in the journal *Chemometrics and Intelligent Laboratory Systems*, vol. 37, pp. 255259, 1997. Each example has measurements at 401 wavelengths. I found this dataset in the R library `pls`. Fit a regression of octane against infrared spectrum using L1 regularized logistic regression.

Solution: I used the `glmnet` package, and I benefited a lot from example code by Trevor Hastie and Junyang Qian and published at https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html. The package will do ridge, lasso and elastic net regressions. One adjusts a parameter in the function call, α , that balances the terms; $\alpha = 0$ is ridge and $\alpha = 1$ is lasso. Not surprisingly, the ridge isn’t great. I tried $\alpha = 0.1$, $\alpha = 0.5$ and $\alpha = 1$. Results in Figure 8.8 suggest fairly strongly that very good predictions should be available with the lasso using quite a small regularization constant; there’s no reason to believe that the best ridge models are better than the best elastic net models, or vice versa. The models are very sparse (look at the number of variables with non-zero weights, plotted on the top).

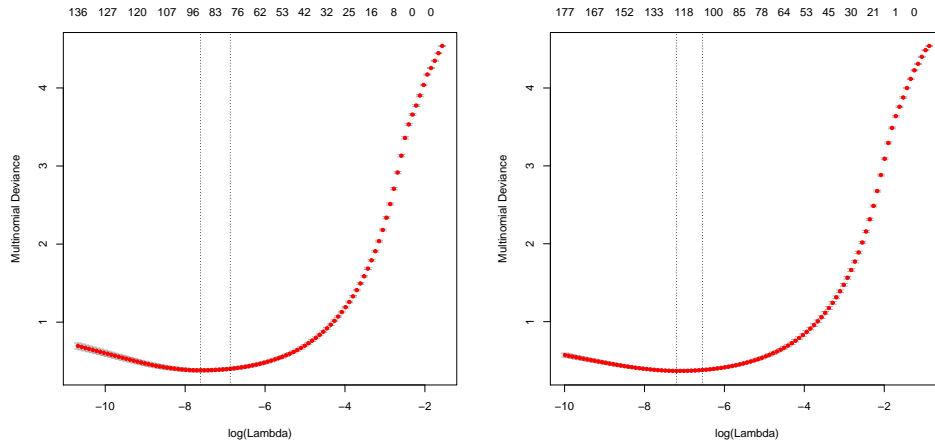


FIGURE 8.9: Multiclass logistic regression on the MNIST data set, using a lasso and elastic net regularizers. On the **left**, deviance of held out data on the digit data set (worked example 8.6), for different values of the log regularization parameter in the lasso case. On the **right**, deviance of held out data on the digit data set (worked example 8.6), for different values of the log regularization parameter in the elastic net case, $\alpha = 0.5$.

8.4.3 Using Sparsity Penalties with Other Models

A really nice feature of using an L1 penalty to enforce sparsity in a model is that it applies to a very wide range of models. For example, we can obtain a sparse SVM by replacing the L2 regularizer with an L1 regularizer. Most SVM packages will do this for you, although I'm not aware of any compelling evidence that this produces an improvement in most cases. All of the generalized linear models I described can be regularized with an L1 regularizer. For these cases, `glmnet` will do the computation required. The worked example shows using a multinomial (i.e. multiclass) logistic regression with an L1 regularizer.

Worked example 8.6 *Multiclass logistic regression with an L1 regularizer*

The MNIST dataset consists of a collection of handwritten digits, which must be classified into 10 classes ($0, \dots, 9$). There is a standard train/test split. This dataset is often called the zip code dataset because the digits come from zip codes, and has been quite widely studied. Yann LeCun keeps a record of the performance of different methods on this dataset at <http://yann.lecun.com/exdb/mnist/>. Obtain the Zip code dataset from <http://statweb.stanford.edu/~tibs/ElemStatLearn/>, and use a multiclass logistic regression with an L1 regularizer to classify it.

Solution: The dataset is rather large, and on my computer the fitting process takes a little time. Figure 8.9 shows what happens with the lasso, and with elasticnet with $\alpha = 0.5$ on the training set, using `glmnet` to predict and cross validation to select λ values. For the lasso, I found an error rate on the held out data of 8.5%, which is OK, but not great compared to other methods. For elastic net, I found a slightly better error rate (8.2%); I believe even lower error rates are possible with these codes.

Listing 8.3: R code used for the digit example of worked example 8.6

```

setwd('/users/daf/Current/courses/Probcourse/Regression/RCode/Digits');
library(glmnet)
digdat<-read.table('zip.train', sep='_', header=FALSE)
y<-as.factor(digdat$V1)
x<-as.matrix(digdat[, -c(1, 258)])
mod<-cv.glmnet(x, y, family="multinomial", alpha=1)
plot(mod)
# ok now we need to predict
digtest<-read.table('zip.test', sep='_', header=FALSE)
ytest<-as.factor(digtest$V1)
xtest<-as.matrix(digtest[, -c(1, 258)])
lmpredn<-predict(mod, xtest, type='class', s='lambda.min')
llpredn<-predict(mod, xtest, type='class', s='lambda.1se')
nmright<-sum(ytest==lmpredn)
errratem<-(1-nmright/dim(lmpredn))
nlright<-sum(ytest==llpredn)
errrate1<-(1-nlright/dim(llpredn))
mod.5<-cv.glmnet(x, y, family="multinomial", alpha=0.5)
plot(mod.5)
# ok now we need to predict
digtest<-read.table('zip.test', sep='_', header=FALSE)
ytest<-as.factor(digtest$V1)
xtest<-as.matrix(digtest[, -c(1, 258)])
lmpredn.5<-predict(mod.5, xtest, type='class', s='lambda.min')
llpredn.5<-predict(mod.5, xtest, type='class', s='lambda.1se')
nmright.5<-sum(ytest==lmpredn.5)
errratem.5<-(1-nmright.5/dim(lmpredn.5))
nlright.5<-sum(ytest==llpredn.5)
errrate1.5<-(1-nlright.5/dim(llpredn.5))

```

8.5 YOU SHOULD

8.5.1 remember:

New term: irreducible error	173
New term: bias	173
New term: variance	173
New term: AIC	175
New term: BIC	176
New term: forward stagewise regression	178
New term: Backward stagewise regression	178
New term: robust regression	180
New term: Huber loss	180
New term: scale	180
New term: inlier	180
New term: iteratively reweighted least squares	182
New term: MAD	183
New term: median absolute deviation	183
New term: Huber's proposal 2	183
New term: link function	183

Definition: Bernoulli random variable 184
New term: logit function 184
New term: logistic regression 184
New term: intensity 186
Definition: Poisson distribution 186
New term: deviance 186
New term: sparse models 187
New term: error cost 187
New term: L2 regularized error 187
New term: lasso 191
New term: elastic net 192

Bumps and Non-Parametric Regression

For many data sets, we expect that points near a given \mathbf{x} will have similar y values. We can exploit this observation in a variety of ways.

9.1 LINEAR MODELS USING BUMPS

Now think about the function y that we wish to interpolate, and assume that \mathbf{x} is “reasonably” scaled, meaning that distances between points are a good guide to their similarity. There are several ways to achieve this. We could whiten the points (section 16.20), or we could use our knowledge of the underlying problem to scale the different features relative to one another. Once we have done this properly, we expect the similarity between $y(\mathbf{u})$ and $y(\mathbf{v})$ to depend on the distance between these points (i.e. $\|\mathbf{u} - \mathbf{v}\|$) rather than on the direction of the vector joining them (i.e. $\mathbf{u} - \mathbf{v}$). Furthermore, we expect that the dependency should decline with increasing $\|\mathbf{x} - \mathbf{x}_i\|$.

This suggests a general approach. Choose a “bump” function, traditionally written $K(u)$ in one variable (u). This is called a bump function because it looks like a bump; it will have a large value at the origin, and fall off reasonably with distance from the origin. Now we can easily make a bump in d dimensions, by looking at $K(\|\mathbf{x} - \mathbf{x}_c\|)$. This will have its bump centered on \mathbf{x}_c , and will fall off with distance from \mathbf{x}_c . Conveniently, the bump falls off in the same way along each of the directions leaving \mathbf{x}_c . We now have a quite general modelling strategy: place a collection of bumps at a variety of points, and build a model of y by adjusting the weights of the bumps.

In most problems, we don’t know how quickly the weights should decline with increasing distance, and it is usual to have a scaling parameter to handle this. The scaling parameter will need to be selected.

9.1.1 Scattered Data: Smoothing and Interpolation

Imagine we have a set of points \mathbf{x}_i on the plane, with a measured height value y_i for each point. We would like to reconstruct a surface from this data. There are two important subcases: **interpolation**, where we want a surface that passes through each value; and **smoothing**, where our surface should be close to the values, but need not pass through them. This case is easily generalised to a larger number of dimensions. Particularly common is to have points in 3D, or in space and time.

Although this problem is very like regression, there is an important difference: we are interested only in the predicted value at each point, rather than in the conditional distribution. Typical methods for dealing with this problem are very like regression methods, but typically the probabilistic infrastructure required to

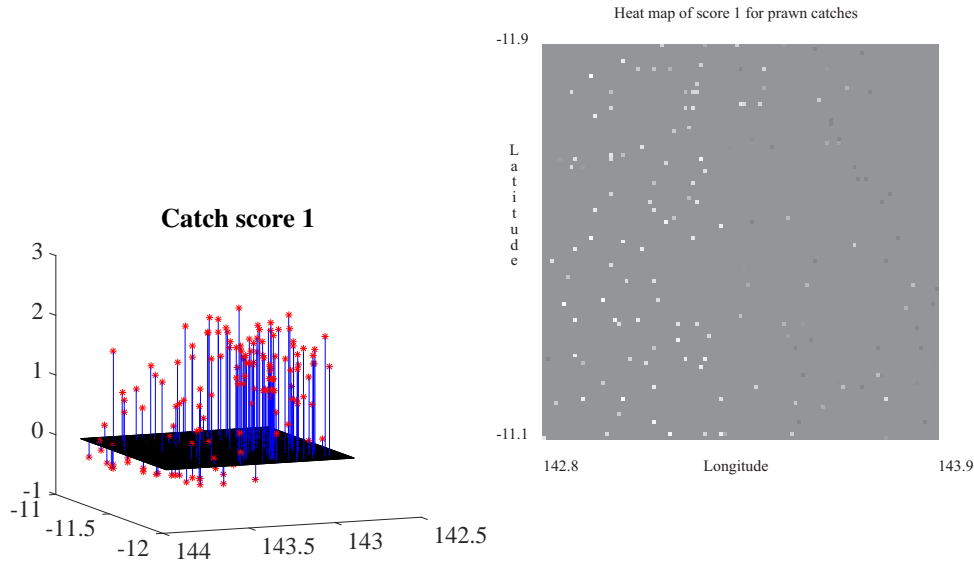


FIGURE 9.1: A dataset recording scores of prawn trawls around the Great Barrier Reef, from <http://www.statsci.org/data/oz/reef.html>. There are two scores; this is score 1. On the **left** I have plotted the data as a 3D scatter plot. This form of plot isn't usually very successful, though it helps to make it slightly easier to read if one supplies vertical lines from each value to zero, and a zero surface. On the **right**, a heat map of this data, made by constructing a fine grid, then computing the average score for each grid box (relatively few boxes get one score, and even fewer get two). The brightest point corresponds to the highest score; mid-grey is zero (most values), and dark points are negative. Notice that the scale is symmetric; the reason there is no very dark point is that the smallest value is considerably larger than the negative of the largest value. The x and y dimensions are longitude and latitude, and I have ignored the curvature of the earth, which is pretty small at this scale.

predict variances, standard errors, and the like are not developed. Interpolation and smoothing problems in one dimension have the remarkable property of being very different to those in two and more dimensions (if you've been through a graphics course, you'll know that, for example, interpolating splines in 2D are very different from those in 1D). We will concentrate on the multidimensional case.

We will formalize the idea of a bump function with a **kernel function**. A kernel function $K(u)$ is a non-negative function such that (a) $K(-u) = K(u)$ and (b) $\int_{-\infty}^{\infty} K(u)du = 1$. Widely used kernel functions are:

- **The Gaussian kernel**, $K(u) = \frac{1}{\sqrt{2\pi}} \exp -\frac{u^2}{2}$. Notice this doesn't have compact support.
- **The Epanechnikov kernel**, $K(u) = \frac{3}{4}(1 - u^2)\mathbb{I}_{[|u| \leq 1]}$. This isn't differentiable at $u = -1$ and at $u = 1$.

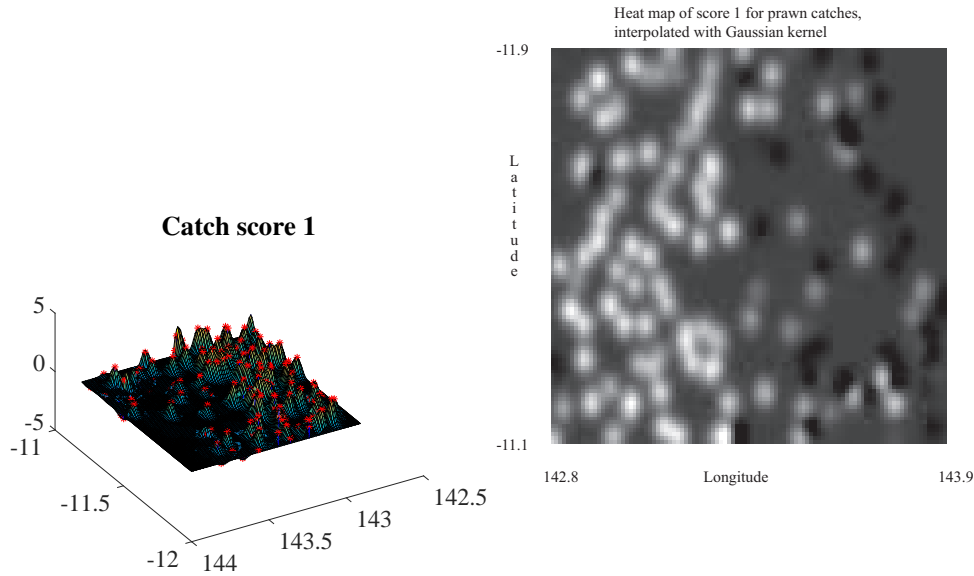


FIGURE 9.2: The prawn data of figure 9.1, interpolated with radial basis functions (in this case, a Gaussian kernel) with scale chosen by cross-validation. On the **left**, a surface shown with the 3D scatter plot. On the **right**, a heat map. I ignored the curvature of the earth, small at this scale, when computing distances between points. This figure is a good example of why interpolation is usually not what one wants to do (text).

- **The Logistic kernel**, $K(u) = \frac{1}{\exp -u + \exp u + 2}$. This doesn't have compact support, either.
- **The Quartic kernel**, $K(u) = \frac{15}{16}(1 - u^2)^2 \mathbb{I}_{|u| \leq 1}$.

You should notice that each is a bump function – it's large at $u = 0$, and falls away as $|u|$ increases. It follows from the two properties above that, for $h > 0$, if $K(u)$ is a kernel function, then $K(u; h) = \frac{1}{h}K(\frac{u}{h})$ is also a kernel function. This means we can vary the width of the bump at the origin in a natural way by choice of h ; this is usually known as the “scale” of the function.

We choose a kernel function K , then build a function

$$y(\mathbf{x}) = \sum_{j=1}^R a_j K\left(\frac{\|\mathbf{x} - \mathbf{b}_j\|}{h}\right)$$

where \mathbf{b}_j is a set of R points *which don't have to be training points*. These are sometimes referred to as **base points**. You can think of this process as placing a weighted bump at each base point. Consider the values that this function takes at

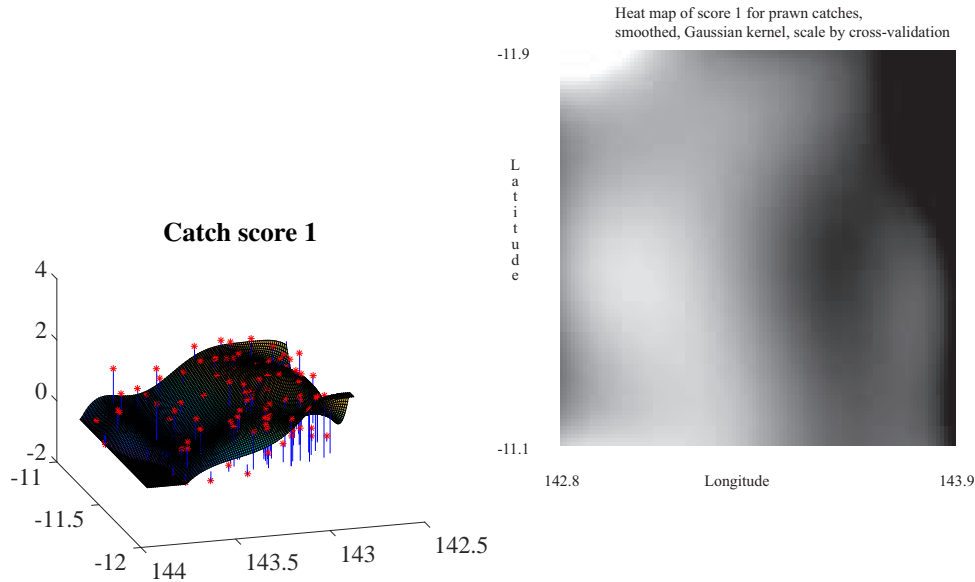


FIGURE 9.3: The prawn data of figure 9.1, smoothed with radial basis functions (in this case, a Gaussian kernel) with scale chosen by cross-validation. On the left, a surface shown with the 3D scatter plot. On the right, a heat map. I ignored the curvature of the earth, small at this scale, when computing distances between points. I used 60 basepoints, constructed by choosing 60 of 155 training points at random, then adding a small random offset.

the training points \mathbf{x}_i . We have

$$y(\mathbf{x}_i) = \sum_{j=1}^R a_j K\left(\frac{\|\mathbf{x}_i - \mathbf{b}_j\|}{h}\right)$$

and we would like to minimize $\sum_i (y_i - y(\mathbf{x}_i))^2$. We can rewrite this with the aid of some linear algebra. Write \mathcal{G} for the **Gram matrix**, whose i, j 'th entry is $K\left(\frac{\|\mathbf{x}_i - \mathbf{b}_j\|}{h}\right)$; write \mathbf{Y} for the vector whose i 'th component is y_i ; and \mathbf{a} for the vector whose j 'th component is a_j . Then we want to minimize

$$(\mathbf{Y} - \mathcal{G}\mathbf{a})^T(\mathbf{Y} - \mathcal{G}\mathbf{a}).$$

There are a variety of cases. For interpolation, we choose the base points to be the same as the training points. A theorem of Micchelli guarantees that for a kernel function that is (a) a function of distance and (b) monotonically decreasing with distance, \mathcal{G} will have full rank. Then we must solve

$$\mathbf{Y} = \mathcal{G}\mathbf{a} \quad (\text{Interpolation})$$

for \mathbf{a} and we will obtain a function that passes through each training point. For smoothing, we may choose any set of $R < N$ basepoints, though it's a good idea to

choose points that are close to the training points. Cluster centers are one useful choice. We then solve the least-squares problem by solving

$$\mathcal{G}^T \mathbf{Y} = \mathcal{G}^T \mathcal{G} \mathbf{a} \quad (\text{Smoothing})$$

for \mathbf{a} . In either case, we choose the scale h by cross-validation. We do this by: selecting a set of scales; holding out some training points, and interpolating (resp. smoothing) the values of the others; then computing the error of the predictions for these held-out points. The error is usually the square of the residual, and it's usually a good idea to average over many points when doing this.

Both interpolation and smoothing can present significant numerical challenges. If the dataset is large, \mathcal{G} will be large. For values of h that are large, \mathcal{G} (resp. $\mathcal{G}^T \mathcal{G}$) is usually very poorly conditioned, so solving the interpolation (resp. smoothing) system accurately is hard. This problem can usually be alleviated by adding a small constant (λ) times an identity matrix (\mathcal{I}) in the appropriate spot. So for interpolation, we solve

$$\mathbf{Y} = (\mathcal{G} + \lambda \mathcal{I}) \mathbf{a} \quad (\text{Interpolation})$$

for \mathbf{a} and we will obtain a function that passes through each training point. Similarly, for smoothing, we solve

$$\mathcal{G}^T \mathbf{Y} = (\mathcal{G}^T \mathcal{G} + \lambda \mathcal{I}) \mathbf{a} \quad (\text{Smoothing})$$

for \mathbf{a} . Usually, we choose λ to be small enough to make the linear algebra work ($1e - 9$ usually works for me), and ignore it.

Remember this: *Interpolation isn't as useful as most people think, because most measurements aren't exactly right or exactly repeatable.*

As Figure 9.2 suggests, interpolation isn't really as useful as you might think. Most measurements aren't exactly right, or exactly repeatable. If you look closely at the figure, you'll see one location where there are two scores; this is entirely to be expected for the score of a prawn trawl at a particular spot in the ocean. This creates problems for interpolation; \mathcal{G} must be rank deficient, because two rows will be the same. Another difficulty is that the scores look "wiggly" — moving a short distance can cause the score to change quite markedly. This is likely the effect of luck in trawling, rather than any real effect. The interpolating method chooses a very short scale, because this causes the least error in cross-validation, caused by predicting zero at the held out point (which is more accurate than any prediction available at any longer scale). The result is an entirely implausible model.

Now look at Figure 9.3. The smoothed surface is a reasonable guide to the scores; in the section of ocean where scores tend to be large and positive, so is the smoothed surface; where they tend to be negative, the smoothed surface is negative, too.

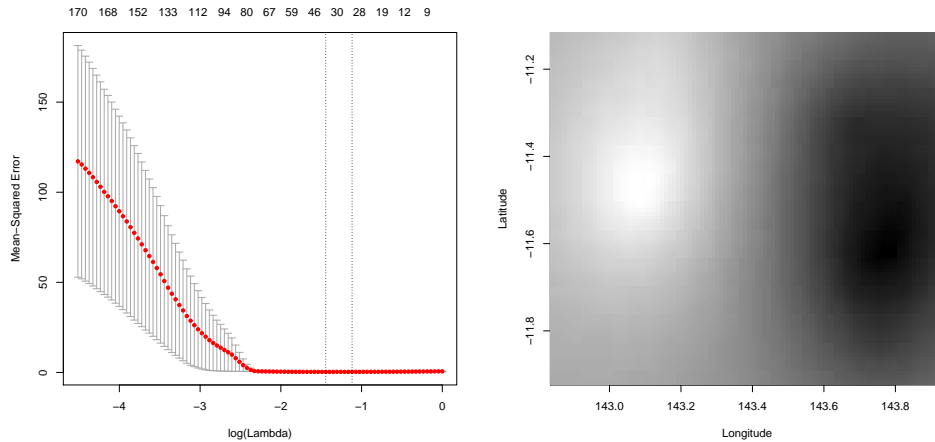


FIGURE 9.4: *The prawn data of figure 9.1, now smoothed with `glmnet`. On the left, the plot of held out error vs log regularization weight produced by `glmnet` for the prawn data, as described in example 9.1. Notice that the best regression uses between 20 and 30 of 930 possible features. On the right, a heat map. case, a Gaussian kernel) with scale chosen by cross-validation.*

9.1.2 Model Selection

We have constructed a model of y as a weighted sum of bump functions. We placed one bump function at each of a set of base points, then solved a linear system to get the weights that best represent y . You could see this as a version of linear regression, where we constructed new independent variables (the bump functions) from the original independent variables (the \mathbf{x}). Apart from the bump functions, we were following the general scheme of linear regression. This is important. It means that we can apply the material of chapter 16.20 to this idea. For example, you could choose the number and location of bumps using AIC or BIC and forward or backward stagewise regression. You could choose the number and location of bumps using a lasso penalty. You could change the squared error cost function to a robust cost function. All of these would be relatively routine (and really useful) combinations of methods we already know. In section 9.2.1, I will show a use of bump functions that is quite different from what we have seen before.

Worked example 9.1 *Using the lasso to predict prawn catches*

One can combine methods, very usefully. It is quite straightforward to persuade `glmnet` to apply the lasso to the prawn catch data.

Solution: I chose a set of base points (in this case, all training data points) and a set of six scales. I then built a bump of each scale at each base point. Write \mathbf{x}_i for each training data point, y_i for the value at the i 'th training point, \mathbf{b}_j for the j 'th base point, s_k for the k 'th scale and $K(\mathbf{x}; \mathbf{b}_j, s_k)$ for the bump of scale s_k centered at \mathbf{b}_j , as a function of position \mathbf{x} . I then constructed a vector of training features for the i 'th training point consisting of

$$\mathbf{f}_i = [K(\mathbf{x}_i; \mathbf{b}_1, s_1), \dots, K(\mathbf{x}_i; \mathbf{b}_1, s_k), \dots, K(\mathbf{x}_i; \mathbf{b}_j, s_1), \dots, K(\mathbf{x}_i; \mathbf{b}_j, s_k), \dots].$$

In this dataset, there are 155 points, and I used 6 scales, so this is a 930 dimensional vector. I then used `glmnet` to predict the values of y_i from \mathbf{f}_i . I used an elastic net regression (i.e. a convex combination of lasso and ridge penalties) because I expected strong correlations between the components of \mathbf{f}_i ; I used $\alpha = 0.5$. Figure 9.4 shows the result, which to me seems a better model than that of Figure 9.3 (there is slightly less wiggling). Notice how good `glmnet` is at reducing the number of variables. The best regression uses between 20 and 30 of 930 possible features (depending on precisely what criterion one uses to choose). This is likely because the method gets to select a scale for each kernel function, and can use more than one kernel at a base point. I have given the R code I used to produce Figure 9.4 in listing 9.1, to give some idea of the procedures involved.

9.2 NON PARAMETRIC BUMP MODELS

All of our regression models to date have been built by choosing some family of parametric functions of the independent variables, then choosing the parameters that give the best prediction. Write θ for a vector of parameters, \mathbf{x} for the independent variables, and $f(\mathbf{x}; \theta)$ for the function in the family chosen by θ ; then we obtain a “good choice” of parameters (write $\hat{\theta}$ for that choice), and our prediction is

$$f(\mathbf{x}; \hat{\theta}).$$

As a concrete example, all the regressions we have discussed so far have used the family of linear functions.

There are tremendous advantages to this approach. By choosing a family of functions with a “small” set of parameters, we can control the variance of the estimate. We can make predictions for points that are some way away from the training data, and they are often quite good. We have a collection of strong algorithms for choosing $\hat{\theta}$, and we can choose algorithms that impose properties on the estimate, too (eg sparseness; robustness to outliers). A model of this form is usually called a **parametric model**.

Listing 9.1: R code used for the elastic net regression example of worked example 9.4

```

setwd('~/.Current/Courses/Probcourse/Regression/Matlabcode/ScatteredData/')
prawns<-read.csv('clean-prawns.csv', header=FALSE)
setwd('~/.Current/Courses/Probcourse/Regression/RCode/Prawns/')
# now comes the good bit

ndat<-dim(prawns)[1]
srange<- c(0.1, 0.15, 0.2, 0.25, 0.3, 0.35)
xmat<-as.matrix(prawns[,c(1, 2)])
spaces<-dist(xmat, method = "euclidean", diag = FALSE, upper = FALSE)
msp<-as.matrix(spaces)
wmat<-exp(-msp/(2*srange[1]^2))
for (i in 2:6)
  {grammmat<-exp(-msp/(2*srange[i]^2))
  wmat<-cbind(wmat, grammmat)}
wmod<-cv.glmnet(wmat, as.vector(prawns[,3]), alpha=0.5)
xmin<-min(xmat[,1])
xmax<-max(xmat[,1])
ymin<-min(xmat[,2])
ymax<-max(xmat[,2])
xvec<-seq(xmin, xmax, length=100)
yvec<-seq(ymin, ymax, length=100)
# these are the points
pmat<-matrix(0, nrow=100*100, ncol=2)
ptr<-1
for (i in 1:100)
  {
    for (j in 1:100)
      {pmat[ptr,1]<-xvec[i]
      pmat[ptr, 2]<-yvec[j]
      ptr<-ptr+1}
  }
diff_ij <- function(i, j) sqrt(rowSums((pmat[i,]-xmat[j,])^2))
distsampletopts <- outer(seq_len(10000), seq_len(dim(xmat)[1]), diff_ij)
wmat<-exp(-distsampletopts/(2*srange[1]^2))
for (i in 2:6)
  {grammmat<-exp(-distsampletopts/(2*srange[i]^2))
  wmat<-cbind(wmat, grammmat)}
preds<-predict.cv.glmnet(wmod, wmat, s='lambda.min')
zmat<-matrix(0, nrow=100, ncol=100)
ptr<-1
for (i in 1:100)
  {for (j in 1:100)
  {zmat[i, j]<-preds[ptr]
  ptr<-ptr+1}
  }
}
wscale=abs(min(preds)), abs(max(preds))
image(yvec, xvec, (t(zmat)+wscale)/(2*wscale),
      xlab='Longitude', ylab='Latitude',
      col=grey(seq(0, 1, length=256)), useRaster=TRUE)
# this gets the heat map
plot(wmod)
# and this the regularization results

```

There are some disadvantages. By forcing the prediction to come from a particular family of functions, we are also forced to make inaccurate predictions for some datasets (those where that family doesn't encode the prediction well). It is hard to build a family of functions that is very flexible in regions where there are lots of data points, and quite inflexible where there are few (you could look at my lasso example on the prawn data as a step in that direction).

An alternative approach is to use a **non-parametric model**. The precise delineation between the classes is difficult to state, but a non-parametric model does not estimate parameter values. Instead, the model finds useful training examples and combines them to make predictions.

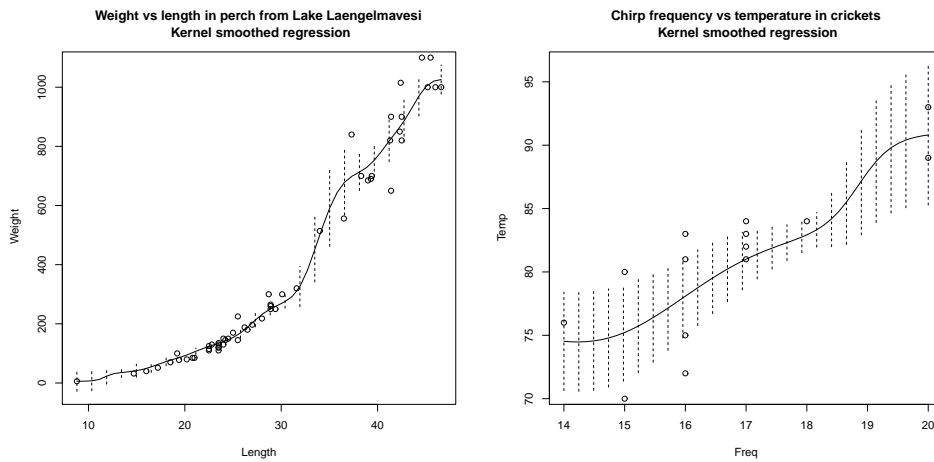


FIGURE 9.5: A non-parametric regression using a kernel smoothing method for the datasets of Figures 7.1 and 7.5. The curve shows the expected value of the independent variable for each value of the explanatory variable. The vertical bars show the standard error of the predicted density for the independent variable, for each value of the explanatory variable. Notice that, as one would expect, the standard deviation is smaller closer to data points, and larger further away. On the **left**, the perch data. On the **right**, the cricket data.

9.2.1 Kernel Smoothing

Bump functions offer one way to build a non-parametric regression. In the previous section, we used a set of bumps, centered on base points, to model a function of \mathbf{x} . The model was a weighted sum of bumps, with the weights adjusted to get the best prediction. This allowed us to use familiar machinery from linear regression to control the number of bumps, the scale, and so on. We did not require that the base points be the same as the example points (but didn't prevent it, either).

An alternative approach is to build functions $W(\mathbf{x}, \mathbf{x}_i)$ that blend together

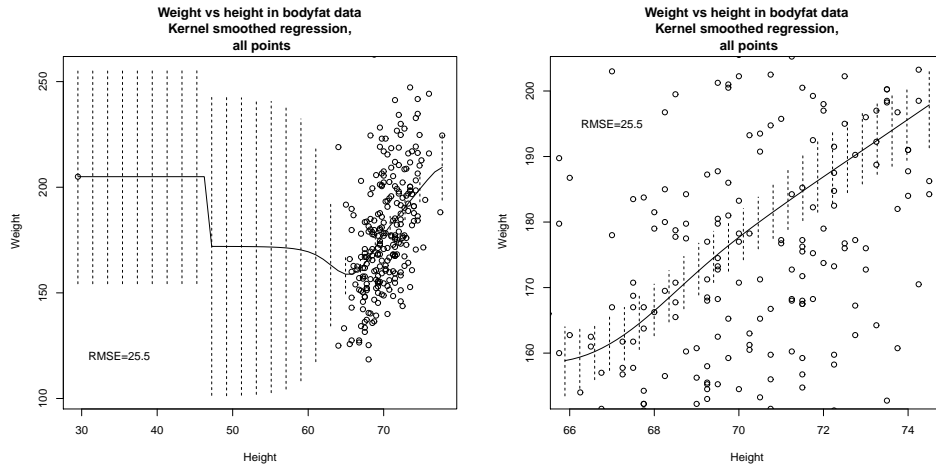


FIGURE 9.6: A non-parametric regression, using kernel smoothing, of weight against height, using the dataset of Figure 16.20. On the **left**, a curve showing the expected value of the independent variable for each value of the explanatory variable, with axes chosen to show the outlying point. The vertical bars show the standard error of the predicted density for the independent variable, for each value of the explanatory variable. Notice the method can make predictions around this point, but that these predictions have very large standard error as one would expect. On the **right**, the same curve, but on axes chosen to display the non-outlying points. The bars are small, because the standard error is small here (this isn't the same as the standard deviation of the predictive distribution).

the values at the training points. Our estimate of $y(\mathbf{x})$ is

$$\sum_{i=1}^N y_i W(\mathbf{x}, \mathbf{x}_i)$$

(i.e. we take each y_i , and blend with a weight function). We need good choices of $W(\mathbf{x}, \mathbf{x}_i)$. There are some simple, natural constraints we can impose. We should like $y(\mathbf{x})$ to be a convex combination of the observed values. This means we want $W(\mathbf{x}, \mathbf{x}_i)$ to be non-negative (so we can think of them as weights) and $\sum_i W(\mathbf{x}, \mathbf{x}_i) = 1$. Generally, we want $W(\mathbf{x}, \mathbf{x}_i)$ to fall off smoothly with distance from \mathbf{x}_i , but we would need to adjust the rate at which it falls off.

There is a natural construction using a kernel function $K(u)$. We use

$$W(\mathbf{x}, \mathbf{x}_i; h_i) = \frac{K\left(\frac{\|\mathbf{x}-\mathbf{x}_i\|}{h_i}\right)}{\sum_{j=1}^k K\left(\frac{\|\mathbf{x}-\mathbf{x}_j\|}{h_i}\right)}$$

where I have expanded the notation for the weight function to keep track of the scaling parameter, h_i , which we will need to select. There is one scaling parameter per training point. This should seem natural to you. If there are few points near a

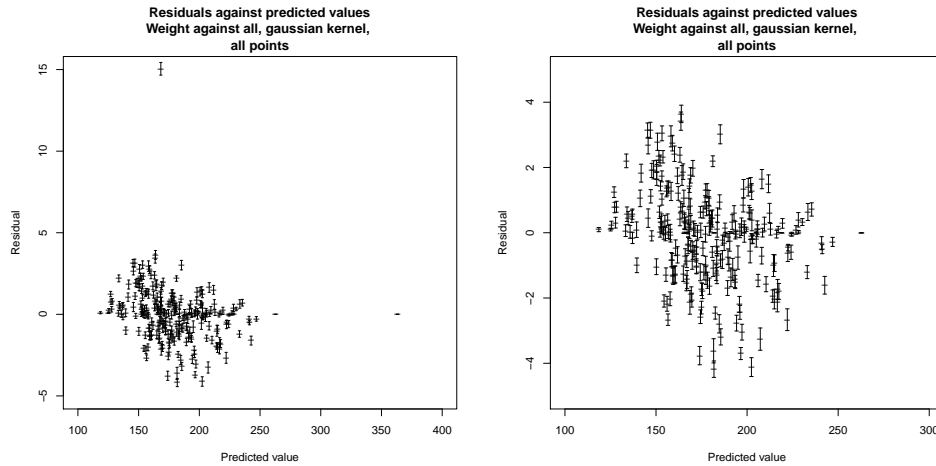


FIGURE 9.7: A non-parametric regression, using kernel smoothing, of weight against all variables, using the bodyfat dataset of Figure 16.20. On the **left**, a scatterplot of the predicted value against residual, on axes that show all points. The residual has not been standardized, because it is unclear how to do so. I have shown a one standard error bar for the residual as a vertical bar, plotted at the mean predicted value. For some points, the standard error is very small, so you can only see a horizontal bar. Notice how some points have a large residual, but most residuals are in a reasonable range. On the **right**, the same plot, but on axes chosen to display the non-outlying points. The bars are small, because the standard error is small here (this isn't the same as the standard deviation of the predictive distribution).

given data point, then we should have weights that fall off slowly, so that we can for a weighted average of those points. But if there are many points nearby, we can have weights that off fast. The weights are non-negative because we are using a non-negative kernel function. The weights sum to 1, because we divide at each point by the sum of all kernels evaluated at that point. For reference, this gives the expression

$$y(x) = \sum_{i=1}^N y_i \left(\frac{K\left(\frac{\|\mathbf{x}-\mathbf{x}_i\|}{h_i}\right)}{\sum_{j=1}^k K\left(\frac{\|\mathbf{x}-\mathbf{x}_j\|}{h_j}\right)} \right)$$

Changing the h_i will change the radius of the bumps, and will cause more (or fewer) points to have more (or less) influence on the shape of $y(\mathbf{x})$. Selecting the h_i is easy in principle. We search for a set of values that minimizes cross-validation error. In practice, this takes an extremely extensive search involving a great deal of computation, particularly if there are lots of points or the dimension is high. For the examples of this section, I used the R package `np`.

Listing 9.2: R code used for the perch example of worked example 9.2

```

setwd('/users/daf/Current/courses/Probcourse/Regression/RCode/FishCatch')
install.packages('np')
library(np)
pod<-read.table('perchonly.txt')
ndat<-dim(pod)
wtab<-matrix(data=NA, nrow=ndat, ncol=2)
wtab[,1]<-pod[,c(3)]
wtab[,2]<-pod[,c(6)]
wframe<-as.data.frame(wtab)
colnames(wframe)<-c('Weight', 'Length')
wm<-npreg(Weight~Length, data=wframe, ckertype='gaussian')
plot(wm, plot.errors.method='bootstrap',
      main="Weight vs length in perch from Lake Laengelmavesi\nKernel smoothed regression")
points(wtab[,2], wtab[, 1])

```

Worked example 9.2 *Using kernel smoothing to predict fish weight*

Use a kernel smoother to predict the weight of a perch from its length, using the dataset of http://www.amstat.org/publications/jse/jse_data_archive.htm (look for “fishcatch” on that page).

Solution: This is mainly an excuse to show you more code. I used the R code of listing 9.2 to produce the perch graph in Figure 9.5

A nice feature of this search is that it generates a great deal of information about the likely accuracy of the model. Each data point is going to be held out multiple times, so that we can estimate the standard error of the prediction at each data point. This is the standard deviation of the error in the estimate of the mean, due to random sampling effects. The details of this estimate are well beyond the scope of these notes. But `np` produces it, and I have shown standard error bars on the plots. You should notice that this is *not* the standard deviation of the predictive distribution (which we haven’t computed, and largely haven’t discussed). You should interpret the bars as the uncertainty in the prediction produced at a point due to the changes from training set to training set. Large bars indicate predictions you should not trust, because resampling the training data set produces large variations in the predictions at that point.

9.2.2 Density Estimation

One specialized application of kernel functions related to smoothing is **density estimation**. Here we have a set of N data points \mathbf{x}_i which we believe to be IID samples from some $p(X)$, and we wish to estimate the probability density function $p(X)$. In the case that we have a parametric model for $p(X)$, we could do so by estimating the parameters with (say) maximum likelihood. But we may not have such a model, or we may have data that does not comfortably conform to any model.

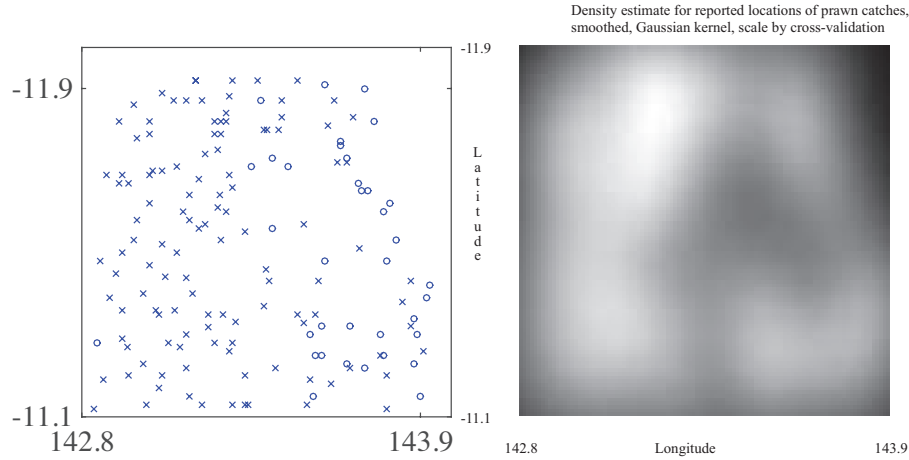


FIGURE 9.8: The prawn data of figure 9.1, now shown on the **left** as a scatter plot of locations from which scores were reported. The 'o's correspond to negative scores, and the '+'s to positive scores. On the **right**, a density estimate of the probability distribution from which the fishing locations were drawn, where lighter pixels correspond to larger density values. I ignored the curvature of the earth, small at this scale, when computing distances between points.

A natural, and important, model is to place probability $1/N$ at each data point, and zero elsewhere. This is sometimes called an **empirical distribution**. However, this model implies that we can only ever see the values we have already seen, which is often implausible or inconvenient. One widely used alternative is to fit a parametric model, using maximum likelihood. Sometimes data doesn't really fit these models well. It is often better to "smooth" the empirical distribution. If the \mathbf{x}_i have low enough dimension, we can construct a non-parametric density estimate in a straightforward way.

Recall that a kernel function is non-negative, and has the property that $\int_{-\infty}^{\infty} K(u)du = 1$. This means that if

$$y(\mathbf{x}) = \sum_{j=1}^R a_j K\left(\frac{\|\mathbf{x} - \mathbf{b}_j\|}{h}\right)$$

we have

$$\int_{-\infty}^{\infty} y(\mathbf{x})d\mathbf{x} = \sum_{j=1}^R a_j.$$

Now imagine we choose a basepoint at each data point, and we choose $a_j = 1/N$ for all j . The resulting function is non-negative, and integrates to one, so can be seen as a probability density function. We are placing a bump function of scale h , weighted by $1/N$ on top of each data point. If there are many data points close together, these bump functions will tend to reinforce one another and the resulting

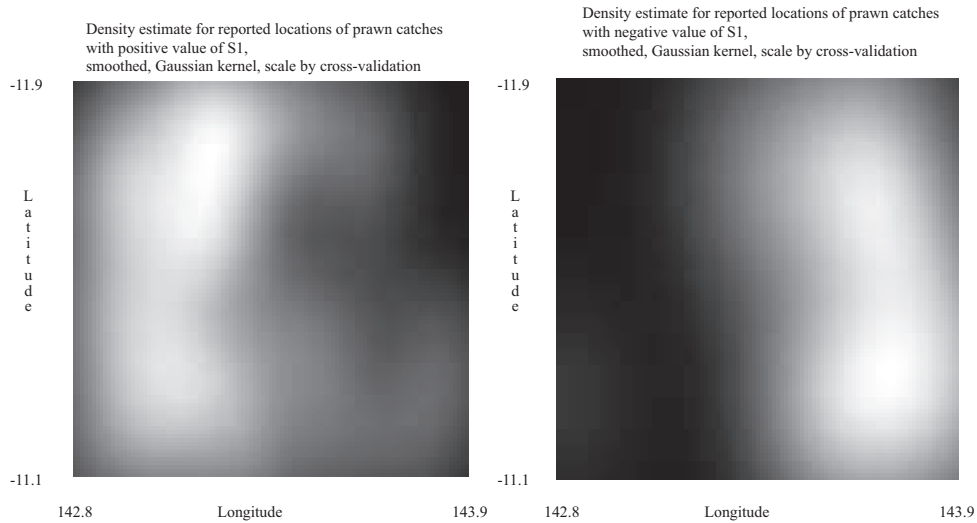


FIGURE 9.9: Further plots of the prawn data of figure 9.1. On the **left**, a density estimate of the probability distribution from which the fishing locations which achieved positive values of the first score were drawn, where lighter pixels correspond to larger density values. On the **right**, a density estimate of the probability distribution from which the fishing locations which achieved negative values of the first score were drawn, where lighter pixels correspond to larger density values. I ignored the curvature of the earth, small at this scale, when computing distances between points.

function will be large in such regions. The function will also be small in gaps between data points that are large compared to h . The resulting model captures the need to (a) represent the data and (b) allow values that have not been observed to have non-zero probability.

Choosing h using cross-validation is straightforward. We choose the h that maximises the log-likelihood of omitted “test” data, averaged over folds. There is one nuisance effect here to be careful of. If you use a kernel that has finite support, you could find that an omitted test item has zero probability; this leads to trouble with logarithms, etc. You could avoid this by obtaining an initial scale estimate with a kernel that has infinite support, then refining this with a kernel with finite support.

Density estimation techniques using kernels are useful in a variety of applications. For example, one can build very good naive Bayes classifiers using a kernel density estimate. This is a good replacement for a normal model of the class-conditional density of feature values. You should be careful about one point. We have seen that, for high dimensional data sets, most pairs of points are far apart. This means that, for high dimensional data sets, kernel density estimates can per-

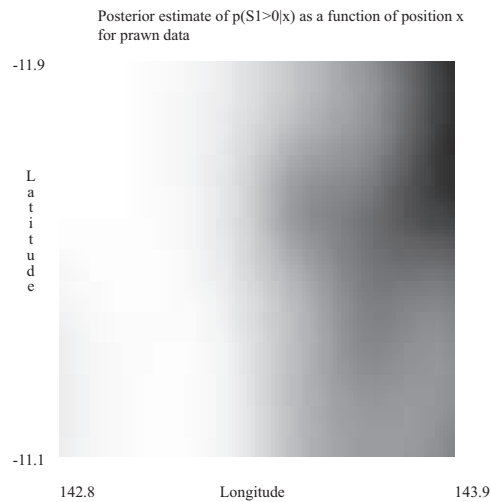


FIGURE 9.10: An estimate of the posterior probability of obtaining a positive value of the first catch score for prawns, as a function of position, from the prawn data of figure 9.2, using the density estimate of figure ???. Lighter pixels correspond to larger density values. I ignored the curvature of the earth, small at this scale, when computing distances between points.

form very poorly. Whether they give good or bad estimates depends quite a lot on how the points are distributed in the feature space, and so on the details of the particular problem.

9.3 EXPLOITING YOUR NEIGHBORS FOR REGRESSION

Nearest neighbors can clearly predict a number for a query example — you find the closest training example, and report its number. This would be one way to use nearest neighbors for regression, but it isn't terribly effective. One important difficulty is that the regression prediction is piecewise constant (Figure 9.11). If there is an immense amount of data, this may not present major problems, because the steps in the prediction will be small and close together. But it's not generally an effective use of data.

A more effective strategy is to find several nearby training examples, and use them to produce an estimate. This approach can produce very good regression estimates, because every prediction is made by training examples that are near to the query example. However, producing a regression estimate is expensive, because for every query one must find the nearby training examples.

Write \mathbf{x} for the query point, and assume that we have already collected the N nearest neighbors, which we write \mathbf{x}_i . Write y_i for the value of the dependent variable for the i 'th of these points. Notice that some of these neighbors could be

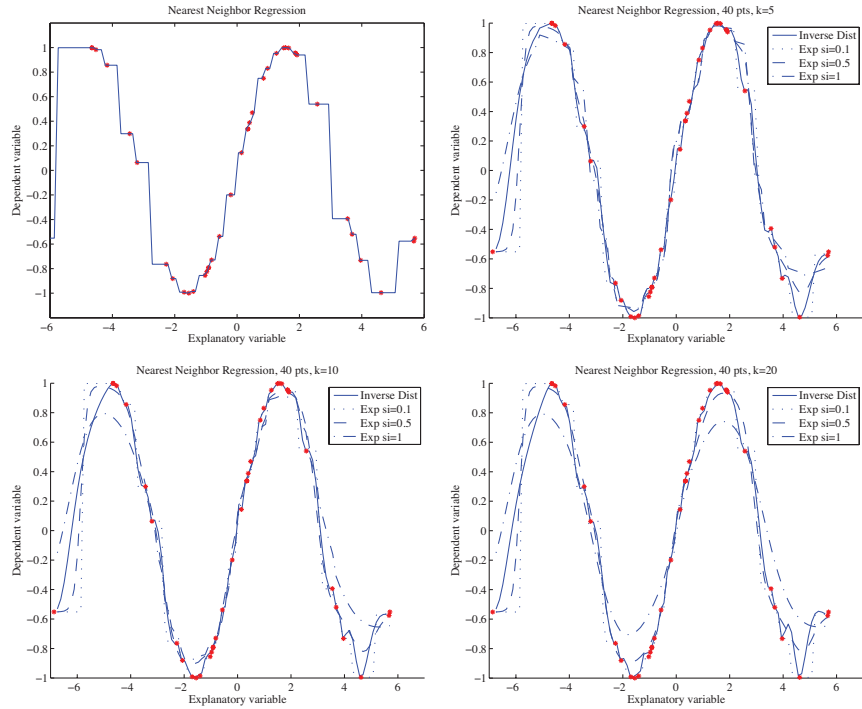


FIGURE 9.11: Different forms of nearest neighbors regression, predicting y from a one-dimensional x , using a total of 40 training points. **Top left:** reporting the nearest neighbor leads to a piecewise constant function. **Top right:** improvements are available by forming a weighted average of the five nearest neighbors, using inverse distance weighting or exponential weighting with three different scales. Notice if the scale is small, then the regression looks a lot like nearest neighbors, and if it is too large, all the weights in the average are nearly the same (which leads to a piecewise constant structure in the regression). **Bottom left and bottom right** show that using more neighbors leads to a smoother regression.

quite far from the query point. We don't want distant points to make as much contribution to the model as nearby points. This suggests forming a weighted average of the predictions of each point. Write w_i for the weight at the i 'th point. Then the estimate is

$$y_{pred} = \frac{\sum_i w_i y_i}{\sum_i w_i}.$$

A variety of weightings are reasonable choices. Write $d_i = \|\mathbf{x} - \mathbf{x}_i\|$ for the distance between the query point and the i 'th nearest neighbor. Then inverse distance weighting uses $w_i = 1/d_i$. Alternatively, we could use an exponential function to strongly weight down more distant points, using

$$w_i = \exp\left(\frac{-d_i^2}{2\sigma^2}\right).$$

We will need to choose a scale σ , which can be done by cross-validation. Hold out some examples, make predictions at the held out examples using a variety of different scales, and choose the scale that gives the best held-out error. Alternatively, if there are enough nearest neighbors, we could form a distance weighted linear regression, then predict the value at the query point from that regression.

Each of these strategies presents some difficulties when \mathbf{x} has high dimension. In that case, it is usual that the nearest neighbor is a lot closer than the second nearest neighbor. If this happens, then each of these weighted averages will boil down to evaluating the dependent variable at the nearest neighbor (because all the others will have very small weight in the average).

Remember this: *Nearest neighbors can be used for regression. In the simplest approach, you find the nearest neighbor to your feature vector, and take that neighbor's number as your prediction. More complex approaches smooth predictions over multiple neighbors.*

9.3.1 Local Polynomial Regression

Imagine wish to predict a value at a point \mathbf{x}_0 . Write y_0 for the value we predict. A reasonable choice would be to choose a number that minimizes a least squares difference to the training values, weighted in some way. We would like weights to be large for training points that are close to the test point, and small for points that are far away. We could achieve this using a kernel function, and choose the y_0 to be the value of α_0 that minimizes

$$\sum_i (y_i - \alpha_0)^2 K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_i\|}{h_i}\right).$$

If you differentiate and set to zero, etc., you will find you have the familiar expression for kernel regression

$$y_0 = \sum_{i=1}^N y_i \left(\frac{K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_i\|}{h_i}\right)}{\sum_{j=1}^k K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_j\|}{h_i}\right)} \right).$$

You can read the $K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_i\|}{h_i}\right)$ terms in the equation

$$\sum_i [y_i - \alpha_0]^2 K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_i\|}{h_i}\right)$$

as weights linking each training point to the test point \mathbf{x}_0 . Then minimizing this expression in α_0 is equivalent to choosing a single value α_0 such that the weighted errors (the terms $[y_i - \alpha_0]^2$) are minimized. There is a new function α_0 for each \mathbf{x}_0 , because when \mathbf{x}_0 changes all the weights change. This suggests something really

fruitful. Rather than thinking of α_0 as a constant placed at each training point, we should think of it as a *function* (that takes a constant value) and is placed at each training point. This is fruitful, because it suggests a way to generalize the model: place more complicated functions at each training point.

I will describe only the case of linear functions, but the others follow fairly naturally. Some, rather fussy, people may prefer me to describe these functions as “affine”. I won’t do this, because if I do, few people will know what I am talking about. We write a linear function around the point \mathbf{x}_i as $\alpha_0 + (\mathbf{x} - \mathbf{x}_i)^T \alpha_1$, where α_1 is a vector of coefficients. You should notice the analogy with a Taylor series.

To compute the value of y_0 corresponding to a test point \mathbf{x}_0 , we first minimize

$$\sum_i [y_i - (\alpha_0 + (\mathbf{x}_0 - \mathbf{x}_i)^T \alpha_1)]^2 K\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_i\|}{h_i}\right)$$

as a function of α_0 and α_1 . You can think about the resulting α_0 , α_1 as the (local) polynomial function of degree d that best approximates the function we are trying to represent. In particular, the weighted errors of trying to predict y_i at \mathbf{x}_i are minimized. We can use this function to predict y_0 at \mathbf{x}_0 ; but $(\mathbf{x}_0 - \mathbf{x}_0) = \mathbf{0}$, so our predicted value is α_0 .

9.3.2 Using your Neighbors to Predict More than a Number

Linear regression takes some features and predicts a number. But in practice, one often wants to predict something more complex than a number. For example, I might want to predict a parse tree (which has combinatorial structure) from a sentence (the explanatory variables). As another example, I might want to predict a map of the shadows in an image (which has spatial structure) against an image (the explanatory variables). As yet another example, I might want to predict which direction to move the controls on a radio-controlled helicopter (which have to be moved together) against a path plan and the current state of the helicopter (the explanatory variables).

Looking at neighbors is a very good way to solve such problems. The general strategy is relatively simple. We find a large collection of pairs of training data. Write \mathbf{x}_i for the explanatory variables for the i ’th example, and \mathbf{y}_i for the dependent variable in the i ’th example. This dependent variable could be anything — it doesn’t need to be a single number. It might be a tree, or a shadow map, or a word, or anything at all. I wrote it as a vector because I needed to choose some notation.

In the simplest, and most general, approach, we obtain a prediction for a new set of explanatory variables \mathbf{x} by (a) finding the nearest neighbor and then (b) producing the dependent variable for that neighbor. We might vary the strategy slightly by using an approximate nearest neighbor. If the dependent variables have enough structure that it is possible to summarize a collection of different dependent variables, then we might recover the k nearest neighbors and summarize their dependent variables. How we summarize rather depends on the dependent variables. For example, it is a bit difficult to imagine the average of a set of trees, but quite straightforward to average images. If the dependent variable was a word, we might not be able to average words, but we can vote and choose the

most popular word. If the dependent variable is a vector, we can compute either distance weighted averages or a distance weighted linear regression.

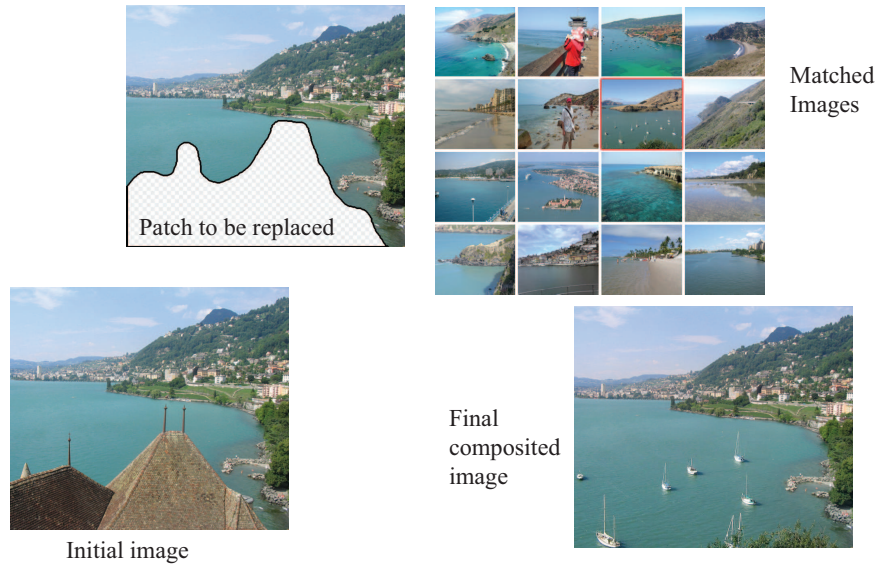


FIGURE 9.12: We can fill large holes in images by matching the image to a collection, choosing one element of the collection, then cutting out an appropriate block of pixels and putting them into the hole in the query image. In this case, the hole has been made by an artist, who wishes to remove the roofline from the view. Notice how there are a range of objects (boats, water) that have been inserted into the hole. These objects are a reasonable choice, because the overall structures of the query and matched image are largely similar — getting an image that matches most of the query image supplies enough context to ensure that the rest “makes sense”.

Example: Filling Large Holes with Whole Images Many different kinds of user want to remove things from images or from video. Art directors might like to remove unattractive telephone wires; restorers might want to remove scratches or marks; there’s a long history of government officials removing people with embarrassing politics from publicity pictures (see the fascinating examples in ?); and home users might wish to remove a relative they dislike from a family picture. All these users must then find something to put in place of the pixels that were removed.

If one has a large hole in a large image, we may not be able to just extend a texture to fill the hole. Instead, entire objects might need to appear in the hole (Figure 9.12). There is a straightforward, and extremely effective, way to achieve this. We match the image to a large collection of images, to find the nearest neighbors (the details of the distance function are below). This yields a set of example images where all the pixels we didn’t want to replace are close to those of the query image. From these, we choose one, and fill in the pixels from that image.

There are several ways to choose. If we wish to do so automatically, we could use the example with the smallest distance to the image. Very often, an artist is involved, and then we could prepare a series of alternatives — using, perhaps, the k closest examples — then show them to the artist, who will choose one. This method, which is very simple to describe, is extremely effective in practice.

It is straightforward to get a useful distance between images. We have an image with some missing pixels, and we wish to find nearby images. We will assume that all images are the same size. If this isn't in fact the case, we could either crop or resize the example images. A good measure of similarity between two images \mathcal{A} and \mathcal{B} can be measured by forming the **sum of squared differences** (or **SSD**) of corresponding pixel values. You should think of an image as an array of pixels. If the images are grey-level images, then each pixel contains a single number, encoding the grey-level. If the images are color images, then each pixel (usually!) contains three values, one encoding the red-level, one encoding the green-level, and one encoding the blue-level. The SSD is computed as

$$\sum_{(i,j)} (\mathcal{A}_{ij} - \mathcal{B}_{ij})^2$$

where i and j range over all pixels. If the images are grey-level images, then by $(\mathcal{A}_{ij} - \mathcal{B}_{ij})^2$, I mean the squared difference between grey levels; if they are color images, then this means the sum of squared differences between red, green and blue values. This distance is small when the images are similar, and large when they are different (it is essentially the length of the difference vector).

Now we don't know some of the pixels in the query image. Write \mathcal{K} for the set of pixels around a point whose values are known, and $\#\mathcal{K}$ for the size of this set. We can now use

$$\frac{1}{\#\mathcal{K}} \sum_{(i,j) \in \mathcal{K}} (\mathcal{A}_{ij} - \mathcal{B}_{ij})^2.$$

Filling in the pixels requires some care. One does not usually get the best results by just copying the missing pixels from the matched image into the hole. Instead, it is better to look for a good **seam**. We search for a curve enclosing the missing pixels which (a) is reasonably close to the boundary of the missing pixels and (b) gives a good boundary between the two images. A good boundary is one where the query image (on one side) is “similar to” the matched image (on the other side). A good sense of similarity requires that pixels match well, and that image gradients crossing the boundary tend to match too.

Remember this: *Nearest neighbors can be used to predict more than numbers. Examples include parse trees, blocks of pixels, and so on.*

CHAPTER 10

Neural Networks

10.1 UNITS AND CLASSIFICATION

We will build complex classification systems out of simple units. A **unit** takes a vector \mathbf{x} of inputs and uses a vector \mathbf{w} of parameters (known as the **weights**), a scalar b (known as the **bias**), and a nonlinear function F to form its output, which is

$$F(\mathbf{w}^T \mathbf{x} + b).$$

Over the years, a wide variety of nonlinear functions have been tried. Current best practice is to use the **RELU** (for rectified linear unit), where

$$F(u) = \max(0, u).$$

For example, if \mathbf{x} was a point on the plane, then a single unit would represent a line, chosen by the choice of \mathbf{w} and b . The output for all points on one side of the line would be zero. The output for points on the other side would be a positive number that is larger for points that are further from the line.

Units are sometimes referred to as **neurons**, and there is a large and rather misty body of vague speculative analogy linking devices built out of units to neuroscience. I deprecate this practice; what we are doing here is quite useful and interesting enough to stand on its own without invoking biological authority. Also, if you want to see a real neuroscientist laugh, explain to them how your neural network is really based on some gobbet of brain tissue or other.

10.1.1 Building a Classifier out of Units: The Cost Function

We will build a multiclass classifier out of units by modelling the class posterior probabilities using the outputs of the units. Each class will get the output of a single unit. Write o_i for the output of the i 'th unit, and θ for all the parameters in all the units. We will organize these units into a vector \mathbf{o} , whose i 'th component is o_i . We want to use that unit to model the probability that the input is of class j , which I will write $p(\text{class} = j | \mathbf{x}, \theta)$. To build this model, I will use the **softmax function**. This is a function that takes a C dimensional vector and returns a C dimensional vector. I will write $\mathbf{s}(\mathbf{u})$ for the softmax function, and the dimension C will always be the number of classes. We have

$$\mathbf{s}(\mathbf{u}) = \left(\frac{1}{\sum_k e^{u_k}} \right) \begin{bmatrix} e^{u_1} \\ e^{u_2} \\ \dots \\ e^{u_C} \end{bmatrix}$$

(recall u_i is the i 'th component of \mathbf{u}). We then use the model

$$p(\text{class} = i | \mathbf{x}, \theta) = s_i(\mathbf{o}(\mathbf{x}, \theta)).$$

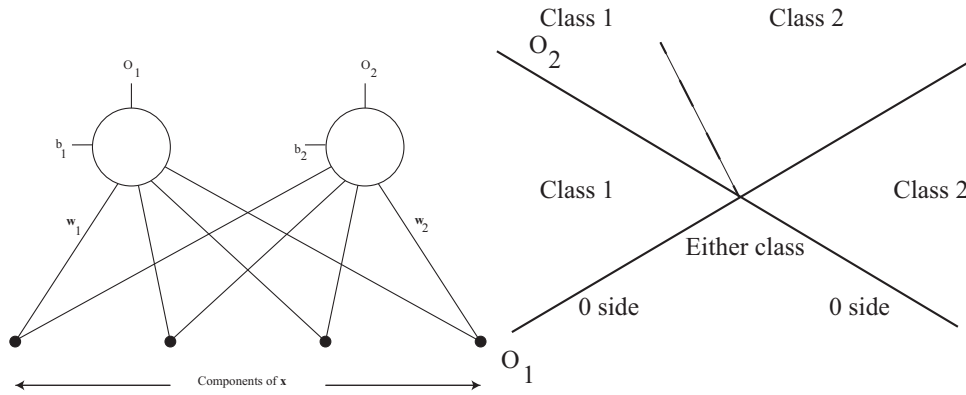


FIGURE 10.1: On the left, two units observing an input vector, and providing outputs. On the right, the decision boundary for two units classifying a point on the plane into one of two classes. The angle of the dashed line depends on the magnitudes of w_1 and w_2 .

Notice that this expression passes important tests for a probability model. Each value is between 0 and 1, and the sum over classes is 1.

In this form, the classifier is not super interesting. For example, imagine that the features x are points on the plane, and we have two classes. Then we have two units, one for each class. There is a line corresponding to each unit; on one side of the line, the unit produces a zero, and on the other side, the unit produces a positive number that increases as with perpendicular distance from the line. We can get a sense of what the decision boundary will be like from this. When a point is on the 0 side of both lines, the class probabilities will be equal (and so both $\frac{1}{2}$ – two classes, remember). When a point is on the positive side of the i 'th line, but the zero side of the other, the class probability for class i will be

$$\frac{e^{o_i(\mathbf{x},\theta)}}{1 + e^{o_i(\mathbf{x},\theta)}}$$

and the point will always be classified in the i 'th class (remember, $o_i \geq 0$). Finally, when a point is on the positive side of both lines, the classifier boils down to choosing the i that has the largest value of $o_i(\mathbf{x},\theta)$. All this leads to the decision boundary shown in figure ???. Notice that this is piecewise linear, and somewhat more complex than the boundary of an SVM. It's quite helpful to try and draw what would happen for three or more classes with x a 2D point.

10.1.2 Building a Classifier out of Units: Strategy

The essential difficulty here is to choose θ that results in the best behavior. We will do so by writing a cost function that estimates the error rate of the classification, then choosing a value $\hat{\theta}$ that minimises that function. We have a set of N examples x_i and for each example we know the class. There are a total of C classes. We encode the class of an example using a **one hot** vector y_i , which is C dimensional.

If the i 'th example is from class j , then the j 'th component of \mathbf{y}_i is 1, and all other components in the vector are 0. I will write y_{ij} for the j 'th component of \mathbf{y}_i .

A natural cost function looks at the log likelihood of the data under the probability model produced from the outputs of the units. If the i 'th example is from class j , we would like $-\log p(\text{class} = j | \mathbf{x}_i, \theta)$ to be small (notice the sign here; it's usual to minimize negative log likelihood). I will write $\log \mathbf{s}$ to mean the vector whose components are the logarithms of the components of \mathbf{s} . This yields a loss function

$$\frac{1}{N} \sum_{i \in \text{data}} \left[\sum_{j \in \text{classes}} \{-\mathbf{y}_i^T \log \mathbf{s}(\mathbf{o}(\mathbf{x}_i, \theta))\} \right].$$

Notice that this loss function is written in a clean way that may lead to a poor implementation. I have used the y_{ij} values as “switches”, as in the discussion of EM. This leads to clean notation, but hides fairly obvious computational efficiencies (when taking the gradient, you need to deal with only one term in the sum over classes). As in the case of the linear SVM (section 16.20), we would like to achieve a low cost with a “small” θ , and so form an overall cost function that will have loss and penalty terms.

There are a variety of possible penalties. For now, we will penalize large sets of weights, but we'll look at other possibilities below. Remember, we have C units (one per class) and so there are C distinct sets of weights. Write the weights for the u 'th unit \mathbf{w}_u . Our penalty becomes

$$\sum_{u \in \text{units}} \mathbf{w}_u^T \mathbf{w}_u.$$

As in the case of the linear SVM (section 16.20), we write λ for a weight applied to the penalty. Our cost function is then

$$S(\theta, \mathbf{x}; \lambda) = \frac{1}{N} \sum_{i \in \text{data}} \left[\sum_{j \in \text{classes}} \{-\mathbf{y}_i^T \log \mathbf{s}(\mathbf{o}(\mathbf{x}_i, \theta))\} \right] + \frac{\lambda}{2} \sum_{u \in \text{units}} \mathbf{w}_u^T \mathbf{w}_u$$

(misclassification loss) (penalty)

10.1.3 Building a Classifier out of Units: Training

I have described a simple classifier built out of units. We must now train this classifier, by choosing a value of θ that results in a small loss. It may be quite hard to get the true minimum, and we may need to settle for a small value. We use stochastic gradient descent, because we have seen it before; because it is effective; and because it is the algorithm of choice when training more complex classifiers built out of units.

For the SVM, we selected one example at random, computed the gradient at that example, updated the parameters, and went again. For neural nets, it is more usual to use **minibatch training**, where we select a subset of the data uniformly and at random, compute a gradient using that subset, update and go again. This is because in the best implementations many operations are vectorized, and using a minibatch can provide a gradient estimate that is clearly better than

that obtained using only one example, but doesn't take longer to compute. The size of the minibatch is usually determined by memory or architectural considerations. It is often a power of two, for this reason.

Now imagine we have chosen a minibatch of M examples. We must compute the gradient of the cost function. This is mainly an exercise in notation, but there's a lot of notation. Write θ_u for a vector containing all the parameters for the u 'th unit, so that $\theta_u = [\mathbf{w}_u, b_u]^T$. Recall $s_k(\mathbf{o}(\mathbf{x}_i, \theta_k))$ is the output of the softmax function for the k 'th unit for input \mathbf{x}_i . This represents the probability that example i is of class k under the current model. Then we must compute

$$\nabla_{\theta_u} \frac{1}{M} \sum_{i \in \text{minibatch}} [\{-\mathbf{y}_i^T \log \mathbf{s}(\mathbf{o}(\mathbf{x}_i, \theta))\}] + \frac{\lambda}{2} \sum_{j \in \text{classes}} \mathbf{w}_j^T \mathbf{w}_j.$$

The gradient is easily computed using the chain rule. The term

$$\frac{\lambda}{2} \sum_{j \in \text{classes}} \mathbf{w}_j^T \mathbf{w}_j$$

presents no challenge, but the other term is more interesting. We must differentiate the softmax function by its inputs, then the units by their parameters. More notation: assume we have a vector valued function of vector inputs, for example, $\mathbf{s}(\mathbf{o})$. Here \mathbf{s} is the function and \mathbf{o} are the inputs. I will write $\#(\mathbf{o})$ to mean the number of components of \mathbf{o} , and o_i for the i 'th component. The matrix of first partial derivatives is extremely important (we will see a lot of these; pay attention). I will write $\mathcal{J}_{\mathbf{s}; \mathbf{o}}$ to mean

$$\begin{pmatrix} \frac{\partial s_1}{\partial o_1} & \cdots & \frac{\partial s_1}{\partial o_{\#(\mathbf{o})}} \\ \cdots & \cdots & \cdots \\ \frac{\partial s_{\#(\mathbf{s})}}{\partial o_1} & \cdots & \frac{\partial s_{\#(\mathbf{s})}}{\partial o_{\#(\mathbf{o})}} \end{pmatrix}$$

and refer to such a matrix of first partial derivatives as a **Jacobian**.

Now we can use the chain rule to write

$$\nabla_{\theta_u} \frac{1}{M} \sum_{i \in \text{minibatch}} [\{-\mathbf{y}_i^T \log \mathbf{s}(\mathbf{x}_i, \theta)\}] = \frac{1}{M} \sum_{i \in \text{minibatch}} [\{-\mathbf{y}_i^T \mathcal{J}_{\log \mathbf{s}; \mathbf{o}} \mathcal{J}_{\mathbf{o}; \theta_u}\}].$$

This isn't particularly helpful without knowing the relevant Jacobians. They're quite straightforward.

Write $\mathbb{I}_{[u=v]}(u, v)$ for the indicator function that is 1 when $u = v$ and zero otherwise. We have

$$\begin{aligned} \frac{\partial \log s_u}{\partial o_v} &= \mathbb{I}_{[u=v]} - \frac{e^{o_v}}{\sum_k e^{o_k}} \\ &= \mathbb{I}_{[u=v]} - s_v. \end{aligned}$$

To get the other Jacobian, we need yet more notation (but this isn't new, it's a reminder). I will write $w_{u,i}$ for the i 'th component of \mathbf{w}_u , and $\mathbb{I}_{[o_u > 0]}(o_u)$ for the indicator function that is 1 if its argument is greater than zero. Then

$$\frac{\partial o_u}{\partial w_{u,i}} = x_i \mathbb{I}_{[o_u > 0]}(o_u)$$

and

$$\frac{\partial o_u}{\partial b_u} = \mathbb{I}_{[o_u > 0]}(o_u).$$

Notice that if $v \neq u$,

$$\frac{\partial o_u}{\partial w_{v,i}} = 0 \text{ and } \frac{\partial o_u}{\partial b_v} = 0.$$

At least in principle, we can build a multiclass classifier in a straightforward way using minibatch gradient descent. We use one unit per class, each one using each component of the feature vector. We obtain training data, and then iterate computing a gradient from a minibatch, and taking a step along the negative of the gradient. If you try, you may run into some of the important small practical problems that cause networks to work badly. Here are some of the ones you may encounter.

Initialization: You need to choose the initial values of all of the parameters. There are many parameters; in our case, with a d dimensional \mathbf{x} and C classes, we have $(d + 1) \times C$ parameters. If you initialize each parameter to zero, you will find that the gradient is also zero, which is not helpful. This occurs because all the o_u will be zero (because the $w_{u,i}$ and the b_u are zero). It is usual to initialize to draw a sample of a zero mean normal random variable for each initial value (appropriate choices of variance get interesting; more below).

Learning rate: Each step will look like $\theta^{(n+1)} = \theta^{(n)} - \eta_n \nabla_{\theta} \text{cost}$. You need to choose η_n for each step. This is widely known as the **learning rate**; an older term is **steplength** (neither term is a super-accurate description). It is not usual for the learning rate to be the same throughout learning. We would like to take “large” steps early, and “small” steps late, in learning, so we would like η_n to be “large” for small n , and “small” for large n . It is tough to be precise about a good choice. As in stochastic gradient descent for a linear SVM, breaking learning into epochs ($e(n)$ is the epoch of the n 'th iteration), then choosing two constants a and b to obtain

$$\eta_n = \frac{1}{a + be(n)}$$

is quite a good choice. The constants, and the epoch size, will need to be chosen by experiment. As we build more complex collections of units, the need for a better process will become pressing; two options appear below.

Ensuring learning is proceeding: We need to keep track of what is going on inside the system as we train it. One way is to plot the loss as a function of the number of steps. These plots can be very informative (Figure ??). If the learning rate is small, the system will make very slow progress but may (eventually) end up in a good state. If the learning rate is large, the system will make fast progress initially, but will then stop improving, because the state will change too quickly to find a good solution. If the learning rate is very large, the system might even diverge. If the learning rate is just right, you should get fast descent to a good value, and then slow but fairly steady improvement. Of course, just as in the case of SVMs, the plot of loss against step isn't a smooth curve, but rather noisy. There is an amusing collection of examples of training problems at lossfunctions.tumblr.com. It is quite usual to plot the error rate, or the accuracy, on a validation dataset

while training. This will allow you to compare the training error with the validation error. If these are very different, you have a problem: the system is overfitting, or not generalizing well. You should increase the regularization constant.

Dead units: Imagine the system gets into a state where for some unit u , $o_u = 0$ for every training data item. This could happen, for example, if the learning rate was too large. Then it can't get out of this state, because the gradient for that unit will be zero for every training data item, too. Such units are referred to as **dead units**. This problem can be contained by keeping the learning rate small enough. In more complex architectures (below), it is also contained by having a large number of units.

Gradient problems: There are a variety of important ways to have gradient problems. By far the most important is making a simple error in code (i.e. you compute the Jacobian elements wrong). This is surprisingly common; everybody does it at least once; and one learns to check gradients. Checking is fairly straightforward. You compute a numerical derivative, and compare that to the exact derivative. If they're too different, you have a gradient problem you need to fix. We will see a second important gradient problem when we see more complex architectures.

Choosing the regularization constant: This follows the recipe we saw for a linear SVM. Hold out a validation dataset. Train for several different values of λ . Evaluate each system on the validation dataset, and choose the best. Notice this involves many rounds of training, which could make things slow.

Does it work? Evaluating the classifier we have described is like evaluating any other classifier. You evaluate the error on a held-out data set that *wasn't* used to choose the regularization constant, or during training.

10.2 LAYERS AND NETWORKS

We have built a multiclass classifier out of units by using one unit per class, then interpreting the outputs of the units as probabilities using a softmax function. This classifier is at best only mildly interesting. The way to get something really interesting is to ask what the features for this classifier should be. To date, we have not looked closely at features. Instead, we've assumed they "come with the dataset" or should be constructed from domain knowledge. Remember that, in the case of regression, we could improve predictions by forming non-linear functions of features. We can do better than that; we could *learn* what non-linear functions to apply, by using the output of one set of units to form the inputs of the next set.

We will focus on systems built by organizing the units into **layers**; these layers form a **neural network** (a term I dislike, for the reasons above, but use because everybody else does). There is an input layer, consisting of the units that receive feature inputs from outside the network. There is an output layer, consisting of units whose outputs are passed outside the network. These two might be the same, as they were in the previous section. The most interesting cases occur when they are not the same. There may be **hidden layers**, whose inputs come from other layers and whose outputs go to other layers. In our case, the layers are ordered, and outputs of a given layer act as inputs to the next layer only (as in Figure 10.2 - we don't allow connections to wander all over the network). For the moment, assume that each unit in a layer receives an input from every unit in the previous layer;

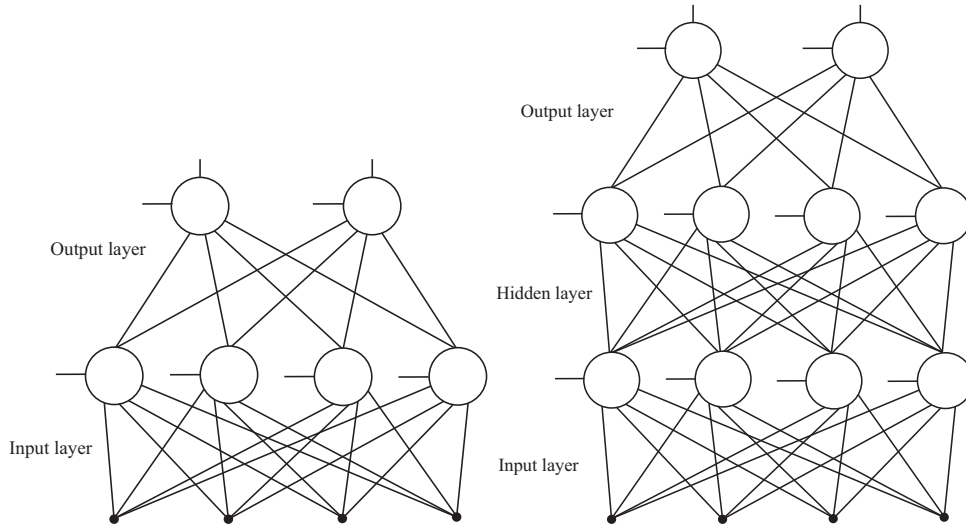


FIGURE 10.2: On the **left**, an input layer connected to an output layer. The units in the input layer take the inputs and compute features; the output layer turns these features into output values that will be turned into class probabilities with the softmax function. On the **right**, there is a hidden layer between input and output layer. This architecture means that the features seen by the output layer can be trained to be a significantly more complex function of the inputs.

this means that our network is **fully connected**. Other architectures are possible, but right now the most important question is how to train the resulting object.

10.2.1 Notation

Inevitably, we need yet more notation. There will be L layers. The input layer is layer 1, and the output layer is L . I will write u_i^l for the i 'th unit in the l 'th layer. This unit has output o_i^l and parameters \mathbf{w}_i^l and b_i^l , which I will stack into a vector θ_i^l . I write θ^l to refer to all the parameters of layer l . If I do not need to identify the layer in which a unit sits (for example, if I am summing over all units) I will drop the superscript. The vector of inputs to this unit is \mathbf{x}_i^l . These inputs are formed by choosing from the outputs of layer $l - 1$. I will write \mathbf{o}^l for all the outputs of the l 'th layer, stacked into a vector. I will represent the connections by a matrix \mathcal{C}_i^l , so that $\mathbf{x}_i^l = \mathcal{C}_i^l \mathbf{o}^{l-1}$. The matrix \mathcal{C}_i^l contains only 1 or 0 entries, and in the case of fully connected layers, it is the identity. Notice that every unit has its own \mathcal{C}_i^l .

I will write $L(\mathbf{y}_i, \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta)))$ for the loss of classifying the i 'th example using softmax. We will continue to use

$$L(\mathbf{y}_i, \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))) = -\mathbf{y}_i^T \log \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))$$

but in other applications, other losses might arise.

Generally, we will train by mini batch gradient descent, though I will describe some tricks that can speed up training and improve results. But we must compute

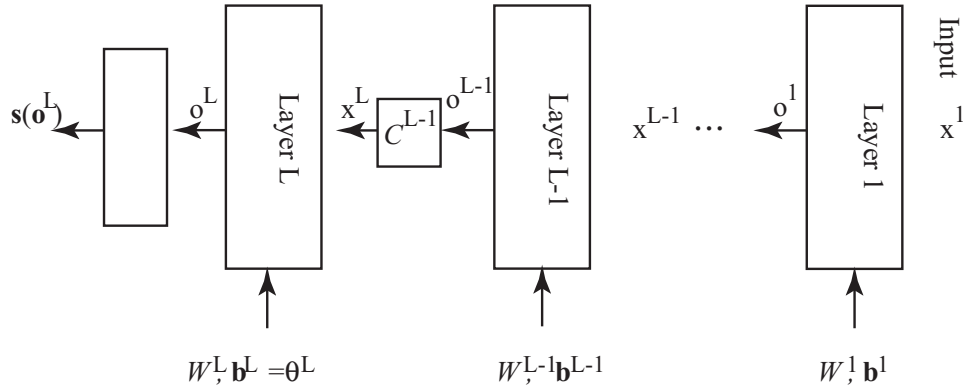


FIGURE 10.3: Notation for layers, inputs, etc.

the gradient. The output layer of our network has C units, one per class. We will apply the softmax to these outputs, as before. Writing E for the cost of error on training examples and R for the regularization term, we can write the cost of using the network as

$$\text{cost} = E + R = (1/N) \sum_{i \in \text{examples}} L(\mathbf{y}_i, \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))) + \frac{\lambda}{2} \sum_{k \in \text{units}} \mathbf{w}_k^T \mathbf{w}_k.$$

You should not let this compactified notation let you lose track of the fact that \mathbf{o}^L depends on \mathbf{x}_i through $\mathbf{o}^{L-1}, \dots, \mathbf{o}^1$. What we really should write is

$$\mathbf{o}^L(\mathbf{o}^{L-1}(\dots(\mathbf{o}^1(\mathbf{x}, \theta^1), \theta^2), \dots), \theta^L).$$

Equivalently, we could stack all the \mathcal{C}_i^l into one linear operator \mathcal{C}^l and write

$$\begin{aligned} \mathbf{o}^L(\mathbf{x}^L, \theta^L) & \quad \text{where} \\ \mathbf{x}^L & = \mathcal{C}^L \mathbf{o}^{L-1}(\mathbf{x}^{L-1}, \theta^{L-1}) \\ \dots & = \dots \\ \mathbf{x}^2 & = \mathcal{C}^2 \mathbf{o}^1(\mathbf{x}^1, \theta^1) \\ \mathbf{x}^1 & = \mathcal{C}^1 \mathbf{x} \end{aligned}$$

This is important, because it allows us to write an expression for the gradient.

10.2.2 Training, Gradients and Backpropagation

Now consider $\nabla_{\theta} E$. We have that E is a sum over examples. The gradient of the loss at a particular example is of most interest, because we will usually train with minibatches. So we are interested in

$$\nabla_{\theta} E_i = \nabla_{\theta} L(\mathbf{y}_i, \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))) = \nabla_{\theta} [-\mathbf{y}_i^T \log \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))]$$

and we can extend our use of the chain rule from section 10.1.3, very aggressively. We have

$$\nabla_{\theta}^L L(\mathbf{y}_i, \mathbf{s}(\mathbf{o}^L(\mathbf{x}_i, \theta))) = -\mathbf{y}_i^T \mathcal{J}_{\log \mathbf{s}; \mathbf{o}^L} \mathcal{J}_{\mathbf{o}^L; \theta^L}$$

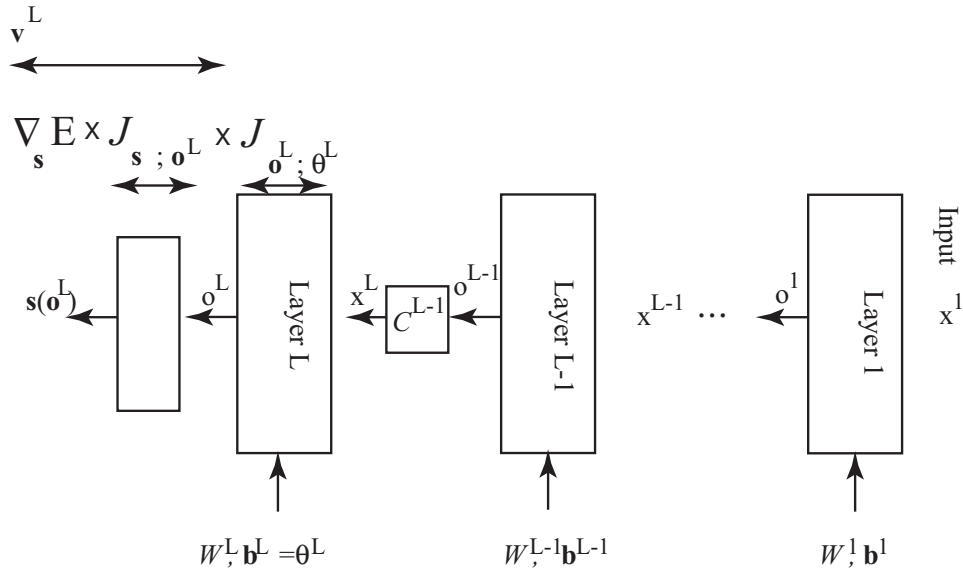


FIGURE 10.4: Constructing the gradient with respect to θ^L .

as in that section. Differentiating with respect to θ^{L-1} is more interesting. Layer L depends on θ^{L-1} in a somewhat roundabout way; layer $L - 1$ uses θ^{L-1} to produce its outputs, and these are fed into layer L as its inputs. So we must have

$$\nabla_{\theta}^{L-1} E_i = -\mathbf{y}_i^T \mathcal{J}_{\log s; \mathbf{o}^L} J_{\mathbf{o}^L; \mathbf{x}^L} J_{\mathbf{x}^L; \theta^{L-1}}$$

(look carefully at the subscripts on the Jacobians). These Jacobians have about the same form as those in section 10.1.3 if you recall that $\mathbf{x}^L = \mathcal{C}^L \mathbf{o}^{L-1}$. In turn, this means that

$$J_{\mathbf{x}^L; \theta^{L-1}} = \mathcal{C}^L J_{\mathbf{o}^{L-1}; \theta^{L-1}}$$

and the form of that Jacobian appears in section 10.1.3. But \mathbf{o}^L depends on θ^{L-2} through its inputs (which are \mathbf{x}^{L-1}), so that

$$\nabla_{\theta}^{L-2} E_i = -\mathbf{y}_i^T \mathcal{J}_{\log s; \mathbf{o}^L} J_{\mathbf{o}^L; \mathbf{x}^L} J_{\mathbf{o}^{L-1}; \mathbf{x}^{L-1}} J_{\mathbf{x}^{L-1}; \theta^{L-2}}$$

(again, look carefully at the subscripts on each of the Jacobians).

We can now get to the point. We have a recursion, which can be made more

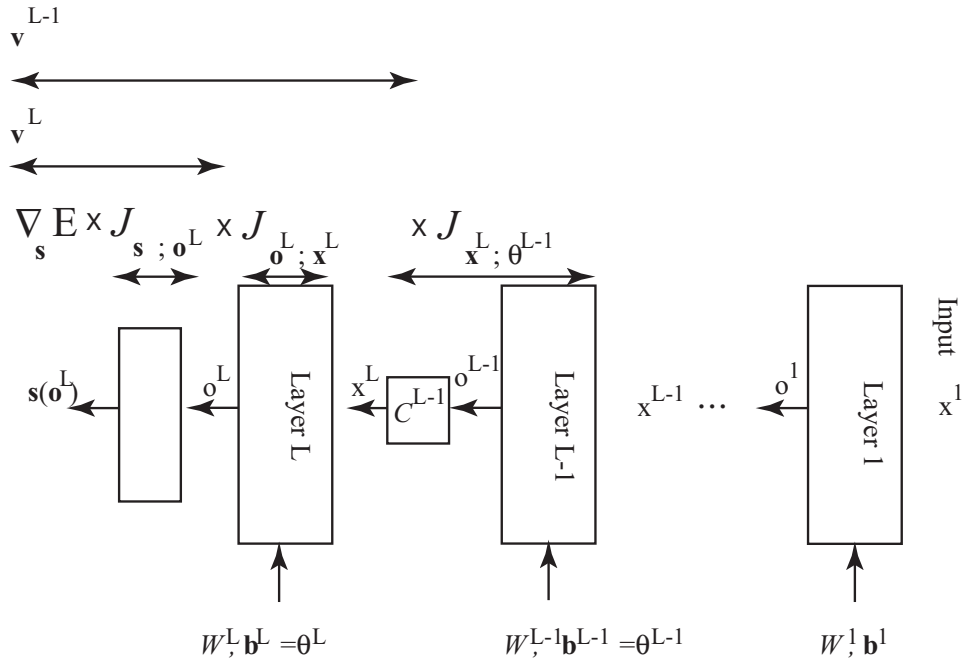


FIGURE 10.5: Constructing the gradient with respect to θ^{L-1} .

obvious with some notation. We have

$$\begin{aligned}
 \mathbf{v}^L &= (\nabla_{\mathbf{o}}^L E_i) \\
 \nabla_{\theta^L} E_i &= \mathbf{v}^L \mathcal{J}_{\mathbf{o}^L; \theta^L} \\
 \mathbf{v}^{L-1} &= \mathbf{v}^L \mathcal{J}_{\mathbf{o}^L; \mathbf{x}^L} \\
 \nabla_{\theta^{L-1}} E &= \mathbf{v}^{L-1} \mathcal{J}_{\mathbf{x}^L; \theta^{L-1}} \\
 &\dots \\
 \mathbf{v}^{i-1} &= \mathbf{v}^i \mathcal{J}_{\mathbf{x}^{i+1}; \mathbf{x}^i} \\
 \nabla_{\theta^{i-1}} E &= \mathbf{v}^{i-1} \mathcal{J}_{\mathbf{x}^i; \theta^{i-1}} \\
 &\dots
 \end{aligned}$$

I have not added notation to keep track of the point at which the partial derivative is evaluated (it should be obvious, and we have quite enough notation already). When you look at this recursion, you should see that, to evaluate \mathbf{v}^{i-1} , you will need to know \mathbf{x}^k for $k \geq i - 1$. This suggests the following strategy. We compute the \mathbf{x} 's (and, equivalently, \mathbf{o} 's) with a “forward pass”, moving from the input layer to the output layer. Then, in a “backward pass” from the output to the input, we compute the gradient. Doing this is often referred to as **backpropagation**.

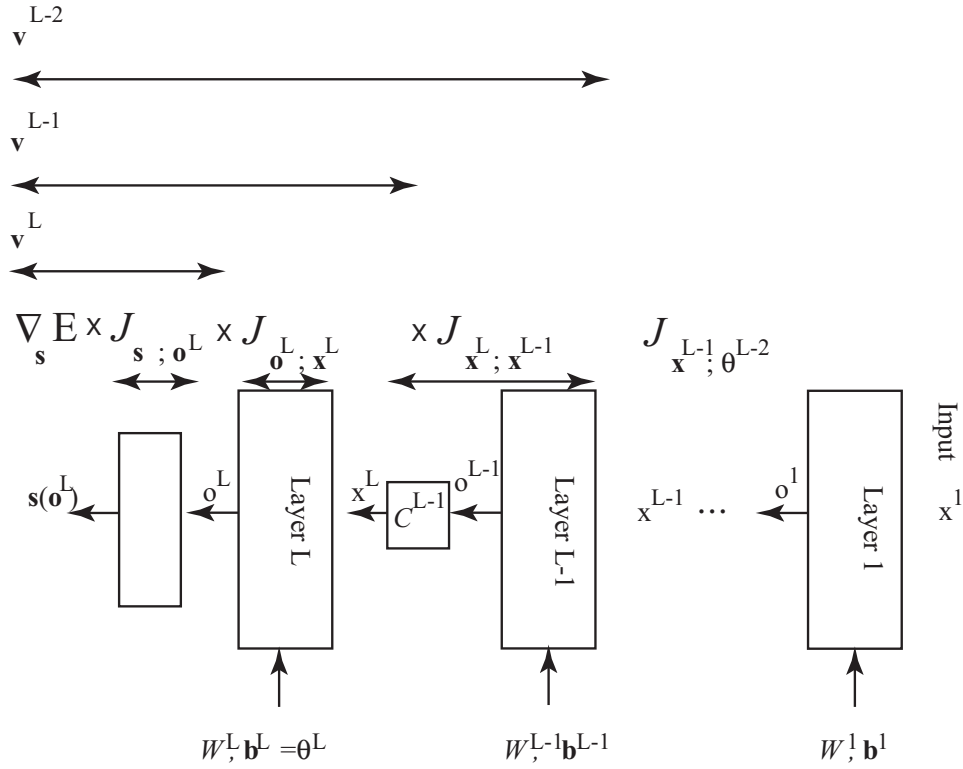


FIGURE 10.6: Constructing the gradient with respect to θ^{L-2} .

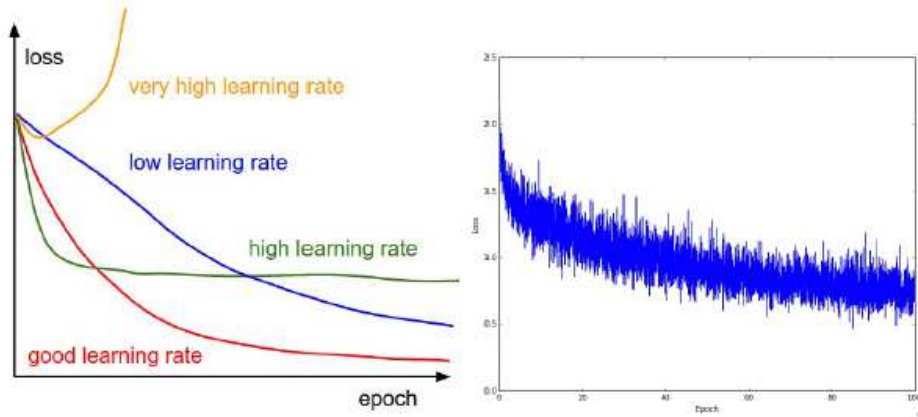


FIGURE 10.7: On the left, a cartoon of the error rate encountered during training a multilayer network and the main phenomena you might observe. On the right, an actual example. Each of these figures is from the excellent course notes for the Stanford class cs231n Convolutional Neural Networks for Visual Recognition, written by Andrej Karpathy. You can find these notes at: <http://cs231n.stanford.edu>.

10.2.3 Training Multiple Layers

A multilayer network represents an extremely complex, highly non-linear function, with an immense number of parameters. Training such networks is not easy. Neural networks are quite an old idea, but have only relatively recently had impact in practical applications. Hindsight suggests the problem is that networks are hard to train successfully. There is now a collection of quite successful tricks — I'll try to describe the most important — but the situation is still not completely clear.

The simplest training strategy is minibatch gradient descent. At round r , we have the set of weights $\theta^{(r)}$. We form the gradient for a minibatch $\nabla_{\theta}E$, and update the weights by taking a small step $\eta^{(r)}$ (usually referred to as the **learning rate**) backwards along the gradient, yielding

$$\theta^{(r+1)} = \theta^{(r)} - \eta^{(r)}\nabla_{\theta}E.$$

The most immediate difficulties are where to start, and what is $\eta^{(r)}$.

Initialization: As for a single layer of units, it is a bad idea to initialize each parameter to zero. It is usual to draw a sample of a zero mean normal random variable for each initial value. However, in a multilayer network, we may well have some units receiving input from more (or fewer) units than others (this is referred to as the **fan in** of the unit). Now assume that we have two units: one with many inputs, and one with few. If we initialize each units weights using the same zero mean normal random variable, the unit with more inputs will have a higher variance output (I'm ignoring the nonlinearity). This tends to lead to problems, because units at the next level will see unbalanced inputs. Experiment has shown that it is a good idea to allow the variance of the random variable you sample to depend on the fan in of the unit whose parameter you are initializing. Write n for the fan in of the unit in question, and ϵ for a small non-negative number. Current best practice appears to be that one initializes each weight with an independent sample of a random variable with mean 0 and *variance*

$$\epsilon \frac{\sqrt{2}}{n}.$$

Choosing ϵ too small or too big can lead to trouble, but I'm not aware of any recipe for coming up with a good choice. Typically, biases are initialized either to 0, or to a small non-negative number; there is mild evidence that 0 is a better choice.

Learning rate: The remarks above about learning rate apply, but for more complicated networks it is usual to apply one of the methods of section 10.2.4, which adjust the gradient to get better optimization behavior.

Ensuring learning is proceeding: We need to keep track of what is going on inside the system as we train it. One way is to plot the loss as a function of the number of steps. These plots can be very informative (Figure 10.7). If the learning rate is small, the system will make very slow progress but may (eventually) end up in a good state. If the learning rate is large, the system will make fast progress initially, but will then stop improving, because the state will change too quickly to find a good solution. If the learning rate is very large, the system might even diverge. If the learning rate is just right, you should get fast descent to a good value, and then slow but fairly steady improvement. Of course, just as in the case

of SVMs, the plot of loss against step isn't a smooth curve, but rather noisy. There is an amusing collection of examples of training problems at lossfunctions.tumblr.com. It is quite usual to plot the error rate, or the accuracy, on a validation dataset while training. This will allow you to compare the training error with the validation error. If these are very different, you have a problem: the system is overfitting, or not generalizing well. You should increase the regularization constant.

Dead units: The remarks above apply.

Gradient problems: The remarks above apply.

Choosing the regularization constant: The remarks above apply, but for more complex networks, it is usual to use the more sophisticated regularization described in section ?? (at considerable training cost).

Does it work? Evaluating the classifier we have described is like evaluating any other classifier. You evaluate the error on a held-out data set that *wasn't* used to choose the regularization constant, or during training.

10.2.4 Gradient Scaling Tricks

Everyone is surprised the first time they learn that the best direction to travel in when you want to minimize a function is not, in fact, backwards down the gradient. The gradient *is* uphill, but repeated downhill steps are often not particularly efficient. An example can help, and we will look at this point several ways because different people have different ways of understanding this point.

We can look at the problem with algebra. Consider $f(x, y) = (1/2)(\epsilon x^2 + y^2)$, where ϵ is a small positive number. The gradient at (x, y) is $[\epsilon x, y]$. For simplicity, use a fixed learning rate η , so we have $[x^{(r)}, y^{(r)}] = [(1 - \epsilon\eta)x^{(r-1)}, (1 - \eta)y^{(r-1)}]$. If you start at, say, $(x^{(0)}, y^{(0)})$ and repeatedly go downhill along the gradient, you will travel very slowly to your destination. You can show that $[x^{(r)}, y^{(r)}] = [(1 - \epsilon\eta)^r x^{(0)}, (1 - \eta)^r y^{(0)}]$. The problem is that the gradient in y is quite large (so y must change quickly) and the gradient in x is small (so x changes slowly). In turn, for steps in y to converge we must have $|1 - \eta| < 1$; but for steps in x to converge, we require only the much weaker constraint $|1 - \epsilon\eta| < 1$. Imagine we choose the largest η we dare for the y constraint. The y value will very quickly have small magnitude, though its sign will change with each step. But the x steps will move you closer to the right spot only extremely slowly.

Another way to see this problem is to reason geometrically. Figure 10.8 shows this effect for this function. The gradient is at right angles to the level curves of the function. But when the level curves form a narrow valley, the gradient points across the valley rather than down it. The effect isn't changed by rotating and translating the function (Figure 10.9).

You may have learned that Newton's method resolves this problem. This is all very well, but to apply Newton's method we would need to know the matrix of second partial derivatives. A network can easily have thousands to millions of parameters, and we simply can't form, store, or work with matrices of these dimensions. Instead, we will need to think more qualitatively about what is causing trouble.

One useful insight into the problem is that fast changes in the gradient vector are worrying. For example, consider $f(x) = (1/2)(x^2 + y^2)$. Imagine you start

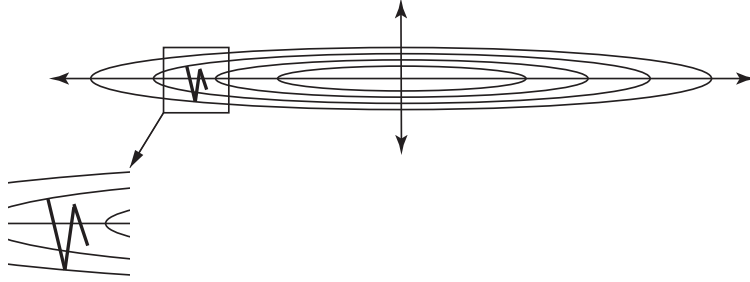


FIGURE 10.8: A plot of the level curves (curves of constant value) of the function $f(x, y) = (1/2)(\epsilon x^2 + y^2)$. Notice that the value changes slowly with large changes in x , and quickly with small changes in y . The gradient points mostly toward the x -axis; this means that gradient descent is a slow zig-zag across the “valley” of the function, as illustrated. We might be able to fix this problem by changing coordinates, if we knew what change of coordinates to use.

far away from the origin. The gradient won’t change much along reasonably sized steps. But now imagine yourself on one side of a valley like the function $f(x) = (1/2)(x^2 + \epsilon y^2)$; as you move along the gradient, the gradient in the x direction gets smaller very quickly, then points back in the direction you came from. You are not justified in taking a large step in this direction, because if you do you will end up at a point with a very different gradient. Similarly, the gradient in the y direction is small, and stays small for quite large changes in y value. You would like to take a small step in the x direction and a large step in the y direction.

You can see that this is the impact of the second derivative of the function (which is what Newton’s method is all about). But we can’t do Newton’s method. We would like to travel further in directions where the gradient doesn’t change much, and less far in directions where it changes a lot. There are several methods for doing so.

Momentum: We should like to discourage parameters from “zig-zagging” as in the example above. In these examples, the problem is caused by components of the gradient changing sign from step to step. It is natural to try and smooth the

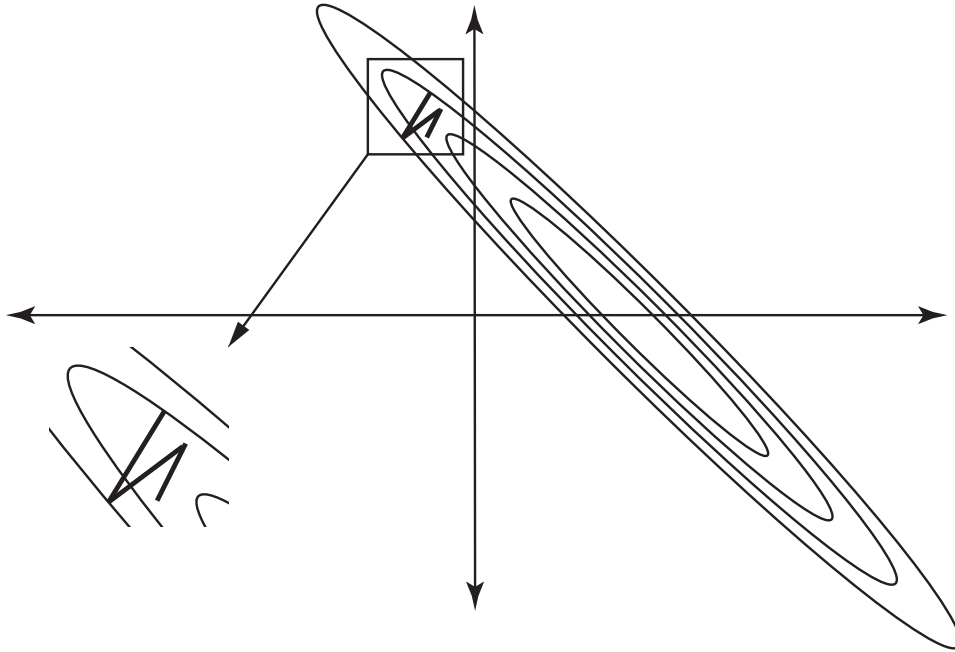


FIGURE 10.9: *Rotating and translating a function rotates and translates the gradient; this is a picture of the function of figure 10.8, but now rotated and translated. The problem of zig-zagging remains. This is important, because it means that we may have serious difficulty choosing a good change of coordinates.*

gradient. We could do so by forming a moving average of the gradient. Construct a vector \mathbf{v} , the same size as the gradient, and initialize this to zero. Choose a positive number $\mu < 1$. Then we iterate

$$\begin{aligned}\mathbf{v}^{(r+1)} &= \mu * \mathbf{v}^{(r)} + \eta \nabla_{\theta} E \\ \theta^{(r+1)} &= \theta^{(r)} - \mathbf{v}^{(r+1)}\end{aligned}$$

Notice that, in this case, the update is an average of all past gradients, each weighted by a power of μ . If μ is small, then only relatively recent gradients will participate in the average, and there will be less smoothing. Larger μ lead to more smoothing. A typical value is $\mu = 0.9$. It is reasonable to make the learning rate go down with epoch when you use momentum, but keep in mind that a very large μ will mean you need to take several steps before the effect of a change in learning rate shows.

Adagrad: We will keep track of the size of each component of the gradient. In particular, we have a running cache \mathbf{c} which is initialized at zero. We choose a small number α (typically 1e-6), and a fixed η . Write $g_i^{(r)}$ for the i 'th component

of the gradient $\nabla_{\theta} E$ computed at the r 'th iteration. Then we iterate

$$\begin{aligned}c_i^{(r+1)} &= c_i^{(r)} + (g_i^{(r)})^2 \\ \theta_i^{(r+1)} &= \theta_i^{(r)} - \eta \frac{g_i^{(r)}}{(c_i^{(r+1)})^{\frac{1}{2}} + \alpha}\end{aligned}$$

Notice that each component of the gradient has its own learning rate, set by the history of previous gradients.

RMSprop: This is a modification of Adagrad, to allow it to “forget” large gradients that occurred far in the past. We choose another number, Δ , (the **decay rate**; typical values might be 0.9, 0.99 or 0.999), and iterate

$$\begin{aligned}c_i^{(r+1)} &= \Delta c_i^{(r)} + (1 - \Delta)(g_i^{(r)})^2 \\ \theta_i^{(r+1)} &= \theta_i^{(r)} - \eta \frac{g_i^{(r)}}{(c_i^{(r+1)})^{\frac{1}{2}} + \alpha}\end{aligned}$$

10.2.5 Dropout

Regularizing by the square of the weights is all very well, but quite quickly we will have problems because there are so many weights. An alternative, and very useful, regularization strategy is to try and ensure that no unit relies too much on the output of any other unit. One can do this as follows. At each training step, randomly select some units, set their outputs to zero (and reweight the inputs of the units receiving input from them), and then take the step. Now units are trained to produce reasonable outputs even if some of their inputs are randomly set to zero — units can't rely too much on one input, because it might be turned off. Notice that this sounds sensible, but it isn't quite a proof that the approach is sound; that comes from experiment. The approach is known as **dropout**.

There are some important details we can't go into. Output units are not subject to dropout, but one can also turn off inputs randomly. At test time, there is no dropout. Every unit computes its usual output in the usual way. This creates an important training issue. Write p for the probability that a unit is dropped out, which will be the same for all units subject to dropout. You should think of the expected output of the i 'th unit at *training* time as $(1 - p)o_i$ (because with probability p , it is zero). But at test time, the next unit will see o_i ; so at training time, you should reweight the inputs by $1/(1 - p)$. In exercises, we will use packages that arrange all the details for us.

10.2.6 It's Still Difficult..

All the tricks above are helpful, but training a multilayer neural network is still difficult. Fully connected layers have many parameters. It's quite natural to take an input feature vector of moderate dimension, build one layer that produces a much higher dimensional vector, then stack a series of quite high dimensional layers on top of that. There is quite good evidence that having many layers can improve practical performance *if* one can train the resulting network. Such an architecture has been known for a long time, but hasn't been particularly successful until recently.

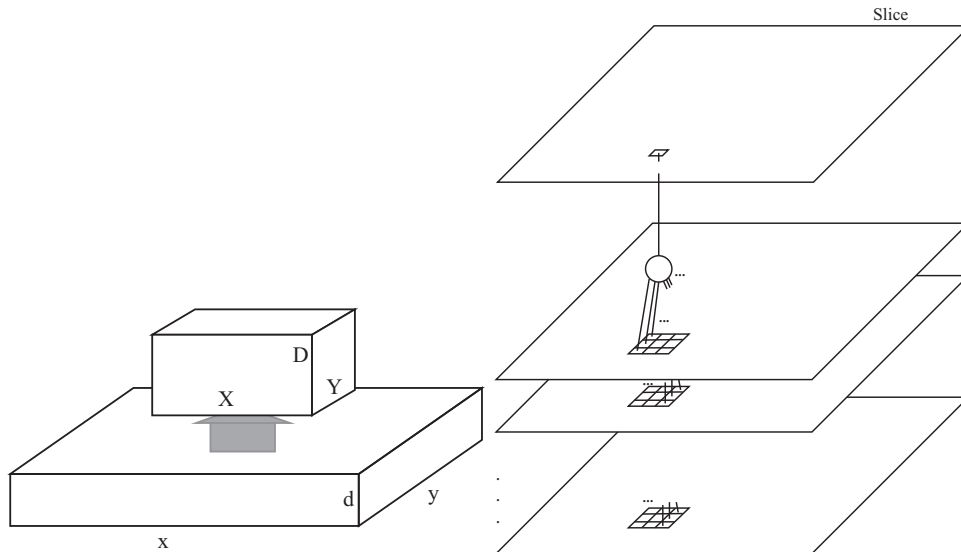


FIGURE 10.10: Terminology for building a convolutional layer. On the **left**, a layer turns one block of data into another. The x and y coordinates index the position of a location in the block, and the d coordinate identifies the data item at that point. A natural example of a block is a color image, which usually has three layers ($d=3$); then x and y identify the pixel, and d chooses the R, G, or B slice. The dimensions x and y are the spatial dimensions. On the **right**, a unit in a convolutional layer forms a weighted sum of set of locations, adds a bias, then applies a RELU. There is one unit for each (X, Y, D) location in the output block. For each (X, Y) in the output block, there is a corresponding window of size (w_x, w_y) in the (x, y) space. Each of the D units whose responses form the D values at the (X, Y) location in the output block forms a weighted sum of all the values covered by that window. These inputs come from each slice in the window below that unit (so the number of inputs is $d \times w_x \times w_y$). Each of the units that feed the a particular slice in the output block has the same set of weights, so you should think of a unit as a form of pattern detector; it will respond strongly if the block below it is “similar” to the weights.

There are several structural obstacles. Without GPU’s, evaluating such a network can be slow, making training slow. The number of parameters in just one fully connected layer is high, meaning that multiple layers will need a lot of data to train, and will take many training batches. There is some reason to believe that multilayer neural networks were discounted in application areas for quite a long time because people underestimated just how much data and how much training was required to make them perform well.

One obstacle that remains technically important has to do with the gradient. Look at the recursion I described for backpropagation. The gradient update at the L ’th (top) layer depends pretty directly on the parameters in that layer. But now consider a layer close to the input end of the network. The gradient update has

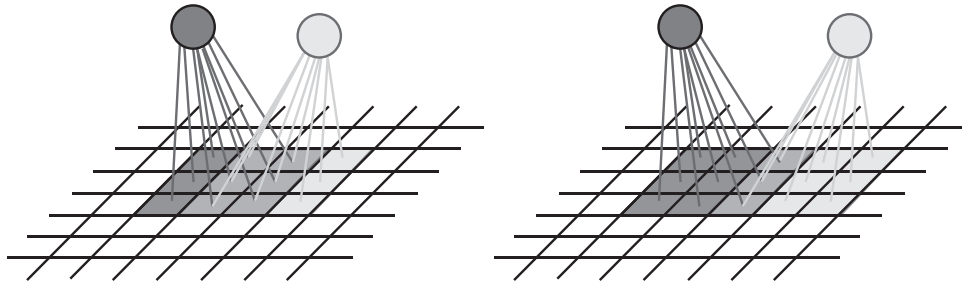


FIGURE 10.11: *Figure 10.10 shows units taking an input block and creating an output block. Each unit is fed by a window on the spatial dimensions of the input block. The window is advanced by the stride to feed the next unit. On the **left**, two units fed by 3×3 windows with a stride of 1. I have shaded the pixels in the window to show which pixels go to which unit; notice the units share 6 pixels. In this case, the spatial dimensions of the output block will be either the same as those of the input block (if we find some values to feed pixels in windows that hang over the edge of the image), or only slightly smaller (if we ignore units whose windows hang over the edges of the image). On the **right**, two units fed by 3×3 windows with a stride of 2. Notice the units share fewer pixels, and the output block will be smaller than the input block.*

been multiplied by several Jacobian matrices. The update may be very small (if these Jacobians shrink their input vectors) or unhelpful (if layers close to the output have poor parameter estimates). For the gradient update to be really helpful, we'd like the layers higher up the network to be right; but we can't achieve this with lower layers that are confused, because they pass their outputs up. If a layer low in the network is in a nonsensical state, it may be very hard to get it out of that state. In turn, this means that adding layers to a network might improve performance, but also might make it worse because the training turns out poorly.

There are a variety of strategies for dealing with this problem. We might just train for a very long time, possibly using gradient rescaling tricks. We might reduce the number of parameters in the layers, by passing to convolutional layers (below) rather than fully connected layers. We might use various tricks to initialize each layer with a good estimate. This is a topic of widespread current interest, but one I can't deal with in any detail here. Finally, we might use architectural tricks (section 16.20) to allow inputs to bypass layers, so that poorly trained layers create fewer difficulties.

10.3 CONVOLUTIONAL NEURAL NETWORKS

One area where neural networks have had tremendous impact is in image understanding. Images have special properties that motivate special constructions of features. These constructions yield a layer architecture that has a significantly reduced number of parameters in the layer. This architecture has proven useful in other applications, but we will work with images.

10.3.1 Images and Convolutional Layers

Using the output of one layer to form features for another layer is attractive. The natural consequence of this idea is that the input to the network would be the image pixels. But this presents some important difficulties. There are an awful lot of pixels in most images, and that means it's likely that there will be an awful lot of parameters.

For a variety of reasons, it doesn't make much sense to have an input layer consisting of units each of which sees all the image pixels. There would be a tremendous number of weights to train. It would be hard to explain why units have different values for the weights. Instead, we can use some intuitions from computer vision.

First, we need to build systems that can handle different versions of what is essentially the same image. For example, imagine you turn the camera slightly to one side, or raise it or lower it when taking the image. If every input unit sees every pixel, this would constrain the form of the weights. Turning the camera up a bit shifts the image down a bit; the representation shouldn't be severely disrupted by this. Second, long experience in computer vision has produced a (very rough) recipe for building image features: you construct features that respond to patterns in small, localized neighborhoods; then other features look at patterns of *those* features; then others look at patterns of those, and so on (*big fleas have little fleas upon their backs to bite 'em; and little fleas have smaller ones, and so ad infinitum*).

We will assume that images are 3D. The first two dimensions will be the x and y dimensions in the image, the third (for the moment!) will identify the color layer of the image (for example, R , G and B). We will build layers that take 3D objects like images (which I will call **blocks**) and make new blocks (Figure 10.10; notice the input block has dimension $x \times y \times d$ and the output block has dimension $X \times Y \times D$). Each block is a stack of **slices**, which — like color layers in an image — have two spatial dimensions.

These layers will draw from a standard recipe for building an image feature that describes a small neighborhood. We construct a **convolution kernel**, which is a small block. This is typically odd sized, and typically from 3×3 to a few tens by a few tens in size along the spatial dimensions, and is always of size d in the other dimension.

Write \mathcal{I} for a block (for example, an image), \mathcal{K} for the kernel, b for a bias term (which might be zero; some, but not all, convolutional layers use a bias term), and \mathcal{I}_{ijk} for the i , j 'th pixel in the k 'th slice of the block. Write F for the function implemented by a RELU, so that $F(x) = \max(0, x)$. Now we form a slice \mathcal{O} , whose u , v 'th entry is

$$\mathcal{O}_{uv} = F(W_{uv} + b) = F\left(\sum_{ijk} \mathcal{I}_{u+i, v+j, k} \mathcal{K}_{ijk} + b\right),$$

where I am assuming the sum goes over all values of i and j , and if the indices to either \mathcal{I} or \mathcal{K} go outside the domain, then the reported value is zero. There is room for some confusion here, because one can use a variety of different indexing schemes, and different authors use different ones (usually for compatibility with the history of convolution); this is of no significance. Figure 10.10 illustrates the

process that produces a slice from a block.

The operation that produces \mathcal{W} from \mathcal{I} and \mathcal{K} is known as **convolution**, and it is usual to write $\mathcal{W} = \mathcal{K} * \mathcal{I}$. We will not go into all the properties of convolution, but you should notice one extremely important property. We obtain the value at a pixel by centering \mathcal{K} on that pixel. We now have a patch sitting over all the layers of the image at some location; we multiply the pixels in that patch (by layer) by the corresponding image pixels (in layers), then accumulate the products — the result goes into \mathcal{O} . This is like a dot-product — we will get a large positive value in \mathcal{O} at that pixel if the image window around that pixel looks like \mathcal{K} , and a small negative value if they're the same up to a sign change. You should think of a convolution kernel as being an example pattern.

Now when you convolve a kernel with an image, you get, at each location, an estimate of how much that image looks like that kernel at that point. The output is the response of a collection of simple pattern detectors, one at each pixel, *for the same pattern*. We may not need every such output; instead, we might look at every second (third, etc.) pixel in each direction. This choice is known as the **stride**. A stride of 1 corresponds to looking at every pixel; of 2, every second pixel; and so on (Figure 10.11)

A slice can be interpreted as a map, giving the response of a local feature detector at every (resp. every second; every third; etc.) pixel. At each pixel of interest (i.e. every pixel; every second, etc. depending on stride), we place a window (which should be odd-sized, to make indexing easier). Every pixel in that window is an input to a unit that corresponds to the window, which multiplies each pixel by a weight, sums all these terms, then applies a RELU. What makes a slice special is that *each unit uses the same set of weights*. The size of this object depends a little on the software package you are using. Assume the input image is of size $n_x \times n_y \times n_z$, and the kernel is of size $2k_x + 1 \times 2k_y + 1 \times n_z$. At least in principle, you cannot place a unit over a pixel that is too close to the edge, because then some of its inputs are outside the image. You could pad the image (either with constants, or by reflecting it, or by attaching copies of the columns/rows at the edge) and supply these inputs; in this case, the output could be $n_x \times n_y \times n_z$. Otherwise, you could place units only over pixels where all of the unit's inputs are inside the image. Then you would have an output of size $n_x - 2k_x \times n_y - 2k_y \times n_z$. The kernel is usually small, so the difference in sizes isn't that great. Most software packages are willing to set up either case.

A slice finds locations in the image where a particular pattern (identified by the weights) occurs. We could attach many slices to the image. They should all have the same stride, so they're all the same size. The output of this collection of slices would be one vector at each pixel location, where the components of the vector represent the similarity between the image patch centered at that location and a particular pattern. A collection of slices is usually referred to as a **convolutional layer**. You should think of a convolutional layer as being like a color image. There are now many different color layers (the slices), so that dimension has been expanded.

10.3.2 Convolutional Layers upon Convolutional Layers

Now the output of the initial convolutional layer is a set of slices, registered to the input image, forming a block of data. That looks like the input of that layer (a set of slices — color layers) forming a **block** of data. This suggests we could use the output of the first convolutional layer could be connected to a second convolutional layer, a second to a third, and so on. Doing so turns out to be an excellent idea.

Think about the output of the first convolutional layer. Each location receives inputs from pixels in a window about that location. Now if we put a second layer on top of the first, each location in the second receives inputs from first layer values in a window about that location. This means that locations in the second layer are affected by a larger window of pixels than those in the first layer. You should think of these as representing “patterns of patterns”. If we place a third layer on top of the second layer, locations in that third layer will depend on an even larger window of pixels. A fourth layer will depend on a yet larger window, and so on.

10.3.3 Pooling

If you have several convolutional layers with stride 1, then each block of data has the same spatial dimensions. This tends to be a problem, because the pixels that feed a unit in the top layer will tend to have a large overlap with the pixels that feed the unit next to it. In turn, the values that the units take will be similar, and so there will be redundant information in the output block. It is usual to try and deal with this by making blocks get smaller. One natural strategy is to occasionally have a layer that has stride 2.

An alternative strategy is to use **max pooling**. A pooling unit reports the largest value of its inputs. In the most usual arrangement, a pooling layer will take an (x, y, d) block to a $(x/2, y/2, d)$ block. For the moment, ignore the entirely minor problems presented by a fractional dimension. The new block is obtained by pooling units that pool a 2x2 window at each slice of the input block to form each slice of the output block. These units are placed so they don’t overlap, so the output block is half the size of the input block (for some reason, this configuration is hard to say but easy to see; Figure 10.12). If x or y or both are odd, there are two options; one could ignore the odd pixel on the boundary, or one could build a row (column; both) of imputed values, most likely by copying the row (column; both) on the edge. These two strategies yield, respectively, $\text{floor}(x/2)$ and $\text{ceil}(x/2)$ for the new dimension. Pooling seems to be falling out of favor, but not so much or so fast that you will not encounter it.

10.4 BUILDING AN IMAGE CLASSIFIER

There are two problems that lie at the core of image understanding. The first is **image classification**, where we decide what class an image of a fixed size belongs to. The taxonomy of classes is provided in advance, but it’s usual to work with a collection of images of objects. These objects will be largely centered in the image, and largely isolated. Each image will have an associated object name. There are many collections with this structure. The best known, by far, is **ImageNet**, which can be found at <http://www.image-net.org>. There is a regular competition to

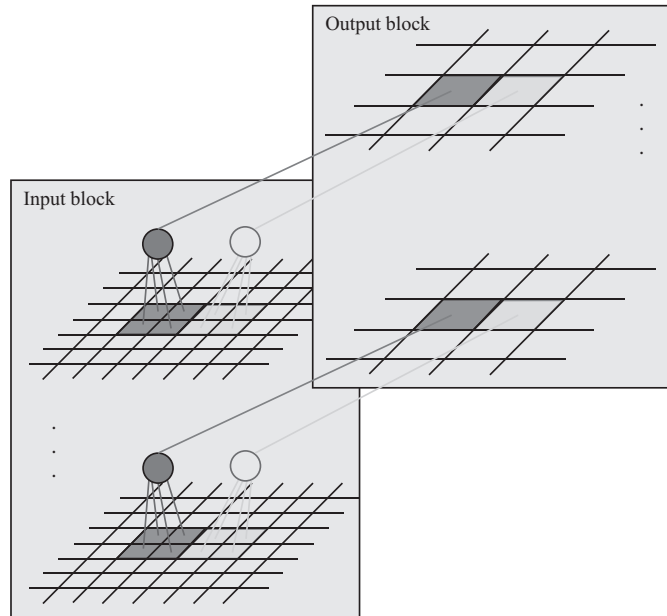


FIGURE 10.12: In a pooling layer, pooling units compute the largest value of their inputs, then pass it on. The most common case is 2×2 , illustrated here. We tile each slice with 2×2 windows that do not overlap. Pooling units compute the max, then pass that on to the corresponding location in the corresponding slice of the output block. As a result, the spatial dimensions of the output block will be about half those of the input block (details depend on how one handles windows that hang over the edge).

classify ImageNet images. Be aware that, while this chapter tries to give a concise description of best practice, you might need to do more than read it to do well in the competition.

The second problem is **object detection**, where we try to find the locations of objects of a set of classes in the image. So we might try to mark all cars, all cats, all camels, and so on. As far as anyone knows, the right way to think about object detection is that we search a collection of windows in an image, apply an image classification method to each window, then resolve disputes between overlapping windows. How windows are to be chosen for this purpose is an active and quickly changing area of research. We will regard image classification as the key building block, and ignore the question of deciding which window to classify.

We have most of the pieces to build an image classifier. Architectural choices will make a difference to its performance. So will a series of tricks.

10.4.1 An Image Classification Architecture

We can now put together an image classifier. A convolutional layer receives image pixel values as input. The output is fed to a stack of convolutional layers, each

feeding the next. The output of the final layer is fed to one or more fully connected layers, with one output per class. The whole is trained by batch gradient descent, or a variant, as above.

There are a number of architectural choices to make, which are typically made by experiment. The main ones are the choice of the number of convolutional layers; the choice of the number of slices in each layer; and the choice of stride for each convolutional layer. There are some constraints on the choice of stride. The first convolutional layer will tend to have stride 1, so that we see all the resolution of the image. But the outputs of that layer are likely somewhat correlated, because they depend on largely the same set of pixels. Later layers might have larger stride for this reason. In turn, this means the spatial dimensions of the representation will get smaller.

Notice that different image classification networks differ by relatively straightforward changes in architectural parameters. Mostly, the same thing will happen to these networks (variants of batch gradient descent on a variety of costs; dropout; evaluation). In turn, this means that we should use some form of specification language to put together a description of the architecture of interest. Ideally, in such an environment, we describe the network architecture, choose an optimization algorithm, and choose some parameters (dropout probability, etc.). Then the environment assembles the net, trains it (ideally, producing log files we can look at) and runs an evaluation. Several such environments exist.

10.4.2 Useful Tricks - 1: Preprocessing Data

It usually isn't possible to simply feed any image into the network. We want each image fed into the network to be the same size. We can achieve this either by resizing the image, or by cropping the image. Resizing might mean we stretch or squash some images, which likely isn't a great idea. Cropping means that we need to make a choice about where the crop box lies in the image. Practical systems quite often apply the same network to different croppings of the same image. For our purposes, we will assume that all the images we deal with have the same size.

It is usually wise to preprocess images before using them. This is because two images with quite similar content might have rather different pixel values. For example, compare image \mathcal{I} and $1.5\mathcal{I}$. One will be brighter than the other, but nothing substantial about the image class will have changed. There is little point in forcing the network to learn something that we know already. There are a variety of preprocessing options, and different options have proven to be best for different problems. I will sketch some of the more useful ones.

You could **whiten pixel values**. You would do this for each pixel in the image grid independently. For each pixel, compute the mean value at that pixel across the training dataset. Subtract this, and divide the result by the standard deviation of the value at that pixel across the training dataset. Each pixel location in the resulting stack of images has mean zero and standard deviation one. Reserve the offset image (the mean at each pixel location) and the scale image (ditto, standard deviation) so that you can normalize test images.

You could **contrast normalize the image** by computing the mean and standard deviation of pixel values in each training (resp. test) image, then subtracting

the mean from the image and dividing the result by the standard deviation.

You could **contrast normalize pixel values locally**. To do so, you compute a smoothed version of the image (convolve with a Gaussian, for insiders; everyone else should skip this paragraph, or perhaps search the internet). You can think of the result as a local estimate of the image mean. At each pixel, you subtract the smoothed value from the image value.

Useful Facts: 10.1 *Whitening a dataset*

For a dataset $\{\mathbf{x}\}$, compute:

- \mathcal{U} , the matrix of eigenvectors of $\text{Covmat}(\{\mathbf{x}\})$;
- and $\text{mean}(\{\mathbf{x}\})$.

Now compute $\{\mathbf{n}\}$ using the rule

$$\mathbf{n}_i = \mathcal{U}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})).$$

Then $\text{mean}(\{\mathbf{n}\}) = \mathbf{0}$ and $\text{Covmat}(\{\mathbf{n}\})$ is diagonal.

Now write Λ for the diagonal matrix of eigenvalues of $\text{Covmat}(\{\mathbf{x}\})$ (so that $\text{Covmat}(\{\mathbf{x}\})\mathcal{U} = \mathcal{U}\Lambda$). Assume that each of the diagonal entries of Λ is greater than zero (otherwise there is a redundant dimension in the data). Write λ_i for the i 'th diagonal entry of Λ , and write $\Lambda^{-(1/2)}$ for the diagonal matrix whose i 'th diagonal entry is $1/\sqrt{\lambda_i}$. Compute $\{\mathbf{z}\}$ using the rule

$$\mathbf{z}_i = \Lambda^{-(1/2)}\mathcal{U}(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})).$$

We have that $\text{mean}(\{\mathbf{z}\}) = \mathbf{0}$ and $\text{Covmat}(\{\mathbf{z}\}) = \mathcal{I}$. The dataset $\{\mathbf{z}\}$ is often known as **whitened data**.

You could whiten the image as in section 10.1. It turns out this doesn't usually help all that much. Instead, you need to use **ZCA-whitening**. I will use the same notation as chapter 16.20, but I reproduce the useful facts box here as a reminder. Notice that, by using the rule

$$\mathbf{z}_i = \Lambda^{-(1/2)}\mathcal{U}(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})),$$

we have rotated the data in the high dimensional space. In the case of images, this means that the image corresponding to \mathbf{z}_i will likely not look like anything coherent. Furthermore, if there are very small eigenvalues, the scaling represented by $\Lambda^{-(1/2)}$ may present serious problems. But notice that the covariance matrix of a dataset is unaffected by rotation. We could choose a small non-negative constant ϵ , and use the rule

$$\mathbf{z}_i = \mathcal{U}^T(\Lambda + \epsilon\mathcal{I})^{-(1/2)}\mathcal{U}(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))$$

instead. The result looks significantly more like an image, and will have a covariance matrix that is the identity (or close, depending on the value of ϵ). This rule is ZCA whitening.

10.4.3 Useful Tricks - 2: Enhancing Training Data

Datasets of images are never big enough to show all effects accurately. This is because an image of a horse is still an image of a horse even if it has been through a small rotation, or has been resized to be a bit bigger or smaller, or has been cropped differently, and so on. There is no way to take account of these effects in the architecture of the network. Generally, a better approach is to expand the training dataset to include different rotations, scalings, and crops of images.

Doing so is relatively straightforward. You take each training image, and generate a collection of extra training images from it. You can obtain this collection by: resizing and then cropping the training image; using different crops of the same training image (assuming that training images are a little bigger than the size of image you will work with); rotating the training image by a small amount, resizing and cropping; and so on. You can't crop too much, because you need to ensure that the modified images are still of the relevant class, and an aggressive crop might cut out the horse, etc. When you rotate then crop, you need to be sure that no "unknown" pixels find their way into the final crop. All this means that only relatively small rescales, crops, rotations, etc. will work. Even so, this approach is an extremely effective way to enlarge the training set.

10.4.4 Useful Tricks - 3: Batch Normalization

There is good experimental evidence that large values of inputs to any layer within a neural network lead to problems. One source of the problem could be this. Imagine some input to some unit has a large absolute value. If the corresponding weight is relatively small, then one gradient step could cause the weight to change sign. In turn, the output of the unit will swing from one side of the RELU's non-linearity to the other. If this happens for too many units, there will be training problems because the gradient is then a poor prediction of what will actually happen to the output. So we should like to ensure that relatively few values at the input of any layer have large absolute values. We will build a new layer, sometimes called a **batch normalization layer**, which can be inserted between two existing layers.

Write \mathbf{x}^b for the input of this layer, and \mathbf{o}^b for its output. The output has the same dimension as the input, and I shall write this dimension d . The layer has two vectors of parameters, γ and β , each of dimension d . Write $\text{diag}(\mathbf{v})$ for the matrix whose diagonal is \mathbf{v} , and with all other entries zero. Assume we know the mean (\mathbf{m}) and standard deviation (\mathbf{s}) of each component of \mathbf{x}^b , where the expectation is taken over all relevant data. The layer forms

$$\begin{aligned}\mathbf{x}^n &= [\text{diag}(\mathbf{s} + \epsilon)]^{-1} (\mathbf{x}^b - \mathbf{m}) \\ \mathbf{o}^b &= [\text{diag}(\gamma)] \mathbf{x}^n + \beta.\end{aligned}$$

Notice that the output of the layer is a differentiable function of γ and β . Notice also that this layer *could* implement the identity transform, if $\gamma = \text{diag}(\mathbf{s} + \epsilon)$ and

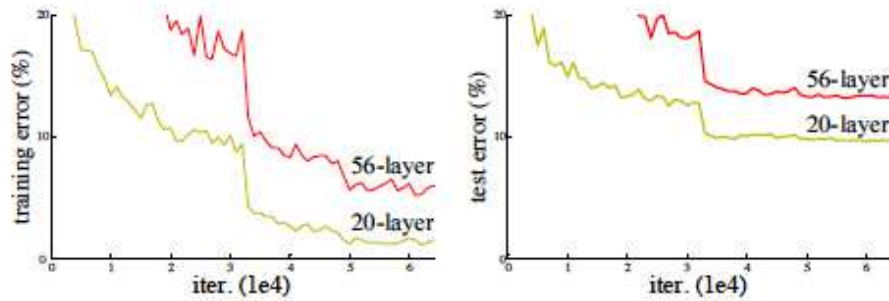


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

FIGURE 10.13: *This figure from Deep Residual Learning for Image Recognition Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun, ICCV 2016, illustrates the difficulties presented by training a deep network.*

$\beta = \mathbf{m}$. We adjust the parameters in training to achieve the best performance. It can be helpful to think about this layer as follows. The layer rescales its input to have zero mean and unit standard deviation, then allows training to readjust the mean and standard deviation as required. In essence, we expect that large values encountered between layers are likely an accident of the difficulty training a network, rather than required for good performance.

The difficulty here is we don’t know either \mathbf{m} or \mathbf{s} , because we don’t know the parameters used for previous layers. Current practice is as follows. First, start with $\mathbf{m} = \mathbf{0}$ and $\mathbf{s} = \mathbf{1}$ for each layer. Now choose a minibatch, and train the network using that minibatch. Once you have taken enough gradient steps and are ready to work on another minibatch, reestimate \mathbf{m} as the mean of values of the inputs to the layer, and \mathbf{s} as the corresponding standard deviations. Now obtain another minibatch, and proceed. Remember, γ and β are parameters that are trained, just like the others (using gradient descent, momentum, adagrad, or whatever). Once the network has been trained, one then takes the mean (resp. standard deviation) of the layer inputs over the training data for \mathbf{m} (resp. \mathbf{s}). Most neural network implementation environments will do all the work for you. It is quite usual to place a batch normalization layer between each layer within the network.

There is a general agreement that batch normalization improves training, but some disagreement about the details. Experiments comparing two networks, one with batch normalization the other without, suggest that the same number of steps tends to produce a lower error rate when batch normalized. Some authors suggest that convergence is faster (which isn’t quite the same thing). Others suggest that larger learning rates can be used.

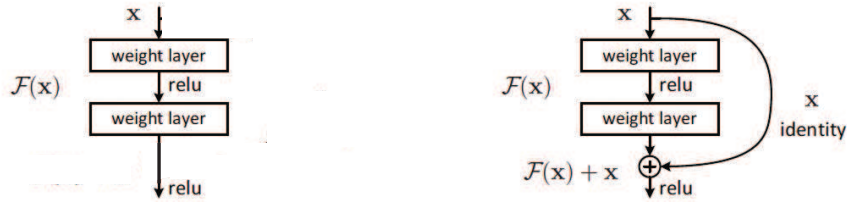


Figure 2. Residual learning: a building block.

FIGURE 10.14: *This figure, which is revised from Deep Residual Learning for Image Recognition Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun, ICCV 2016, conveys the intention of a residual network.*

10.4.5 Useful Tricks - 4: Residual Networks

A randomly initialized deep network can so severely mangle its inputs that only a wholly impractical amount of training will cause the latest layers to do anything useful. As a result, there have been practical limits on the number of layers that can be stacked (Figure 10.13). One recent strategy for avoiding this difficulty is to build a **residual layer**. Figure 10.14 sketches the idea in the form currently best understood. Remember, $F(x)$ is a RELU. Our usual layers produce $\mathbf{x}^{l+1} = \mathcal{F}(\mathbf{x}^l; \theta) = F(\mathcal{W}\mathbf{x}^l + \mathbf{b})$ as its output. This layer could be anything, but is most likely a fully connected or a convolutional layers. Then we can replace this layer with one that produces

$$\begin{aligned}\mathbf{x}^{l+1} &= F(\mathbf{x}^l + \mathcal{W}_1\mathbf{q} + \mathbf{b}_1) \\ \mathbf{q} &= F(\mathcal{W}_2\mathbf{x}^l + \mathbf{b}_2).\end{aligned}$$

It is usual, if imprecise, to think of this as producing an output that is $\mathbf{x} + \mathcal{F}(\mathbf{x}; \theta)$ — the layer passes on its input with a residual added to it. The point of all this is that, at least in principle, this residual layer can represent its output as a small offset on its input. If it is presented with large inputs, it can produce large outputs by passing on the input. Its output is also significantly less mangled by stacking layers, because its output is largely given by its input plus a non-linear function.

Very recently, an improvement on this strategy has surfaced, in *Identity Mappings in Deep Residual Networks* by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (which you can find on ArXiv using a search engine). Rather than use the expression above (corresponding to Figure 10.14), we use a layer that produces

$$\begin{aligned}\mathbf{x}^{l+1} &= \mathbf{x}^l + \mathcal{W}_1\mathbf{q} + \mathbf{b}_1 \\ \mathbf{q} &= F(\mathcal{W}_2F(\mathbf{x}^l) + \mathbf{b}_2).\end{aligned}$$

It is rather more informative to think of this as producing an output that is $\mathbf{x} + \mathcal{F}(\mathbf{x}; \theta)$ — the layer passes on its input with a residual added to it. There is good

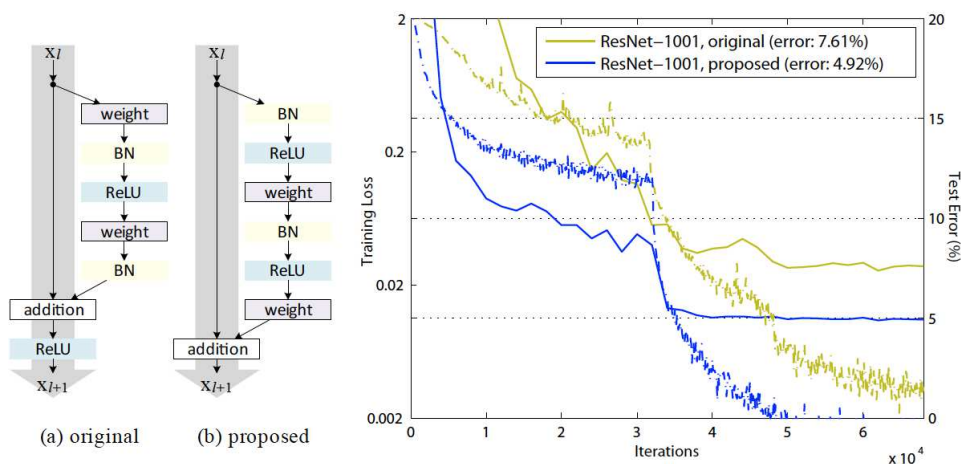


FIGURE 10.15: *This figure is revised from Identity Mappings in Deep Residual Networks by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (which you can find on ArXiv using a search engine). On the left, (a) shows the original residual network architecture, also described in Figure 10.14; (b) shows the current best architecture. On the right, train (dashed) and test (full) curves for the old and new architectures on CIFAR-10. Notice the significant improvement in performance.*

evidence that such layers can be stacked very deeply indeed (the paper I described uses 1001 layers to get under 5% error on CIFAR-10; beat that if you can!). One reason is that there are useful components to the gradient for each layer that do not get mangled by previous layers. You can see this by considering the Jacobian of such a layer with respect to its inputs. You will see that this Jacobian will have the form

$$\mathcal{J}_{x^{t+1};x^t} = (\mathcal{I} + \mathcal{M}_l)$$

where \mathcal{I} is the identity matrix and \mathcal{M}_l is a set of terms that depend on \mathcal{W} and \mathbf{b} . Now remember that, when we construct the gradient at the k 'th layer, we evaluate by multiplying a set of Jacobians corresponding to the layers above. This product in turn must look like

$$\mathcal{J}_{o;\theta^k} = \mathcal{J}_{o;x^{k+1}} \mathcal{J}_{x^{k+1};\theta^k} = (\mathcal{I} + \mathcal{M}_1 + \dots) \mathcal{J}_{x^{k+1};\theta^k}$$

which means that some components of the gradient at that layer do not get mangled by being passed through a sequence of poorly estimated Jacobians. One reason I am having trouble making a compelling argument explaining why this architecture is better is that the argument doesn't seem to be known in any tighter form (it certainly isn't to me). There is overwhelming evidence that the architecture is, in practice, better; it has produced networks that are (a) far deeper and (b) far more accurate than anything produced before. But why it works remains somewhat veiled.

Structured Models, Learning and Inference

* general comments. example cases * Inference * Learning * Structure Learning

11.1 HMM'S, DYNAMIC PROGRAMMING, AND TREE-STRUCTURED MODELS

We will need models that can represent time sequences with a fair amount of structure, and that are relatively easy to work with. Hidden Markov models (HMM's; Section 11.1.1) are a useful choice of probabilistic model, with straightforward procedures for inference (Section 11.1.2) and learning (Section 11.1.3). It turns out that inference methods appropriate to hidden Markov models apply to a broader class of model, which doesn't have to be probabilistic and isn't purely temporal. This class of model, which we call a tree-structured model (Section 11.1.4), is the mainstay of the parsing community.

11.1.1 Hidden Markov Models

A program that reads American Sign Language from a video sequence of someone signing must infer a state, internal to the user, for each sign. The program will infer state from measurements of hand position that are unlikely to be accurate, but will depend—hopefully quite strongly—on the state. The signs change state in a random (but quite orderly) fashion. In particular, some sequences of states occur very seldom (e.g., a sequence of letter signs for the sequence “wkwk” is extremely unlikely). This means that both the measurements *and* the relative probabilities of different sequences of signs can be used to determine what actually happened.

The elements of this kind of problem are:

- there is a sequence of random variables (in our example, the signs), each of which is conditionally independent of all others given its predecessor and
- each random variable generates a measurement (the measurements of hand position) whose probability distribution depends on the state.

Similar elements are to be found in examples such as interpreting the movement of dancers or of martial artists. An extremely useful formal model, known as a **hidden Markov model**, corresponds to these elements.

A sequence of random variables X_n is said to be a **Markov chain** if

$$P(\mathbf{X}_n = \mathbf{a} | \mathbf{X}_{n-1} = \mathbf{b}, \mathbf{X}_{n-2} = \mathbf{c}, \dots, \mathbf{X}_0 = \mathbf{x}) = P(\mathbf{X}_n = \mathbf{a} | \mathbf{X}_{n-1} = \mathbf{b})$$

and a **homogeneous Markov chain** if this probability does not depend on n . Markov chains can be thought of as sequences with very little memory; the new

state depends on the previous state, but not on the whole history. It turns out that this property is surprisingly useful in modeling, both because many physical variables appear to have it and because it enables a variety of simple inference algorithms. There are very slightly different notations for Markov chains on discrete and continuous state spaces; we shall discuss only the discrete case.

Assume that we have a finite discrete state space. Write the elements of the space as s_i and assume that there are k elements. Assume that we have a sequence of random variables taking values in that state space that forms a homogeneous Markov chain. Now we write

$$P(X_n = s_j | X_{n-1} = s_i) = p_{ij},$$

and because the chain is independent of n , so is p_{ij} . We can write a matrix \mathcal{P} with i, j th element p_{ij} which describes the behavior of the chain; this matrix is called the **state transition matrix**. Assume that X_0 has probability distribution $P(X_0 = s_i) = \pi_i$, and we will write π as a vector with i th element π_i . This means that

$$\begin{aligned} P(X_1 = s_j) &= \sum_{i=1}^k P(X_1 = s_j | X_0 = s_i) P(X_0 = s_i) \\ &= \sum_{i=1}^k P(X_1 = s_j | X_0 = s_i) \pi_i \\ &= \sum_{i=1}^k p_{ij} \pi_i. \end{aligned}$$

and so the probability distribution for the state of X_1 is given by $\mathcal{P}^T \pi$. By a similar argument, the probability distribution for the state of X_n is given by $(\mathcal{P}^T)^n \pi$. For all Markov chains, there is at least one distribution π^s such that $\pi^s = \mathcal{P}^T \pi^s$; this is known as the **stationary distribution**. Markov chains allow quite simple and informative pictures. We can draw a weighted, directed graph with a node for each state and the weight on each edge, indicating the probability of a state transition (Figure 11.1).

If we observe the random variable X_n , then inference is easy—we know what state the chain is in. This is a poor observation model, however. A much better model is to say that, for each element of the sequence, we observe *another* random variable, whose probability distribution depends on the state of the chain. That is, we observe some Y_n , where the probability distribution is some $P(Y_n | X_n = s_i) = q_i(Y_n)$. We can arrange these elements into a matrix \mathcal{Q} . Specifying a hidden Markov model requires providing the state transition process, the relationship between state and the probability distribution on Y_n , and the initial distribution on states, i.e., the model is given by $(\mathcal{P}, \mathcal{Q}, \pi)$. We will assume that the state space has k elements.

We will assume that we are dealing with a hidden Markov model on a discrete state space; this simplifies computation considerably, usually at no particular cost. There are two important problems:

- **Inference:** We need to determine what underlying set of states gave rise to

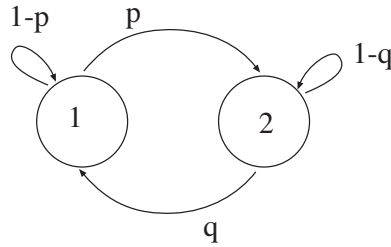


FIGURE 11.1: A simple, two-state Markov chain. In this chain, the probability of going from state one to state two is p ; from state one to state one is $1 - p$; etc. Its stationary distribution is $(q/(p+q), p/(p+q))$. This makes sense; for example, if p is very small and q is close to one, the chain will spend nearly all its time in state one. Notice that, if p and q are both very small, the chain will stay in one state for a long time, and then flip to the other state, where it will stay for a long time.

our observations. This will make it possible to, for example, infer what the dancer is doing or the signer is saying.

- **Fitting:** We need to choose a hidden Markov model that represents a sequence of past observations well.

Each has an efficient, standard solution.

11.1.2 Inference for an HMM

Assume that we have a series of N measurements \mathbf{Y}_i that we believe to be the output of a hidden Markov model. We can set up these measurements in a structure called a **trellis**. This is a weighted, directed graph consisting of N copies of the state space, which we arrange in columns. There is a column corresponding to each measurement. We weight the node representing the case that state $\mathbf{X}_j = s_k$ in the column corresponding to \mathbf{Y}_j with $\log q_k(\mathbf{Y}_j)$.

We join up the elements from column to column as follows. Consider the column corresponding to the \mathbf{Y}_j ; we join the element in this column representing state $\mathbf{X}_j = s_k$ to the element in the column corresponding to \mathbf{Y}_{j+1} representing state $\mathbf{X}_{j+1} = s_l$ if p_{kl} is nonzero. This arc represents the fact that there is a possible transition between these states. This arc is weighted with $\log p_{kl}$. Figure 11.2 shows a trellis constructed from an HMM.

The trellis has the following interesting property: each (directed) path through the trellis represents a legal sequence of states. Now because each node of the trellis is weighted with the log of the emission probability and each arc is weighted with the log of the transition probability, the joint probability of a sequence of states with the measurements can be obtained by identifying the path corresponding to this sequence, and summing the weights (of arcs and nodes) along the path. This yields an extremely effective algorithm for finding the best path, known as either **dynamic programming** or the **Viterbi algorithm**.

We start at the *final* column of the tellis. We know the log-likelihood of a one-state path, ending at each node, as this is just the weight of that node. Now

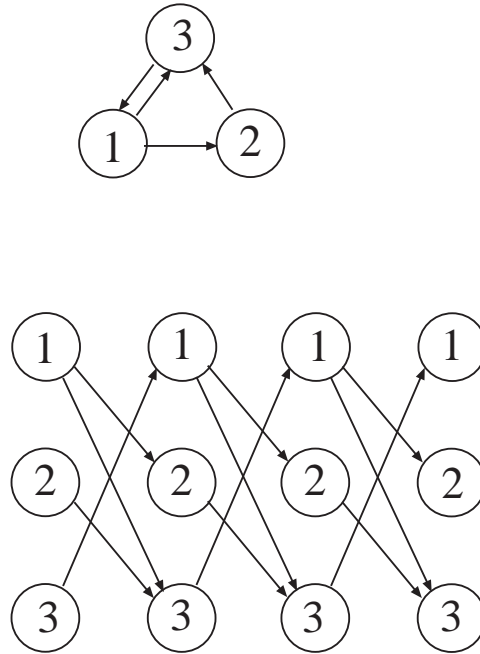


FIGURE 11.2: At the **top**, a simple state transition model. Each outgoing edge has some probability, though the topology of the model forces two of these probabilities to be 1. Below, the trellis corresponding to that model. Notice that each path through the trellis corresponds to a legal sequence of states, for a sequence of four measurements. We weight the arcs with the log of the transition probabilities, and the nodes with the log of the emission probabilities; weights are not shown here, to reduce the complexity of the drawing.

consider a two-state path, which will start at the second last column of the trellis. We can easily obtain the best path leaving each node in this column. Consider a node: we know the weight of each arc leaving the node and the weight of the node at the far end of the arc, so we can choose the path segment with the largest value of the sum; this arc is the best we can do leaving that node. Now, for each node, we add the weight at the node to the value of the best path segment leaving that node (i.e., the arc weight plus the weight of the node at the far end). This sum is the best value obtainable on reaching that node—which we'll call the **node value**.

Now, because we know the best value obtainable on reaching each node in the second-last column, we can figure out the best value obtainable on reaching each node in the third-last column. At each node in the third-last column, we have a choice of arcs, each reaching a node whose value we know. We choose the arc with the largest value of (arc weight plus node value), then add this value to the weight at the starting node in the third last column, and this yields the value of the starting node. We can repeat this process, until we have a value for each of the nodes in the first column; the largest value is the maximum likelihood.

We can also get the path with the maximum likelihood value. When we

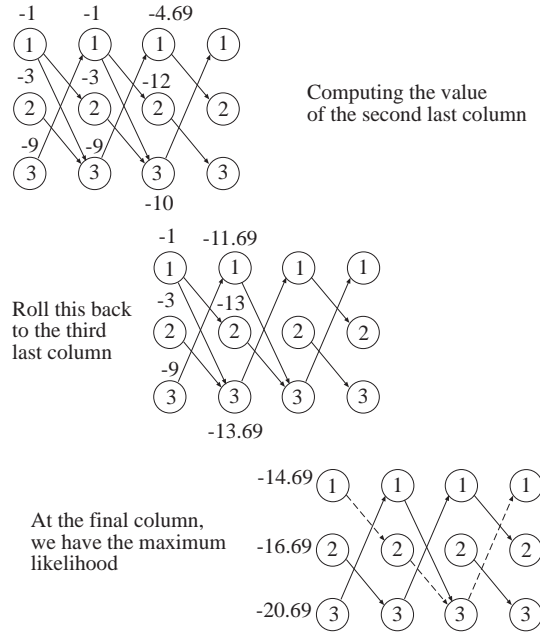


FIGURE 11.3: *It is a simple matter to find the best path through the trellis of Figure 11.2 (or any other trellis, for that matter!). We assume that each 1 node has log-probability -1 , each 2 node has log-probability -3 , and each 3 node has log-probability -9 . We also assume that the probabilities of leaving a node are uniform (check our numbers!). Now the value of each node in the second-last column is the value of the node plus the best value to be obtained by leaving that node. This is easily computed. The algorithm involves computing the value of each node in the second-last column; then of each node in the third-last column, etc., as described in the text. Once we get to the start of the trellis, the largest weight is the maximum of the log-likelihood; because we erased all but the best path segments, we have the best path, too (indicated by a dashed line).*

compute the value of a node, we erase all but the best arc leaving that node. Once we reach the first column, we simply follow the path from the node with the best value. Figure 11.3 illustrates this extremely simple and very powerful algorithm.

We can write all this out more formally. For inference, we have a series of observations $\{Y_0, Y_1, \dots, Y_n\}$, and we would like to obtain the sequence of $n + 1$ states $\mathbf{S} = \{S_0, S_1, \dots, S_n\}$ that maximizes

$$P(\mathbf{S} | \{Y_0, Y_1, \dots, Y_n\}, (\mathcal{P}, \mathcal{Q}, \pi)),$$

which is the same as maximizing the joint distribution

$$P(\mathbf{S}, \{Y_0, Y_1, \dots, Y_n\} | (\mathcal{P}, \mathcal{Q}, \pi)).$$

There is a standard algorithm for this purpose, the **Viterbi algorithm**. We seek an $n + 1$ element path through the states (from S_0 to S_n). There are k^{n+1}

This algorithm yields the path through an HMM that maximizes the joint, and the value of the joint at this path. Here δ and ψ are convenient bookkeeping variables (as in the text); p^* is the maximum value of the joint; and q_t^* is the t th state in the optimal path.

1. Initialization:

$$\begin{aligned}\delta_1(j) &= \pi_j q_j(Y_1) \quad 1 \leq j \leq N \\ \psi_1(j) &= 0\end{aligned}$$

2. Recursion:

$$\begin{aligned}\delta_n(j) &= \left(\max_i \delta_{n-1}(i) P_{ij} \right) q_j(Y_n) \\ \psi_n(j) &= \arg \max_i (\delta_{n-1}(i) P_{ij})\end{aligned}$$

3. Termination:

$$\begin{aligned}p^* &= \max_i (\delta_N(i)) \\ q_N^* &= \arg \max_i (\delta_N(i))\end{aligned}$$

4. Path backtracking:

$$q_t^* = \psi_{t+1}(q_{t+1}^*)$$

Algorithm 11.1: *The Viterbi Algorithm*

We fit a model to a data sequence \mathbf{Y} is achieved by a version of EM. We assume a model $(\mathcal{P}, \mathcal{Q}, \pi)_i$, and then compute the coefficients of a new model; this iteration is guaranteed to converge to a local maximum of $P(\mathbf{Y} | (\mathcal{P}, \mathcal{Q}, \pi))$.

Until $(\mathcal{P}, \mathcal{Q}, \pi)_{i+1}$ is the same as $(\mathcal{P}, \mathcal{Q}, \pi)_i$
 compute the forward variables α and β
 using the procedures of algorithms 11.3 and 11.4

$$\text{compute } \xi_t(i, j) = \frac{\alpha_t(i) p_{ij} q_j(Y_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) p_{ij} q_j(Y_{t+1}) \beta_{t+1}(j)}$$

compute the updated parameters using the procedures of Algorithm 11.5

These values are the elements of $(\mathcal{P}, \mathcal{Q}, \pi)_{i+1}$
 end

Algorithm 11.2: *Fitting Hidden Markov Models with EM*

such paths, because we could choose from each state for every element of the path (assuming that there are no zeros in \mathcal{P} ; in most cases, there are certainly $O(k^{n+1})$ paths). We can't look at every path, but in fact we don't have to. The approach is as follows: assume that, for each possible state s_l , we know the value of the joint for the best n -step path that ends in $S_{n-1} = s_l$; then, the path that maximizes the joint for an $n + 1$ step path must consist of one of these paths, combined with another step. All we have to do is find the missing step.

$$\begin{aligned}\alpha_0(j) &= \pi_j q_j(Y_0) \\ \alpha_{t+1}(j) &= \left[\sum_{l=1}^k \alpha_t(l) p_{lj} \right] q_j(Y_{t+1}) \quad 1 \leq t \leq n-1\end{aligned}$$

Algorithm 11.3: *Computing the Forward Variable for Fitting an HMM*

$$\begin{aligned}\beta_N(j) &= 1 \\ \beta_t(j) &= \left[\sum_{l=1}^k p_{jl} q_l(Y_{t+1}) \right] \beta_{t+1}(j) \quad 1 \leq t \leq k-1\end{aligned}$$

Algorithm 11.4: *Computing the Backward Variable for Fitting an HMM*

We can approach finding the path with the maximum value of the joint as an induction problem. Assume that, for each value j of S_{n-1} , we know the value of the joint for the best path that ends in $S_{n-1} = j$, which we write as

$$\delta_{n-1}(j) = \max_{S_0, S_1, \dots, S_{n-2}} P(\{S_0, S_1, \dots, S_{n-1} = j\}, \{Y_0, Y_1, \dots, Y_{n-1}\} | (\mathcal{P}, \mathcal{Q}, \pi)).$$

Now we have that

$$\delta_n(j) = \left(\max_i \delta_{n-1}(i) P_{ij} \right) q_j(Y_n).$$

We need not only the maximum *value*, but also the path that gave rise to this value. We define another variable

$$\psi_n(j) = \arg \max (\delta_{n-1}(i) P_{ij})$$

(i.e., the best path that ends in $S_n = j$). This gives us an inductive algorithm for getting the best path.

The reasoning is as follows: I know the best path to each state for the $n - 1$ th measurement; for each state for the n th measurement, I can look backward and

$$\begin{aligned}
\overline{\pi}_i &= \text{expected frequency of being in state } s_i \text{ at time 1} \\
&= \sum_{j=1}^N \xi_1(i, j) \\
\overline{p}_{ij} &= \frac{\text{expected number of transitions from } s_i \text{ to } s_j}{\text{expected number of transitions from state } s_i} \\
&= \frac{\sum_{t=1}^T \xi_t(i, j)}{\sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j)} \\
\overline{q}_i(k) &= \frac{\text{expected number of times in } s_i \text{ and observing } Y = y_k}{\text{expected number of times in state } s_i} \\
&= \frac{\sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j) \delta(Y_t, y_k)}{\sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j)}
\end{aligned}$$

here $\delta(u, v)$ is one if its arguments are equal and zero otherwise.

Algorithm 11.5: *Updating Parameters for Fitting an HMM*

choose the best state for the $n - 1$ th measurement; but I know the best path from there, so I have the best path to each state for the n th measurements. We have put everything together in Algorithm 11.1.

11.1.3 Fitting an HMM with EM

We have a dataset \mathbf{Y} for which we believe a hidden Markov model is an appropriate model, but which hidden Markov model should we use? We wish to choose a model that best represents a set of data. To do this, we will use a version of the Expectation-Maximization algorithm of Section ???. In this algorithm, we assume that we have an HMM, $(\mathcal{P}, \mathcal{Q}, \pi)$; we now want to use this model and our dataset to estimate a new set of values for these parameters. We now estimate $(\overline{\mathcal{P}}, \overline{\mathcal{Q}}, \overline{\pi})$ using a procedure that is given below. There will be two possibilities (a fact that we won't prove). Either $P(\mathbf{Y} | (\overline{\mathcal{P}}, \overline{\mathcal{Q}}, \overline{\pi})) > P(\mathbf{Y} | (\mathcal{P}, \mathcal{Q}, \pi))$, or $(\overline{\mathcal{P}}, \overline{\mathcal{Q}}, \overline{\pi}) = (\mathcal{P}, \mathcal{Q}, \pi)$.

The updated values of the model parameters will have the form:

$$\begin{aligned}
\overline{\pi}_i &= \text{expected frequency of being in state } s_i \text{ at time 1} \\
\overline{p}_{ij} &= \frac{\text{expected number of transitions from } s_i \text{ to } s_j}{\text{expected number of transitions from state } s_i} \\
\overline{q}_j(k) &= \frac{\text{expected number of times in } s_j \text{ and observing } Y = y_k}{\text{expected number of times in state } s_j}
\end{aligned}$$

We need to be able to evaluate these expressions. In particular, we need to be able to determine

$$P(X_t = s_i, X_{t+1} = s_j | \mathbf{Y}, (\mathcal{P}, \mathcal{Q}, \pi)),$$

which we shall write as $\xi_t(i, j)$. If we know $\xi_t(i, j)$, we have

$$\begin{aligned} \text{expected number of transitions from } s_i \text{ to } s_j &= \sum_{t=1}^T \xi_t(i, j) \\ \text{expected number of transitions from } s_i &= \sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j) \\ \text{expected number of times in } s_i &= \sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j) \\ \text{expected frequency of being in } s_i \text{ at time 1} &= \sum_{j=1}^N \xi_1(i, j) \\ \text{expected number of times in } s_i \text{ and observing } Y = y_k &= \sum_{t=1}^T \sum_{j=1}^N \xi_t(i, j) \delta(Y_t, y_k) \end{aligned}$$

where $\delta(u, v)$ is one if its arguments are equal and zero otherwise.

To evaluate $\xi_t(i, j)$, we need two intermediate variables: a **forward variable** and a **backward variable**.

The forward variable is $\alpha_n(j) = P(Y_0, Y_1, \dots, Y_n, X_n = s_j | (\mathcal{P}, \mathcal{Q}, \pi))$.

The backward variable is $\beta_t(j) = P(\{Y_{t+1}, Y_{t+2}, \dots, Y_n\} | X_t = s_j, (\mathcal{P}, \mathcal{Q}, \pi))$.

If we assume that we know the values of these variables, we have that

$$\begin{aligned} \xi_t(i, j) &= P(X_t = s_i, X_{t+1} = s_j | \mathbf{Y}, (\mathcal{P}, \mathcal{Q}, \pi)) \\ &= \frac{P(\mathbf{Y}, X_t = s_i, X_{t+1} = s_j | (\mathcal{P}, \mathcal{Q}, \pi))}{P(\mathbf{Y} | (\mathcal{P}, \mathcal{Q}, \pi))} \\ &= \frac{\left\{ \begin{array}{l} P(Y_0, Y_1, \dots, Y_t, X_t = s_i | (\mathcal{P}, \mathcal{Q}, \pi)) \\ \times P(Y_{t+1} | X_{t+1} = s_j, (\mathcal{P}, \mathcal{Q}, \pi)) \\ \times P(X_{t+1} = s_j | X_t = s_i, (\mathcal{P}, \mathcal{Q}, \pi)) \\ \times P(Y_{t+2}, \dots, Y_n | X_{t+1} = s_j, (\mathcal{P}, \mathcal{Q}, \pi)) \end{array} \right\}}{P(\mathbf{Y} | (\mathcal{P}, \mathcal{Q}, \pi))} \\ &= \frac{\alpha_t(i) p_{ij} q_j(Y_{t+1}) \beta_{t+1}(j)}{P(\mathbf{Y} | (\mathcal{P}, \mathcal{Q}, \pi))} \\ &= \frac{\alpha_t(i) p_{ij} q_j(Y_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) p_{ij} q_j(Y_{t+1}) \beta_{t+1}(j)}. \end{aligned}$$

Both the forward and backward variables can be evaluated by induction. We

get $\alpha_n(j)$, by observing that:

$$\begin{aligned}
\alpha_0(j) &= P(Y_0, X_0 = s_j | (\mathcal{P}, \mathcal{Q}, \pi)) \\
&= \pi_j q_j(Y_0) \\
\alpha_{t+1}(j) &= P(Y_0, Y_1, \dots, Y_{t+1}, X_{t+1} = s_j | (\mathcal{P}, \mathcal{Q}, \pi)) \\
&= P(Y_0, Y_1, \dots, Y_t, X_{t+1} = s_j | (\mathcal{P}, \mathcal{Q}, \pi)) P(Y_{t+1} | X_{t+1} = s_j) \\
&= \sum_{l=1}^k [P(Y_0, Y_1, \dots, Y_t, X_t = s_l, X_{t+1} = s_j | (\mathcal{P}, \mathcal{Q}, \pi)) P(Y_{t+1} | X_{t+1} = s_j)] \\
&= \left(\sum_{l=1}^k \left[P(Y_0, Y_1, \dots, Y_t, X_t = s_l | (\mathcal{P}, \mathcal{Q}, \pi)) \right. \right. \\
&\quad \left. \left. \times P(X_{t+1} = s_j | X_t = s_l) \right] P(Y_{t+1} | X_{t+1} = s_j) \right) \\
&= \left[\sum_{l=1}^k \alpha_t(l) p_{lj} \right] q_j(Y_{t+1}) \quad 1 \leq t \leq n-1.
\end{aligned}$$

This backward variable can also be obtained by induction as:

$$\begin{aligned}
\beta_N(j) &= P(\text{no further output} | X_n = s_j, (\mathcal{P}, \mathcal{Q}, \pi)) \\
&= 1 \\
\beta_t(j) &= P(\{Y_{t+1}, Y_{t+2}, \dots, Y_n\} | X_t = s_j, (\mathcal{P}, \mathcal{Q}, \pi)) \\
&= \sum_{l=1}^k [P(\{Y_{t+1}, Y_{t+2}, \dots, Y_n\}, X_t = s_l | X_{t+1} = s_j, (\mathcal{P}, \mathcal{Q}, \pi))] \\
&= \left(\left[\sum_{l=1}^k P(X_t = s_l, Y_{t+1} | X_{t+1} = s_j) \right] \right. \\
&\quad \left. \times P(\{Y_{t+2}, \dots, Y_n\} | X_{t+1} = s_j, (\mathcal{P}, \mathcal{Q}, \pi)) \right) \\
&= \left[\sum_{l=1}^k p_{jl} q_l(Y_{t+1}) \right] \beta_{t+1}(j) \quad 1 \leq t \leq k-1.
\end{aligned}$$

As a result, we have a simple fitting algorithm, collected in Algorithm 11.2.

11.1.4 Tree-Structured Energy Models

A major attraction of hidden Markov models is the relative simplicity of inference and learning. This is a consequence of the combinatorial structure of the model; in fact, our algorithms apply to somewhat richer combinatorial structures, and to models that are not intrinsically probabilistic. We now consider a situation where we will choose values for a set discrete variables, traditionally written X_i for $i = 1, \dots, n$, to maximize the value of an objective function $f(X_1, \dots, X_n)$. This objective function is a sum of **unary terms** (i.e., functions that take one argument), which we write $u_j(X_i)$, and **binary terms** (i.e., functions that take two arguments), which we write $b_k(X_i, X_j)$. This is a relatively general model of parts and relations. There is one score (which would be the u_k) associated with each

part (which would be the X_i) and another associated with some of the relations (i.e., the b_k). For example, in the case of the HMM, the variables would be the hidden states, the unary terms would be the logs of emission probabilities, and binary terms would be the logs of transition probabilities. It is natural to think of the unary terms as nodes in a graph, and the binary terms as edges. However, it is not required to think of the unary or binary terms as log probabilities; instead, you could think of them as negative energies (because we are maximizing), or you could minimize, and think of them as energies or costs. For this kind of model, maximization is straightforward if the graph we have described is a forest.

HMM's are a special case, because the graph in that case is a chain. We will redescribe inference for an HMM in this more general setting, because it will then follow easily that the method applies to a forest. We can write the objective function as

$$f_{\text{chain}}(X_1, \dots, X_n) = \sum_{i=1}^{i=n} u_i(X_i) + \sum_{i=1}^{i=n-1} b_i(X_i, X_{i+1})$$

and we wish to maximize this function (you should check that the terms match terms in the expression for the joint for an HMM; a strategically placed logarithm will help). Now we define a new function, the **cost-to-go function**, with a recursive definition. Write

$$f_{\text{cost-to-go}}^{(n-1)}(X_{n-1}) = \max_{X_n} b_{n-1}(X_{n-1}, X_n) + u_n(X_n),$$

and notice that we have

$$\operatorname{argmax}_{X_1, \dots, X_n} f_{\text{chain}}(X_1, \dots, X_n)$$

is equal to

$$\operatorname{argmax}_{X_1, \dots, X_{n-1}} \left(f_{\text{chain}}(X_1, \dots, X_{n-1}) + f_{\text{cost-to-go}}^{(n-1)}(X_{n-1}) \right),$$

which means that we can eliminate the n th variable from the optimization by replacing the term $b_{n-1}(X_{n-1}, X_n) + u_n(X_n)$ with a function of X_{n-1} . This function is obtained by maximizing this term with respect to X_n . Equivalently, assume we must choose a value for X_{n-1} . The cost-to-go function tells us the value of $b_{n-1}(X_{n-1}, X_n) + u_n(X_n)$ obtained by making the best choice of X_n conditioned on our choice of X_{n-1} . Because any other choice would not lead to a maximum, if we know the cost-to-go function at X_{n-1} , we can now compute the best choice of X_{n-1} conditioned on our choice of X_{n-2} . This yields that

$$\max_{X_{n-1}, X_n} [b_{n-2}(X_{n-2}, X_{n-1}) + u_{n-1}(X_{n-1}) + b_{n-1}(X_{n-1}, X_n) + u_n(X_n)]$$

is equal to

$$\max_{X_{n-1}} \left[b_{n-2}(X_{n-2}, X_{n-1}) + u_{n-1}(X_{n-1}) + \left(\max_{X_n} b_{n-1}(X_{n-1}, X_n) + u_n(X_n) \right) \right].$$

But all this can go on recursively, yielding

$$f_{\text{cost-to-go}}^{(k)}(X_k) = \max_{X_{k+1}} b_k(X_k, X_{k+1}) + u_k(X_k) + f_{\text{cost-to-go}}^{(k+1)}(X_{k+1}).$$

We can expand this to describe our use of the trellis in Section 11.1.2. Notice that

$$\operatorname{argmax}_{X_1, \dots, X_n} f_{\text{chain}}(X_1, \dots, X_n)$$

is equal to

$$\operatorname{argmax}_{X_1, \dots, X_{n-1}} \left(f_{\text{chain}}(X_1, \dots, X_{n-1}) + f_{\text{cost-to-go}}^{(n-1)}(X_{n-1}) \right)$$

which is equal to

$$\operatorname{argmax}_{X_1, \dots, X_{n-2}} \left(f_{\text{chain}}(X_1, \dots, X_{n-2}) + f_{\text{cost-to-go}}^{(n-2)}(X_{n-2}) \right),$$

and we can apply the recursive definition of the cost-to-go function to get

$$\operatorname{argmax}_{X_1, \dots, X_n} f_{\text{chain}}(X_1, \dots, X_n) = \operatorname{argmax}_{X_1} \left(f_{\text{chain}}(X_1) + f_{\text{cost-to-go}}^1(X_1) \right),$$

which yields an extremely powerful maximization strategy. We start at X_n , and construct $f_{\text{cost-to-go}}^{(n-1)}(X_{n-1})$. We can represent this function as a table, giving the value of the cost-to-go function for each possible value of X_{n-1} . We build a second table giving the optimum X_n for each possible value of X_{n-1} . From this, we can build $f_{\text{cost-to-go}}^{(n-2)}(X_{n-2})$, again as a table, and also the best X_{n-1} as a function of X_{n-2} , again as a table, and so on. Now we arrive at X_1 . We obtain the solution for X_1 by choosing the X_1 that yields the best value of $\left(f_{\text{chain}}(X_1) + f_{\text{cost-to-go}}^2(X_2) \right)$. But from this solution, we can obtain the solution for X_2 by looking in the table that gives the best X_2 as a function of X_1 ; and so on. It should be clear that this process yields a solution in polynomial time; in the exercises, you will show that, if each X_i can take one of k values, then the time is $O(nK^2)$.

This strategy will work for a model with the structure of a forest. The proof is an easy induction. If the forest has no edges (i.e., consists entirely of nodes), then it is obvious that a simple strategy applies (choose the best value for each X_i independently). This is clearly polynomial. Now assume that the algorithm yields a result in polynomial time for a forest with e edges, and show that it works for a forest with $e + 1$ edges. There are two cases. The new edge could link two existing trees, in which case we could re-order the trees so the nodes that are linked are roots, construct a cost-to-go function for each root, and then choose the best pair of states for these roots from the cost-to-go functions. Otherwise, one tree had a new edge added, joining the tree to an isolated node. In this case, we reorder the tree so that this new node is the root and build a cost-to-go function from the leaves to the root. The fact that the algorithm works is a combinatorial insight, but many kinds of model have a tree structure. Models of this form are particularly important in cases of tracking and of parsing.

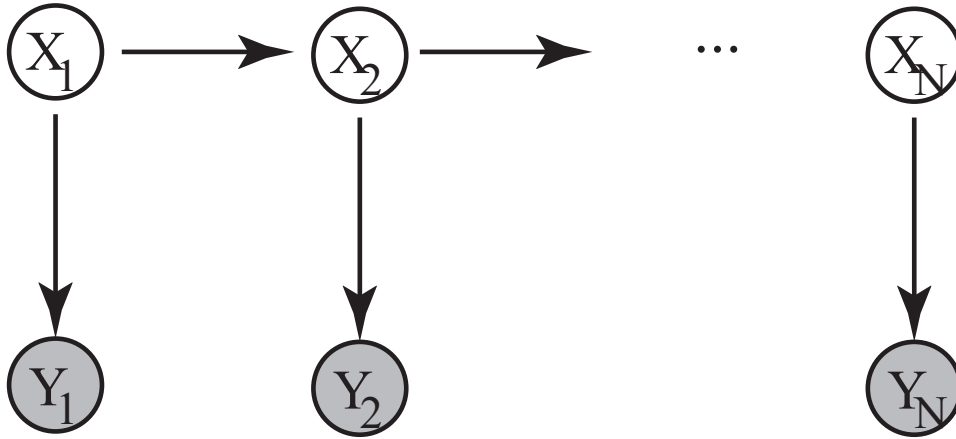


FIGURE 11.4:

11.2 EXAMPLES

11.2.1 Transcribing Speech

11.2.2 Tracking a Blob

11.2.3 Reading Text

11.2.4 Human Parsing

11.3 SIMPLE CONDITIONAL RANDOM FIELD MODELS

HMM models have been widely used, but have one odd feature that is inconsistent with practical experience. Recall X_i are the hidden variables, and Y_i are the observations. HMM's model

$$P(Y_1, \dots, Y_n | X_1, \dots, X_n) \propto P(Y_1, \dots, Y_n, X_1, \dots, X_n),$$

which is the probability of observations given the hidden variables. This is modelled using the factorization

$$P(Y_1, Y_2, \dots, Y_N, X_1, X_2, \dots, X_N) = P(X_1)P(Y_1|X_1)P(X_2|X_1)P(Y_2|X_2) \dots P(X_N|X_{N-1})P(Y_N|X_N).$$

In much of what we will do, this seems unnatural. For example, in the case of reading written text, we would be modelling the probability of the observed ink given the original text. But we would not attempt to find a single character by modelling the probability of the observed ink given the character (I will call this a **generative** strategy). Instead, we would search using a classifier, which is a model of the probability of the character conditioned on the observed ink (I will call this a **discriminative** strategy). The two strategies are quite different in practice. A generative strategy would need to explain all possible variants of the ink that a character could produce, but a discriminative strategy just needs to judge whether the ink observed is the character or not.

As another example, in the case of parsing people, we would be modelling the probability of the observed pixels given the configuration of the body. Again, this

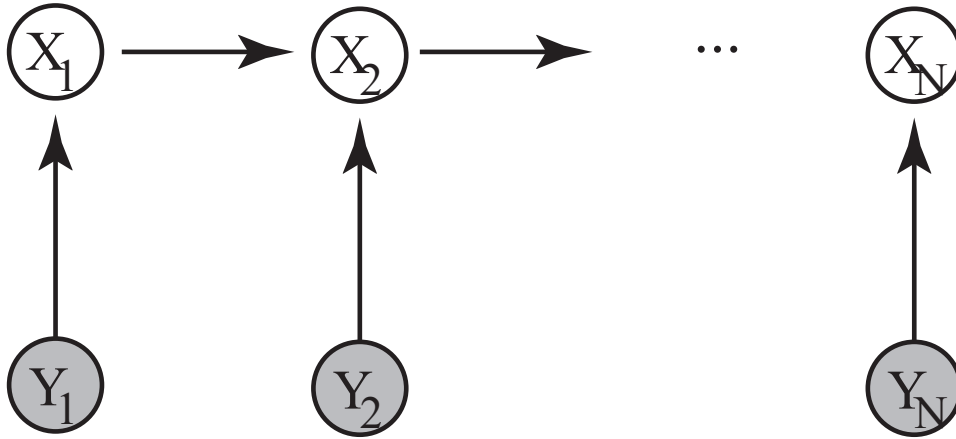


FIGURE 11.5:

is an odd thing to do, because we may need quite sophisticated models to predict any set of pixels that could be produced by a body part. The alternative is to use a body part detector – a model of the probability the part is present, given the observed pixels. This reasoning applies to each of our examples.

11.3.1 Drawing a Model

We now adopt a convention that allows us to draw a model as a graph. Assume we have a probability distribution over a collection of variables. In the HMM example, these would be the hidden and the observed states in the HMM example. We will draw each variable as a circle. We will fill the circle if the variable’s value is observed. We will add edges to this drawing according to the following rules. Assume we know a factorization of the probability distribution we are modelling into a collection of terms. There will be an edge between each pair of variables that appears together in a at least once in any factor. If two variables, say X_i and Y_i , only appear together in a factor in the form $P(Y_i|X_i)$, then the edge will be directed from X_i to Y_i (the arrow points toward the variable that is “generated”).

TODO: FIX

These rules allow us to draw an HMM as Figure 11.4. This drawing corresponds to the factorization above.

11.3.2 MEMM’s and Label Bias

An alternative that might seem plausible is to try and build a model with the structure suggested by Figure 11.5. This would correspond to

$$P(X_1, X_2, \dots, X_N | Y_1, Y_2, \dots, Y_N) = P(X_1 | Y_1) P(X_2 | Y_2, X_1) P(X_3 | X_2, Y_2), \dots P(X_N | X_{N-1}, Y_N).$$

Such models are known as **maximum entropy markov models** or **MEMM’s**. It turns out these models present practical difficulties, which are worth understanding. Assume we have fitted a model, and wish to recover the best sequence of X_i

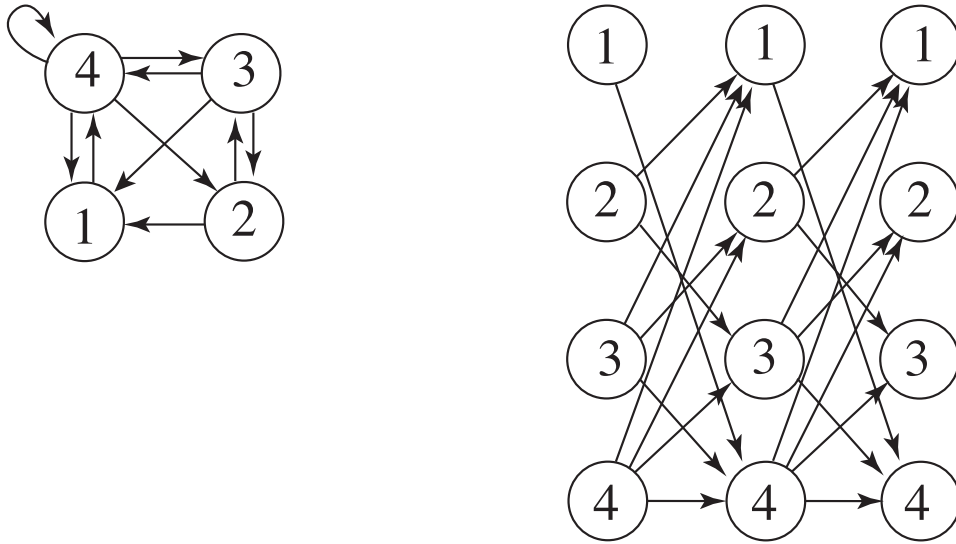


FIGURE 11.6:

corresponding to a given sequence of observations Y_i . We could minimize

$$-\log P(X_1, X_2, \dots, X_N | Y_1, Y_2, \dots, Y_N) = -\log P(X_1 | Y_1) - \log P(X_2 | Y_2, X_1) - \log P(X_3 | X_2, Y_2), \dots - \log P(X_N | X_{N-1}, Y_N)$$

by choice of X_1, \dots, X_N . We can represent this cost function on a trellis, as for the HMM, but now notice that the costs on the trellis behave differently. For an HMM, each state (circle) in a trellis had a cost, corresponding to $-\log P(Y_i | X_i)$, and each edge had a cost ($-\log P(X_{i+1} | X_i)$), and the cost of a particular sequence was the sum of the costs along the implied path. But for an MEMM, the representation is slightly different. There is no term associated with each state in the trellis; instead, we associate the edge going from the state $X_i = U$ to the state $X_{i+1} = V$ with the cost $-\log P(X_{i+1} = V | X_i = U, Y_i)$. Again, the cost of a sequence of states is represented by the sum of costs along the corresponding path. This may look to you like a subtle change, but it has nasty effects.

Look at the example of Figure 11.6. Notice that when the model is in state 1, it can only transition to state 2. In turn, this means that $-\log P(X_{i+1} = 2 | X_i = 1, Y_i) = 0$ *whatever* the measurement Y_i is. Furthermore, either $P(X_{i+1} = 3 | X_i = 2, Y_i) \geq 0.5$ or $P(X_{i+1} = 1 | X_i = 2, Y_i) \geq 0.5$ (because there are only two options leaving state 2). Here the measurement can determine which of the two options has higher probability. That figure shows a trellis corresponding to three measurements. In this trellis, the path 2 1 4 will be the lowest cost path *unless* the first measurement overwhelmingly disfavors the transition $2 \rightarrow 1$. This is because most other paths must share weights between many outgoing edges; but $1 \rightarrow 4$ is very cheap, and $2 \rightarrow 1$ will be cheap unless there is an unusual measurement. Paths which pass through many states with few outgoing edges are strongly favored. This is known as the **label bias problem**. There are some fixes that can be applied, but it is better to reengineer the model.

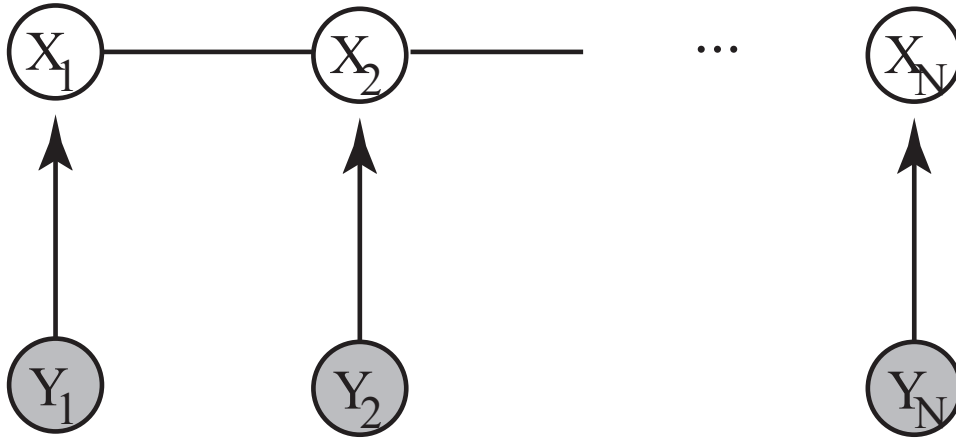


FIGURE 11.7:

11.3.3 Conditional Random Field Models

We want a model of sequences that is discriminative, but doesn't have the label bias problem. We'd also like that model to be as tractable as an HMM for sequences. Notice that what made an HMM easy to deal with was that we could solve for the best sequence of states using dynamic programming. We can use a trellis to represent the inference problem for any model where transitions between the hidden states are Markov. We can certainly use dynamic programming to solve for the best sequence of states if the cost model consists of per-edge and per-state costs (as it does in the case of the HMM). But to make the model discriminative, we need to change the probabilistic interpretation of these costs.

We start with the costs, written on the trellis. Write $E_i(a, b)$ for the cost of the edge from $X_i = a$ to $X_{i+1} = b$, and write $V_i(a)$ for the cost of assigning $X_i = a$. We will assume that $E_i(a, b) \geq 0$ and $V_i(a) \geq 0$. Note that for a useful model, we would like $V_i(a)$ to depend on Y_i . With appropriate rescaling, etc., we can interpret $V_i(a)$ as $-\log P(X_i = a | Y_i)$ (exercise: show that this doesn't affect the best path). We now apply a probabilistic interpretation to this model as

$$P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N | Y_1, Y_2, \dots, Y_N) = \frac{1}{K} \exp [V_1(x_1) + E_1(x_1, x_2) + V_2(x_2) + E_2(x_2, x_3) + \dots V_N(x_N)]$$

where K is some normalizing constant, chosen to ensure the probability distribution sums to 1. Notice the crucial difference with the MEMM; there are now node as well as edge costs *but* we can't interpret the edge costs as transition probabilities. A model with this structure is known as a **conditional random field**. We can expand our graphical vocabulary to include undirected edges, corresponding to the $E_i(a, b)$ terms. The result is the model depicted in Figure 11.7; this figure is interpreted as yielding the probabilistic model above.

Assume we have known E_i and V_i terms. Then inference proceeds by dynamic programming, as above. What is much more interesting is *learning* a CRF. We can likely come up with a reasonable scheme for estimating the V_i terms, but we don't

know K , and K depends on the E_i and V_i terms in a complicated way. In particular, we have

$$K = \sum_{x_1, x_2, \dots, x_n} P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_n | Y_1, Y_2, \dots, Y_N)$$

which means that choosing E_i values by maximizing the posterior probability of observed examples is going to be hard.

11.4 DISCRIMINATIVE LEARNING OF CRFS

A really powerful strategy for learning a CRF follows from ignoring the probabilistic semantics of the model and thinking about the model as representing an energy or cost. We write the cost of a sequence x_1, x_2, \dots, x_n as

$$\begin{aligned} C(x_1, x_2, \dots, x_n) &= V_1(x_1) + E_1(x_1, x_2) + V_2(x_2) + E_2(x_2, x_3) + \dots + V_N(x_n) \\ &= V_1(x_1) + \sum_{i=1}^{n-1} [E_i(x_i, x_{i+1}) + V_{i+1}(x_{i+1})] \\ &= \sum_{i=1}^n V_i(x_i) + \sum_{i=1}^{n-1} E_i(x_i, x_{i+1}). \end{aligned}$$

Notice that the slightly funny form of the sum in the second line occurs because when there are N vertical arrows in Figure ??, there are $N - 1$ horizontal edges. The notation in the third line separates edge costs from node costs. Each notation appears on occasion. You should think of each V_i as a vector of costs (one per state in the model) and each E_i as a table of costs (one per pair of states in the model). Inference involves finding the sequence of states that minimizes this cost. Learning will involve building the V_i and E_i from a collection of examples. Each example will be a pair of sequences, one of observations and the other of hidden states. These sequences might not be of the same length (think about learning to read words; they're not all the same length).

11.4.1 Representing the Model

We need a parametric representation of V_i and E_i ; we'll then search for the parameters that yield the right model. For the V , we will assume that V_i and V_j differ only if $Y_i \neq Y_j$. We then construct a set of functions $\phi_j^{(v)}(X, Y)$, and a vector of parameters $\theta_j^{(v)}$. Finally, we assert that $V_i(x) = \sum_j \theta_j^{(v)} \phi_j^{(v)}(x, Y_i)$. Similarly, for E_i we will construct a set of functions $\phi_j^{(v)}(U, V)$, and have $E_i(U, V) = \sum_j \theta_j^{(v)} \phi_j^{(v)}(U, V)$.

I give some sample constructions below, but you may find them somewhat involved at first glance. What is important is that (a) we have some parameters θ so that, for different choices of parameter, we get different cost functions; and (b) for any choice of parameters (say $\hat{\theta}$) and any sequence of Y_i we can label each vertex and each edge on the trellis with a number representing the cost of that vertex or edge. Assume we have these properties. Then for any particular $\hat{\theta}$ we can construct the best sequence of X_i using dynamic programming (as above). Furthermore, we

can try to adjust the choice of $\hat{\theta}$ so that for the i 'th training sequence, $\mathbf{y}^{(i)}$, inference yields $\mathbf{x}^{(i)}$ or something close by.

Building V : The choice of $\phi_j^{(v)}(X, Y)$ will naturally depend on the application. We can get some guidance by noticing we would like V to behave like a classifier — we'd like smaller values for $V(x, Y)$ when the x is a “good” interpretation of the observed state, and larger values when it isn't. For example, assume that the observed states are inked numerals (like MNIST), and that they appear in a window of fixed size. Multinomial logistic regression works quite well on MNIST using just the pixel values as features. This means that I can compute 10 linear functions of the pixel values (one for each numeral) such that the linear function corresponding to the right numeral is smaller than any other of the functions, at least most of the time. Now for each possible true numeral x and each pixel location p we can build a feature function $\phi(X, Y) = \mathbb{I}_{[X=x]} \mathbb{I}_{[Y(p)=0]}$. This is 1 if $X = x$ (i.e. for a particular numeral, x) and the ink at pixel location p is dark, and otherwise zero. We index these feature functions in any way that seems convenient to get $\phi_j^{(v)}(X, Y)$, and stack them into a vector. Now $\sum_j \theta_j^{(v)} \phi_j^{(v)}(X, Y)$ is (a) linear in the ink for any fixed value of X and (b) different for different values of X and the same ink. With a good choice of $\theta^{(v)}$, we could produce a function that was small for the right numeral and large for the wrong numeral.

Another strategy would be to build the $\phi(X, Y)$ out of classifiers. For each class x , we will build several classifiers each of which can tell that class from all the others. We obtain different classifiers by using different training sets; or different features; or different classifier architectures; or all of these. Write $g_{i,x}(Y)$ for the i 'th classifier for class x . We ensure that $g_{i,x}(Y)$ is small if Y is of class x , and large otherwise. Then for each classifier and each class we can build a feature function by $\phi(X, Y) = g_{i,x}(Y)$. We index these feature functions in any way that seems convenient. Now, with appropriate choice of $\theta_j^{(v)}$, we have that $\sum_j \theta_j^{(v)} \phi_j^{(v)}(X, Y)$ is a weighted sum of classifiers that produces a small value if X has the same class as Y and a large value otherwise.

Building E : We must now build $E_i(a, b)$. There are several possibilities. I will use U and V as dummy variables; each could take the value of any state. I will assume the states are labelled with counting numbers, without any loss of generality, and will write a, b for particular values of the state. One simple construction is to build a set of feature functions $\phi_{ab}^{(e)}(U, V) = \mathbb{I}_{[U=a]} \mathbb{I}_{[V=b]}$. There is one feature function for each possible pair of states, so if there are S possible states, there will be S^2 of these feature functions. Each one takes the value 1 when U and V take the corresponding state values, otherwise is 0. Now we construct a vector of S^2 parameters $\theta_{ab}^{(e)}$, and construct $E(U, V) = \sum_{a,b} \theta_{ab}^{(e)} \phi_{ab}^{(e)}(U, V)$. This construction will allow us to represent any possible cost for any transitions, as long as this doesn't depend on the observations.

We can build a set of feature functions that depend on the observed states, which I will write $\phi_{ab}^{(e)}(U, V; Y_i, Y_{i+1})$. The construction would depend on the application. For example, in the case of reading numerals, we could build a collection of classifiers that look at the ink corresponding to a pair of numerals together. Write $g_{j,ab}(Y_i, Y_{i+1})$ for these classifiers. Here j is an index that identifies the particular classifier. We require that $g_{j,ab}(Y_i, Y_{i+1})$ give a small response for examples

where Y_i is the ink for numeral a and Y_{i+1} is the ink for numeral b ; otherwise, the response should be large. Then we construct a vector of parameters $\theta_j^{(e)}$, and construct $E(U, V) = \sum_j \theta_j^{(e)} g_{j,UV}(Y_i, Y_{i+1})$. This yields a function that is small when the ink is consistent with the two states in the argument, and large otherwise.

General notation: We now have a model of the cost. I will write sequences like vectors, so \mathbf{x} is a sequence, and x_i is the i 'th element of that sequence. Write $C(\mathbf{x}; \mathbf{y}, \theta)$ for the cost of a sequence \mathbf{x} of hidden variables, conditioned on observed values \mathbf{y} and parameters θ . I'm suppressing the number of items in this sequence for conciseness, but will use N if I need to represent it. We have

$$C(\mathbf{x}; \mathbf{y}, \theta) = \sum_{i=1}^N \left[\sum_j \theta_j^{(v)} \phi_j^{(v)}(x_i, y_i) \right] + \sum_{i=1}^{N-1} \left[\left(\sum_j \theta_j^{(e)} \phi_j^{(e)}(x_i, x_{i+1}) \right) \right].$$

Notice that this cost function is linear in θ . We will use this to build a search for the best setting of θ .

11.4.2 Setting Up the Learning Problem

I will write $\mathbf{x}^{(i)}$ for the i 'th training sequence of hidden states, and $\mathbf{y}^{(i)}$ for the i 'th training sequence of observations. I will write $x_j^{(i)}$ for the hidden state at step j in the i 'th training sequence, etc. The general principle we will adopt is that we should train a model by choosing θ such that, if we apply inference to $\mathbf{y}^{(i)}$, we will recover $\mathbf{x}^{(i)}$ (or something very similar).

For *any* sequence \mathbf{x} , we would like to have $C(\mathbf{x}^{(i)}; \mathbf{y}^{(i)}, \theta) \leq C(\mathbf{x}; \mathbf{y}^{(i)}, \theta)$. This inequality is much more general than it seems, because it covers any available sequence. Now assume we engage in inference on the model represented by θ , using $\mathbf{y}^{(i)}$ as observed variables. Write $\mathbf{x}^{*,i}$ for the sequence recovered by inference, so that

$$\mathbf{x}^{*,i} = \underset{\mathbf{x}}{\operatorname{argmin}} C(\mathbf{x}; \mathbf{y}^{(i)}, \theta)$$

which means that $\mathbf{x}^{*,i}$ is the sequence recovered from the model by inference if the parameters take the value θ . In turn, we want to have

$$C(\mathbf{x}^{(i)}; \mathbf{y}^{(i)}, \theta) \leq C(\mathbf{x}^{*,i}; \mathbf{y}^{(i)}, \theta).$$

This constraint is likely to be violated in practice. Assume that

$$\xi_i = C(\mathbf{x}^{(i)}; \mathbf{y}^{(i)}, \theta) - C(\mathbf{x}^{*,i}; \mathbf{y}^{(i)}, \theta)$$

so that ξ_i measures the extent to which the constraint is violated. Notice that $\xi_i \geq 0$, because we cannot have $C(\mathbf{x}^{(i)}; \mathbf{y}^{(i)}, \theta)$ smaller than $C(\mathbf{x}^{*,i}; \mathbf{y}^{(i)}, \theta)$ (by definition).

We would like to choose θ so that we have the smallest possible set of constraint violations. It is natural to want to minimize the sum of ξ_i over all training data. But we also want to ensure that θ is not "too large", for the same reasons we regularized

Background: First Tools for Looking at Data

The single most important question for a working scientist — perhaps the single most useful question anyone can ask — is: “what’s going on here?” Answering this question requires creative use of different ways to make pictures of datasets, to summarize them, and to expose whatever structure might be there. This is an activity that is sometimes known as “Descriptive Statistics”. There isn’t any fixed recipe for understanding a dataset, but there is a rich variety of tools we can use to get insights.

12.1 DATASETS

A dataset is a collection of descriptions of different instances of the same phenomenon. These descriptions could take a variety of forms, but it is important that they are descriptions of the same thing. For example, my grandfather collected the daily rainfall in his garden for many years; we could collect the height of each person in a room; or the number of children in each family on a block; or whether 10 classmates would prefer to be “rich” or “famous”. There could be more than one description recorded for each item. For example, when he recorded the contents of the rain gauge each morning, my grandfather could have recorded (say) the temperature and barometric pressure. As another example, one might record the height, weight, blood pressure and body temperature of every patient visiting a doctor’s office.

The descriptions in a dataset can take a variety of forms. A description could be **categorical**, meaning that each data item can take a small set of prescribed values. For example, we might record whether each of 100 passers-by preferred to be “Rich” or “Famous”. As another example, we could record whether the passers-by are “Male” or “Female”. Categorical data could be **ordinal**, meaning that we can tell whether one data item is larger than another. For example, a dataset giving the number of children in a family for some set of families is categorical, because it uses only non-negative integers, but it is also ordinal, because we can tell whether one family is larger than another.

Some ordinal categorical data appears not to be numerical, but can be assigned a number in a reasonably sensible fashion. For example, many readers will recall being asked by a doctor to rate their pain on a scale of 1 to 10 — a question that is usually relatively easy to answer, but is quite strange when you think about it carefully. As another example, we could ask a set of users to rate the usability of an interface in a range from “very bad” to “very good”, and then record that using -2 for “very bad”, -1 for “bad”, 0 for “neutral”, 1 for “good”, and 2 for “very good”.

Many interesting datasets involve **continuous** variables (like, for example,

height or weight or body temperature) when you could reasonably expect to encounter any value in a particular range. For example, we might have the heights of all people in a particular room; or the rainfall at a particular place for each day of the year; or the number of children in each family on a list.

You should think of a dataset as a collection of d -tuples (a d -tuple is an ordered list of d elements). Tuples differ from vectors, because we can always add and subtract vectors, but we cannot necessarily add or subtract tuples. We will always write N for the number of tuples in the dataset, and d for the number of elements in each tuple. The number of elements will be the same for every tuple, though sometimes we may not know the value of some elements in some tuples (which means we must figure out how to predict their values, which we will do much later).

Index	net worth	Index	Taste score	Index	Taste score
1	100, 360	1	12.3	11	34.9
2	109, 770	2	20.9	12	57.2
3	96, 860	3	39	13	0.7
4	97, 860	4	47.9	14	25.9
5	108, 930	5	5.6	15	54.9
6	124, 330	6	25.9	16	40.9
7	101, 300	7	37.3	17	15.9
8	112, 710	8	21.9	18	6.4
9	106, 740	9	18.1	19	18
10	120, 170	10	21	20	38.9

TABLE 12.1: *On the left, net worths of people you meet in a bar, in US \$; I made this data up, using some information from the US Census. The index column, which tells you which data item is being referred to, is usually not displayed in a table because you can usually assume that the first line is the first item, and so on. On the right, the taste score (I'm not making this up; higher is better) for 20 different cheeses. This data is real (i.e. not made up), and it comes from <http://lib.stat.cmu.edu/DASL/Datafiles/Cheese.html>.*

Each element of a tuple has its own type. Some elements might be categorical. For example, one dataset we shall see several times has entries for Gender; Grade; Age; Race; Urban/Rural; School; Goals; Grades; Sports; Looks; and Money for 478 children, so $d = 11$ and $N = 478$. In this dataset, each entry is categorical data. Clearly, these tuples are not vectors because one cannot add or subtract (say) Gender, or add Age to Grades.

Most of our data will be vectors. We use the same notation for a tuple and for a vector. We write a vector in bold, so \mathbf{x} could represent a vector or a tuple (the context will make it obvious which is intended).

The entire data set is $\{\mathbf{x}\}$. When we need to refer to the i 'th data item, we write \mathbf{x}_i . Assume we have N data items, and we wish to make a new dataset out of them; we write the dataset made out of these items as $\{\mathbf{x}_i\}$ (the i is to suggest you are taking a set of items and making a dataset out of them).

In this chapter, we will work mainly with continuous data. We will see a

variety of methods for plotting and summarizing 1-tuples. We can build these plots from a dataset of d -tuples by extracting the r 'th element of each d -tuple. All through the book, we will see many datasets downloaded from various web sources, because people are so generous about publishing interesting datasets on the web. In the next chapter, we will look at 2-dimensional data, and we look at high dimensional data in chapter 4.

12.2 WHAT'S HAPPENING? - PLOTTING DATA

The very simplest way to present or visualize a dataset is to produce a table. Tables can be helpful, but aren't much use for large datasets, because it is difficult to get any sense of what the data means from a table. As a continuous example, table 12.1 gives a table of the net worth of a set of people you might meet in a bar (I made this data up). You can scan the table and have a rough sense of what is going on; net worths are quite close to \$ 100, 000, and there aren't any very big or very small numbers. This sort of information might be useful, for example, in choosing a bar.

People would like to measure, record, and reason about an extraordinary variety of phenomena. Apparently, one can score the goodness of the flavor of cheese with a number (bigger is better); table 12.1 gives a score for each of thirty cheeses (I did not make up this data, but downloaded it from <http://lib.stat.cmu.edu/DASL/Datafiles/Cheese.html>). You should notice that a few cheeses have very high scores, and most have moderate scores. It's difficult to draw more significant conclusions from the table, though.

Gender	Goal	Gender	Goal
boy	Sports	girl	Sports
boy	Popular	girl	Grades
girl	Popular	boy	Popular
girl	Popular	boy	Popular
girl	Popular	boy	Popular
girl	Popular	girl	Grades
girl	Popular	girl	Sports
girl	Grades	girl	Popular
girl	Sports	girl	Grades
girl	Sports	girl	Sports

TABLE 12.2: *Chase and Dunner (?) collected data on what students thought made other students popular. As part of this effort, they collected information on (a) the gender and (b) the goal of students. This table gives the gender (“boy” or “girl”) and the goal (to make good grades — “Grades”; to be popular — “Popular”; or to be good at sports — “Sports”). The table gives this information for the first 20 of 478 students; the rest can be found at <http://lib.stat.cmu.edu/DASL/Datafiles/PopularKids.html>. This data is clearly categorical, and not ordinal.*

Table 12.2 shows a table for a set of categorical data. Psychologists collected data from students in grades 4-6 in three school districts to understand what factors students thought made other students popular. This fascinating data set can be found at <http://lib.stat.cmu.edu/DASL/Datafiles/PopularKids.html>, and was pre-

pared by Chase and Dunner (?). Among other things, for each student they asked whether the student's goal was to make good grades ("Grades", for short); to be popular ("Popular"); or to be good at sports ("Sports"). They have this information for 478 students, so a table would be very hard to read. Table 12.2 shows the gender and the goal for the first 20 students in this group. It's rather harder to draw any serious conclusion from this data, because the full table would be so big. We need a more effective tool than eyeballing the table.

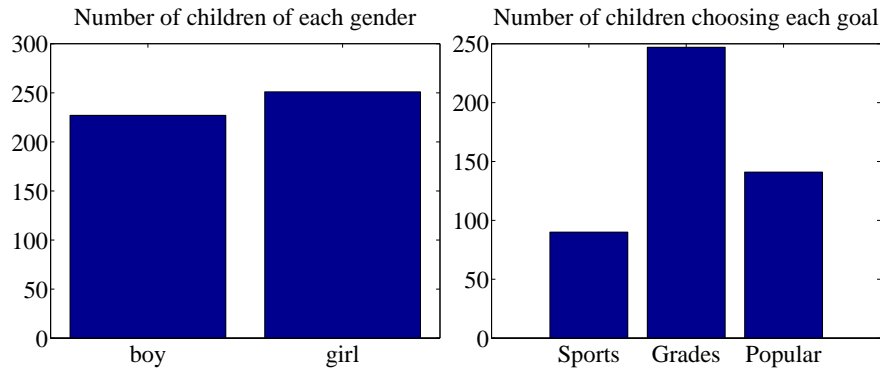


FIGURE 12.1: *On the left*, a bar chart of the number of children of each gender in the Chase and Dunner study (). Notice that there are about the same number of boys and girls (the bars are about the same height). *On the right*, a bar chart of the number of children selecting each of three goals. You can tell, at a glance, that different goals are more or less popular by looking at the height of the bars.

12.2.1 Bar Charts

A **bar chart** is a set of bars, one per category, where the height of each bar is proportional to the number of items in that category. A glance at a bar chart often exposes important structure in data, for example, which categories are common, and which are rare. Bar charts are particularly useful for categorical data. Figure 12.1 shows such bar charts for the genders and the goals in the student dataset of Chase and Dunner (). You can see at a glance that there are about as many boys as girls, and that there are more students who think grades are important than students who think sports or popularity is important. You couldn't draw either conclusion from Table 12.2, because I showed only the first 20 items; but a 478 item table is very difficult to read.

12.2.2 Histograms

Data is continuous when a data item could take any value in some range or set of ranges. In turn, this means that we can reasonably expect a continuous dataset contains few or no pairs of items that have *exactly* the same value. Drawing a bar chart in the obvious way — one bar per value — produces a mess of unit height bars, and seldom leads to a good plot. Instead, we would like to have fewer bars, each representing more data items. We need a procedure to decide which data

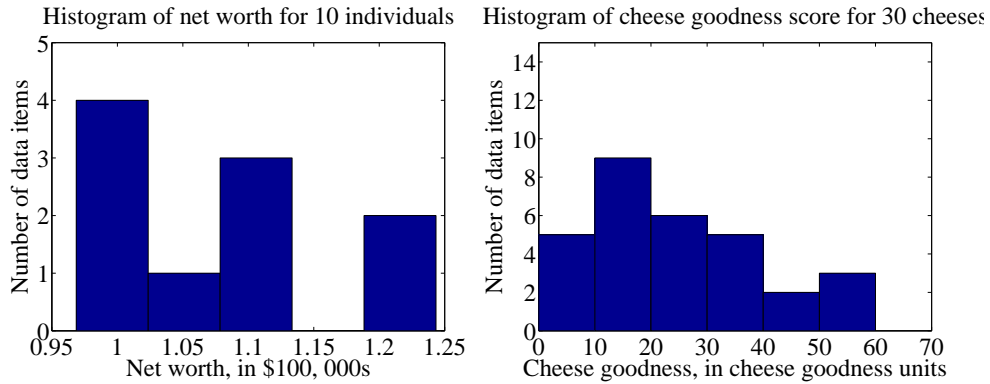


FIGURE 12.2: On the **left**, a histogram of net worths from the dataset described in the text and shown in table 12.1. On the **right**, a histogram of cheese goodness scores from the dataset described in the text and shown in table 12.1.

items count in which bar.

A simple generalization of a bar chart is a **histogram**. We divide the range of the data into intervals, which do not need to be equal in length. We think of each interval as having an associated pigeonhole, and choose one pigeonhole for each data item. We then build a set of boxes, one per interval. Each box sits on its interval on the horizontal axis, and its height is determined by the number of data items in the corresponding pigeonhole. In the simplest histogram, the intervals that form the bases of the boxes are equally sized. In this case, the height of the box is given by the number of data items in the box.

Figure 12.2 shows a histogram of the data in table 12.1. There are five bars — by my choice; I could have plotted ten bars — and the height of each bar gives the number of data items that fall into its interval. For example, there is one net worth in the range between \$102,500 and \$107,500. Notice that one bar is invisible, because there is no data in that range. This picture suggests conclusions consistent with the ones we had from eyeballing the table — the net worths tend to be quite similar, and around \$100,000.

Figure 12.2 also shows a histogram of the data in table 12.1. There are six bars (0-10, 10-20, and so on), and the height of each bar gives the number of data items that fall into its interval — so that, for example, there are 9 cheeses in this dataset whose score is greater than or equal to 10 and less than 20. You can also use the bars to estimate other properties. So, for example, there are 14 cheeses whose score is less than 20, and 3 cheeses with a score of 50 or greater. This picture is much more helpful than the table; you can see at a glance that quite a lot of cheeses have relatively low scores, and few have high scores.

12.2.3 How to Make Histograms

Usually, one makes a histogram by finding the appropriate command or routine in your programming environment. I use Matlab, and chapter ?? sketches some useful Matlab commands. However, it is useful to understand the procedures used.

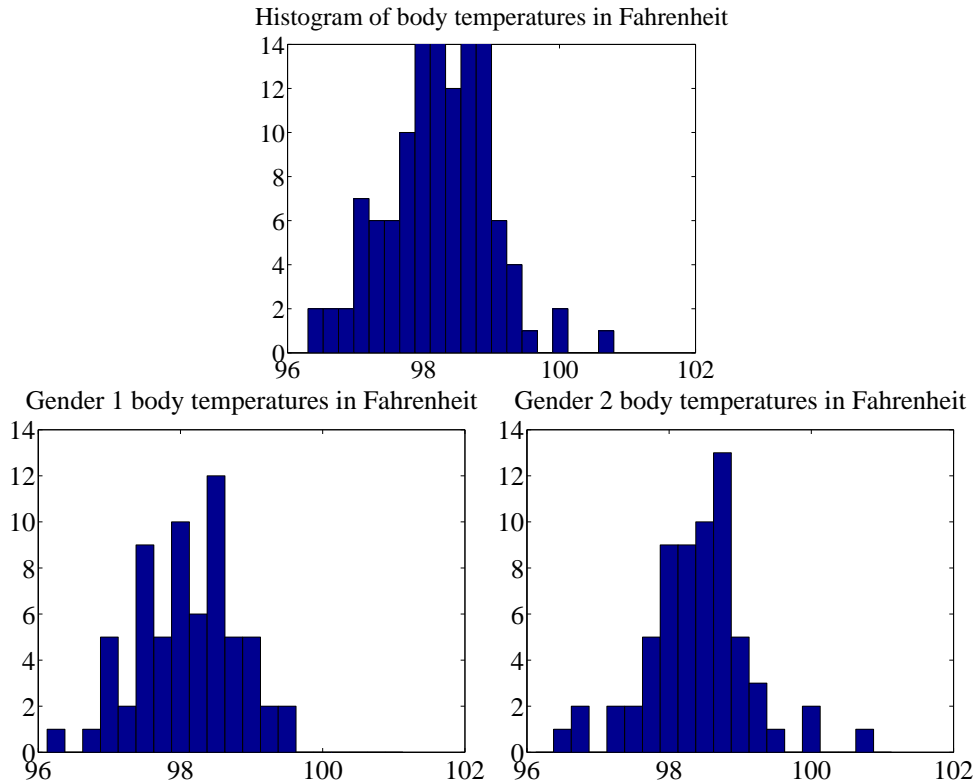


FIGURE 12.3: *On top*, a histogram of body temperatures, from the dataset published at <http://www2.stetson.edu/~jrasp/data.htm>. These seem to be clustered fairly tightly around one value. The **bottom row** shows histograms for each gender (I don't know which is which). It looks as though one gender runs slightly cooler than the other.

Histograms with Even Intervals: The easiest histogram to build uses equally sized intervals. Write x_i for the i 'th number in the dataset, x_{\min} for the smallest value, and x_{\max} for the largest value. We divide the range between the smallest and largest values into n intervals of even width $(x_{\max} - x_{\min})/n$. In this case, the height of each box is given by the number of items in that interval. We could represent the histogram with an n -dimensional vector of counts. Each entry represents the count of the number of data items that lie in that interval. Notice we need to be careful to ensure that each point in the range of values is claimed by exactly one interval. For example, we could have intervals of $[0 - 1)$ and $[1 - 2)$, or we could have intervals of $(0 - 1]$ and $(1 - 2]$. We could *not* have intervals of $[0 - 1]$ and $[1 - 2]$, because then a data item with the value 1 would appear in two boxes. Similarly, we could not have intervals of $(0 - 1)$ and $(1 - 2)$, because then a data item with the value 1 would not appear in any box.

Histograms with Uneven Intervals: For a histogram with even intervals, it is natural that the height of each box is the number of data items in that box.

But a histogram with even intervals can have empty boxes (see figure 12.2). In this case, it can be more informative to have some larger intervals to ensure that each interval has some data items in it. But how high should we plot the box? Imagine taking two consecutive intervals in a histogram with even intervals, and fusing them. It is natural that the height of the fused box should be the average height of the two boxes. This observation gives us a rule.

Write dx for the width of the intervals; n_1 for the height of the box over the first interval (which is the number of elements in the first box); and n_2 for the height of the box over the second interval. The height of the fused box will be $(n_1 + n_2)/2$. Now the *area* of the first box is $n_1 dx$; of the second box is $n_2 dx$; and of the fused box is $(n_1 + n_2) dx$. For each of these boxes, the *area* of the box is proportional to the number of elements in the box. This gives the correct rule: plot boxes such that the area of the box is proportional to the number of elements in the box.

12.2.4 Conditional Histograms

Most people believe that normal body temperature is 98.4° in Fahrenheit. If you take other people's temperatures often (for example, you might have children), you know that some individuals tend to run a little warmer or a little cooler than this number. I found data giving the body temperature of a set of individuals at <http://www2.stetson.edu/~jrasp/data.htm>. As you can see from the histogram (figure 12.3), the body temperatures cluster around a small set of numbers. But what causes the variation?

One possibility is gender. We can investigate this possibility by comparing a histogram of temperatures for males with histogram of temperatures for females. Such histograms are sometimes called **conditional histograms** or **class-conditional histograms**, because each histogram is conditioned on something (in this case, the histogram uses only data that comes from gender).

The dataset gives genders (as 1 or 2 - I don't know which is male and which female). Figure 12.3 gives the class conditional histograms. It does seem like individuals of one gender run a little cooler than individuals of the other, although we don't yet have mechanisms to test this possibility in detail (chapter 16.20).

12.3 SUMMARIZING 1D DATA

For the rest of this chapter, we will assume that data items take values that are continuous real numbers. Furthermore, we will assume that values can be added, subtracted, and multiplied by constants in a meaningful way. Human heights are one example of such data; you can add two heights, and interpret the result as a height (perhaps one person is standing on the head of the other). You can subtract one height from another, and the result is meaningful. You can multiply a height by a constant — say, $1/2$ — and interpret the result (A is half as high as B).

12.3.1 The Mean

One simple and effective summary of a set of data is its **mean**. This is sometimes known as the **average** of the data.

Definition: 12.1 *Mean*

Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . Their mean is

$$\text{mean}(\{x\}) = \frac{1}{N} \sum_{i=1}^{i=N} x_i.$$

For example, assume you're in a bar, in a group of ten people who like to talk about money. They're average people, and their net worth is given in table 12.1 (you can choose who you want to be in this story). The mean of this data is \$107,903.

Properties of the Mean The mean has several important properties you should remember:

- Scaling data scales the mean: or $\text{mean}(\{kx_i\}) = k\text{mean}(\{x_i\})$.
- Translating data translates the mean: or $\text{mean}(\{x_i + c\}) = \text{mean}(\{x_i\}) + c$.
- The sum of signed differences from the mean is zero. This means that

$$\sum_{i=1}^N (x_i - \text{mean}(\{x_i\})) = 0.$$

- Choose the number μ such that the sum of squared distances of data points to μ is minimized. That number is the mean. In notation

$$\arg \min_{\mu} \sum_i (x_i - \mu)^2 = \text{mean}(\{x_i\})$$

These properties are easy to prove (and so easy to remember). I have broken these out into a box of useful facts below, to emphasize them. All but one proof is relegated to the exercises.

Proposition: $\arg \min_{\mu} \sum_i (x_i - \mu)^2 = \text{mean}(\{x\})$

Proof: Choose the number μ such that the sum of squared distances of data points to μ is minimized. That number is the mean. In notation:

$$\arg \min_{\mu} \sum_i (x_i - \mu)^2 = \text{mean}(\{x\})$$

We can show this by actually minimizing the expression. We must have that the derivative of the expression we are minimizing is zero at the value of μ we are seeking. So we have

$$\begin{aligned} \frac{d}{d\mu} \sum_{i=1}^N (x_i - \mu)^2 &= \sum_{i=1}^N 2(x_i - \mu) \\ &= 2 \sum_{i=1}^N (x_i - \mu) \\ &= 0 \end{aligned}$$

so that $2N\text{mean}(\{x\}) - 2N\mu = 0$, which means that $\mu = \text{mean}(\{x\})$.

Property 12.1: The Average Squared Distance to the Mean is Minimized

Now this result means that the mean is the single number that is closest to all the data items. The mean tells you where the overall blob of data lies. For this reason, it is often referred to as a “location parameter”. If you choose to summarize the dataset with a number that is as close as possible to each data item, the mean is the number to choose. The mean is also a guide to what new values will look like, if you have no other information. For example, in the case of the bar, a new person walks in, and I must guess that person’s net worth. Then the mean is the best guess, because it is closest to all the data items we have already seen. In the case of the bar, if a new person walked into this bar, and you had to guess that person’s net worth, you should choose \$107, 903.

Useful Facts: 12.1 *Mean*

- $\text{mean}(\{kx_i\}) = k\text{mean}(\{x_i\})$.
- $\text{mean}(\{x_i + c\}) = \text{mean}(\{x_i\}) + c$.
- $[\sum_{i=1}^N (x_i - \text{mean}(\{x_i\}))] = 0$.

-

$$\arg \min_{\mu} \sum_i (x_i - \mu)^2 = \text{mean}(\{x_i\})$$

The mean is a location parameter; it tells you where the data lies along a number line.

12.3.2 Standard Deviation and Variance

We would also like to know the extent to which data items are close to the mean. This information is given by the **standard deviation**, which is the root mean square of the offsets of data from the mean.

Definition: 12.2 *Standard deviation*

Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . The standard deviation of this dataset is:

$$\text{std}(\{x_i\}) = \sqrt{\frac{1}{N} \sum_{i=1}^{i=N} (x_i - \text{mean}(\{x\}))^2} = \sqrt{\text{mean}(\{(x_i - \text{mean}(\{x\}))^2\})}$$

You should think of the standard deviation as a scale. It measures the size of the average deviation from the mean for a dataset, or how wide the spread of data is. For this reason, it is often referred to as a scale parameter. When the standard deviation of a dataset is large, there are many items with values much larger than, or much smaller than, the mean. When the standard deviation is small, most data items have values close to the mean. This means it is helpful to talk about how many standard deviations away from the mean a particular data item is. Saying that data item x_j is “within k standard deviations from the mean” means that

$$\text{abs}(x_j - \text{mean}(\{x\})) \leq k\text{std}(\{x_i\}).$$

Similarly, saying that data item x_j is “more than k standard deviations from the mean” means that

$$\text{abs}(x_i - \text{mean}(\{x\})) > k\text{std}(\{x\}).$$

As I will show below, there must be some data at least one standard deviation away from the mean, and there can be very few data items that are many standard deviations away from the mean.

Properties of the Standard Deviation Standard deviation has very important properties:

- Translating data does not change the standard deviation, i.e. $\text{std}(\{x_i + c\}) = \text{std}(\{x_i\})$.
- Scaling data scales the standard deviation, i.e. $\text{std}(\{kx_i\}) = k\text{std}(\{x_i\})$.
- For any dataset, there can be only a few items that are many standard deviations away from the mean. In particular, assume we have N data items, x_i , whose standard deviation is σ . Then there are at most $\frac{1}{k^2}$ data points lying k or more standard deviations away from the mean.
- For any dataset, there must be at least one data item that is at least one standard deviation away from the mean.

The first two properties are easy to prove, and are relegated to the exercises. I prove the others below. Again, for emphasis, I have broken these properties out in a box below.

Proposition: Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . Assume the standard deviation of this dataset is $\text{std}(\{x\}) = \sigma$. Then there are at most $\frac{1}{k^2}$ data points lying k or more standard deviations away from the mean.

Proof: Assume the mean is zero. There is no loss of generality here, because translating data translates the mean, but doesn't change the standard deviation. Now we must construct a dataset with the largest possible fraction r of data points lying k or more standard deviations from the mean. To achieve this, our data should have $N(1 - r)$ data points each with the value 0, because these contribute 0 to the standard deviation. It should have Nr data points with the value $k\sigma$; if they are further from zero than this, each will contribute more to the standard deviation, so the fraction of such points will be fewer. Because

$$\text{std}(\{x\}) = \sigma = \sqrt{\frac{\sum_i x_i^2}{N}}$$

we have that, for this rather specially constructed dataset,

$$\sigma = \sqrt{\frac{Nrk^2\sigma^2}{N}}$$

so that

$$r = \frac{1}{k^2}.$$

We constructed the dataset so that r would be as large as possible, so

$$r \leq \frac{1}{k^2}$$

for any kind of data at all.

Property 12.2: For any dataset, it is hard for data items to get many standard deviations away from the mean.

This bound (box 12.2) is true for *any kind of data*. This bound implies that, for example, at most 100% of *any* dataset could be one standard deviation away from the mean, 25% of *any* dataset is 2 standard deviations away from the mean and at most 11% of *any* dataset could be 3 standard deviations away from the mean. But the configuration of data that achieves this bound is very unusual. This means the bound tends to wildly *overstate* how much data is far from the mean for most practical datasets. Most data has more random structure, meaning that we expect to see very much *less* data far from the mean than the bound predicts. For example, much data can reasonably be modelled as coming from a normal distribution (a topic we'll go into later). For such data, we expect that about 68% of the data is within one standard deviation of the mean, 95% is within two standard deviations of the mean, and 99.7% is within three standard deviations of the mean, and the percentage of data that is within ten standard deviations of the mean is essentially indistinguishable from 100%. This kind of behavior is quite common; the crucial point about the standard deviation is that you won't see much

data that lies many standard deviations from the mean, because you can't.

Proposition:

$$(\text{std}(\{x\}))^2 \leq \max_i (x_i - \text{mean}(\{x\}))^2.$$

Proof: You can see this by looking at the expression for standard deviation. We have

$$\text{std}(\{x\}) = \sqrt{\frac{1}{N} \sum_{i=1}^{i=N} (x_i - \text{mean}(\{x\}))^2}.$$

Now, this means that

$$N(\text{std}(\{x\}))^2 = \sum_{i=1}^{i=N} (x_i - \text{mean}(\{x\}))^2.$$

But

$$\sum_{i=1}^{i=N} (x_i - \text{mean}(\{x\}))^2 \leq N \max_i (x_i - \text{mean}(\{x\}))^2$$

so

$$(\text{std}(\{x\}))^2 \leq \max_i (x_i - \text{mean}(\{x\}))^2.$$

Property 12.3: For any dataset, there must be at least one data item that is at least one standard deviation away from the mean.

Boxes 12.2 and 12.3 mean that the standard deviation is quite informative. Very little data is many standard deviations away from the mean; similarly, at least some of the data should be one or more standard deviations away from the mean. So the standard deviation tells us how data points are scattered about the mean.

Useful Facts: 12.2 *Standard deviation*

- $\text{std}(\{x + c\}) = \text{std}(\{x\})$.
- $\text{std}(\{kx\}) = k \text{std}(\{x\})$.
- For N data items, x_i , whose standard deviation is σ , there are at most $\frac{1}{k^2}$ data points lying k or more standard deviations away from the mean.
- $(\text{std}(\{x\}))^2 \leq \max_i (x_i - \text{mean}(\{x\}))^2$.

The standard deviation is often referred to as a scale parameter; it tells you how broadly the data spreads about the mean.

There is an ambiguity that comes up often here because two (very slightly) different numbers are called the standard deviation of a dataset. One — the one we use in this chapter — is an estimate of the scale of the data, as we describe it. The other differs from our expression very slightly; one computes

$$\sqrt{\frac{\sum_i (x_i - \text{mean}(\{x\}))^2}{N - 1}}$$

(notice the $N - 1$ for our N). If N is large, this number is basically the same as the number we compute, but for smaller N there is a difference that can be significant. Irritatingly, this number is also called the standard deviation; even more irritatingly, we will have to deal with it, but not yet. I mention it now because you may look up terms I have used, find this definition, and wonder whether I know what I'm talking about. In this case, I do (although I would say that).

The confusion arises because sometimes the datasets we see are actually samples of larger datasets. For example, in some circumstances you could think of the net worth dataset as a sample of all the net worths in the USA. In such cases, we are often interested in the standard deviation *of the dataset that was sampled*. The second number is a slightly better way to estimate this standard deviation than the definition we have been working with. Don't worry - the N in our expressions is the right thing to use for what we're doing.

12.3.3 Variance

It turns out that thinking in terms of the square of the standard deviation, which is known as the **variance**, will allow us to generalize our summaries to apply to higher dimensional data.

Definition: 12.3 *Variance*

Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . where $N > 1$. Their variance is:

$$\text{var}(\{x\}) = \frac{1}{N} \left(\sum_{i=1}^{i=N} (x_i - \text{mean}(\{x\}))^2 \right) = \text{mean}(\{(x_i - \text{mean}(\{x\}))^2\}).$$

One good way to think of the variance is as the mean-square error you would incur if you replaced each data item with the mean. Another is that it is the square of the standard deviation.

Properties of the Variance The properties of the variance follow from

the fact that it is the square of the standard deviation. I have broken these out in a box, for emphasis.

Useful Facts: 12.3 *Variance*

- $\text{var}(\{x + c\}) = \text{var}(\{x\})$.
- $\text{var}(\{kx\}) = k^2 \text{var}(\{x\})$.

While one could restate the other two properties of the standard deviation in terms of the variance, it isn't really natural to do so. The standard deviation is in the same units as the original data, and should be thought of as a scale. Because the variance is the square of the standard deviation, it isn't a natural scale (unless you take its square root!).

12.3.4 The Median

One problem with the mean is that it can be affected strongly by extreme values. Go back to the bar example, of section 12.3.1. Now Warren Buffett (or Bill Gates, or your favorite billionaire) walks in. What happened to the average net worth?

Assume your billionaire has net worth \$ 1, 000, 000, 000. Then the mean net worth suddenly has become

$$\frac{10 \times \$107,903 + \$1,000,000,000}{11} = \$91,007,184$$

But this mean isn't a very helpful summary of the people in the bar. It is probably more useful to think of the net worth data as ten people together with one billionaire. The billionaire is known as an **outlier**.

One way to get outliers is that a small number of data items are very different, due to minor effects you don't want to model. Another is that the data was misrecorded, or mistranscribed. Another possibility is that there is just too much variation in the data to summarize it well. For example, a small number of extremely wealthy people could change the average net worth of US residents dramatically, as the example shows. An alternative to using a mean is to use a **median**.

Definition: 12.4 *Median*

The median of a set of data points is obtained by sorting the data points, and finding the point halfway along the list. If the list is of even length, it's usual to average the two numbers on either side of the middle. We write

$$\text{median}(\{x_i\})$$

for the operator that returns the median.

For example,

$$\text{median}(\{3, 5, 7\}) = 5,$$

$$\text{median}(\{3, 4, 5, 6, 7\}) = 5,$$

and

$$\text{median}(\{3, 4, 5, 6\}) = 4.5.$$

For much, but not all, data, you can expect that roughly half the data is smaller than the median, and roughly half is larger than the median. Sometimes this property fails. For example,

$$\text{median}(\{1, 2, 2, 2, 2, 2, 2, 3\}) = 2.$$

With this definition, the median of our list of net worths is \$107,835. If we insert the billionaire, the median becomes \$108,930. Notice by how little the number has changed — it remains an effective summary of the data.

Properties of the median You can think of the median of a dataset as giving the “middle” or “center” value. It is another way of estimating where the dataset lies on a number line (and so is another location parameter). This means it is rather like the mean, which also gives a (slightly differently defined) “middle” or “center” value. The mean has the important properties that if you translate the dataset, the mean translates, and if you scale the dataset, the mean scales. The median has these properties, too, which I have broken out in a box. Each is easily proved, and proofs are relegated to the exercises.

Useful Facts: 12.4 *Median*

- $\text{median}(\{x + c\}) = \text{median}(\{x\}) + c.$
- $\text{median}(\{kx\}) = k\text{median}(\{x\}).$

12.3.5 Interquartile Range

Outliers are a nuisance in all sorts of ways. Plotting the histogram of the net worth data with the billionaire included will be tricky. Either you leave the billionaire out of the plot, or all the histogram bars are tiny. Visualizing this plot shows outliers can affect standard deviations severely, too. For our net worth data, the standard deviation without the billionaire is \$9265, but if we put the billionaire in there, it is $\$3.014 \times 10^8$. When the billionaire is in the dataset, the mean is about 91M\$ and the standard deviation is about 300M\$ — so all but one of the data items lie about a third of a standard deviation away from the mean on the small side. The other data item (the billionaire) is about three standard deviations away from the mean on the large side. In this case, the standard deviation has done its work of informing us that there are huge changes in the data, but isn't really helpful as a description of the data.

The problem is this: describing the net worth data with billionaire as a having a mean of $\$9.101 \times 10^7$ with a standard deviation of $\$3.014 \times 10^8$ really isn't terribly helpful. Instead, the data really should be seen as a clump of values that are near \$100,000 and moderately close to one another, and one massive number (the billionaire outlier).

One thing we could do is simply remove the billionaire and compute mean and standard deviation. This isn't always easy to do, because it's often less obvious which points are outliers. An alternative is to follow the strategy we did when we used the median. Find a summary that describes scale, but is less affected by outliers than the standard deviation. This is the **interquartile range**; to define it, we need to define percentiles and quartiles, which are useful anyway.

Definition: 12.5 *Percentile*

The k 'th percentile is the value such that $k\%$ of the data is less than or equal to that value. We write $\text{percentile}(\{x\}, k)$ for the k 'th percentile of dataset $\{x\}$.

Definition: 12.6 *Quartiles*

The first quartile of the data is the value such that 25% of the data is less than or equal to that value (i.e. $\text{percentile}(\{x\}, 25)$). The second quartile of the data is the value such that 50% of the data is less than or equal to that value, which is usually the median (i.e. $\text{percentile}(\{x\}, 50)$). The third quartile of the data is the value such that 75% of the data is less than or equal to that value (i.e. $\text{percentile}(\{x\}, 75)$).

Definition: 12.7 *Interquartile Range*

The interquartile range of a dataset $\{x\}$ is $\text{iqr}\{x\} = \text{percentile}(\{x\}, 75) - \text{percentile}(\{x\}, 25)$.

Like the standard deviation, the interquartile range gives an estimate of how widely the data is spread out. But it is quite well-behaved in the presence of outliers. For our net worth data without the billionaire, the interquartile range is \$12350; with the billionaire, it is \$17710.

Properties of the interquartile range You can think of the interquartile range of a dataset as giving an estimate of the scale of the difference from the mean. This means it is rather like the standard deviation, which also gives a (slightly differently defined) scale. The standard deviation has the important properties that if you translate the dataset, the standard deviation translates, and if you scale the dataset, the standard deviation scales. The interquartile range has these properties, too, which I have broken out into a box. Each is easily proved, and proofs are relegated to the exercises.

Useful Facts: 12.5 *Interquartile range*

- $\text{iqr}\{x + c\} = \text{iqr}\{x\}$.
- $\text{iqr}\{kx\} = k^2 \text{iqr}\{x\}$.

For most datasets, interquartile ranges tend to be somewhat larger than standard deviations. This isn't really a problem. Each is a method for estimating the scale of the data — the range of values above and below the mean that you are likely to see. It is neither here nor there if one method yields slightly larger estimates than another, as long as you don't compare estimates across methods.

12.3.6 Using Summaries Sensibly

One should be careful how one summarizes data. For example, the statement that “the average US family has 2.6 children” invites mockery (the example is from Andrew Vickers' book *What is a p-value anyway?*), because you can't have fractions of a child — no family has 2.6 children. A more accurate way to say things might

be “the average of the number of children in a US family is 2.6”, but this is clumsy. What is going wrong here is the 2.6 is a mean, but the number of children in a family is a categorical variable. Reporting the mean of a categorical variable is often a bad idea, because you may never encounter this value (the 2.6 children). For a categorical variable, giving the median value and perhaps the interquartile range often makes much more sense than reporting the mean.

For continuous variables, reporting the mean is reasonable because you could expect to encounter a data item with this value, even if you haven’t seen one in the particular data set you have. It is sensible to look at both mean and median; if they’re significantly different, then there is probably something going on that is worth understanding. You’d want to plot the data using the methods of the next section before you decided what to report.

You should also be careful about how precisely numbers are reported (equivalently, the number of significant figures). Numerical and statistical software will produce very large numbers of digits freely, but not all are always useful. This is a particular nuisance in the case of the mean, because you might add many numbers, then divide by a large number; in this case, you will get many digits, but some might not be meaningful. For example, Vickers (ibid) describes a paper reporting the mean length of pregnancy as 32.833 weeks. That fifth digit suggests we know the mean length of pregnancy to about 0.001 weeks, or roughly 10 minutes. Neither medical interviewing nor people’s memory for past events is that detailed. Furthermore, when you interview them about embarrassing topics, people quite often lie. There is no prospect of knowing this number with this precision.

People regularly report silly numbers of digits because it is easy to miss the harm caused by doing so. But the harm is there: you are implying to other people, and to yourself, that you know something more accurately than you do. At some point, someone will suffer for it.

12.4 PLOTS AND SUMMARIES

Knowing the mean, standard deviation, median and interquartile range of a dataset gives us some information about what its histogram might look like. In fact, the summaries give us a language in which to describe a variety of characteristic properties of histograms that are worth knowing about (Section 12.4.1). Quite remarkably, many different datasets have histograms that have about the same shape (Section 12.4.2). For such data, we know roughly what percentage of data items are how far from the mean.

Complex datasets can be difficult to interpret with histograms alone, because it is hard to compare many histograms by eye. Section 12.4.3 describes a clever plot of various summaries of datasets that makes it easier to compare many cases.

12.4.1 Some Properties of Histograms

The **tails** of a histogram are the relatively uncommon values that are significantly larger (resp. smaller) than the value at the peak (which is sometimes called the **mode**). A histogram is **unimodal** if there is only one peak; if there are more than one, it is **multimodal**, with the special term **bimodal** sometimes being used for the case where there are two peaks (Figure 12.4). The histograms we have seen

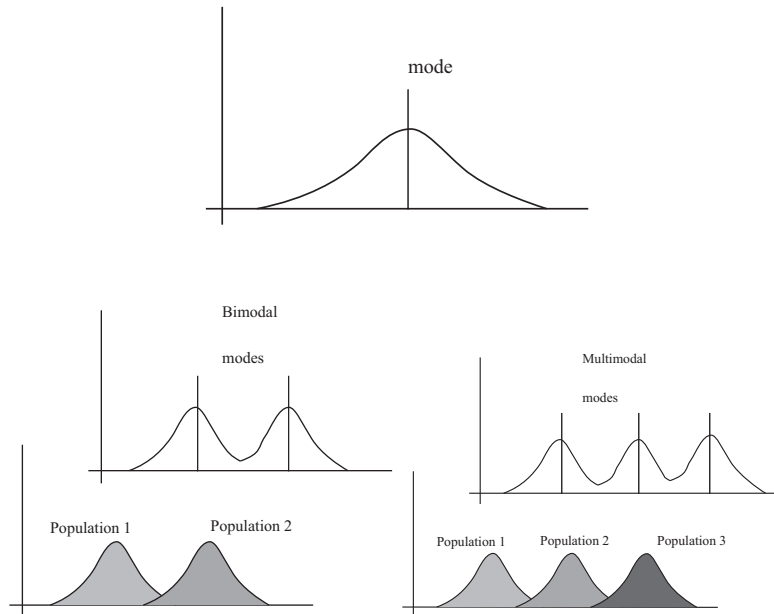


FIGURE 12.4: Many histograms are unimodal, like the example on the **top**; there is one peak, or mode. Some are bimodal (two peaks; **bottom left**) or even multimodal (two or more peaks; **bottom right**). One common reason (but not the only reason) is that there are actually two populations being conflated in the histograms. For example, measuring adult heights might result in a bimodal histogram, if male and female heights were slightly different. As another example, measuring the weight of dogs might result in a multimodal histogram if you did not distinguish between breeds (eg chihuahua, terrier, german shepherd, pyrenean mountain dog, etc.).

have been relatively symmetric, where the left and right tails are about as long as one another. Another way to think about this is that values a lot larger than the mean are about as common as values a lot smaller than the mean. Not all data is symmetric. In some datasets, one or another tail is longer (figure 12.5). This effect is called **skew**.

Skew appears often in real data. SOCR (the Statistics Online Computational Resource) publishes a number of datasets. Here we discuss a dataset of citations to faculty publications. For each of five UCLA faculty members, SOCR collected the number of times each of the papers they had authored had been cited by other authors (data at http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_Dinov_072108_H_Index_Pubs). Generally, a small number of papers get many citations, and many papers get few citations. We see this pattern in the histograms of citation numbers (figure 12.6). These are very different from (say) the body temperature pictures. In the citation histograms, there are many data items that have very few citations, and few that have many citations. This means that the right tail of the histogram is longer, so the histogram is skewed to the right.

One way to check for skewness is to look at the histogram; another is to

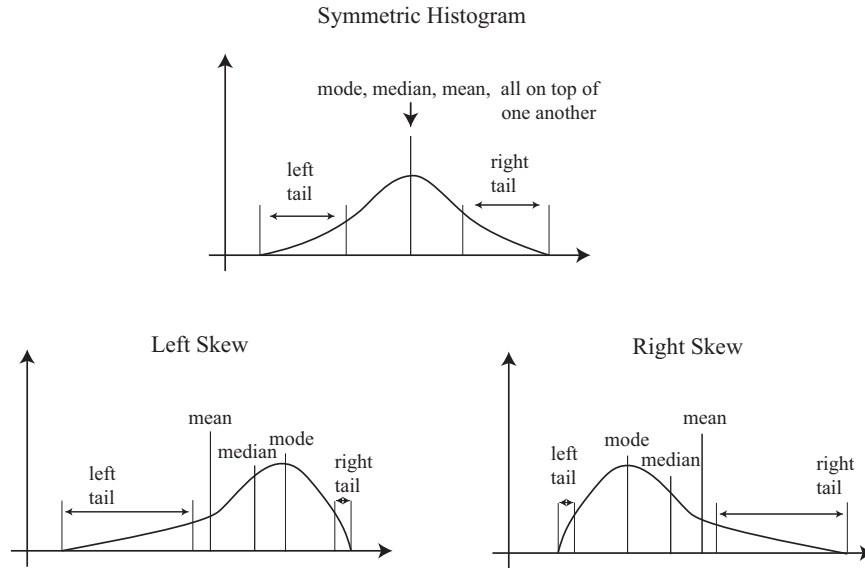


FIGURE 12.5: On the **top**, an example of a symmetric histogram, showing its tails (relatively uncommon values that are significantly larger or smaller than the peak or mode). **Lower left**, a sketch of a left-skewed histogram. Here there are few large values, but some very small values that occur with significant frequency. We say the left tail is “long”, and that the histogram is left skewed. You may find this confusing, because the main bump is to the right — one way to remember this is that the left tail has been stretched. **Lower right**, a sketch of a right-skewed histogram. Here there are few small values, but some very large values that occur with significant frequency. We say the right tail is “long”, and that the histogram is right skewed.

compare mean and median (though this is not foolproof). For the first citation histogram, the mean is 24.7 and the median is 7.5; for the second, the mean is 24.4, and the median is 11. In each case, the mean is a lot bigger than the median. Recall the definition of the median (form a ranked list of the data points, and find the point halfway along the list). For much data, the result is larger than about half of the data set and smaller than about half the dataset. So if the median is quite small compared to the mean, then there are many small data items and a small number of data items that are large — the right tail is longer, so the histogram is skewed to the right.

Left-skewed data also occurs; figure 12.6 shows a histogram of the birth weights of 44 babies born in Brisbane, in 1997 (from http://www.amstat.org/publications/jse/jse_data_archive.htm). This data appears to be somewhat left-skewed, as birth weights can be a lot smaller than the mean, but tend not to be much larger than the mean.

Skewed data is often, but not always, the result of constraints. For example, good obstetrical practice tries to ensure that very large birth weights are rare (birth

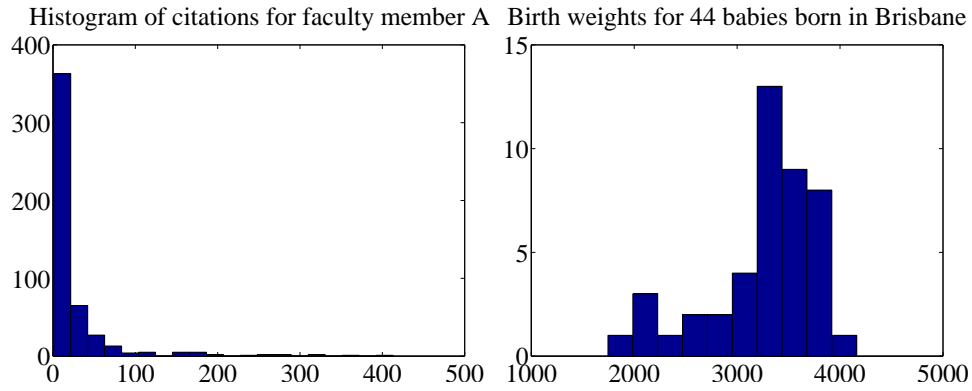


FIGURE 12.6: *On the left*, a histogram of citations for a faculty member, from data at http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_Dinov_072108_H_Index_Pubs. Very few publications have many citations, and many publications have few. This means the histogram is strongly right-skewed. *On the right*, a histogram of birth weights for 44 babies borne in Brisbane in 1997. This histogram looks slightly left-skewed.

is typically induced before the baby gets too heavy), but it may be quite hard to avoid some small birth weights. This could skew birth weights to the left (because large babies will get born, but will not be as heavy as they could be if obstetricians had not interfered). Similarly, income data can be skewed to the right by the fact that income is always positive. Test mark data is often skewed — whether to right or left depends on the circumstances — by the fact that there is a largest possible mark and a smallest possible mark.

12.4.2 Standard Coordinates and Normal Data

It is useful to look at lots of histograms, because it is often possible to get some useful insights about data. However, in their current form, histograms are hard to compare. This is because each is in a different set of units. A histogram for length data will consist of boxes whose horizontal units are, say, metres; a histogram for mass data will consist of boxes whose horizontal units are in, say, kilograms. Furthermore, these histograms typically span different ranges.

We can make histograms comparable by (a) estimating the “location” of the plot on the horizontal axis and (b) estimating the “scale” of the plot. The location is given by the mean, and the scale by the standard deviation. We could then normalize the data by subtracting the location (mean) and dividing by the standard deviation (scale). The resulting values are unitless, and have zero mean. They are often known as **standard coordinates**.

Definition: 12.8 *Standard coordinates*

Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . We represent these data items in standard coordinates by computing

$$\hat{x}_i = \frac{(x_i - \text{mean}(\{x\}))}{\text{std}(\{x\})}.$$

We write $\{\hat{x}\}$ for a dataset that happens to be in standard coordinates.

Standard coordinates have some important properties. Assume we have N data items. Write x_i for the i 'th data item, and \hat{x}_i for the i 'th data item in standard coordinates (I sometimes refer to these as “normalized data items”). Then we have

$$\text{mean}(\{\hat{x}\}) = 0.$$

We also have that

$$\text{std}(\{\hat{x}\}) = 1.$$

An extremely important fact about data is that, for many kinds of data, histograms of these standard coordinates look the same. Many completely different datasets produce a histogram that, in standard coordinates, has a very specific appearance. It is symmetric, unimodal; it looks like a narrow bump. If there were enough data points and the histogram boxes were small enough, the curve would look like the curve in figure 12.7. This phenomenon is so important that data of this form has a special name.

Definition: 12.9 *Standard normal data*

Data is **standard normal data** if, when we have a great deal of data, the histogram is a close approximation to the **standard normal curve**. This curve is given by

$$y(x) = \frac{1}{\sqrt{2\pi}} e^{(-x^2/2)}$$

(which is shown in figure 12.7).

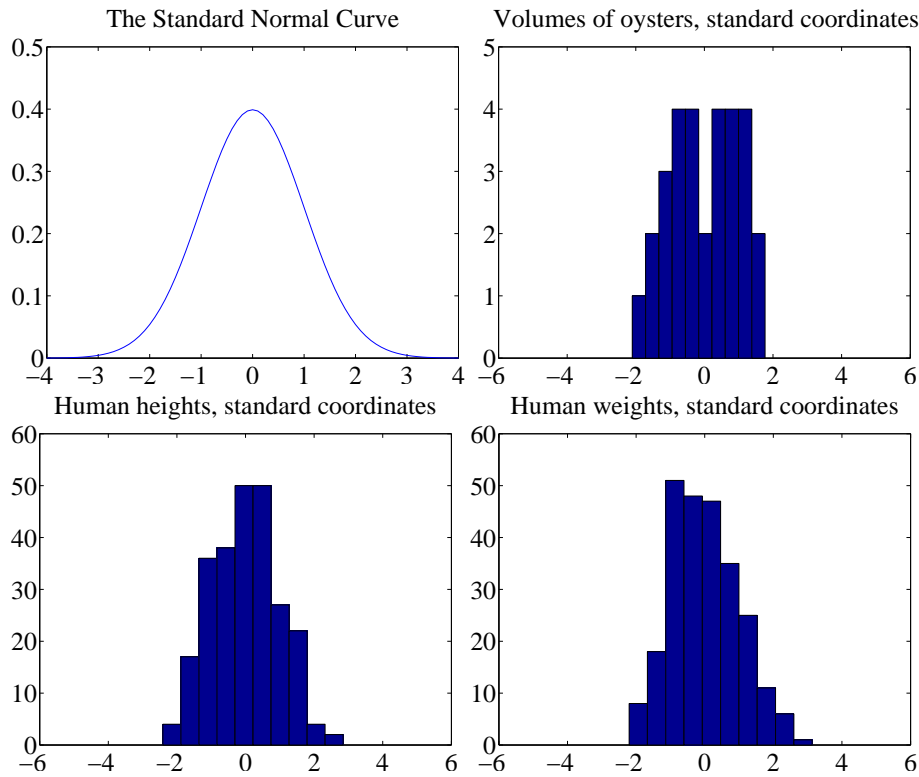


FIGURE 12.7: Data is standard normal data when its histogram takes a stylized, bell-shaped form, plotted above. One usually requires a lot of data and very small histogram boxes for this form to be reproduced closely. Nonetheless, the histogram for normal data is unimodal (has a single bump) and is symmetric; the tails fall off fairly fast, and there are few data items that are many standard deviations from the mean. Many quite different data sets have histograms that are similar to the normal curve; I show three such datasets here.

Definition: 12.10 *Normal data*

Data is **normal data** if, when we subtract the mean and divide by the standard deviation (i.e. compute standard coordinates), it becomes standard normal data.

It is not always easy to tell whether data is normal or not, and there are a variety of tests one can use, which we discuss later. However, there are many examples of normal data. Figure 12.7 shows a diverse variety of data sets, plotted

as histograms in standard coordinates. These include: the volumes of 30 oysters (from http://www.amstat.org/publications/jse/jse_data_archive.htm; look for 30oysters.dat.txt); human heights (from <http://www2.stetson.edu/~jrasp/data.htm>; look for bodyfat.xls, with two outliers removed); and human weights (from <http://www2.stetson.edu/~jrasp/data.htm>; look for bodyfat.xls, with two outliers removed).

Properties of normal data For the moment, assume we know that a dataset is normal. Then we expect it to have the following properties:

- If we normalize it, its histogram will be close to the standard normal curve. This means, among other things, that the data is not significantly skewed.
- About 68% of the data lie within one standard deviation of the mean. We will prove this later.
- About 95% of the data lie within two standard deviations of the mean. We will prove this later.
- About 99% of the data lie within three standard deviations of the mean. We will prove this later.

In turn, these properties imply that data that contains outliers (points many standard deviations away from the mean) is not normal. This is usually a very safe assumption. It is quite common to model a dataset by excluding a small number of outliers, then modelling the remaining data as normal. For example, if I exclude two outliers from the height and weight data from <http://www2.stetson.edu/~jrasp/data.htm>, the data looks pretty close to normal.

12.4.3 Boxplots

It is usually hard to compare multiple histograms by eye. One problem with comparing histograms is the amount of space they take up on a plot, because each histogram involves multiple vertical bars. This means it is hard to plot multiple overlapping histograms cleanly. If you plot each one on a separate figure, you have to handle a large number of separate figures; either you print them too small to see enough detail, or you have to keep flipping over pages.

A **boxplot** is a way to plot data that simplifies comparison. A boxplot displays a dataset as a vertical picture. There is a vertical box whose height corresponds to the interquartile range of the data (the width is just to make the figure easy to interpret). Then there is a horizontal line for the median; and the behavior of the rest of the data is indicated with whiskers and/or outlier markers. This means that each dataset makes is represented by a vertical structure, making it easy to show multiple datasets on one plot *and interpret the plot* (Figure 12.8).

To build a boxplot, we first plot a box that runs from the first to the third quartile. We then show the median with a horizontal line. We then decide which data items should be outliers. A variety of rules are possible; for the plots I show, I used the rule that data items that are larger than $q_3 + 1.5(q_3 - q_1)$ or smaller than $q_1 - 1.5(q_3 - q_1)$, are outliers. This criterion looks for data items that are more than one and a half interquartile ranges above the third quartile, or more than one and a half interquartile ranges below the first quartile.

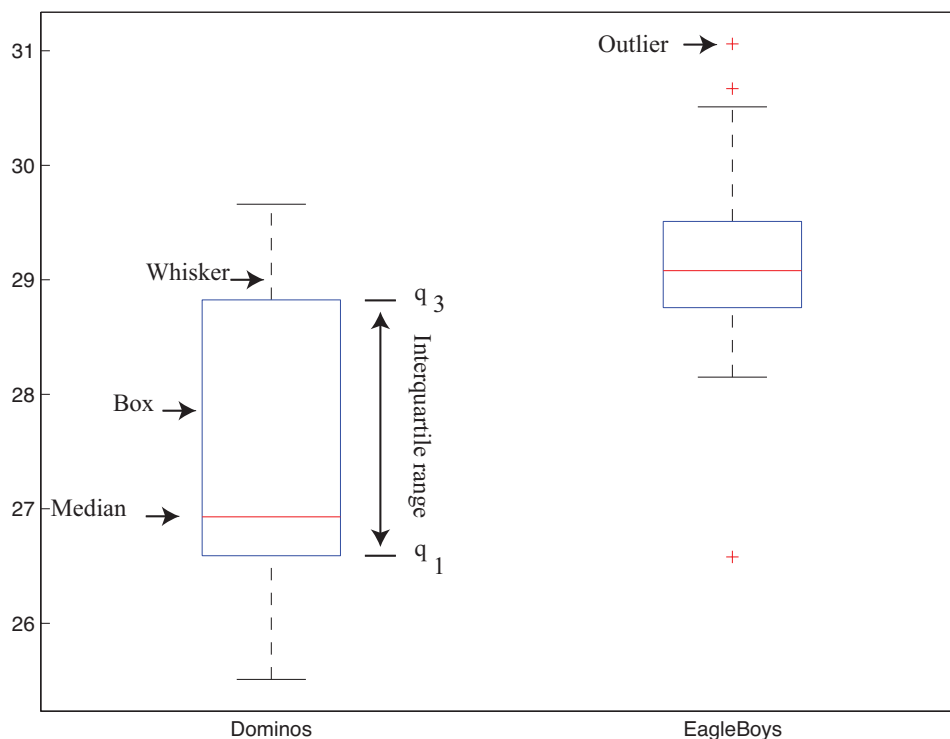


FIGURE 12.8: A boxplot showing the box, the median, the whiskers and two outliers. Notice that we can compare the two datasets rather easily; the next section explains the comparison.

Once we have identified outliers, we plot these with a special symbol (crosses in the plots I show). We then plot whiskers, which show the range of non-outlier data. We draw a whisker from q_1 to the smallest data item that is not an outlier, and from q_3 to the largest data item that is not an outlier. While all this sounds complicated, any reasonable programming environment will have a function that will do it for you. Figure 12.8 shows an example boxplot. Notice that the rich graphical structure means it is quite straightforward to compare two histograms.

12.5 WHOSE IS BIGGER? INVESTIGATING AUSTRALIAN PIZZAS

At http://www.amstat.org/publications/jse/jse_data_archive.htm), you will find a dataset giving the diameter of pizzas, measured in Australia (search for the word “pizza”). This website also gives the backstory for this dataset. Apparently, EagleBoys pizza claims that their pizzas are always bigger than Dominos pizzas, and published a set of measurements to support this claim (the measurements were available at <http://www.eagleboys.com.au/realsizepizza> as of Feb 2012, but seem not to be there anymore).

Whose pizzas are bigger? and why? A histogram of all the pizza sizes appears

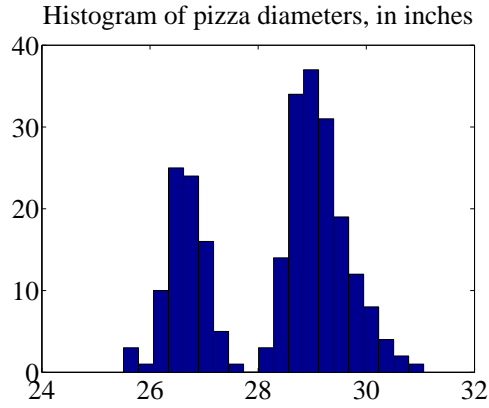


FIGURE 12.9: A histogram of pizza diameters from the dataset described in the text. Notice that there seem to be two populations.

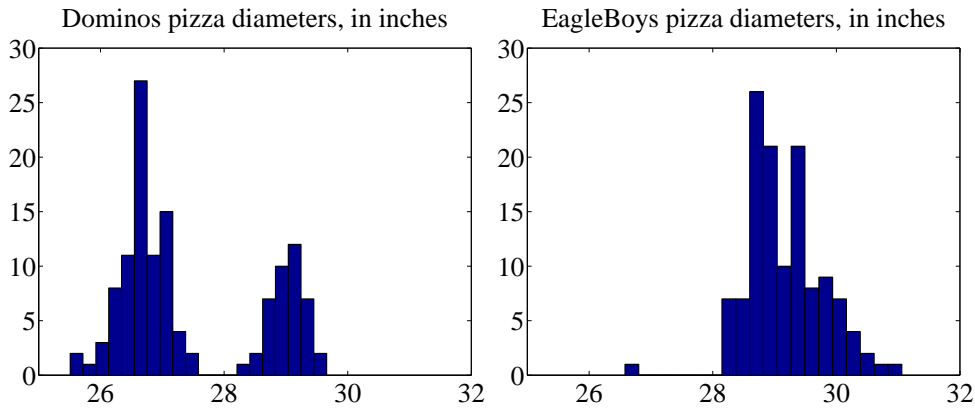


FIGURE 12.10: On the **left**, the class-conditional histogram of Dominos pizza diameters from the pizza data set; on the **right**, the class-conditional histogram of EagleBoys pizza diameters. Notice that EagleBoys pizzas seem to follow the pattern we expect — the diameters are clustered tightly around a mean, and there is a small standard deviation — but Dominos pizzas do not seem to be like that. There is more to understand about this data.

in figure 12.9. We would not expect every pizza produced by a restaurant to have exactly the same diameter, but the diameters are probably pretty close to one another, and pretty close to some standard value. This would suggest that we'd expect to see a histogram which looks like a single, rather narrow, bump about a mean. This is not what we see in figure 12.9 — instead, there are two bumps, which suggests two populations of pizzas. This isn't particularly surprising, because we know that some pizzas come from EagleBoys and some from Dominos.

If you look more closely at the data in the dataset, you will notice that each data item is tagged with the company it comes from. We can now easily plot con-

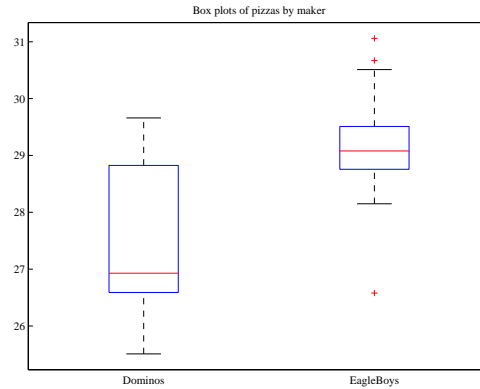


FIGURE 12.11: *Boxplots of the pizza data, comparing EagleBoys and Dominos pizza. There are several curiosities here: why is the range for Dominos so large (25.5-29)? EagleBoys has a smaller range, but has several substantial outliers; why? One would expect pizza manufacturers to try and control diameter fairly closely, because pizzas that are too small present risks (annoying customers; publicity; hostile advertising) and pizzas that are too large should affect profits.*

ditional histograms, conditioning on the company that the pizza came from. These appear in figure 12.10. Notice that EagleBoys pizzas seem to follow the pattern we expect — the diameters are clustered tightly around one value — but Dominos pizzas do not seem to be like that. This is reflected in a boxplot (figure 12.11), which shows the range of Dominos pizza sizes is surprisingly large, and that EagleBoys pizza sizes have several large outliers. There is more to understand about this data. The dataset contains labels for the type of crust and the type of topping — perhaps these properties affect the size of the pizza?

EagleBoys produces DeepPan, MidCrust and ThinCrust pizzas, and Dominos produces DeepPan, ClassicCrust and ThinNCrispy pizzas. This may have something to do with the observed patterns, but comparing six histograms by eye is unattractive. A boxplot is the right way to compare these cases (figure 12.12). The boxplot gives some more insight into the data. Dominos thin crust appear to have a narrow range of diameters (with several outliers), where the median pizza is rather larger than either the deep pan or the classic crust pizza. EagleBoys pizzas all have a range of diameters that is (a) rather similar across the types and (b) rather a lot like the Dominos thin crust. There are outliers, but few for each type.

Another possibility is that the variation in size is explained by the topping. We can compare types and toppings by producing a set of conditional boxplots (i.e. the diameters for each type and each topping). This leads to rather a lot of boxes (figure 12.13), but they’re still easy to compare by eye. The main difficulty is that the labels on the plot have to be shortened. I made labels using the first letter from the manufacturer (“D” or “E”); the first letter from the crust type (previous paragraph); and the first and last letter of the topping. Toppings for Dominos are: Hawaiian; Supreme; BBQMeatlovers. For EagleBoys, toppings are: Hawaiian; SuperSupremo; and BBQMeatlovers. This gives the labels: ‘DCBs’; (Dominos; Clas-

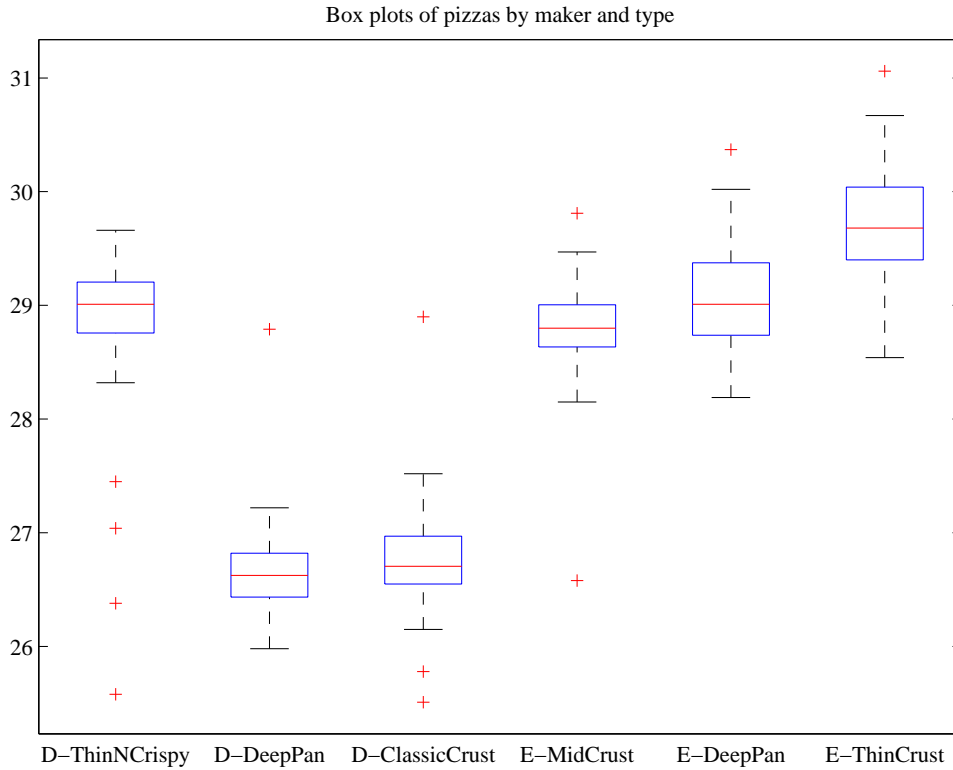


FIGURE 12.12: *Boxplots for the pizza data, broken out by type (thin crust, etc.).*

sicCrust; BBQMeatlovers); 'DCHn'; 'DCSe'; 'DDBs'; 'DDHn'; 'DDSe'; 'DTBs'; 'DTHn'; 'DTSe'; 'EDBs'; 'EDHn'; 'EDSo'; 'EMBs'; 'EMHn'; 'EMSo'; 'ETBs'; 'ETHn'; 'ETSo'. Figure 12.13 suggests that the topping isn't what is important, but the crust (group the boxplots by eye).

What could be going on here? One possible explanation is that EagleBoys have tighter control over the size of the final pizza. One way this could happen is that all EagleBoys pizzas start the same size and shrink the same amount in baking, whereas all Dominos pizzas start a standard diameter, but different Dominos crusts shrink differently in baking. Another way is that Dominos makes different size crusts for different types, but that the cooks sometimes get confused. Yet another possibility is that Dominos controls portions by the mass of dough (so thin crust diameters tend to be larger), but EagleBoys controls by the diameter of the crust.

You should notice that this is more than just a fun story. If you were a manager at a pizza firm, you'd need to make choices about how to control costs. Labor costs, rent, and portion control (i.e. how much pizza, topping, etc. a customer gets for their money) are the main thing to worry about. If the same kind of pizza has a wide range of diameters, you have a problem, because some customers are getting too much (which affects your profit) or too little (which means they might call someone

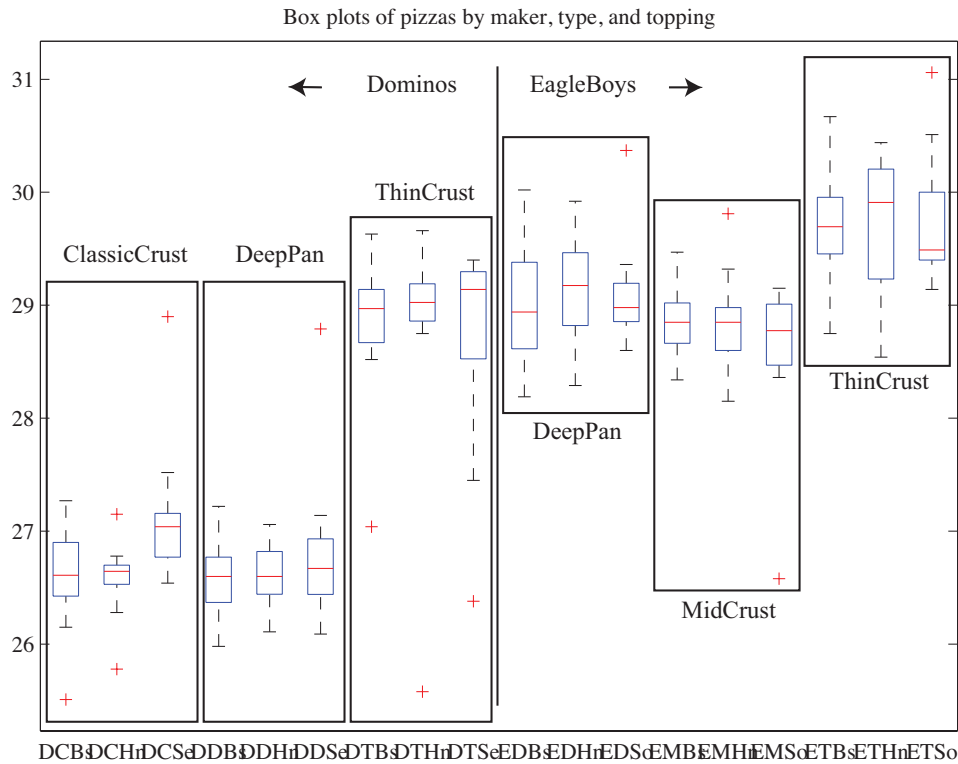


FIGURE 12.13: The pizzas are now broken up by topping as well as crust type (look at the source for the meaning of the names). I have separated Dominos from EagleBoys with a vertical line, and grouped each crust type with a box. It looks as though the issue is not the type of topping, but the crust. EagleBoys seems to have tighter control over the size of the final pizza.

else next time). But making more regular pizzas might require more skilled (and so more expensive) labor. The fact that Dominos and EagleBoys seem to be following different strategies successfully suggests that more than one strategy might work. But you can't choose if you don't know what's happening. As I said at the start, "what's going on here?" is perhaps the single most useful question anyone can ask.

12.6 YOU SHOULD

12.6.1 be able to:

- Plot a bar chart for a dataset.
- Plot a histogram for a dataset.
- Tell whether the histogram is skewed or not, and in which direction.
- Plot a box plot for one or several datasets.
- Interpret a box plot.

12.6.2 remember:

New term: categorical	266
New term: ordinal	266
New term: continuous	266
New term: bar chart	269
New term: histogram	270
New term: conditional histograms	272
New term: class-conditional histograms	272
New term: mean	272
Definition: Mean	273
Useful facts: Mean	275
New term: standard deviation	275
Definition: Standard deviation	275
Useful facts: Standard deviation	278
New term: variance	279
Definition: Variance	279
Useful facts: Variance	280
New term: outlier	280
New term: median	280
Definition: Median	281
Useful facts: Median	281
New term: interquartile range	282
Definition: Percentile	282
Definition: Quartiles	283
Definition: Interquartile Range	283
Useful facts: Interquartile range	283
New term: tails	284
New term: mode	284
New term: unimodal	284
New term: multimodal	284
New term: bimodal	284
New term: skew	285
New term: standard coordinates	287
Definition: Standard coordinates	288
New term: standard normal data	288

New term: standard normal curve	288
Definition: Standard normal data	288
New term: normal data	289
Definition: Normal data	289
New term: boxplot	290

Background: Looking at Relationships

We think of a dataset as a collection of d -tuples (a d -tuple is an ordered list of d elements). For example, the Chase and Dunner dataset had entries for Gender; Grade; Age; Race; Urban/Rural; School; Goals; Grades; Sports; Looks; and Money (so it consisted of 11-tuples). The previous chapter explored methods to visualize and summarize a set of values obtained by extracting a single element from each tuple. For example, I could visualize the heights or the weights of a population (as in Figure 12.7). But I could say nothing about the relationship between the height and weight. In this chapter, we will look at methods to visualize and summarize the relationships between pairs of elements of a dataset.

13.1 PLOTTING 2D DATA

We take a dataset, choose two different entries, and extract the corresponding elements from each tuple. The result is a dataset consisting of 2-tuples, and we think of this as a two dimensional dataset. The first step is to plot this dataset in a way that reveals relationships. The topic of how best to plot data fills many books, and we can only scratch the surface here. Categorical data can be particularly tricky, because there are a variety of choices we can make, and the usefulness of each tends to depend on the dataset and to some extent on one's cleverness in graphic design (section 13.1.1).

For some continuous data, we can plot the one entry as a function of the other (so, for example, our tuples might consist of the date and the number of robberies; or the year and the price of lynx pelts; and so on, section 13.1.2).

Mostly, we use a simple device, called a scatter plot. Using and thinking about scatter plots will reveal a great deal about the relationships between our data items (section 13.1.3).

13.1.1 Categorical Data, Counts, and Charts

Categorical data is a bit special. Assume we have a dataset with several categorical descriptions of each data item. One way to plot this data is to think of it as belonging to a richer set of categories. Assume the dataset has categorical descriptions, which are not ordinal. Then we can construct a new set of categories by looking at each of the cases for each of the descriptions. For example, in the Chase and Dunner data of table 12.2, our new categories would be: “boy-sports”; “girl-sports”; “boy-popular”; “girl-popular”; “boy-grades”; and “girl-grades”. A large set of categories like this can result in a poor bar chart, though, because there may be too many bars to group the bars successfully. Figure 13.1 shows such a bar chart. Notice that it is hard to group categories by eye to compare; for example,

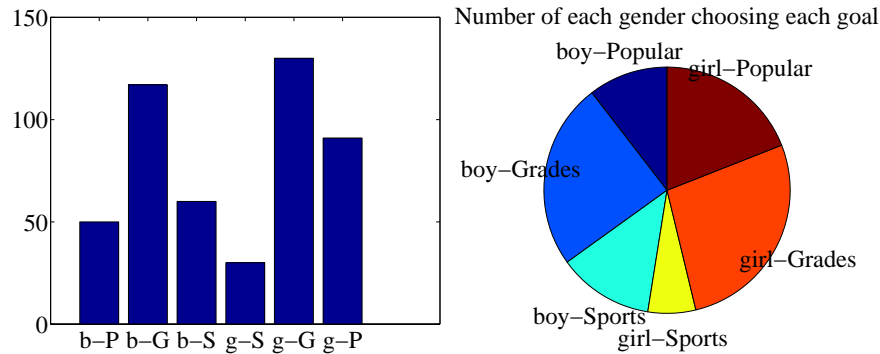


FIGURE 13.1: I sorted the children in the Chase and Dunner study into six categories (two genders by three goals), and counted the number of children that fell into each cell. I then produced the bar chart on the **left**, which shows the number of children of each gender, selecting each goal. On the **right**, a pie chart of this information. I have organized the pie chart so it is easy to compare boys and girls by eye — start at the top; going down on the left side are boy goals, and on the right side are girl goals. Comparing the size of the corresponding wedges allows you to tell what goals (resp. girls) identify with more or less often.

you can see that slightly more girls think grades are important than boys do, but to do so you need to compare two bars that are separated by two other bars. An alternative is a **pie chart**, where a circle is divided into sections whose angle is proportional to the size of the data item. You can think of the circle as a pie, and each section as a slice of pie. Figure 13.1 shows a pie chart, where each section is proportional to the number of students in its category. In this case, I’ve used my judgement to lay the categories out in a way that makes comparisons easy. I’m not aware of any tight algorithm for doing this, though.

Pie charts have problems, because it is hard to judge small differences in area accurately by eye. For example, from the pie chart in figure 13.1, it’s hard to tell that the “boy-sports” category is slightly bigger than the “boy-popular” category (try it; check using the bar chart). For either kind of chart, it is quite important to think about *what* you plot. For example, the plot of figure 13.1 shows the total number of respondents, and if you refer to figure 12.1, you will notice that there are slightly more girls in the study. Is the *percentage* of boys who think grades are important smaller (or larger) than the *percentage* of girls who think so? you can’t tell from these plots, and you’d have to plot the percentages instead.

An alternative is to use a **stacked bar chart**. You can (say) regard the data as of two types, “Boys” and “Girls”. Within those types, there are subtypes (“Popularity”, “Grades” and “Sport”). The height of the bar is given by the number of elements in the type, and the bar is divided into sections corresponding to the number of elements of that subtype. Alternatively, if you want the plot to show relative frequencies, the bars could all be the same height, but the shading corresponds to the fraction of elements of that subtype. This is all much harder to say than to see or to do (Figure 13.2).

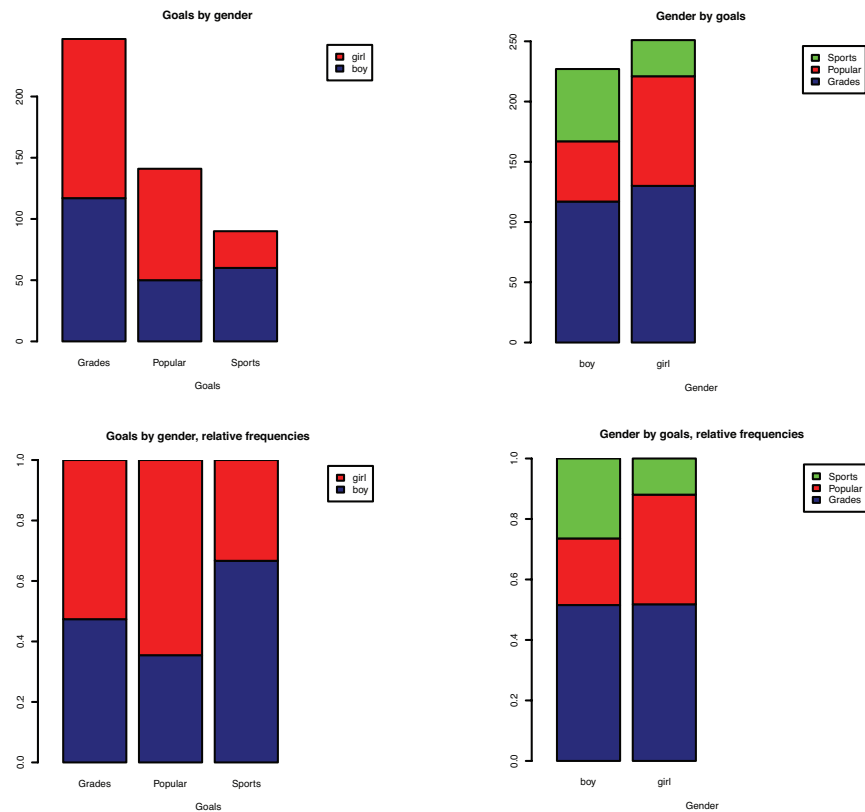


FIGURE 13.2: These bar charts use stacked bars. In the **top row**, the overall height of the bar is given by the number of elements of that type but each different subtype is identified by shading, so you can tell by eye, for example, how many of the “Grades” in the study were “Boys”. This layout makes it hard to tell what fraction of, say, “Boys” aspire to “Popularity”. In the **bottom row**, all bars have the same height, but the shading of the bar identifies the fraction of that type that has a corresponding subtype. This means you can tell by eye what fraction of “Girls” aspire to “Sports”.

An alternative to a pie chart that is very useful for two dimensional data is a **heat map**. This is a method of displaying a matrix as an image. Each entry of the matrix is mapped to a color, and the matrix is represented as an image. For the Chase and Dunner study, I constructed a matrix where each row corresponds to a choice of “sports”, “grades”, or “popular”, and each column corresponds to a choice of “boy” or “girl”. Each entry contains the count of data items of that type. Zero values are represented as white; the largest values as red; and as the value increases, we use an increasingly saturated pink. This plot is shown in figure 13.3

If the categorical data is ordinal, the ordering offers some hints for making a good plot. For example, imagine we are building a user interface. We build an initial version, and collect some users, asking each to rate the interface on scales for

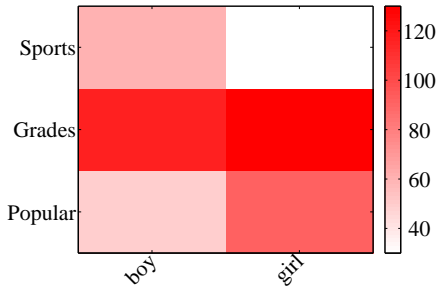


FIGURE 13.3: A heat map of the Chase and Dunner data. The color of each cell corresponds to the count of the number of elements of that type. The colorbar at the side gives the correspondence between color and count. You can see at a glance that the number of boys and girls who prefer grades is about the same; that about the same number of boys prefer sports and popularity, with sports showing a mild lead; and that more girls prefer popularity to sports.

	-2	-1	0	1	2
-2	24	5	0	0	1
-1	6	12	3	0	0
0	2	4	13	6	0
1	0	0	3	13	2
2	0	0	0	1	5

TABLE 13.1: I simulated data representing user evaluations of a user interface. Each cell in the table on the **left** contains the count of users rating “ease of use” (horizontal, on a scale of -2 -very bad- to 2 -very good) vs. “enjoyability” (vertical, same scale). Users who found the interface hard to use did not like using it either. While this data is categorical, it’s also ordinal, so that the order of the cells is determined. It wouldn’t make sense, for example, to reorder the columns of the table or the rows of the table.

“ease of use” (-2, -1, 0, 1, 2, running from bad to good) and “enjoyability” (again, -2, -1, 0, 1, 2, running from bad to good). It is natural to build a 5x5 table, where each cell represents a pair of “ease of use” and “enjoyability” values. We then count the number of users in each cell, and build graphical representations of this table. One natural representation is a **3D bar chart**, where each bar sits on its cell in the 2D table, and the height of the bars is given by the number of elements in the cell. Table 13.1 shows a table and figure 13.4 shows a 3D bar chart for some simulated data. The main difficulty with a 3D bar chart is that some bars are hidden behind others. This is a regular nuisance. You can improve things by using an interactive tool to rotate the chart to get a nice view, but this doesn’t always work. Heatmaps don’t suffer from this problem (Figure 13.4), another reason they are a good choice.

Counts of user responses for a user interface

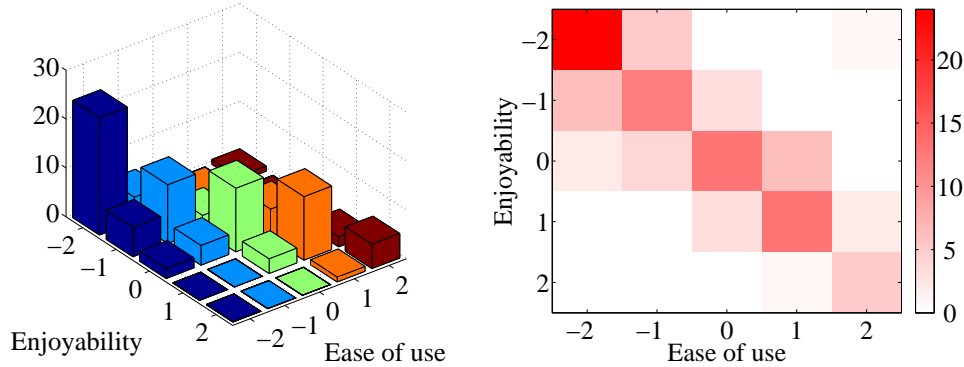


FIGURE 13.4: On the left, a 3D bar chart of the data. The height of each bar is given by the number of users in each cell. This figure immediately reveals that users who found the interface hard to use did not like using it either. However, some of the bars at the back are hidden, so some structure might be hard to infer. On the right, a heat map of this data. Again, this figure immediately reveals that users who found the interface hard to use did not like using it either. It's more apparent that everyone disliked the interface, though, and it's clear that there is no important hidden structure.

Remember this: There are a variety of tools for plotting categorical data. It's difficult to give strict rules for which to use when, but usually one tries to avoid pie charts (angles are hard to judge by eye) and 3D bar charts (where occlusion can hide important effects).

13.1.2 Series

Sometimes one component of a dataset gives a natural ordering to the data. For example, we might have a dataset giving the maximum rainfall for each day of the year. We could record this either by using a two-dimensional representation, where one dimension is the number of the day and the other is the temperature, or with a convention where the i 'th data item is the rainfall on the i 'th day. For example, at <http://lib.stat.cmu.edu/DASL/Datafiles/timeseriesdat.html>, you can find four datasets indexed in this way. It is natural to plot data like this as a function of time. From this dataset, I extracted data giving the number of burglaries each month in a Chicago suburb, Hyde Park. I have plotted part this data in Figure 13.5 (I left out the data to do with treatment effects). It is natural to plot a graph of the burglaries as a function of time (in this case, the number of the month). The plot shows each data point explicitly. I also told the plotting software to draw lines joining data points, because burglaries do not all happen on a specific day.

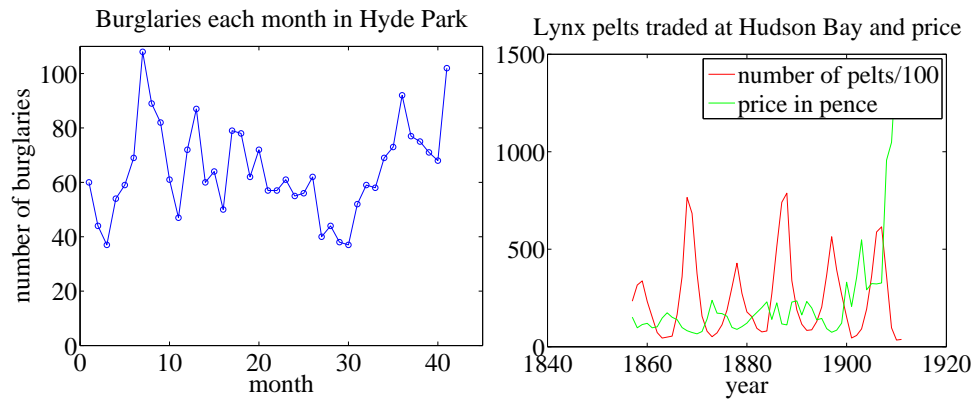


FIGURE 13.5: **Left**, the number of burglaries in Hyde Park, by month. **Right**, a plot of the number of lynx pelts traded at Hudson Bay and of the price paid per pelt, as a function of the year. Notice the scale, and the legend box (the number of pelts is scaled by 100).

The lines suggest, reasonably enough, the rate at which burglaries are happening between data points.

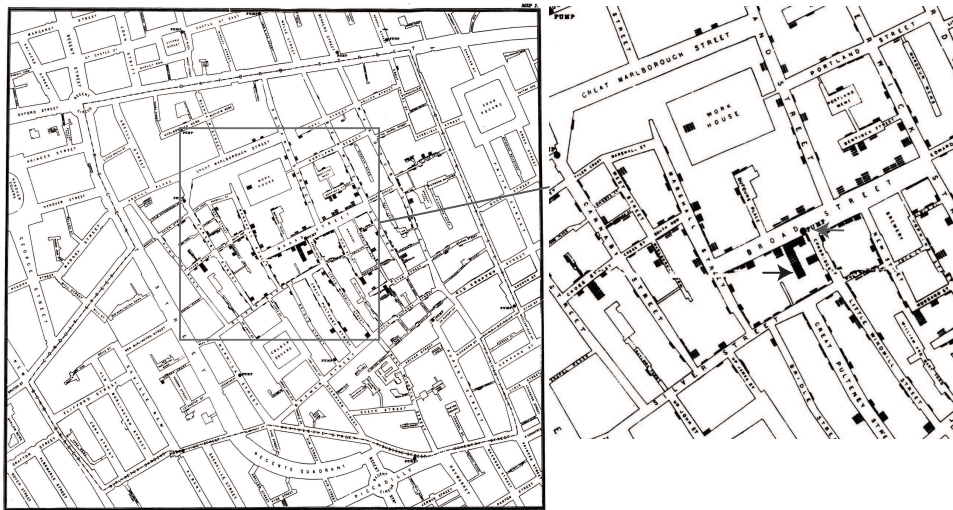


FIGURE 13.6: Snow's scatter plot of cholera deaths on the **left**. Each cholera death is plotted as a small bar on the house in which the bar occurred (for example, the black arrow points to one stack of these bars, indicating many deaths, in the detail on the **right**). Notice the fairly clear pattern of many deaths close to the Broad street pump (grey arrow in the detail), and fewer deaths further away (where it was harder to get water from the pump).

As another example, at <http://lib.stat.cmu.edu/datasets/Andrews/> you can

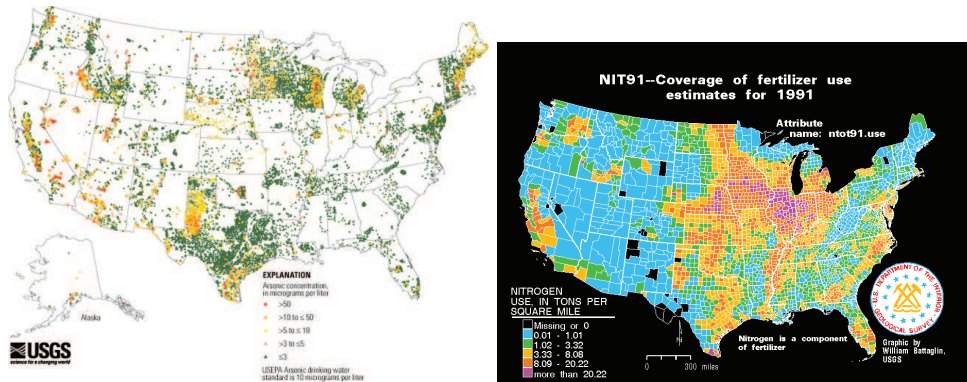


FIGURE 13.7: **Left**, a scatter plot of arsenic levels in US groundwater, prepared by the US Geological Survey (you can find the data at http://water.usgs.gov/GIS/metadata/usgswrd/XML/arsenic_map.xml). Here the shape and color of each marker shows the amount of arsenic, and the spatial distribution of the markers shows where the wells were sampled. **Right**, the usage of Nitrogen (a component of fertilizer) by US county in 1991, prepared by the US Geological Survey (you can find the data at <http://water.usgs.gov/GIS/metadata/usgswrd/XML/nit91.xml>). In this variant of a scatter plot (which usually takes specialized software to prepare) one fills each region with a color indicating the data in that region.

find a dataset that records the number of lynx pelts traded to the Hudson's Bay company and the price paid for each pelt. This version of the dataset appeared first in table 3.2 of *Data: a Collection of Problems from many Fields for the Student and Research Worker* by D.F. Andrews and A.M. Herzberg, published by Springer in 1985. I have plotted it in figure 13.5. The dataset is famous, because it shows a periodic behavior in the number of pelts (which is a good proxy for the number of lynx), which is interpreted as a result of predator-prey interactions. Lynx eat rabbits. When there are many rabbits, lynx kittens thrive, and soon there will be many lynx; but then they eat most of the rabbits, and starve, at which point the rabbit population rockets. You should also notice that after about 1900, prices seem to have gone up rather quickly. I don't know why this is. There is also some suggestion, as there should be, that prices are low when there are many pelts, and high when there are few.

13.1.3 Scatter Plots for Spatial Data

It isn't always natural to plot data as a function. For example, in a dataset containing the temperature and blood pressure of a set of patients, there is no reason to believe that temperature is a function of blood pressure, or the other way round. Two people could have the same temperature, and different blood pressures, or vice-versa. As another example, we could be interested in what causes people to die of cholera. We have data indicating *where* each person died in a particular outbreak. It isn't helpful to try and plot such data as a function.

The **scatter plot** is a powerful way to deal with this situation. In the first

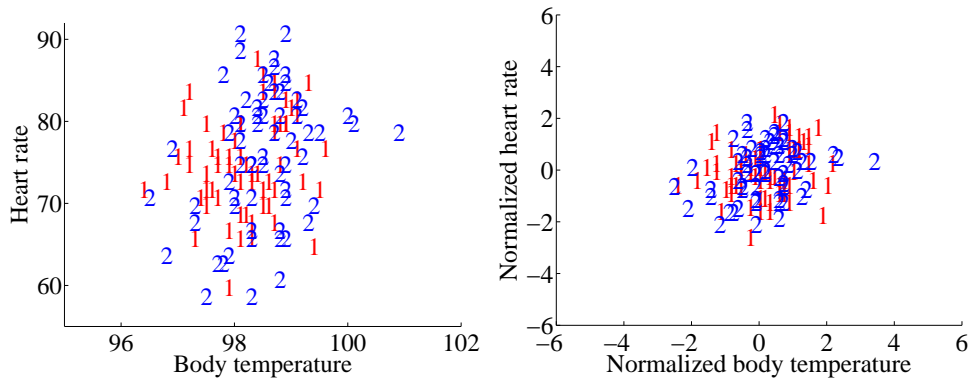


FIGURE 13.8: A scatter plot of body temperature against heart rate, from the dataset at <http://www2.stetson.edu/~jrasp/data.htm>; *normtemp.xls*. I have separated the two genders by plotting a different symbol for each (though I don't know which gender is indicated by which letter); if you view this in color, the differences in color makes for a greater separation of the scatter. This picture suggests, but doesn't conclusively establish, that there isn't much dependence between temperature and heart rate, and any dependence between temperature and heart rate isn't affected by gender.

instance, assume that our data points actually describe points on the a real map. Then, to make a scatter plot, we make a mark on the map at a place indicated by each data point. What the mark looks like, and how we place it, depends on the particular dataset, what we are looking for, how much we are willing to work with complex tools, and our sense of graphic design.

Figure 13.6 is an extremely famous scatter plot, due to John Snow. Snow — one of the founders of epidemiology — used a scatter plot to reason about a cholera outbreak centered on the Broad Street pump in London in 1854. At that time, the mechanism that causes cholera was not known. Snow plotted cholera deaths as little bars (more bars, more deaths) on the location of the house where the death occurred. More bars means more deaths, fewer bars means fewer deaths. There are more bars per block close to the pump, and few far away. This plot offers quite strong evidence of an association between the pump and death from cholera. Snow used this scatter plot as evidence that cholera was associated with water, and that the Broad Street pump was the source of the tainted water.

Figure 13.7 shows a scatter plot of arsenic levels in groundwater for the United States, prepared by the US Geological Survey. The data set was collected by Focazio and others in 2000; by Welch and others in 2000; and then updated by Ryker 2001. It can be found at http://water.usgs.gov/GIS/metadata/usgswrd/XML/arsenic_map.xml. One variant of a scatter plot that is particularly useful for geographic data occurs when one fills regions on a map with different colors, following the data in that region. Figure 13.7 shows the nitrogen usage by US county in 1991; again, this figure was prepared by the US Geological Survey.

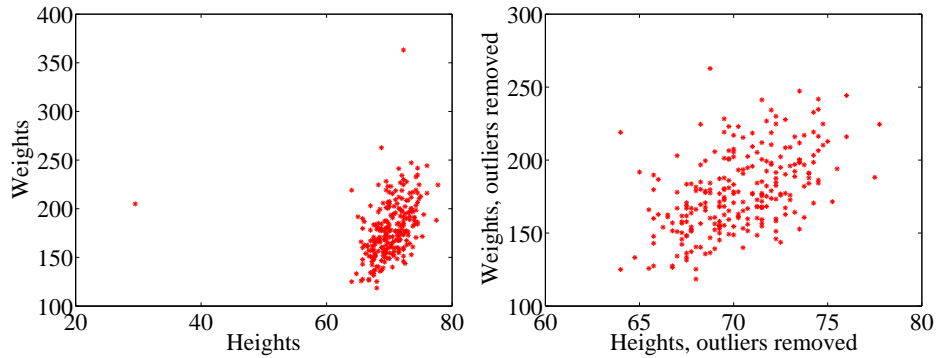


FIGURE 13.9: A scatter plots of weight against height, from the dataset at <http://www2.stetson.edu/~jrasp/data.htm>. **Left:** Notice how two outliers dominate the picture, and to show the outliers, the rest of the data has had to be bunched up. **Right** shows the data with the outliers removed. The structure is now somewhat clearer.

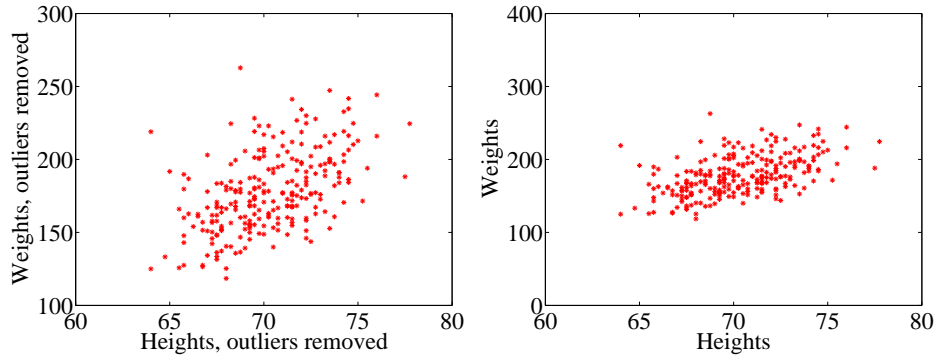


FIGURE 13.10: Scatter plots of weight against height, from the dataset at <http://www2.stetson.edu/~jrasp/data.htm>. **Left:** data with two outliers removed, as in figure 13.9. **Right:** this data, rescaled slightly. Notice how the data looks less spread out. But there is no difference between the datasets. Instead, your eye is easily confused by a change of scale.

Remember this: Scatter plots are a most effective tool for geographic data and 2D data in general. A scatter plot should be your first step with a new 2D dataset.

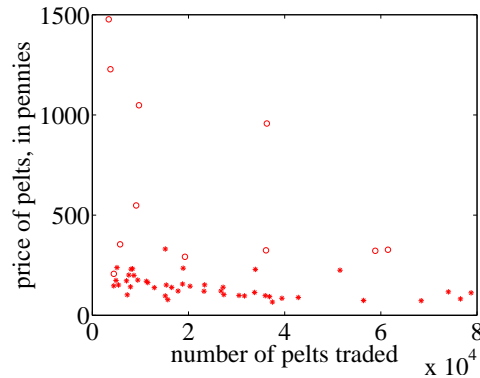


FIGURE 13.11: A scatter plot of the price of lynx pelts against the number of pelts. I have plotted data for 1901 to the end of the series as circles, and the rest of the data as *'s. It is quite hard to draw any conclusion from this data, because the scale is confusing. Furthermore, the data from 1900 on behaves quite differently from the other data.

13.1.4 Exposing Relationships with Scatter Plots

Scatter plots are natural for geographic data, but a scatter plot is a useful, simple tool for ferreting out associations in other kinds of data as well. Now we need some notation. Assume we have a dataset $\{x\}$ of N data items, x_1, \dots, x_N . Each data item is a d dimensional vector (so its components are numbers). We wish to investigate the relationship between two components of the dataset. For example, we might be interested in the 7'th and the 13'th component of the dataset. We will produce a two-dimensional plot, one dimension for each component. It does not really matter which component is plotted on the x -coordinate and which on the y -coordinate (though it will be some pages before this is clear). But it is very difficult to write sensibly without talking about the x and y coordinates.

We will make a two-dimensional dataset out of the components that interest us. We must choose which component goes first in the resulting 2-vector. We will plot this component on the x -coordinate (and we refer to it as the x -coordinate), and to the other component as the y -coordinate. This is just to make it easier to describe what is going on; there's no important idea here. It really will not matter which is x and which is y . The two components make a dataset $\{\mathbf{x}_i\} = \{(x_i, y_i)\}$. To produce a scatter plot of this data, we plot a small shape at the location of each data item.

Such scatter plots are very revealing. For example, figure 13.8 shows a scatter plot of body temperature against heart rate for humans. In this dataset, the gender of the subject was recorded (as "1" or "2" — I don't know which is which), and so I have plotted a "1" at each data point with gender "1", and so on. Looking at the data suggests there isn't much difference between the blob of "1" labels and the blob of "2" labels, which suggests that females and males are about the same in this respect.

The scale used for a scatter plot matters. For example, plotting lengths in

meters gives a very different scatter from plotting lengths in millimeters. Figure 13.9 shows two scatter plots of weight against height. Each plot is from the same dataset, but one is scaled so as to show two outliers. Keeping these outliers means that the rest of the data looks quite concentrated, just because the axes are in large units. In the other plot, the axis scale has changed (so you can't see the outliers), but the data looks more scattered. This may or may not be a misrepresentation. Figure 13.10 compares the data with outliers removed, with the same plot on a somewhat different set of axes. One plot looks as though increasing height corresponds to increasing weight; the other looks as though it doesn't. This is purely due to deceptive scaling — each plot shows the same dataset.

Dubious data can also contribute to scaling problems. Recall that, in figure 13.5, price data before and after 1900 appeared to behave differently. Figure 13.11 shows a scatter plot of the lynx data, where I have plotted number of pelts against price. I plotted the post-1900 data as circles, and the rest as asterisks. Notice how the circles seem to form a quite different figure, which supports the suggestion that something interesting happened around 1900. We can reasonably choose to analyze data after 1900 separately from before 1900. A choice like this should be made with care. If you exclude every data point that might disagree with your hypothesis, you may miss the fact that you are wrong. Leaving out data is an essential component of many kinds of fraud. You should always reveal whether you have excluded data, and why, to allow the reader to judge the evidence.

When you look at Figure 13.11, you should notice the scatter plot does not seem to support the idea that prices go up when supply goes down. This is puzzling because it's generally a pretty reliable idea. In fact, the plot is just hard to interpret because it is poorly scaled. Scale is an important nuisance, and it's easy to get misled by scale effects.

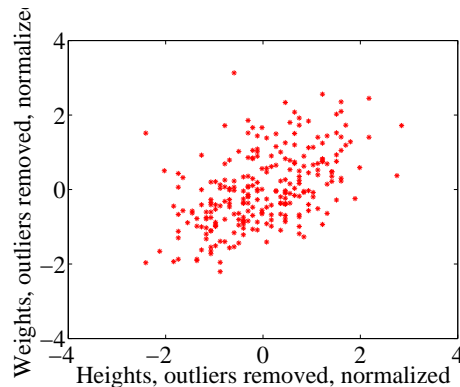


FIGURE 13.12: A normalized scatter plot of weight against height, from the dataset at <http://www2.stetson.edu/~jrasp/data.htm>. Now you can see that someone who is a standard deviation taller than the mean will tend to be somewhat heavier than the mean too.

The way to avoid the problem is to plot in standard coordinates. We can normalize without worrying about the dimension of the data — we normalize each

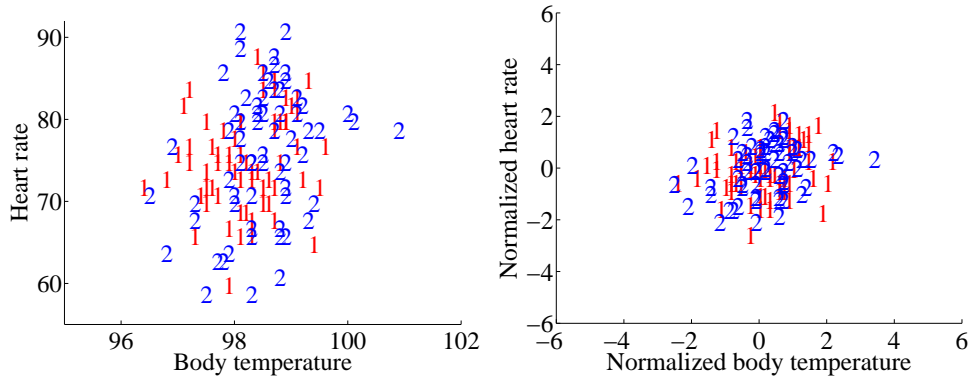


FIGURE 13.13: **Left:** A scatter plot of body temperature against heart rate, from the dataset at <http://www2.stetson.edu/~jrasp/data.htm; normtemp.xls>. I have separated the two genders by plotting a different symbol for each (though I don't know which gender is indicated by which letter); if you view this in color, the differences in color makes for a greater separation of the scatter. This picture suggests, but doesn't conclusively establish, that there isn't much dependence between temperature and heart rate, and any dependence between temperature and heart rate isn't affected by gender. The scatter plot of the normalized data, in standard coordinates, on the right supports this view.

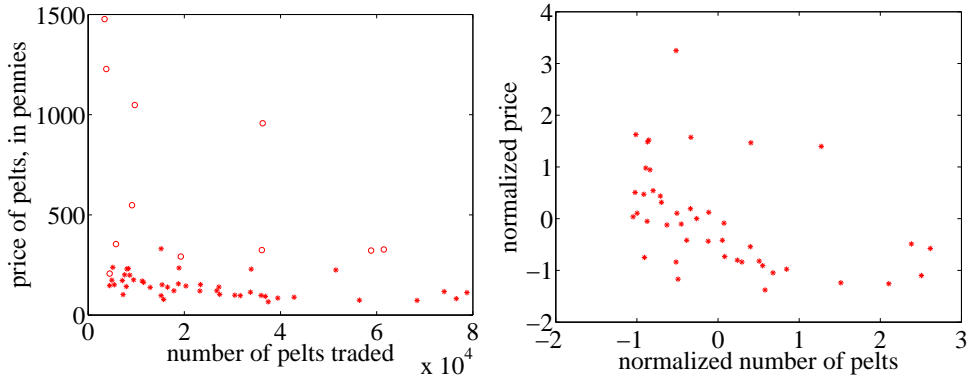


FIGURE 13.14: **Left:** A scatter plot of the price of lynx pelts against the number of pelts (this is a repeat of figure 13.11 for reference). I have plotted data for 1901 to the end of the series as circles, and the rest of the data as *'s. It is quite hard to draw any conclusion from this data, because the scale is confusing. **Right:** A scatter plot of the price of pelts against the number of pelts for lynx pelts. I excluded data for 1901 to the end of the series, and then normalized both price and number of pelts. Notice that there is now a distinct trend; when there are fewer pelts, they are more expensive, and when there are more, they are cheaper.

dimension independently by subtracting the mean of that dimension and dividing by the standard deviation of that dimension. This means we can normalize the x and y coordinates of the two-dimensional data separately. We continue to use the convention of writing the normalized x coordinate as \hat{x} and the normalized y coordinate as \hat{y} . So, for example, we can write $\hat{x}_j = (x_j - \text{mean}(\{x\}))/\text{std}(\{x\})$ for the \hat{x} value of the j 'th data item in normalized coordinates. Normalizing shows us the dataset on a standard scale. Once we have done this, it is quite straightforward to read off simple relationships between variables from a scatter plot.

Remember this: *The plot scale can mask effects in scatter plots, and it's usually a good idea to plot in standard coordinates.*

13.2 CORRELATION

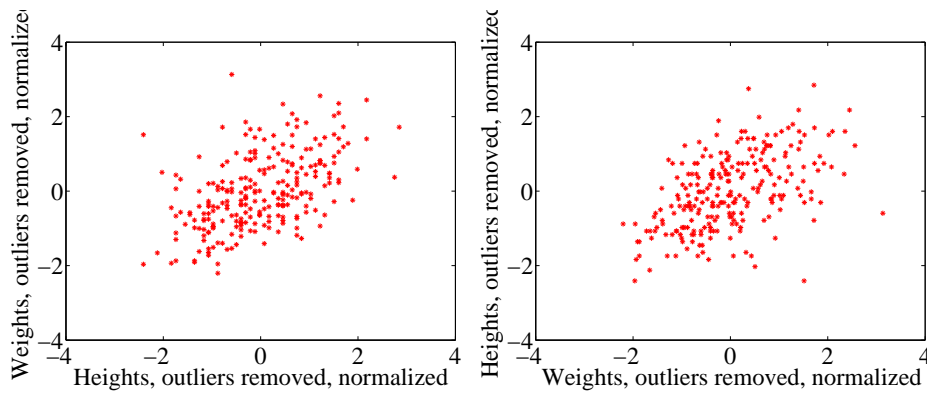


FIGURE 13.15: *On the left, a normalized scatter plot of weight (y -coordinate) against height (x -coordinate). On the right, a scatter plot of height (y -coordinate) against weight (x -coordinate). I've put these plots next to one another so you don't have to mentally rotate (which is what you should usually do).*

The simplest, and most important, relationship to look for in a scatter plot is this: when \hat{x} increases, does \hat{y} tend to increase, decrease, or stay the same? This is straightforward to spot in a normalized scatter plot, because each case produces a very clear shape on the scatter plot. Any relationship is called **correlation** (we will see later how to measure this), and the three cases are: positive correlation, which means that larger \hat{x} values tend to appear with larger \hat{y} values; zero correlation, which means no relationship; and negative correlation, which means that larger \hat{x} values tend to appear with smaller \hat{y} values. I have shown these cases together in one figure using a real data example (Figure 13.16), so you can compare the appearance of the plots.

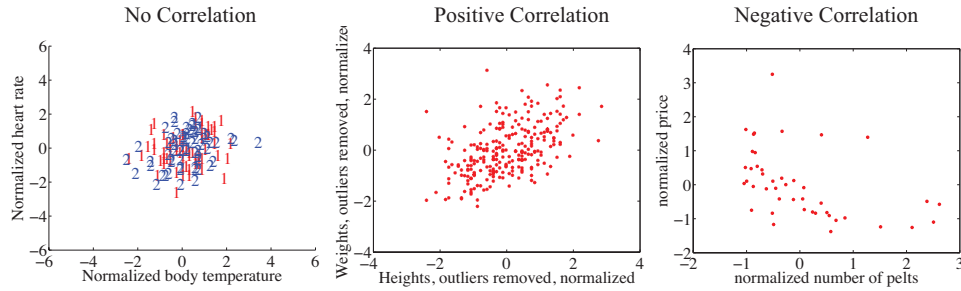


FIGURE 13.16: The three kinds of scatter plot are less clean for real data than for our idealized examples. Here I used the body temperature vs heart rate data for the zero correlation; the height-weight data for positive correlation; and the lynx data for negative correlation. The pictures aren't idealized — real data tends to be messy — but you can still see the basic structures.

Positive correlation occurs when larger \hat{x} values tend to appear with larger \hat{y} values. This means that data points with with small (i.e. negative with large magnitude) \hat{x} values must have small \hat{y} values, otherwise the mean of \hat{x} (resp. \hat{y}) would be too big. In turn, this means that the scatter plot should look like a “smear” of data from the bottom left of the graph to the top right. The smear might be broad or narrow, depending on some details we'll discuss below. Figure 13.12 shows normalized scatter plots of weight against height, and of body temperature against heart rate. In the weight-height plot, you can clearly see that individuals who are higher tend to weigh more. The important word here is “tend” — taller people could be lighter, but mostly they tend not to be. Notice, also, that I did NOT say that they weighed more *because* they were taller, but only that they tend to be heavier.

Negative correlation occurs when larger \hat{x} values tend to appear with smaller \hat{y} values. This means that data points with with small \hat{x} values must have large \hat{y} values, otherwise the mean of \hat{x} (resp. \hat{y}) would be too big. In turn, this means that the scatter plot should look like a “smear” of data from the top left of the graph to the bottom right. The smear might be broad or narrow, depending on some details we'll discuss below. Figure 13.14 shows a normalized scatter plot of the lynx pelt-price data, where I have excluded the data from 1901 on. I did so because there seemed to be some other effect operating to drive prices up, which was inconsistent with the rest of the series. This plot suggests that when there were more pelts, prices were lower, as one would expect.

Zero correlation occurs when there is no relationship. This produces a characteristic shape in a scatter plot, but it takes a moment to understand why. If there really is no relationship, then knowing \hat{x} will tell you nothing about \hat{y} . All we know is that $\text{mean}(\{\hat{y}\}) = 0$, and $\text{var}(\{\hat{y}\}) = 1$. This is enough information to predict what the plot will look like. We know that $\text{mean}(\{\hat{x}\}) = 0$ and $\text{var}(\{\hat{x}\}) = 1$; so there will be many data points with \hat{x} value close to zero, and few with a much larger or much smaller \hat{x} value. The same applies to \hat{y} . Now consider the data points in a strip of \hat{x} values. If this strip is far away from the origin, there will

be few data points in the strip, because there aren't many big \hat{x} values. If there is no relationship, we don't expect to see large or small \hat{y} values in this strip, because there are few data points in the strip and because large or small \hat{y} values are uncommon — we see them only if there are many data points, and then seldom. So for a strip with \hat{x} close to zero, we might see some \hat{y} values that are far from zero because we will see many \hat{y} values. For a strip with \hat{x} that is far from zero, we expect to see few \hat{y} values that are far from zero, because we see few points in this strip. This reasoning means the data should form a round blob, centered at the origin. In the temperature-heart rate plot of figure 13.13, it looks as though nothing of much significance is happening. The average heart rate seems to be about the same for people who run warm or who run cool. There is probably not much relationship here.

The correlation is not affected by which variable is plotted on the x -axis and which is plotted on the y -axis. Figure 13.15 compares a plot of height against weight to one of weight against height. Usually, one just does this by rotating the page, or by imagining the new picture. The left plot tells you that data points with higher height value tend to have higher weight value; the right plot tells you that data points with higher weight value tend to have higher height value — i.e. the plots tell you the same thing. It doesn't really matter which one you look at. Again, the important word is “tend” — the plot doesn't tell you anything about *why*, it just tells you that when one variable is larger the other tends to be, too.

13.2.1 The Correlation Coefficient

Consider a normalized data set of N two-dimensional vectors. We can write the i 'th data point *in standard coordinates* (\hat{x}_i, \hat{y}_i) . We already know many important summaries of this data, because it is in standard coordinates. We have $\text{mean}(\{\hat{x}\}) = 0$; $\text{mean}(\{\hat{y}\}) = 0$; $\text{std}(\{\hat{x}\}) = 1$; and $\text{std}(\{\hat{y}\}) = 1$. Each of these summaries is itself the mean of some monomial. So $\text{std}(\{\hat{x}\})^2 = \text{mean}(\{\hat{x}^2\}) = 1$; $\text{std}(\{\hat{y}\})^2 = \text{mean}(\{\hat{y}^2\})$ (the other two are easy). We can rewrite this information in terms of means of monomials, giving $\text{mean}(\{\hat{x}\}) = 0$; $\text{mean}(\{\hat{y}\}) = 0$; $\text{mean}(\{\hat{x}^2\}) = 1$; and $\text{mean}(\{\hat{y}^2\}) = 1$. There is one monomial missing here, which is $\hat{x}\hat{y}$.

The term $\text{mean}(\{\hat{x}\hat{y}\})$ captures correlation between x and y . The term is known as the **correlation coefficient** or **correlation**.

Definition: 13.1 *Correlation coefficient*

Assume we have N data items which are 2-vectors $(x_1, y_1), \dots, (x_N, y_N)$, where $N > 1$. These could be obtained, for example, by extracting components from larger vectors. We compute the correlation coefficient by first normalizing the x and y coordinates to obtain $\hat{x}_i = \frac{(x_i - \text{mean}(\{x\}))}{\text{std}(x)}$, $\hat{y}_i = \frac{(y_i - \text{mean}(\{y\}))}{\text{std}(y)}$. The correlation coefficient is the mean value of $\hat{x}\hat{y}$, and can be computed as:

$$\text{corr}(\{(x, y)\}) = \frac{\sum_i \hat{x}_i \hat{y}_i}{N}$$

Correlation is a measure of our ability to predict one value from another. The correlation coefficient takes values between -1 and 1 (we'll prove this below). If the correlation coefficient is close to 1 , then we are likely to predict very well. Small correlation coefficients (under about 0.5 , say, but this rather depends on what you are trying to achieve) tend not to be all that interesting, because (as we shall see) they result in rather poor predictions. Figure 13.17 gives a set of scatter plots of different real data sets with different correlation coefficients. These all come from data set of age-height-weight, which you can find at <http://www2.stetson.edu/~jrasp/data.htm> (look for bodyfat.xls). In each case, two outliers have been removed. Age and height are hardly correlated, as you can see from the figure. Younger people do tend to be slightly taller, and so the correlation coefficient is -0.25 . You should interpret this as a small correlation. However, the variable called “adiposity” (which isn't defined, but is presumably some measure of the amount of fatty tissue) is quite strongly correlated with weight, with a correlation coefficient is 0.86 . Average tissue density is quite strongly negatively correlated with adiposity, because muscle is much denser than fat, so these variables are negatively correlated — we expect high density to appear with low adiposity, and vice versa. The correlation coefficient is -0.86 . Finally, density is very strongly correlated with body weight. The correlation coefficient is -0.98 .

It's not always convenient or a good idea to produce scatter plots in standard coordinates (among other things, doing so hides the units of the data, which can be a nuisance). Fortunately, scaling or translating data does not change the value of the correlation coefficient (though it can change the sign if one scale is negative). This means that it's worth being able to spot correlation in a scatter plot that isn't in standard coordinates (even though correlation is always *defined* in standard coordinates). Figure 13.18 shows different correlated datasets plotted in their original units. These data sets are the same as those used in figure 13.17

Properties of the Correlation Coefficient

You should memorize the following properties of the correlation coefficient:

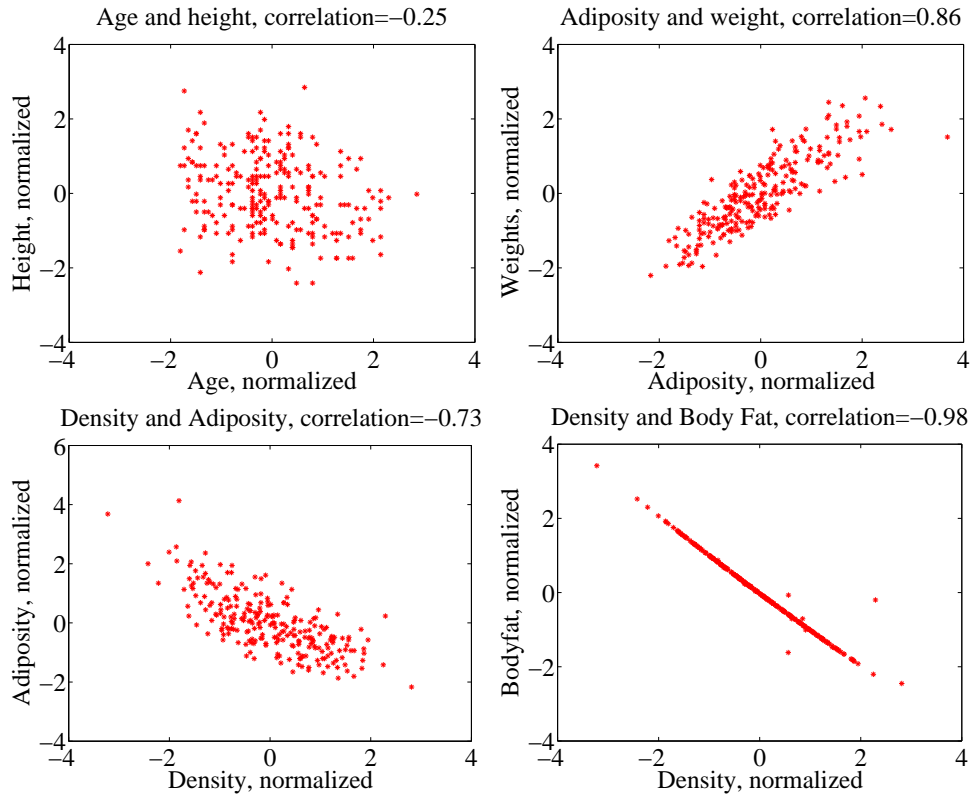


FIGURE 13.17: Scatter plots for various pairs of variables for the age-height-weight dataset from <http://www2.stetson.edu/~jrasp/data.htm>; *bodyfat.xls*. In each case, two outliers have been removed, and the plots are in standard coordinates (compare to figure 13.18, which shows these data sets plotted in their original units). The legend names the variables.

- The correlation coefficient is symmetric (it doesn't depend on the order of its arguments), so

$$\text{corr}(\{(x, y)\}) = \text{corr}(\{(y, x)\})$$

- The value of the correlation coefficient is not changed by translating the data. Scaling the data can change the sign, but not the absolute value. For constants $a \neq 0$, b , $c \neq 0$, d we have

$$\text{corr}(\{(ax + b, cx + d)\}) = \text{sign}(ac)\text{corr}(\{(x, y)\})$$

- If \hat{y} tends to be large (resp. small) for large (resp. small) values of \hat{x} , then the correlation coefficient will be positive.
- If \hat{y} tends to be small (resp. large) for large (resp. small) values of \hat{x} , then the correlation coefficient will be negative.

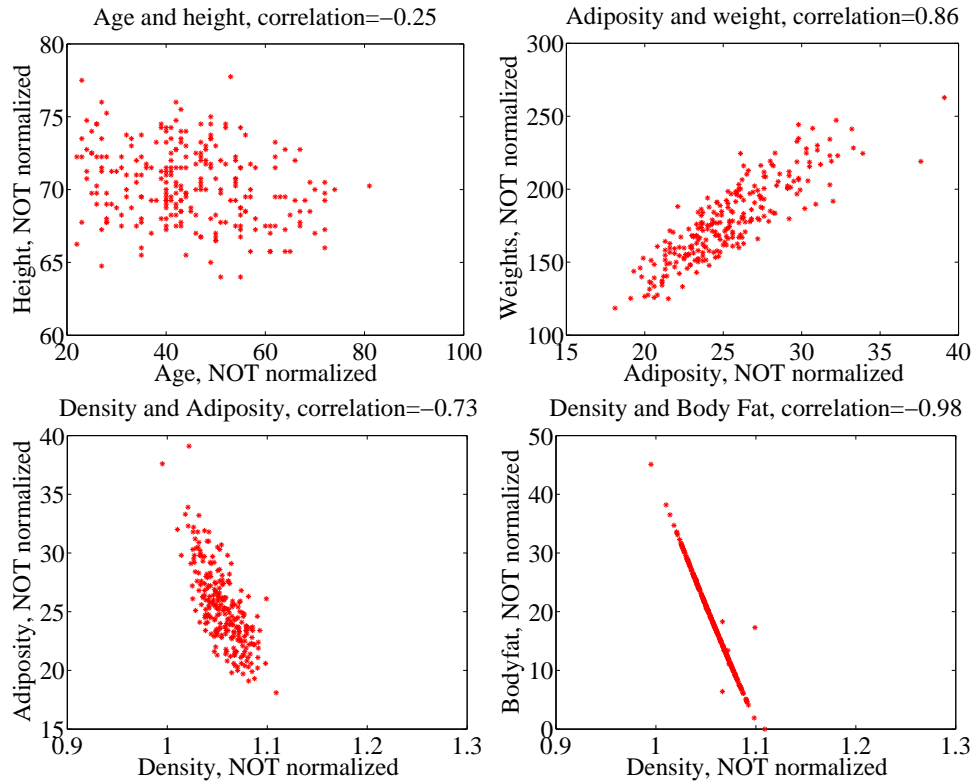


FIGURE 13.18: Scatter plots for various pairs of variables for the age-height-weight dataset from <http://www2.stetson.edu/~jrjasp/data.htm>; *bodyfat.xls*. In each case, two outliers have been removed, and the plots are NOT in standard coordinates (compare to figure 13.17, which shows these data sets plotted in normalized coordinates). The legend names the variables.

- If \hat{y} doesn't depend on \hat{x} , then the correlation coefficient is zero (or close to zero).
- The largest possible value is 1, which happens when $\hat{x} = \hat{y}$.
- The smallest possible value is -1, which happens when $\hat{x} = -\hat{y}$.

The first property is easy, and we relegate that to the exercises. One way to see that the correlation coefficient isn't changed by translation or scale is to notice that it is defined in standard coordinates, and scaling or translating data doesn't change those. Another way to see this is to scale and translate data, then write out the equations; notice that taking standard coordinates removes the effects of the scale and translation. In each case, notice that if the scale is negative, the sign of the correlation coefficient changes.

The property that, if \hat{y} tends to be large (resp. small) for large (resp. small) values of \hat{x} , then the correlation coefficient will be positive, doesn't really admit

a formal statement. But it's relatively straightforward to see what's going on. Because $\text{mean}(\{\hat{x}\}) = 0$, small values of $\text{mean}(\{\hat{x}\})$ must be negative and large values must be positive. But $\text{corr}(\{(x, y)\}) = \frac{\sum_i \hat{x}_i \hat{y}_i}{N}$; and for this sum to be positive, it should contain mostly positive terms. It can contain few or no hugely positive (or hugely negative) terms, because $\text{std}(\hat{x}) = \text{std}(\hat{y}) = 1$ so there aren't many large (or small) numbers. For the sum to contain mostly positive terms, then the sign of \hat{x}_i should be the same as the sign \hat{y}_i for most data items. Small changes to this argument work to show that if \hat{y} tends to be small (resp. large) for large (resp. small) values of \hat{x} , then the correlation coefficient will be negative.

Showing that no relationship means zero correlation requires slightly more work. Divide the scatter plot of the dataset up into thin vertical strips. There are S strips. Each strip is narrow, so the \hat{x} value does not change much for the data points in a particular strip. For the s 'th strip, write $N(s)$ for the number of data points in the strip, $\hat{x}(s)$ for the \hat{x} value at the center of the strip, and $\bar{\hat{y}}(s)$ for the mean of the \hat{y} values within that strip. Now the strips are narrow, so we can approximate all data points within a strip as having the same value of \hat{x} . This yields

$$\text{mean}(\{\hat{x}\hat{y}\}) \approx \frac{1}{S} \sum_{s \in \text{strips}} \hat{x}(s) [N(s)\bar{\hat{y}}(s)]$$

(where you could replace \approx with $=$ if the strips were narrow enough). Now assume that $\bar{\hat{y}}(s)$ does not change from strip to strip, meaning that there is no relationship between \hat{x} and \hat{y} in this dataset (so the picture is like the left hand side in figure 13.16). Then each value of $\bar{\hat{y}}(s)$ is the same — we write $\bar{\hat{y}}$ — and we can rearrange to get

$$\text{mean}(\{\hat{x}\hat{y}\}) \approx \bar{\hat{y}} \frac{1}{S} \sum_{s \in \text{strips}} \hat{x}(s).$$

Now notice that

$$0 = \text{mean}(\{\hat{y}\}) \approx \frac{1}{S} \sum_{s \in \text{strips}} N(s)\bar{\hat{y}}(s)$$

(where again you could replace \approx with $=$ if the strips were narrow enough). This means that if every strip has the same value of $\bar{\hat{y}}(s)$, then that value must be zero. In turn, if there is no relationship between \hat{x} and \hat{y} , we must have $\text{mean}(\{\hat{x}\hat{y}\}) = 0$.

Proposition:

$$-1 \leq \text{corr}(\{(x, y)\}) \leq 1$$

Proof: Writing \hat{x} , \hat{y} for the normalized coefficients, we have

$$\text{corr}(\{(x, y)\}) = \frac{\sum_i \hat{x}_i \hat{y}_i}{N}$$

and you can think of the value as the inner product of two vectors. We write

$$\mathbf{x} = \frac{1}{\sqrt{N}} [\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N] \quad \text{and} \quad \mathbf{y} = \frac{1}{\sqrt{N}} [\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N]$$

and we have $\text{corr}(\{(x, y)\}) = \mathbf{x}^T \mathbf{y}$. Notice $\mathbf{x}^T \mathbf{x} = \text{std}(x)^2 = 1$, and similarly for \mathbf{y} . But the inner product of two vectors is at its maximum when the two vectors are the same, and this maximum is 1. This argument is also sufficient to show that smallest possible value of the correlation is -1 , and this occurs when $\hat{x}_i = -\hat{y}_i$ for all i .

Property 13.1: The largest possible value of the correlation is 1, and this occurs when $\hat{x}_i = \hat{y}_i$ for all i . The smallest possible value of the correlation is -1 , and this occurs when $\hat{x}_i = -\hat{y}_i$ for all i .

13.2.2 Using Correlation to Predict

Assume we have N data items which are 2-vectors $(x_1, y_1), \dots, (x_N, y_N)$, where $N > 1$. These could be obtained, for example, by extracting components from larger vectors. As usual, we will write \hat{x}_i for x_i in normalized coordinates, and so on. Now assume that we know the correlation coefficient is r (this is an important, traditional notation). What does this mean?

One (very useful) interpretation is in terms of prediction. Assume we have a data point $(x_0, ?)$ where we know the x -coordinate, but not the y -coordinate. We can use the correlation coefficient to predict the y -coordinate. First, we transform to standard coordinates. Now we must obtain the best \hat{y}_0 value to predict, using the \hat{x}_0 value we have.

We want to construct a prediction function which gives a prediction for any value of \hat{x} . This predictor should behave as well as possible on our existing data. For each of the (\hat{x}_i, \hat{y}_i) pairs in our data set, the predictor should take \hat{x}_i and produce a result as close to \hat{y}_i as possible. We can choose the predictor by looking at the errors it makes at each data point.

We write \hat{y}_i^p for the value of \hat{y}_i predicted at \hat{x}_i . The simplest form of predictor is linear. If we predict using a linear function, then we have, for some unknown a, b , that $\hat{y}_i^p = a\hat{x}_i + b$. Now think about $u_i = \hat{y}_i - \hat{y}_i^p$, which is the error in our prediction. We would like to have $\text{mean}(\{u\}) = 0$ (otherwise, we could reduce the

error of the prediction just by subtracting a constant).

$$\begin{aligned}
 \text{mean}(\{u\}) &= \text{mean}(\{\hat{y} - \hat{y}^p\}) \\
 &= \text{mean}(\{\hat{y}\}) - \text{mean}(\{a\hat{x}_i + b\}) \\
 &= \text{mean}(\{\hat{y}\}) - a\text{mean}(\{\hat{x}\}) + b \\
 &= 0 - a0 + b \\
 &= 0.
 \end{aligned}$$

This means that we must have $b = 0$.

To estimate a , we need to think about $\text{var}(\{u\})$. We should like $\text{var}(\{u\})$ to be as small as possible, so that the errors are as close to zero as possible (remember, small variance means small standard deviation which means the data is close to the mean). We have

$$\begin{aligned}
 \text{var}(\{u\}) &= \text{var}(\{\hat{y} - \hat{y}^p\}) \\
 &= \text{mean}(\{(\hat{y} - a\hat{x})^2\}) \quad \text{because } \text{mean}(\{u\}) = 0 \\
 &= \text{mean}(\{(\hat{y})^2 - 2a\hat{x}\hat{y} + a^2(\hat{x})^2\}) \\
 &= \text{mean}(\{(\hat{y})^2\}) - 2a\text{mean}(\{\hat{x}\hat{y}\}) + a^2\text{mean}(\{(\hat{x})^2\}) \\
 &= 1 - 2ar + a^2,
 \end{aligned}$$

which we want to minimize by choice of a . At the minimum, we must have

$$\frac{d\text{var}(\{u_i\})}{da} = 0 = -2r + 2a$$

so that $a = r$ and the correct prediction is

$$\hat{y}_0^p = r\hat{x}_0$$

You can use a version of this argument to establish that if we have $(?, \hat{y}_0)$, then the best prediction for \hat{x}_0 (*which is in standard coordinates*) is $r\hat{y}_0$. It is important to notice that the coefficient of \hat{y}_i is NOT $1/r$; you should work this example, which appears in the exercises. We now have a prediction procedure, outlined below.

Procedure: 13.1 *Predicting a value using correlation*

Assume we have N data items which are 2-vectors $(x_1, y_1), \dots, (x_N, y_N)$, where $N > 1$. These could be obtained, for example, by extracting components from larger vectors. Assume we have an x value x_0 for which we want to give the best prediction of a y value, based on this data. The following procedure will produce a prediction:

- Transform the data set into standard coordinates, to get

$$\hat{x}_i = \frac{1}{\text{std}(x)}(x_i - \text{mean}(\{x\}))$$

$$\hat{y}_i = \frac{1}{\text{std}(y)}(y_i - \text{mean}(\{y\}))$$

$$\hat{x}_0 = \frac{1}{\text{std}(x)}(x_0 - \text{mean}(\{x\})).$$

- Compute the correlation

$$r = \text{corr}(\{(x, y)\}) = \text{mean}(\{\hat{x}\hat{y}\}).$$

- Predict $\hat{y}_0 = r\hat{x}_0$.
- Transform this prediction into the original coordinate system, to get

$$y_0 = \text{std}(y)r\hat{x}_0 + \text{mean}(\{y\})$$

Now assume we have a y value y_0 , for which we want to give the best prediction of an x value, based on this data. The following procedure will produce a prediction:

- Transform the data set into standard coordinates.
- Compute the correlation.
- Predict $\hat{x}_0 = r\hat{y}_0$.
- Transform this prediction into the original coordinate system, to get

$$x_0 = \text{std}(x)r\hat{y}_0 + \text{mean}(\{x\})$$

There is another way of thinking about this prediction procedure, which is often helpful. Assume we need to predict a value for x_0 . In normalized coordinates, our prediction is $\hat{y}^p = r\hat{x}_0$; if we revert back to the original coordinate system, the

prediction becomes

$$\frac{(y^p - \text{mean}(\{y\}))}{\text{std}(y)} = r \left(\frac{(x_0 - \text{mean}(\{x\}))}{\text{std}(x)} \right).$$

This gives a really useful rule of thumb, which I have broken out in the box below.

Procedure: 13.2 *Predicting a value using correlation: Rule of thumb - 1*

If x_0 is k standard deviations from the mean of x , then the predicted value of y will be rk standard deviations away from the mean of y , and the sign of r tells whether y increases or decreases.

An even more compact version of the rule of thumb is in the following box.

Procedure: 13.3 *Predicting a value using correlation: Rule of thumb - 2*

The predicted value of y goes up by r standard deviations when the value of x goes up by one standard deviation.

We can compute the average root mean square error that this prediction procedure will make. The square of this error must be

$$\begin{aligned} \text{mean}(\{u^2\}) &= \text{mean}(\{y^2\}) - 2r\text{mean}(\{xy\}) + r^2\text{mean}(\{x^2\}) \\ &= 1 - 2r^2 + r^2 \\ &= 1 - r^2 \end{aligned}$$

so the root mean square error will be $\sqrt{1 - r^2}$. This is yet another interpretation of correlation; if x and y have correlation close to one, then predictions could have very small root mean square error, and so might be very accurate. In this case, knowing one variable is about as good as knowing the other. If they have correlation close to zero, then the root mean square error in a prediction might be as large as the root mean square error in \hat{y} — which means the prediction is nearly a pure guess.

The prediction argument means that we can spot correlations for data in other kinds of plots — one doesn't have to make a scatter plot. For example, if we were to observe a child's height from birth to their 10'th year (you can often find these observations in ballpen strokes, on kitchen walls), we could plot height as a function of year. If we also had their weight (less easily found), we could plot weight as a function of year, too. The prediction argument above says that, if you can predict the weight from the height (or vice versa) then they're correlated. One way to spot this is to look and see if one curve goes up when the other does (or

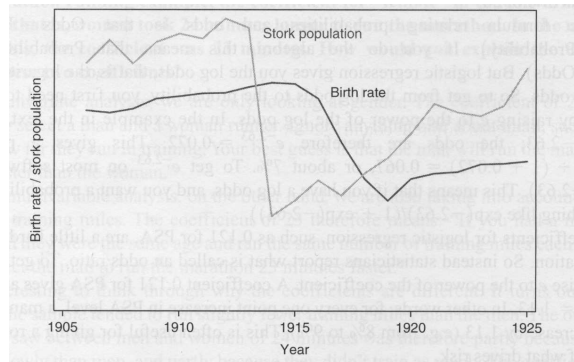


FIGURE 13.19: *This figure, from Vickers (ibid, p184) shows a plot of the stork population as a function of time, and the human birth rate as a function of time, for some years in Germany. The correlation is fairly clear; but this does not mean that reducing the number of storks means there are fewer able to bring babies. Instead, this is the impact of the first world war — a hidden or latent variable.*

goes down when the other goes up). You can see this effect in figure 13.5, where (before 19h00), prices go down when the number of pelts goes up, and vice versa. These two variables are negatively correlated.

13.2.3 Confusion caused by correlation

There is one very rich source of potential (often hilarious) mistakes in correlation. When two variables are correlated, they change together. If the correlation is positive, that means that, in typical data, if one is large then the other is large, and if one is small the other is small. In turn, this means that one can make a reasonable prediction of one from the other. However, correlation DOES NOT mean that changing one variable causes the other to change (sometimes known as causation).

Two variables in a dataset could be correlated for a variety of reasons. One important reason is pure accident. If you look at enough pairs of variables, you may well find a pair that appears to be correlated just because you have a small set of observations. Imagine, for example, you have a dataset consisting of only two vectors — there is a pretty good chance that there is some correlation between the coefficients. Such accidents can occur in large datasets, particularly if the dimensions are high.

Another reason variables could be correlated is that there is some causal relationship — for example, pressing the accelerator tends to make the car go faster, and so there will be some correlation between accelerator position and car acceleration. As another example, adding fertilizer does tend to make a plant grow bigger. Imagine you record the amount of fertilizer you add to each pot, and the size of the resulting potplant. There should be some correlation.

Yet another reason variables could be correlated is that there is some other background variable — often called a **latent variable** — linked causally to each of

the observed variables. For example, in children (as Freedman, Pisani and Purves note in their excellent *Statistics*), shoe size is correlated with reading skills. This DOES NOT mean that making your feet grow will make you read faster, or that you can make your feet shrink by forgetting how to read. The real issue here is the age of the child. Young children tend to have small feet, and tend to have weaker reading skills (because they've had less practice). Older children tend to have larger feet, and tend to have stronger reading skills (because they've had more practice). You can make a reasonable prediction of reading skills from foot size, because they're correlated, even though there is no direct connection.

This kind of effect can mask correlations, too. Imagine you want to study the effect of fertilizer on potplants. You collect a set of pots, put one plant in each, and add different amounts of fertilizer. After some time, you record the size of each plant. You expect to see correlation between fertilizer amount and plant size. But you might not if you had used a different species of plant in each pot. Different species of plant can react quite differently to the same fertilizer (some plants just die if over-fertilized), so the species could act as a latent variable. With an unlucky choice of the different species, you might even conclude that there was a negative correlation between fertilizer and plant size. This example illustrates why you need to take great care in setting up experiments and interpreting their results.

This sort of thing happens often, and it's an effect you should look for. Another nice example comes from Vickers (ibid). The graph, shown in Figure 13.19, shows a plot of (a) a dataset of the stork population in Europe over a period of years and (b) a dataset of the birth rate over those years. This isn't a scatter plot; instead, the data has been plotted on a graph. You can see by eye that these two datasets are quite strongly correlated. Even more disturbing, the stork population dropped somewhat before the birth rate dropped. Is this evidence that storks brought babies in Europe during those years? No (the usual arrangement seems to have applied). For a more sensible explanation, look at the dates. The war disturbed both stork and human breeding arrangements. Storks were disturbed immediately by bombs, etc., and the human birth rate dropped because men died at the front.

13.3 STERILE MALES IN WILD HORSE HERDS

Large herds of wild horses are (apparently) a nuisance, but keeping down numbers by simply shooting surplus animals would provoke outrage. One strategy that has been adopted is to sterilize males in the herd; if a herd contains sufficient sterile males, fewer foals should result. But catching stallions, sterilizing them, and reinserting them into a herd is a performance — does this strategy work?

We can get some insight by plotting data. At <http://lib.stat.cmu.edu/DASL/Datafiles/WildHorses.html>, you can find a dataset covering herd management in wild horses. I have plotted part of this dataset in figure 13.20. In this dataset, there are counts of all horses, sterile males, and foals made on each of a small number of days in 1986, 1987, and 1988 for each of two herds. I extracted data for one herd. I have plotted this data as a function of the count of days since the first data point, because this makes it clear that some measurements were taken at about the same time, but there are big gaps in the measurements. In this plot, the data

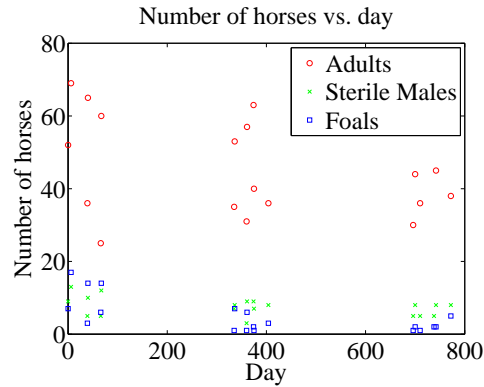


FIGURE 13.20: A plot of the number of adult horses, sterile males, and foals in horse herds over a period of three years. The plot suggests that introducing sterile males might cause the number of foals to go down. Data from <http://lib.stat.cmu.edu/DASL/Datafiles/WildHorses.html>.

points are shown with a marker. Joining them leads to a confusing plot because the data points vary quite strongly. However, notice that the size of the herd drifts down slowly (you could hold a ruler against the plot to see the trend), as does the number of foals, when there is a (roughly) constant number of sterile males.

Does sterilizing males result in fewer foals? This is likely hard to answer for this dataset, but we could ask whether herds with more sterile males have fewer foals. A scatter plot is a natural tool to attack this question. However, the scatter plots of figure 13.21 suggest, rather surprisingly, that when there are more sterile males there are more adults (and vice versa), and when there are more sterile males there are more foals (and vice versa). This is borne out by a correlation analysis. The correlation coefficient between foals and sterile males is 0.74, and the correlation coefficient between adults and sterile males is 0.68. You should find this very surprising — how do the horses know how many sterile males there are in the herd? You might think that this is an effect of scaling the plot, but there is a scatter plot in normalized coordinates in figure 13.21 that is entirely consistent with the conclusions suggested by the unnormalized plot. What is going on here?

The answer is revealed by the scatter plots of figure 13.22. Here, rather than plotting a '*' at each data point, I have plotted the day number of the observation. This is in days from the first observation. You can see that the whole herd is shrinking — observations where there are many adults (resp. sterile adults, foals) occur with small day numbers, and observations where there are few have large day numbers. Because the whole herd is shrinking, it is true that when there are more adults and more sterile males, there are also more foals. Alternatively, you can see the plots of figure 13.20 as a scatter plot of herd size (resp. number of foals, number of sterile males) against day number. Then it becomes clear that the whole herd is shrinking, as is the size of each group. To drive this point home, we can look at the correlation coefficient between adults and days (-0.24), between sterile adults and days (-0.37), and between foals and days (-0.61). We can use the rule of thumb in

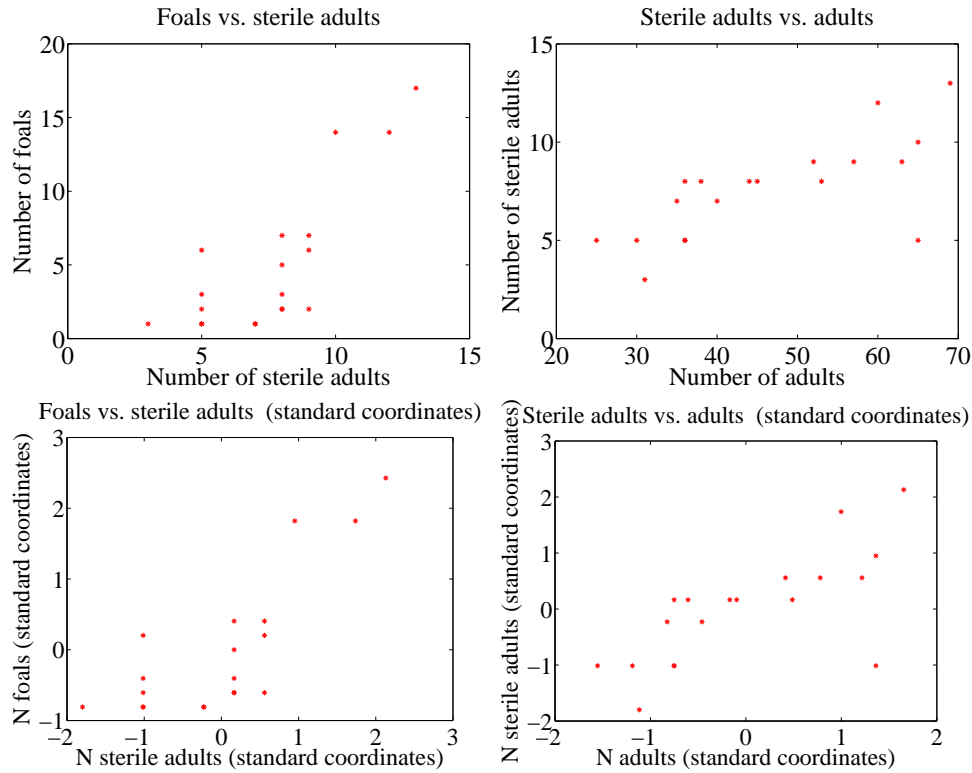


FIGURE 13.21: Scatter plots of the number of sterile males in a horse herd against the number of adults, and the number of foals against the number of sterile males, from data of <http://lib.stat.cmu.edu/DASL/Datafiles/WildHorses.html>. **Top:** unnormalized; **bottom:** standard coordinates.

box 13.3 to interpret this. This means that every 282 days, the herd loses about three adults; about one sterile adult; and about three foals. For the herd to have a stable size, it needs to gain by birth as many foals as it loses both to growing up and to death. If the herd is losing three foals every 282 days, then if they all grow up to replace the missing adults, the herd will be shrinking slightly (because it is losing four adults in this time); but if it loses foals to natural accidents, etc., then it is shrinking rather fast.

The message of this example is important. To understand a simple dataset, you might need to plot it several ways. You should make a plot, look at it and ask what it says, and then try to use another type of plot to confirm or refute what you think might be going on.

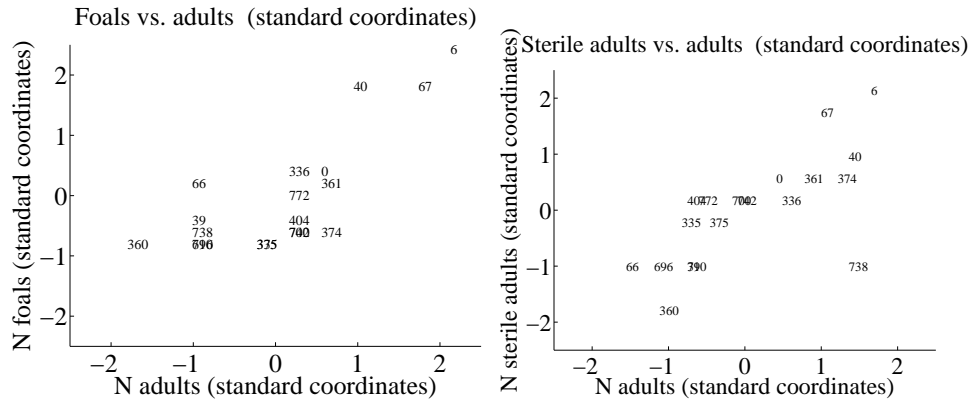


FIGURE 13.22: Scatter plots of the number of foals vs. the number of adults and the number of adults vs. the number of sterile adults for the wild horse herd, from <http://lib.stat.cmu.edu/DASL/Datafiles/WildHorses.html>. Rather than plot data points as dots, I have plotted the day on which the observation was made. Notice how the herd starts large, and then shrinks.

13.4 YOU SHOULD

13.4.1 be able to:

- Plot a bar chart, a heat map, and a pie chart for a categorical dataset.
- Plot a dataset as a graph, making sensible choices about markers, lines and the like.
- Plot a scatter plot for a dataset.
- Plot a normalized scatter plot for a dataset.
- Interpret the scatter plot to tell the sign of the correlation between two variables, and estimate the size of the correlation coefficient.
- Compute a correlation coefficient.
- Interpret a correlation coefficient.
- Use correlation to make predictions.

13.4.2 remember:

New term: pie chart	298
New term: stacked bar chart	298
New term: heat map	299
New term: 3D bar chart	300
There are a variety of tools for plotting categorical data.	301
New term: scatter plot	303
A scatter plot should be your first step with a new 2D dataset.	305

It's usually a good idea to plot in standard coordinates. 309
New term: correlation 309
New term: correlation coefficient 311
New term: correlation 311
Definition: Correlation coefficient 312
New term: latent variable 320

Background: Useful Probability Distributions

We will use probability as a tool to resolve practical questions about data. I describe some forms these questions take below, for concreteness. Generally, resolving these questions requires some form of model. This model gives an abstract representation of the problem that is useful for problem solving, and (typically) comes with recipes for attacking the major types of problem.

We could ask **what process produced the data?** For example, I observe a set of independent coin flips. I would now like to know the probability of observing a head when the coin is flipped. This might seem a bit empty as a problem, but an analogous problem is: are male children or female children more frequent? Notice that the answers to these questions are typically not exact. Instead, they are estimates. We will see a variety of methods to estimate the probability of observing a head from a sequence of independent coin flips, but none is guaranteed to return the “right” answer (because you can’t). Instead, we have to live with information about how accurate the estimate is.

We could ask **what sort of data can we expect in the future?** For example, we could ask: is gender assigned independently? equivalently, can you predict the gender of the next child born to a couple more accurately if you look at the genders of the previous children? In reliability engineering, one asks: how long will it be until this product fails? One version of this question that occupies many people is: how long until I die? By the way, the best answer seems to be subtract your age from a number that seems to be close to 85. Again, these are questions which don’t lend themselves to the “right” answer, as opposed to the best possible estimate. You might get hit by a truck tomorrow.

We could ask **what labels should we attach to unlabelled data?** For example, we might see a large number of credit card transactions, some known to be legitimate and others known to be fraudulent. We now see a new transaction: is it legitimate? You can see that versions of this question appear in many applications. As another example, you see many programs downloaded from the web, some known to be legitimate and others known to be malware. You now see a new program: is it safe to run? It may be possible in some circumstances to know the “right” answer to this question, but we are usually stuck with the best answer.

We could ask **is an effect easily explained by chance variations, or is it real?** For example, you believe that the average weight of a mouse is 15 grams. You could test this by catching 100 mice (a sample) and weighing them; but, when you do this, the answer won’t be 15 grams *even if* the average weight of a mouse is 15 grams. This is because you have a random selection of mice. But we will be able to build a model of the random variations in the sample average. We can use

this model to tell whether the difference between sample average and 15 grams is easily explained by random variations in the sample, or is significant. This allows us to tell how forcefully the evidence contradicts your original belief.

Building a model requires understanding the specific problem you want to solve, then choosing from a vocabulary of many different models that might apply to the problem. Long experience shows that even if the model you choose does not match the problem exactly, it can still be useful. In this chapter, I describe the properties of some probability distributions that are used again and again in model building.

14.1 DISCRETE DISTRIBUTIONS

14.1.1 The Discrete Uniform Distribution

If every value of a discrete random variable has the same probability, then the probability distribution is the discrete uniform distribution. We have seen this distribution before, numerous times. For example, I define a random variable by the number that shows face-up on the throw of a fair die. This has a uniform distribution. As another example, write the numbers 1-52 on the face of each card of a standard deck of playing cards. The number on the face of the first card drawn from a well-shuffled deck is a random variable with a uniform distribution.

One can construct expressions for the mean and variance of a discrete uniform distribution, but they're not usually much use (too many terms, not often used). Keep in mind that if two random variables have a uniform distribution, their sum and difference will not (recall example ??).

14.1.2 Bernoulli Random Variables

A Bernoulli random variable models a biased coin with probability p of coming up heads in any one flip.

Definition: 14.1 *Bernoulli random variable*

A Bernoulli random variable takes the value 1 with probability p and 0 with probability $1 - p$. This is a model for a coin toss, among other things

Useful Facts: 14.1 *Bernoulli random variable*

1. A Bernoulli random variable has mean p .
2. A Bernoulli random variable has variance $p(1 - p)$.

Proofs are easy, and in the exercises.

14.1.3 The Geometric Distribution

We have a biased coin. The probability it will land heads up, $P(\{H\})$ is given by p . We flip this coin until the first head appears. The number of flips required is a discrete random variable which takes integer values greater than or equal to one, which we shall call X . To get n flips, we must have $n - 1$ tails followed by 1 head. This event has probability $(1 - p)^{(n-1)}p$. We can now write out the probability distribution that n flips are required.

Definition: 14.2 *Geometric distribution*

We have an experiment with a binary outcome (i.e. heads or tails; 0 or 1; and so on), with $P(H) = p$ and $P(T) = 1 - p$. We repeat this experiment until the first head occurs. The probability distribution for n , the number of repetitions, is the geometric distribution. It has the form

$$P(\{X = n\}) = (1 - p)^{(n-1)}p.$$

for $0 \leq p \leq 1$ and $n \geq 1$; for other n it is zero. p is called the **parameter** of the distribution.

Notice that the geometric distribution is non-negative everywhere. It is straightforward to show that it sums to one, and so is a probability distribution (exercises).

Useful Facts: 14.2 *Geometric distribution*

1. The mean of the geometric distribution is $\frac{1}{p}$.
2. The variance of the geometric distribution is $\frac{1-p}{p^2}$.

The proof of these facts requires some work with series, and is relegated to the exercises.

It should be clear that this model isn't really about coins, but about repeated trials. The trial could be anything that has some probability of failing. Each trial is independent, and the rule for repeating is that you keep trying until the first success. Textbooks often set exercises involving missiles and aircraft; I'll omit these on the grounds of taste.

14.1.4 The Binomial Probability Distribution

Assume we have a biased coin with probability p of coming up heads in any one flip. The binomial probability distribution gives the probability that it comes up heads h times in N flips.

Worked example ?? yields one way of deriving this distribution. In that example, I showed that there are

$$N!/(h!(N-h)!)$$

outcomes of N coin flips that have h heads. These outcomes are disjoint, and each has probability $p^h(1-p)^{(N-h)}$. As a result, we must have the probability distribution below.

Definition: 14.3 *Binomial distribution*

In N independent repetitions of an experiment with a binary outcome (ie heads or tails; 0 or 1; and so on) with $P(H) = p$ and $P(T) = 1 - p$, the probability of observing a total of h H 's and $N - h$ T 's is

$$P_b(h; N, p) = \binom{N}{h} p^h (1-p)^{(N-h)}$$

as long as $0 \leq h \leq N$; in any other case, the probability is zero.

Useful Fact: 14.3 *Binomial distribution*

Write $P_b(i; N, p)$ for the binomial distribution that one observes i H 's in N trials.

$$\sum_{i=0}^N P_b(i; N, p) = (p + (1 - p))^N = (1)^N = 1$$

by pattern matching to the binomial theorem. As a result,

$$\sum_{i=0}^N P_b(i; N, p) = 1$$

The binomial distribution satisfies a recurrence relation. We must have that

$$P_b(h; N, p) = pP_b(h - 1; N - 1, p) + (1 - p)P_b(h; N - 1, p).$$

This is because can get h heads in N flips either by having $h - 1$ heads in $N - 1$ flips, then flipping another, or by having h heads in N flips then flipping a tail. You can verify by induction that the binomial distribution satisfies this recurrence relation.

Useful Facts: 14.4 *Binomial distribution*

1. The mean of $P_b(i; N, p)$ is Np .
2. The variance of $P_b(i; N, p)$ is $Np(1 - p)$

The proofs are informative, and so are not banished to the exercises.

Proof: 14.1 *The binomial distribution*

Notice that the number of heads in N coin tosses is can be obtained by adding the number of heads in each toss. Write Y_i for the Bernoulli random variable representing the i 'th toss. If the coin comes up heads, $Y_i = 1$, otherwise $Y_i = 0$. The Y_i are independent. Now

$$\begin{aligned}\mathbb{E}[X] &= \mathbb{E}\left[\sum_{j=1}^N Y_j\right] \\ &= \sum_{j=1}^N \mathbb{E}[Y_j] \\ &= N\mathbb{E}[Y_1] && \text{because the } Y_i \text{ are independent} \\ &= Np.\end{aligned}$$

The variance is easy, too. Each coin toss is independent, so the variance of the sum of coin tosses is the sum of the variances. This gives

$$\begin{aligned}\text{var}[X] &= \text{var}\left[\sum_{j=1}^N Y_j\right] \\ &= N\text{var}[Y_1] \\ &= Np(1-p)\end{aligned}$$

14.1.5 Multinomial probabilities

The binomial distribution describes what happens when a coin is flipped multiple times. But we could toss a die multiple times too. Assume this die has k sides, and we toss it N times. The distribution of outcomes is known as the **multinomial distribution**.

We can guess the form of the multinomial distribution in rather a straightforward way. The die has k sides. We toss the die N times. This gives us a sequence of N numbers. Each toss of the die is independent. Assume that side 1 appears n_1 times, side 2 appears n_2 times, ... side k appears n_k times. Any single sequence with this property will appear with probability $p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$, because the tosses are independent. However, there are

$$\frac{N!}{n_1! n_2! \dots n_k!}$$

such sequences. Using this reasoning, we arrive at the distribution below

Definition: 14.4 *Multinomial distribution*

I perform N independent repetitions of an experiment with k possible outcomes. The i 'th such outcome has probability p_i . I see outcome 1 n_1 times, outcome 2 n_2 times, etc. Notice that $n_1 + n_2 + n_3 + \dots + n_k = N$. The probability of observing this set of outcomes is

$$P_m(n_1, \dots, n_k; N, p_1, \dots, p_k) = \frac{N!}{n_1!n_2!\dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}.$$

Worked example 14.1 *Dice*

I throw five fair dice. What is the probability of getting two 2's and three 3's?

Solution: $\frac{5!}{2!3!} (\frac{1}{6})^2 (\frac{1}{6})^3$

14.1.6 The Poisson Distribution

Assume we are interested in counts that occur in an interval of time (e.g. within a particular hour). Because they are counts, they are non-negative and integer valued. We know these counts have two important properties. First, they occur with some fixed average rate. Second, an observation occurs independent of the interval since the last observation. Then the Poisson distribution is an appropriate model.

There are numerous such cases. For example, the marketing phone calls you receive during the day time are likely to be well modelled by a Poisson distribution. They come at some average rate — perhaps 5 a day as I write, during the last phases of an election year — and the probability of getting one clearly doesn't depend on the time since the last one arrived. Classic examples include the number of Prussian soldiers killed by horse-kicks each year; the number of calls arriving at a call center each minute; the number of insurance claims occurring in a given time interval (outside of a special event like a hurricane, etc.).

Definition: 14.5 *Poisson distribution*

A non-negative, integer valued random variable X has a Poisson distribution when its probability distribution takes the form

$$P(\{X = k\}) = \frac{\lambda^k e^{-\lambda}}{k!},$$

where $\lambda > 0$ is a parameter often known as the **intensity** of the distribution.

Notice that the Poisson distribution is a probability distribution, because it is non-negative and because

$$\sum_{i=0}^{\infty} \frac{\lambda^i}{i!} = e^{\lambda}$$

so that

$$\sum_{k=0}^{\infty} \frac{\lambda^k e^{-\lambda}}{k!} = 1$$

Useful Facts: 14.5 *Poisson distribution*

1. The mean of a Poisson distribution with intensity λ is λ .
2. The variance of a Poisson distribution with intensity λ is λ (no, that's not an accidentally repeated line or typo).

The proof of these facts requires some work with series, and is relegated to the exercises.

I described the Poisson distribution as a natural model for counts of randomly distributed points along a time axis. But it doesn't really matter that this is a time axis — it could be a space axis instead. For example, you could take a length of road, divide it into even intervals, then count the number of road-killed animals in each interval. If the location of each animal is independent of the location of any other animal, then you could expect a Poisson model to apply to the count data. Assume that the Poisson model that best describes the data has parameter λ . One property of such models is that if you doubled the length of the intervals, then the resulting dataset would be described by a Poisson model with parameter 2λ ; similarly, if you halved the length of the intervals, the best model would have parameter $\lambda/2$. This corresponds to our intuition about such data; roughly, the

number of road-killed animals in two miles of road should be twice the number in one mile of road. This property means that no pieces of the road are “special” — each behaves the same as the other.

We can build a really useful model of spatial randomness by observing this fact and generalizing very slightly. A **Poisson point process** with intensity λ is a set of random points with the property that the number of points in an interval of length s is a Poisson random variable with parameter λs . Notice how this captures our intuition that if points are “very randomly” distributed, there should be twice as many of them in an interval that is twice as long.

This model is easily, and very usefully, extended to points on the plane, on surfaces, and in 3D. In each case, the process is defined on a domain D (which has to meet some very minor conditions that are of no interest to us). The number of points in any subset s of D is a Poisson random variable, with intensity $\lambda m(s)$, where $m(s)$ is the area (resp. volume) of s . These models are useful, because they capture the property that (a) the points are random and (b) the probability you find a point doesn’t depend on where you are. You could reasonably believe models like this apply to, say, dead flies on windcreens; the places where you find acorns at the foot of an oak tree; the distribution of cowpats in a field; the distribution of cherries in a fruitcake; and so on.

14.2 CONTINUOUS DISTRIBUTIONS

14.2.1 The Continuous Uniform Distribution

Some continuous random variables have a natural upper bound and a natural lower bound but otherwise we know nothing about them. For example, imagine we are given a coin of unknown properties by someone who is known to be a skillful maker of unfair coins. The manufacturer makes no representations as to the behavior of the coin. The probability that this coin will come up heads is a random variable, about which we know nothing except that it has a lower bound of zero and an upper bound of one.

If we know nothing about a random variable apart from the fact that it has a lower and an upper bound, then a **uniform distribution** is a natural model. Write l for the lower bound and u for the upper bound. The probability density function for the uniform distribution is

$$p(x) = \begin{cases} 0 & x < l \\ 1/(u-l) & l \leq x \leq u \\ 0 & x > u \end{cases}$$

A continuous random variable whose probability distribution is the uniform distribution is often called a **uniform random variable**.

14.2.2 The Beta Distribution

It’s hard to explain now why the Beta (or β) distribution is useful, but it will come in useful later (section 16.20). The Beta distribution is a probability distribution for a continuous random variable x in the range $0 \leq x \leq 1$. There are two parameters, $\alpha > 0$ and $\beta > 0$. Recall the definition of the Γ function from section 2.2. We have

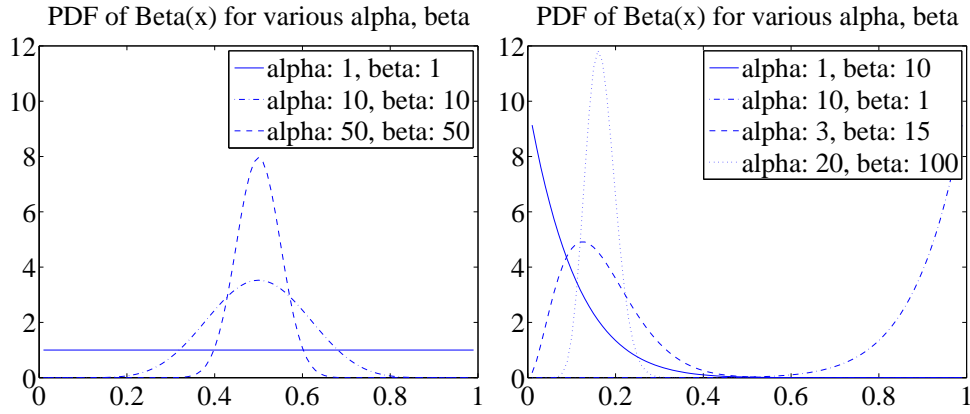


FIGURE 14.1: Probability density functions for the Beta distribution with a variety of different choices of α and β .

that

$$P_{\beta}(x|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{(\alpha-1)}(1-x)^{(\beta-1)}.$$

From this expression, you can see that:

- $P_{\beta}(x|1, 1)$ is a uniform distribution on the unit interval.
- $P_{\beta}(x|\alpha, \beta)$ has a single maximum at $x = (\alpha - 1)/(\alpha + \beta - 2)$ for $\alpha > 1, \beta > 1$ (differentiate and set to zero).
- Generally, as α and β get larger, this peak gets narrower.
- For $\alpha = 1, \beta > 1$ the largest value of $P_{\beta}(x|\alpha, \beta)$ is at $x = 0$.
- For $\alpha > 1, \beta = 1$ the largest value of $P_{\beta}(x|\alpha, \beta)$ is at $x = 1$.

Figure 14.1 shows plots of the probability density function of the Beta distribution for a variety of different values of α and β .

Useful Facts: 14.6 Beta distribution

For a Beta distribution with parameters α, β

1. The mean is $\frac{\alpha}{\alpha + \beta}$.
2. The variance is $\frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$.

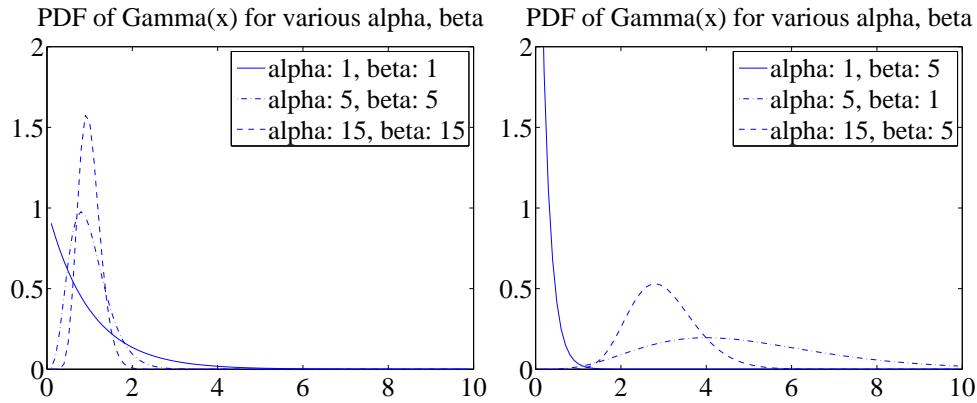


FIGURE 14.2: Probability density functions for the Gamma distribution with a variety of different choices of α and β .

14.2.3 The Gamma Distribution

The Gamma (or γ) distribution will also come in useful later on (section 16.20). The Gamma distribution is a probability distribution for a non-negative continuous random variable $x \geq 0$. There are two parameters, $\alpha > 0$ and $\beta > 0$. The probability density function is

$$P_{\gamma}(x|\alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{(\alpha-1)} e^{-\beta x}.$$

Figure 14.2 shows plots of the probability density function of the Gamma distribution for a variety of different values of α and β .

Useful Facts: 14.7 Gamma distribution

For a Gamma distribution with parameters α, β

1. The mean is $\frac{\alpha}{\beta}$.
2. The variance is $\frac{\alpha}{\beta^2}$.

14.2.4 The Exponential Distribution

Assume we have an infinite interval of time or space, with points distributed on it. Assume these points form a Poisson point process, as above. For example, we might consider the times at which email arrives; or the times at which phone calls arrive at a large telephone exchange; or the locations of roadkill on a road. The distance (or span of time) between two consecutive points is a random variable X . This

random variable takes an **exponential distribution**. There is a single parameter, λ . We have that

$$P_{\text{exp}}(x|\lambda) = \begin{cases} \lambda \exp^{-\lambda x} & \text{for } x \geq 0 \\ 0 & \text{otherwise} \end{cases} .$$

This distribution is often useful in modelling the failure of objects. We assume that failures form a Poisson process in time; then the time to the next failure is exponentially distributed.

Useful Facts: 14.8 *Exponential distribution*

For an exponential distribution with parameter λ

1. The mean is $\frac{1}{\lambda}$.
2. The variance is $\frac{1}{\lambda^2}$.

Notice the relationship between this parameter and the parameter of the Poisson distribution. If (say) the phone calls are distributed with Poisson distribution with intensity λ (per hour), then your expected number of calls per hour is λ . The time between calls will be exponentially distributed with parameter λ , and the expected time to the next call is $1/\lambda$ (in hours).

14.3 THE NORMAL DISTRIBUTION

14.3.1 The Standard Normal Distribution

Definition: 14.6 *Standard Normal distribution*

The probability density function

$$p(x) = \left(\frac{1}{\sqrt{2\pi}} \right) \exp \left(\frac{-x^2}{2} \right) .$$

is known as the **standard normal distribution**

The first step is to plot this probability density function (Figure 14.3). You should notice it is quite familiar from work on histograms, etc. in Chapter 16.20. It has the shape of the histogram of standard normal data, or at least the shape that the histogram of standard normal data aspires to.

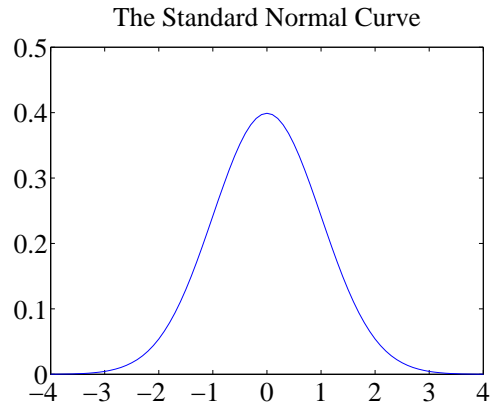


FIGURE 14.3: A plot of the probability density function of the standard normal distribution. Notice how probability is concentrated around zero, and how there is relatively little probability density for numbers with large absolute values.

Useful Facts: 14.9 *standard normal distribution*

1. The mean of the standard normal distribution is 0.
2. The variance of the standard normal distribution is 1.

These results are easily established by looking up (or doing!) the relevant integrals; they are relegated to the exercises.

A continuous random variable is a **standard normal random variable** if its probability density function is a standard normal distribution.

14.3.2 The Normal Distribution

Any probability density function that is a standard normal distribution *in standard coordinates* is a **normal distribution**. Now write μ for the mean of a random variable and σ for its standard deviation; we are saying that, if

$$\frac{x - \mu}{\sigma}$$

has a standard normal distribution, then $p(x)$ is a normal distribution. We can work out the form of the probability density function of a general normal distribution in two steps: first, we notice that for any normal distribution, we must have

$$p(x) \propto \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right].$$

But, for this to be a probability density function, we must have $\int_{-\infty}^{\infty} p(x)dx = 1$. This yields the constant of proportionality, and we get

Definition: 14.7 *Normal distribution*

The probability density function

$$p(x) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right) \exp \left(\frac{-(x - \mu)^2}{2\sigma^2} \right).$$

is a normal distribution.

Useful Facts: 14.10 *normal distribution*

The probability density function

$$p(x) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right) \exp \left(\frac{-(x - \mu)^2}{2\sigma^2} \right).$$

has

1. mean μ
2. and variance σ .

These results are easily established by looking up (or doing!) the relevant integrals; they are relegated to the exercises.

A continuous random variable is a **normal random variable** if its probability density function is a **normal distribution**. Notice that it is quite usual to call normal distributions **gaussian distributions**.

14.3.3 Properties of the Normal Distribution

Normal distributions are important, because one often runs into data that is well described by a normal distribution. It turns out that anything that behaves like a binomial distribution with a lot of trials — for example, the number of heads in many coin tosses; as another example, the percentage of times you get the outcome of interest in a simulation in many runs — should produce a normal distribution (Section 14.4). For this reason, pretty much any experiment where you perform a simulation, then count to estimate a probability or an expectation, should give you an answer that has a normal distribution.

It is a remarkable and deep fact, known as the **central limit theorem**, that adding many independent random variables produces a normal distribution pretty much *whatever* the distributions of those random variables. I've not shown this in detail because it's a nuisance to prove. However, if you add together many random variables, each of pretty much any distribution, then the answer has a distribution close to the normal distribution. It turns out that many of the processes we observe add up subsidiary random variables. This means that you will see normal distributions very often in practice.

A normal random variable tends to take values that are quite close to the mean, measured in standard deviation units. We can demonstrate this important fact by computing the probability that a standard normal random variable lies between u and v . We form

$$\int_u^v \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du.$$

It turns out that this integral can be evaluated relatively easily using a special function. The **error function** is defined by

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2) dt$$

so that

$$\frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) = \int_0^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du.$$

Notice that $\operatorname{erf}(x)$ is an odd function (i.e. $\operatorname{erf}(-x) = -\operatorname{erf}(x)$). From this (and tables for the error function, or Matlab) we get that, for a standard normal random variable

$$\frac{1}{\sqrt{2\pi}} \int_{-1}^1 \exp\left(-\frac{x^2}{2}\right) dx \approx 0.68$$

and

$$\frac{1}{\sqrt{2\pi}} \int_{-2}^2 \exp\left(-\frac{x^2}{2}\right) dx \approx 0.95$$

and

$$\frac{1}{\sqrt{2\pi}} \int_{-3}^3 \exp\left(-\frac{x^2}{2}\right) dx \approx 0.99.$$

These are very strong statements. They measure how often a standard normal random variable has values that are in the range $-1, 1$, $-2, 2$, and $-3, 3$ respectively. But these measurements apply to normal random variables if we recognize that they now measure how often the normal random variable is some number of standard deviations away from the mean. In particular, it is worth remembering that:

Useful Facts: 14.11 *Normal Random Variables*

- About 68% of the time, a normal random variable takes a value within one standard deviation of the mean.
- About 95% of the time, a normal random variable takes a value within one standard deviation of the mean.
- About 99% of the time, a normal random variable takes a value within one standard deviation of the mean.

14.4 APPROXIMATING BINOMIALS WITH LARGE N

The Binomial distribution appears to be a straightforward thing. We assume we flip a coin N times, where N is a very large number. The coin has probability p of coming up heads, and so probability $q = 1 - p$ of coming up tails. The number of heads h follows the binomial distribution, so

$$P(h) = \frac{N!}{h!(N-h)!} p^h q^{(N-h)}$$

The mean of this distribution is Np , the variance is Npq , and the standard deviation is \sqrt{Npq} .

Evaluating this probability distribution for large N is very difficult, because factorials grow fast. We will construct an approximation to the binomial distribution for large N that allows us to evaluate the probability that h lies in some range. This approximation will show that the probability that h is within one standard deviation of the mean is approximately 68%.

This is important, because it shows that our model of probability as frequency is consistent. Consider the probability that the number of heads you see lies within one standard deviation of the mean. The size of that interval is $2\sqrt{Npq}$. As N gets bigger, the size of that interval, *relative to the total number of flips*, gets smaller. If I flip a coin N times, in principle I could see a number of heads h that ranges from 0 to N . However, we will establish that about 68% of the time, h will lie in the interval within one standard deviation of the mean. The size of this interval, *relative to the total number of flips* is

$$2\frac{\sqrt{Npq}}{N} = 2\sqrt{\frac{pq}{N}}.$$

As a result, as $N \rightarrow \infty$,

$$\frac{h}{N} \rightarrow p$$

because h will tend to land in an interval around pN that gets narrower as N gets larger.

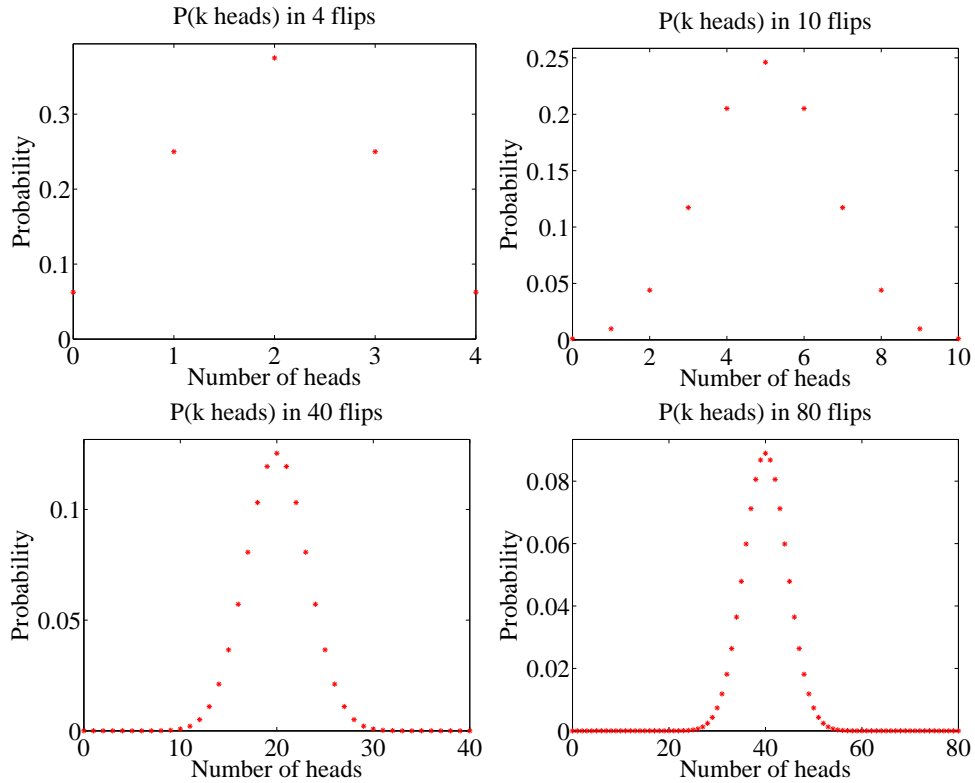


FIGURE 14.4: Plots of the binomial distribution for $p = q = 0.5$ for different values of N . You should notice that the set of values of h (the number of heads) that have substantial probability is quite narrow compared to the range of possible values. This set gets narrower as the number of flips increases. This is because the mean is pN and the standard deviation is \sqrt{Npq} — so the fraction of values that is within one standard deviation of the mean is $O(1/\sqrt{N})$.

The main difficulty with Figure 14.4 (and with the argument above) is that the mean and standard deviation of the binomial distribution tends to infinity as the number of coin flips tends to infinity. This can confuse issues. For example, the plots of Figure 14.4 show narrowing probability distributions — but is this because the scale is compacted, or is there a real effect? It turns out there is a real effect, and a good way to see it is to consider the normalized number of heads.

14.4.1 Large N

Recall that to normalize a dataset, you subtract the mean and divide the result by the standard deviation. We can do the same for a random variable. We now consider

$$x = \frac{h - Np}{\sqrt{Npq}}.$$

The probability distribution of x can be obtained from the probability distribution for h , because $h = Np + x\sqrt{Npq}$, so

$$P(x) = \left(\frac{N!}{(Np + x\sqrt{Npq})!(Nq - x\sqrt{Npq})!} \right) p^{(Np + x\sqrt{Npq})} q^{(Nq - x\sqrt{Npq})}.$$

I have plotted this probability distribution for various values of N in Figure 14.5.

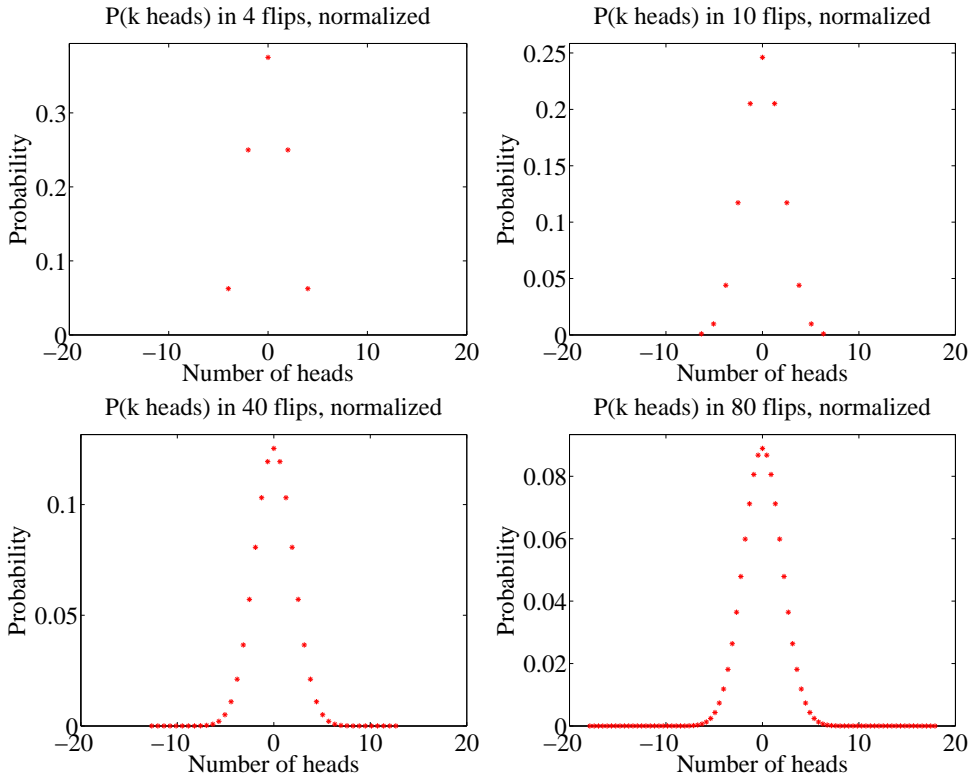


FIGURE 14.5: Plots of the distribution for the normalized variable x , with $P(x)$ given in the text, obtained from the binomial distribution with $p = q = 0.5$ for different values of N . These distributions are normalized (mean 0, variance 1). They look increasingly like a standard normal distribution EXCEPT that the value at their mode gets smaller as N gets bigger (there are more possible outcomes). In the text, we will establish that the standard normal distribution is a limit, in a useful sense.

But it is hard to work with this distribution for very large N . The factorials become very difficult to evaluate. Second, it is a discrete distribution on N points, spaced $1/\sqrt{Npq}$ apart. As N becomes very large, the number of points that have non-zero probability becomes very large, and x can be very large, or very small. For example, there is some probability, though there may be very little indeed, on the point where $h = N$, or, equivalently, $x = N(p + \sqrt{Npq})$. For sufficiently large N ,

we think of this probability distribution as a probability density function. We can do so, for example, by spreading the probability for x_i (the i 'th value of x) evenly over the interval between x_i and x_{i+1} . We then have a probability density function that looks like a histogram, with bars that become narrower as N increases. But what is the limit?

14.4.2 Getting Normal

To proceed, we need Stirling's approximation, which says that, for large N ,

$$N! \approx \sqrt{2\pi} \sqrt{N} \left(\frac{N}{e}\right)^N.$$

This yields

$$P(h) \approx \left(\frac{Np}{h}\right)^h \left(\frac{Nq}{N-h}\right)^{(N-h)} \sqrt{\frac{N}{2\pi h(N-h)}}$$

Recall we used the normalized variable

$$x = \frac{h - Np}{\sqrt{Npq}}.$$

We will encounter the term \sqrt{Npq} often, and we use $\sigma = \sqrt{Npq}$ as a shorthand. We can compute h and $N - h$ from x by the equalities

$$h = Np + \sigma x \qquad N - h = Nq - \sigma x.$$

So the probability distribution written in this new variable x is

$$P(x) \approx \left(\frac{Np}{Np + \sigma x}\right)^{(Np + \sigma x)} \left(\frac{Nq}{Nq - \sigma x}\right)^{(Nq - \sigma x)} \sqrt{\frac{N}{2\pi(Np + \sigma x)(Nq - \sigma x)}}$$

There are three terms to deal with here. It is easiest to work with $\log P$. Now

$$\log(1 + x) = x - \frac{1}{2}x^2 + O(x^3)$$

so we have

$$\begin{aligned} \log\left(\frac{Np}{Np + \sigma x}\right) &= -\log\left(1 + \frac{\sigma x}{Np}\right) \\ &\approx -\frac{\sigma x}{Np} + \left(\frac{1}{2}\right)\left(\frac{\sigma x}{Np}\right)^2 \end{aligned}$$

and

$$\log\left(\frac{Nq}{Nq - \sigma x}\right) \approx \frac{\sigma x}{Nq} + \left(\frac{1}{2}\right)\left(\frac{\sigma x}{Nq}\right)^2.$$

From this, we have that

$$\begin{aligned} \log \left[\left(\frac{Np}{Np + \sigma x} \right)^{(Np + \sigma x)} \left(\frac{Nq}{Nq - \sigma x} \right)^{(Nq - \sigma x)} \right] &\approx [Np + \sigma x] \left[-\frac{\sigma x}{Np} + \left(\frac{1}{2} \right) \left(\frac{\sigma x}{Np} \right)^2 \right] + \\ &\quad [Nq - \sigma x] \left[\frac{\sigma x}{Nq} + \left(\frac{1}{2} \right) \left(\frac{\sigma x}{Nq} \right)^2 \right] \\ &= -\left(\frac{1}{2} \right) x^2 + O((\sigma x)^3) \end{aligned}$$

(recall $\sigma = \sqrt{Npq}$ if you're having trouble with the last step). Now we look at the square-root term. We have

$$\begin{aligned} \log \sqrt{\frac{N}{2\pi(Np + \sigma x)(Nq - \sigma x)}} &= -\frac{1}{2} (\log [Np + \sigma x] + \log [Nq - \sigma x] - \log N + \log 2\pi) \\ &= -\frac{1}{2} \left(\begin{array}{l} \log Np + O\left(\left(\frac{\sigma x}{Np}\right)\right) \\ + \log Nq - O\left(\left(\frac{\sigma x}{Nq}\right)\right) \\ - \log N + \log 2\pi \end{array} \right) \end{aligned}$$

but, since N is very large compared to σx , we can ignore the $O\left(\left(\frac{\sigma x}{Np}\right)\right)$ terms. Then this term is not a function of x . So we have

$$\log P(x) \approx \frac{-x^2}{2} + \text{constant}.$$

Now because N is very large, our probability distribution P limits to a probability density function p , with

$$p(x) \propto \exp\left(\frac{-x^2}{2}\right).$$

We can get the constant of proportionality from integrating, to

$$p(x) = \left(\frac{1}{\sqrt{2\pi}} \right) \exp\left(\frac{-x^2}{2}\right).$$

This constant of proportionality deals with the effect in figure 14.5, where the mode of the distribution gets smaller as N gets bigger. It does so because there are more points with non-zero probability to be accounted for. But we are interested in the limit where N tends to infinity. This must be a probability density function, so it must integrate to one.

Review this blizzard of terms. We started with a binomial distribution, but standardized the variables so that the mean was zero and the standard deviation was one. We then assumed there was a very large number of coin tosses, so large that that the distribution started to look like a continuous function. The function we get is the standard normal distribution.

14.4.3 So What?

I have proven an extremely useful fact, which I shall now put in a box.

Useful Fact: 14.12 *Binomial distribution for large N*

Assume h follows the binomial distribution with parameters p and q . Write $x = \frac{h - Np}{\sqrt{Npq}}$. Then, for sufficiently large N , the probability distribution $P(x)$ can be approximated by the probability density function

$$\left(\frac{1}{\sqrt{2\pi}}\right) \exp\left(\frac{-x^2}{2}\right)$$

in the sense that

$$P(\{x \in [a, b]\}) \approx \int_a^b \left(\frac{1}{\sqrt{2\pi}}\right) \exp\left(\frac{-u^2}{2}\right) du$$

This justifies our model of probability as frequency. I interpreted an event having probability p to mean that, if I had a large number N of independent repetitions of the experiment, the number that produced the event would be close to Np , and would get closer as N got larger. We know that, for example, 68% of the time a standard normal random variable takes a value between 1 and -1 . In this case, the standard normal random variable is

$$\frac{h - (Np)}{\sqrt{Npq}}$$

so that 68% of the time, h must take a value in the range $[Np - \sqrt{Npq}, Np + \sqrt{Npq}]$. Equivalently, the relative frequency h/N must take a value in the range

$$\left[p - \frac{pq}{\sqrt{N}}, p + \frac{pq}{\sqrt{N}}\right]$$

but as $N \rightarrow \infty$ this range gets smaller and smaller, and h/N limits to p . So our view of probability as a frequency is consistent.

To obtain h , we added N independent Bernoulli random variables. So you can interpret the box as saying that the sum of many independent Bernoulli random variables has a probability distribution that limits to the normal distribution as the number added together gets larger. Remember that I have stated, though not precisely, but not proved the deep and useful fact that the sum of pretty much any independent random variables has a distribution that gets closer to a normal distribution as the number added together gets larger.

14.5 YOU SHOULD

14.5.1 remember:

Definition: Bernoulli random variable 326
 Useful facts: Bernoulli random variable 327
 Definition: Geometric distribution 327
 Useful facts: Geometric distribution 328
 Definition: Binomial distribution 328
 Useful fact: Binomial distribution 329
 Useful facts: Binomial distribution 329
 Definition: Multinomial distribution 331
 New term: intensity 332
 Definition: Poisson distribution 332
 Useful facts: Poisson distribution 332
 New term: Poisson point process 333
 New term: uniform distribution 333
 New term: uniform random variable 333
 Useful facts: Beta distribution 334
 Useful facts: Gamma distribution 335
 New term: exponential distribution 336
 Useful facts: Exponential distribution 336
 New term: standard normal distribution 336
 Definition: Standard Normal distribution 336
 Useful facts: standard normal distribution 337
 New term: standard normal random variable 337
 New term: normal distribution 337
 Definition: Normal distribution 338
 Useful facts: normal distribution 338
 New term: normal random variable 338
 New term: normal distribution 338
 New term: gaussian distributions 338
 New term: central limit theorem 339
 New term: error function 339
 Useful facts: Normal Random Variables 340
 Useful fact: Binomial distribution for large N 345

Background: Inference: Making Point Estimates

Inference is the process of drawing conclusions from data. One form of inference is to estimate a number, or set of numbers, that describes a dataset. The result is known as a **point estimate**. An alternative is to estimate an interval within which a number lies, with some degree of certainty. Such estimates are known as **interval estimates**. Finally, one might wish to assess the extent to which a body of evidence supports rejecting an hypothesis — known as **hypothesis testing**. In this chapter, we deal with point estimates. In the following chapter, we deal with interval estimates and hypothesis testing, which require an understanding of point estimates. There are two, somewhat distinct, situations in which we could make point estimates.

In the first, we have a dataset $\{\mathbf{x}\}$, and a probability model we believe applies to that dataset. But we need to select appropriate values of the parameters to ensure that the model describes the data. For example, we might have a set of N coin flips which we believe to be independent and identically distributed. Of these, k flips came up H . We know that a binomial distribution with $p(H) = p$ is a good model — but what value of p should we use? Your intuition is likely to suggest using k/N , but we'd like a more robust procedure than guessing. We need an inference procedure to obtain the unknown parameter from the data. Notice that this will be an estimate, rather than the “true” value. As we shall see, there is more than one possible procedure to apply, depending to some extent on the problem. In some cases (section 15.1), we estimate parameter values based solely on data; in others (section 15.2), we are able to use prior information about the parameters to affect the estimate.

In the second situation, we want to know some property of a population. For example, we may wish to know the mean weight of a person, or the mean response of a mouse to a drug. It would be a stretch to describe this population with one of the probability models that we have seen. In principle, the number we want is not even necessarily random; in principle, we could measure everyone on the planet and average the weights. In practice, this doesn't make sense, for quite straightforward reasons. You can't really weigh every person, or dose every mouse, on the planet. Instead, to estimate this property, we obtain a sample (some people; some mice; etc.) of the population, and estimate the property from the sample. There is now an important problem. Different samples lead to different estimates of the property. We will arrange to have the sample drawn randomly from the population, so the sample we see represents the value of a set of random variables. If you followed the proof of the weak law of large numbers, you should suspect that the mean of this sample could be a good estimate of the population mean. This turns out to be the

case (section 16.20). However, there is some random error in the estimate, and we can tell (on average) how large the error caused by random sampling could be.

15.1 ESTIMATING MODEL PARAMETERS WITH MAXIMUM LIKELIHOOD

Assume we have a dataset $\mathcal{D} = \{\mathbf{x}\}$, and a probability model we believe applies to that dataset. Generally, application logic suggests the type of model (i.e. normal probability density; Poisson probability; geometric probability; and so on). But usually, we do not know the parameters of the model — for example, the mean and standard deviation of a normal distribution; the intensity of a poisson distribution; and so on. Notice that this situation is unlike what we have seen to date. In chapter 14, we assumed that we knew θ , and could then use the model to assign a probability to a set of data items \mathcal{D} . Here we *know* the value of \mathcal{D} , but don't know θ . Our model will be better or worse depending on how well we choose the parameters. We need a strategy to estimate the parameters of a model from a sample dataset. Notice how each of the following examples fits this pattern.

Example: 15.1 *Inferring p from repeated flips — binomial*

We could flip the coin N times, and count the number of heads k . We know that an appropriate probability model for a set of independent coin flips is the binomial model $P(k; N, p)$. But we do not know p , which is the parameter — we need a strategy to extract a value of p from the data.

Example: 15.2 *Inferring p from repeated flips — geometric*

We could flip the coin repeatedly until we see a head. We know that, in this case, the number of flips has the geometric distribution with parameter p . In this case, the data is a sequence of T 's with a final H from the coin flips. There are N flips (or terms) and the last flip is a head. We know that an appropriate probability model is the geometric distribution $P_g(N; p)$. But we do not know p , which is the parameter — we need a strategy to extract a value of p from the data.

Example: 15.3 *Inferring the intensity of spam — poisson*

It is reasonable to assume that the number of spam emails one gets in an hour has a Poisson distribution. But what is the intensity parameter λ ? We could count the number of spam emails that arrive in each of a set of distinct hours, giving a dataset of counts \mathcal{D} . We need a strategy to wrestle an estimate of λ from this dataset.

Example: 15.4 *Inferring the mean and standard deviation of normal data*

Imagine we know for some reason that our data is well described by a normal distribution. We could ask what is the mean and standard deviation of the normal distribution that best represents the data?

Write θ for the parameters of our model. If we knew θ , then the probability of observing the data \mathcal{D} would be $P(\mathcal{D}|\theta)$. We know \mathcal{D} , and we don't know θ , so the value of $P(\mathcal{D}|\theta)$ is a function of θ .

Definition: 15.1 *Likelihood*

The function $P(\mathcal{D}|\theta)$, which is a function of θ , is known as the **likelihood** of the data D , and is often written $\mathcal{L}(\theta)$ (or $\mathcal{L}(\theta; \mathcal{D})$ if you want to remember that data is involved).

15.1.1 The Maximum Likelihood Principle

We need a “reasonable” procedure to choose a value of θ to report. One — and I stress this is not the only one — is the **maximum likelihood principle**.

Definition: 15.2 *The maximum likelihood principle*

Choose θ such that $\mathcal{L}(\theta) = P(\mathcal{D}|\theta)$ is maximised, as a function of θ .

In words, this means: Choose the parameter such that the probability of observing the data you actually see, is maximised. This should strike you as being a reasonable choice. You should also be aware that this is not the only possible choice (we'll see another one in section 15.2).

For the examples we work with, the data will be **independent and identically distributed** or **IID**. This means that each data item is an independently obtained sample from the same probability distribution (see section ??). In turn, this means that the likelihood is a product of terms, one for each data item, which we can write as

$$\mathcal{L}(\theta) = P(\mathcal{D}|\theta) = \prod_{i \in \text{dataset}} P(d_i|\theta).$$

It is traditional to write θ for any set of parameters that are unknown. There are two, distinct, important concepts we must work with. One is the unknown parameter(s), which we will write θ . The other is the *estimate* of the value of that parameter, which we will write $\hat{\theta}$. This estimate is the best we can do — it may not be the “true” value of the parameter.

Worked example 15.1 *Inferring $p(H)$ for a coin from flips using a binomial model*

In N independent coin flips, you observe k heads. Use the maximum likelihood principle to infer $p(H)$.

Solution: The coin has $\theta = p(H)$, which is the unknown parameter. We know that an appropriate probability model is the binomial model $P(k; N, \theta)$. We have that

$$\mathcal{L}(\theta) = P(\mathcal{D}|\theta) = P_b(k; N, \theta) = \binom{N}{k} \theta^k (1 - \theta)^{(N-k)}$$

which is a function of θ — the unknown probability that a coin comes up heads; k and N are known. We must find the value of θ that maximizes this expression. Now the maximum occurs when

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = 0.$$

We have

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \binom{N}{k} \left(k\theta^{k-1}(1 - \theta)^{(N-k)} - \theta^k(N - k)(1 - \theta)^{(N-k-1)} \right)$$

and this is zero when

$$k\theta^{k-1}(1 - \theta)^{(N-k)} = \theta^k(N - k)(1 - \theta)^{(N-k-1)}$$

so the maximum occurs when

$$k(1 - \theta) = \theta(N - k).$$

This means the maximum likelihood estimate is

$$\hat{\theta} = \frac{k}{N}$$

which is what we guessed would happen, but now we know why that guess “makes sense”.

Worked example 15.2 *Inferring $p(H)$ from coin flips using a geometric model*

You flip a coin N times, stopping when you see a head. Use the maximum likelihood principle to infer $p(H)$ for the coin.

Solution: The coin has $\theta = p(H)$, which is the unknown parameter. We know that an appropriate probability model is the geometric model $P_g(N; \theta)$. We have that

$$\mathcal{L}(\theta) = P(\mathcal{D}|\theta) = P_g(N; \theta) = (1 - \theta)^{(N-1)}\theta$$

which is a function of θ — the unknown probability that a coin comes up heads; N is known. We must find the value of θ that maximizes this expression. Now the maximum occurs when

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = 0 = ((1 - \theta)^{(N-1)} - (N - 1)(1 - \theta)^{(N-2)}\theta)$$

So the maximum likelihood estimate is

$$\hat{\theta} = \frac{1}{N}.$$

We didn't guess this.

Worked example 15.3 *Inferring die probabilities from multiple rolls and a multinomial distribution*

You throw a die N times, and see n_1 ones, ... and n_6 sixes. Write p_1, \dots, p_6 for the probabilities that the die comes up one, ..., six. Use the maximum likelihood principle to estimate p_1, \dots, p_6 .

Solution: The data are n, n_1, \dots, n_6 . The parameters are $\theta = (p_1, \dots, p_6)$. $P(\mathcal{D}|\theta)$ comes from the multinomial distribution. In particular,

$$\mathcal{L}(\theta) = P(\mathcal{D}|\theta) = \frac{n!}{n_1! \dots n_6!} p_1^{n_1} p_2^{n_2} \dots p_6^{n_6}$$

which is a function of $\theta = (p_1, \dots, p_6)$. Now we want to maximize this function by choice of θ . Notice that we could do this by simply making all p_i very large — but this omits a fact, which is that $p_1 + p_2 + p_3 + p_4 + p_5 + p_6 = 1$. So we substitute using $p_6 = 1 - p_1 - p_2 - p_3 - p_4 - p_5$ (there are other, neater, ways of dealing with this issue, but they take more background knowledge). At the maximum, we must have that for all i ,

$$\frac{\partial \mathcal{L}(\theta)}{\partial p_i} = 0$$

which means that, for p_i , we must have

$$n_i p_i^{(n_i-1)} (1 - p_1 - p_2 - p_3 - p_4 - p_5)^{n_6} - p_i^{n_i} n_6 (1 - p_1 - p_2 - p_3 - p_4 - p_5)^{(n_6-1)} = 0$$

so that, for each p_i , we have

$$n_i (1 - p_1 - p_2 - p_3 - p_4 - p_5) - n_6 p_i = 0$$

or

$$\frac{p_i}{1 - p_1 - p_2 - p_3 - p_4 - p_5} = \frac{n_i}{n_6}.$$

You can check that this equation is solved by

$$\hat{\theta} = \frac{1}{(n_1 + n_2 + n_3 + n_4 + n_5 + n_6)} (n_1, n_2, n_3, n_4, n_5, n_6)$$

The logarithm is a monotonic function (i.e. if $x > 0, y > 0, x > y$, then $\log(x) > \log(y)$). This means that the values of θ that maximise the log-likelihood are the same as the values that maximise the likelihood. This observation is very useful, because it allows us to transform a product into a sum. The derivative of a product involves numerous terms; the derivative of a sum is easy to take. We have

$$\log P(\mathcal{D}|\theta) = \log \prod_{i \in \text{dataset}} P(d_i|\theta) = \sum_{i \in \text{dataset}} \log P(d_i|\theta)$$

and in some cases, $\log P(d_i|\theta)$ takes a convenient, easy form.

Definition: 15.3 *The log-likelihood of a dataset under a model*

is a function of the unknown parameters, and you will often see it written as

$$\log \mathcal{L}(\theta) = \sum_{i \in \text{dataset}} \log P(d_i|\theta).$$

Worked example 15.4 *Poisson distributions*

You observe N intervals, each of the same, fixed length (in time, or space). You know that, in these intervals, events occur with a Poisson distribution (for example, you might be observing Prussian officers being kicked by horses, or telemarketer calls...). You know also that the intensity of the Poisson distribution is the same for each observation. The number of events you observe in the i 'th interval is n_i . What is the intensity, λ ?

Solution: The likelihood is

$$\mathcal{L}(\theta) = \prod_{i \in \text{intervals}} P(\{n_i \text{ events}\}|\theta) = \prod_{i \in \text{intervals}} \frac{\theta^{n_i} e^{-\theta}}{n_i!}.$$

It will be easier to work with logs. The log-likelihood is

$$\log \mathcal{L}(\theta) = \sum_i (n_i \log \theta - \theta - \log n_i!)$$

so that we must solve

$$\frac{\partial \log \mathcal{L}(\theta)}{\partial \theta} = \sum_i \left(\frac{n_i}{\theta} - 1 \right) = 0$$

which yields a maximum likelihood estimate of

$$\hat{\theta} = \frac{\sum_i n_i}{N}$$

Worked example 15.5 *The intensity of swearing*

A famously swearsome politician gives a talk. You listen to the talk, and for each of 30 intervals 1 minute long, you record the number of swearwords. You record this as a histogram (i.e. you count the number of intervals with zero swear words, with one, etc.). For the first 10 intervals, you see

no. of swear words	no. of intervals
0	5
1	2
2	2
3	1
4	0

and for the following 20 intervals, you see

no. of swear words	no. of intervals
0	9
1	5
2	3
3	2
4	1

Assume that the politician's use of swearwords is Poisson. What is the intensity using the first 10 intervals? the second 20 intervals? all the intervals? why are they different?

Solution: Use the expression from worked example 15.4 to find

$$\begin{aligned}\hat{\lambda}_{10} &= \frac{\text{total number of swearwords}}{\text{number of intervals}} \\ &= \frac{9}{10}\end{aligned}$$

$$\begin{aligned}\hat{\lambda}_{20} &= \frac{\text{total number of swearwords}}{\text{number of intervals}} \\ &= \frac{21}{20}\end{aligned}$$

$$\begin{aligned}\hat{\lambda}_{30} &= \frac{\text{total number of swearwords}}{\text{number of intervals}} \\ &= \frac{30}{30}.\end{aligned}$$

These are different because the maximum likelihood estimate is an *estimate* — we can't expect to recover the exact value from a dataset. Notice, however, that the estimates are quite close.

Worked example 15.6 *Normal distributions*

Assume we have x_1, \dots, x_N , and we wish to model these data with a normal distribution. Use the maximum likelihood principle to estimate the mean of that normal distribution.

Solution: The likelihood of a set of data values under the normal distribution with unknown mean θ and standard deviation σ is

$$\begin{aligned}\mathcal{L}(\theta) &= P(x_1, \dots, x_N | \theta, \sigma) \\ &= P(x_1 | \theta, \sigma) P(x_2 | \theta, \sigma) \dots P(x_N | \theta, \sigma) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \theta)^2}{2\sigma^2}\right)\end{aligned}$$

and this expression is a moderate nuisance to work with. The log of the likelihood is

$$\log \mathcal{L}(\theta) = \left(\sum_{i=1}^N -\frac{(x_i - \theta)^2}{2\sigma^2} \right) + \text{term not depending on } \theta.$$

We can find the maximum by differentiating wrt θ and setting to zero, which yields

$$\begin{aligned}\frac{\partial \log \mathcal{L}(\theta)}{\partial \theta} &= \sum_{i=1}^N \frac{2(x_i - \theta)}{2\sigma^2} \\ &= 0 \\ &= \frac{1}{\sigma^2} \left(\sum_{i=1}^N x_i - N\theta \right)\end{aligned}$$

so the maximum likelihood estimate is

$$\hat{\theta} = \frac{\sum_{i=1}^N x_i}{N}$$

which probably isn't all that surprising. Notice we did not have to pay attention to σ in this derivation — we did not assume it was known, it just doesn't do anything.

Worked example 15.7 *Normal distributions -II*

Assume we have x_1, \dots, x_N which are data that can be modelled with a normal distribution. Use the maximum likelihood principle to estimate the standard deviation of that normal distribution.

Solution: Now we have to write out the log of the likelihood in more detail. Write μ for the mean of the normal distribution and θ for the unknown standard deviation of the normal distribution. We get

$$\log \mathcal{L}(\theta) = \left(\sum_{i=1}^N -\frac{(x_i - \mu)^2}{2\theta^2} \right) - N \log \theta + \text{Term not depending on } \theta$$

We can find the maximum by differentiating wrt θ and setting to zero, which yields

$$\frac{\partial \log \mathcal{L}(\theta)}{\partial \theta} = \frac{-2}{\theta^3} \sum_{i=1}^N \frac{-(x_i - \mu)^2}{2} - \frac{N}{\theta} = 0$$

so the maximum likelihood estimate is

$$\hat{\theta} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

which probably isn't all that surprising, either.

You should notice that one could maximize the likelihood of a normal distribution with respect to mean and standard deviation in one go (i.e. I could have done worked examples 15.6 and 15.7 in one worked example, instead of two). I did this example in two parts because I felt it was more accessible that way; if you object, you're likely to be able to fill in the details yourself very easily.

The maximum likelihood principle has a variety of neat properties we cannot expound. One worth knowing about is **consistency**; for our purposes, this means that the maximum likelihood estimate of parameters can be made arbitrarily close to the right answer by having a sufficiently large dataset. Now assume that our data doesn't actually come from the underlying model. This is the usual case, because we can't usually be sure that, say, the data truly is normal or truly comes from a Poisson distribution. Instead we choose a model that we think will be useful. When the data doesn't come from the model, maximum likelihood produces an estimate of θ that corresponds to the model that is (in quite a strong sense, which we can't explore here) the closest to the source of the data. Maximum likelihood is very widely used because of these neat properties. But there are some difficulties.

15.1.2 Cautions about Maximum Likelihood

One important problem is that it might be hard to find the maximum of the likelihood exactly. There are strong numerical methods for maximizing functions, and

these are very helpful, but even today there are likelihood functions where it is very hard to find the maximum.

The second is that small amounts of data can present nasty problems. For example, in the binomial case, if we have only one flip we will estimate p as either 1 or 0. We should find this report unconvincing. In the geometric case, with a fair coin, there is a probability 0.5 that we will perform the estimate and then report that the coin has $p = 1$. This should also worry you. As another example, if we throw a die only a few times, we could reasonably expect that, for some i , $n_i = 0$. This doesn't necessarily mean that $p_i = 0$, though that's what the maximum likelihood inference procedure will tell us.

This creates a very important technical problem — how can I estimate the probability of events that haven't occurred? This might seem like a slightly silly question to you, but it isn't. Solving this problem has really significant practical consequences. As one example, consider a biologist trying to count the number of butterfly species on an island. The biologist catches and classifies a lot of butterflies, then leaves. But are there more butterfly species on the island? To get some sense that we can reason successfully about this problem, compare two cases. In the first, the biologist catches many individuals of each of the species observed. In this case, you should suspect that catching more butterflies is unlikely to yield more species. In the second case, there are many species where the biologist sees only one individual of that species. In this case, you should suspect that catching more butterflies might very well yield new species.

As another example, a really important part of natural language processing involves estimating the probability of groups of three words. These groups are usually known as “trigrams”. People typically know an awful lot of words (tens to hundreds of thousands, depending on what you mean by a word). This means that there are a tremendous number of trigrams, and you can expect that any real dataset lacks almost all of them, because the dataset isn't big enough for there to be even just one of each trigram. Some are missing because they don't occur in real life, but others are not there simply because they are unusual (eg “Atom Heart Mother” actually occurs in real life, but you may not have seen it — try a web search if the phrase doesn't ring a bell). Modern speech recognition systems need to know how probable *every* trigram is. If the speech system thinks a trigram has zero probability and the trigram actually occurs, the system will make a mistake. We can't solve this problem just by giving each trigram a very small non-zero probability, because there are too many trigrams — it is important to distinguish between rare ones, and ones that don't ever occur. But what probability should I use for a rare trigram? Maximum likelihood would say use zero, but this would generate problems. Formalizing all this gets difficult quite quickly,

15.2 INCORPORATING PRIORS WITH BAYESIAN INFERENCE

Another important issue with maximum likelihood is that there is no mechanism to incorporate prior beliefs. For example, imagine you get a new die from a reliable store, roll it six times and see a one once. You would be happy to believe that $p(6) = 1/6$ for this die. Now imagine you borrow a die from a friend with a long history of making weighted dice. Your friend tells you this die is weighted so that

$p(1) = 1/2$. You roll the die six times and see a one once; in this case, you might worry that $p(6)$ isn't $1/6$, and you just happened to get a slightly unusual set of rolls. You'd worry because you have good reason to believe the die isn't fair, and you'd want more evidence to believe $p(6) = 1/6$. Maximum likelihood can't distinguish between these two cases.

The difference lies in prior information — information we possess before we look at the data. We would like to take this information into account when we estimate the model. One way to do so is to place a **prior probability distribution** $p(\theta)$ on the parameters θ . Then, rather than working with the likelihood $p(\mathcal{D}|\theta)$, we could apply Bayes' rule, and form the **posterior** $p(\theta|\mathcal{D})$. This posterior represents the probability that θ takes various values, given the data \mathcal{D} .

Definition: 15.4 *Bayesian inference*

Extracting information from the posterior $p(\theta|\mathcal{D})$ is usually called **Bayesian inference**

Definition: 15.5 *MAP estimate*

A natural estimate of θ is the value that maximizes the posterior $p(\theta|\mathcal{D})$. This estimate is known as a **maximum a posteriori estimate** or **MAP estimate**.

15.2.1 Constructing the Posterior

Bayes' rule tells us that

$$p(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$$

but (as we shall see) it can be hard to work out $P(\mathcal{D})$. For some problems, we might not need to know it.

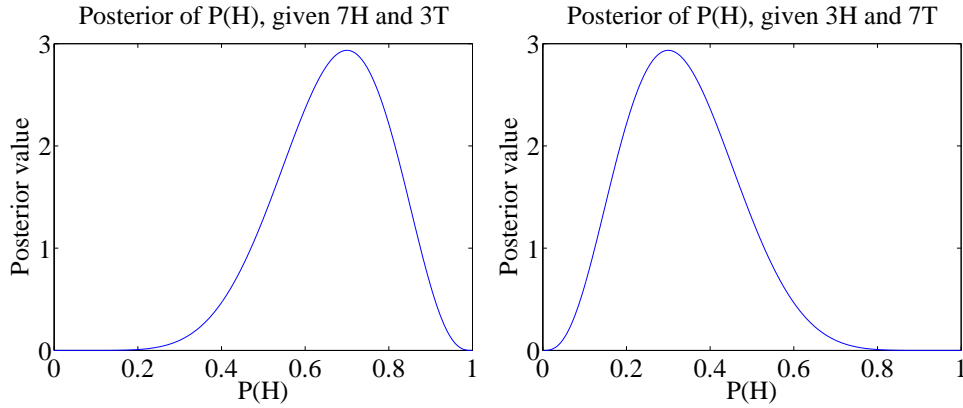


FIGURE 15.1: The curves show a function proportional to the posterior on θ , for the two cases of example ???. Notice that this information is rather richer than the single value we would get from maximum likelihood inference.

Worked example 15.8 *Flipping a coin*

We have a coin with probability θ of coming up heads when flipped. We start knowing nothing about θ . We then flip the coin 10 times, and see 7 heads (and 3 tails). Plot a function proportional to $p(\theta | \{7 \text{ heads and } 3 \text{ tails}\})$. What happens if there are 3 heads and 7 tails?

Solution: We know nothing about p , except that $0 \leq \theta \leq 1$, so we choose a uniform prior on p . We have that $p(\{7 \text{ heads and } 3 \text{ tails}\} | \theta)$ is binomial. The joint distribution is $p(\{7 \text{ heads and } 3 \text{ tails}\} | \theta) \times p(\theta)$ but $p(\theta)$ is uniform, so doesn't depend on θ . So the posterior is *proportional to*: $\theta^7(1 - \theta)^3$ which is graphed in figure 15.1. The figure also shows $\theta^3(1 - \theta)^7$ which is *proportional to* the posterior for 3 heads and 7 tails. In each case, the evidence does not rule out the possibility that $\theta = 0.5$, but tends to discourage the conclusion. Maximum likelihood would give $\theta = 0.7$ or $\theta = 0.3$, respectively.

In Example 15.8, it is interesting to follow how the posterior on p changes as evidence come in, which is easy to do because the posterior is proportional to a binomial distribution. Figure 15.2 shows a set of these posteriors for different sets of evidence.

For other problems, we will need to marginalize out θ , by computing

$$P(\mathcal{D}) = \int_{\theta} P(\mathcal{D}|\theta)P(\theta)d\theta.$$

It is usually impossible to do this in closed form, so we would have to use a numerical integral or a trick. The next section expounds one useful trick.

Conjugate priors

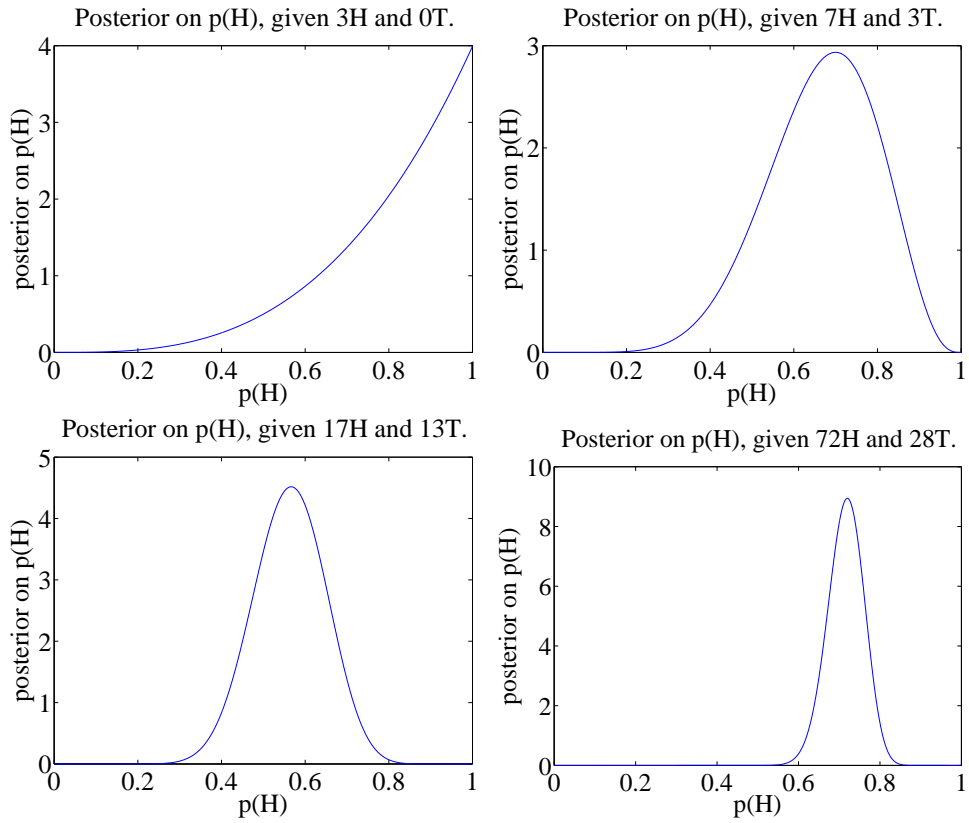


FIGURE 15.2: The probability that an unknown coin will come up heads when flipped is $p(H)$. For these figures, I simulated coin flips from a coin with $p = 0.75$. I then plotted the posterior for various data. Notice how, as we see more flips, we get more confident about p . The vertical axis changes substantially between plots in this figure.

In some cases, $P(\theta)$ and $P(\mathcal{D}|\theta)$, when multiplied together, take a familiar form. This happens when $P(\mathcal{D}|\theta)$ and $P(\theta)$ each belong to parametric families where there is a special relationship between the families. When a prior has this property, it is called a **conjugate** prior. There are some cases worth knowing, given in the worked examples.

Worked example 15.9 *Flipping a coin - II*

We have a coin with probability θ of coming up heads when flipped. We model the prior on θ with a Beta distribution, with parameters $\alpha > 0$, $\beta > 0$. We then flip the coin N times, and see h heads. What is $P(\theta|N, h, \alpha, \beta)$?

Solution: We have that $P(N, h|\theta)$ is binomial, and that $P(\theta|N, h, \alpha, \beta) \propto P(N, h|\theta)P(\theta|\alpha, \beta)$. This means that

$$P(\theta|N, h, \alpha, \beta) \propto \binom{N}{h} \theta^h (1-\theta)^{(N-h)} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}.$$

and we can write

$$P(\theta|N, h, \alpha, \beta) \propto \theta^{(\alpha+h-1)} (1-\theta)^{(\beta+N-h-1)}.$$

Notice this has the form of a Beta distribution, so it is easy to recover the constant of proportionality. We have

$$P(\theta|N, h, \alpha, \beta) = \frac{\Gamma(\alpha+\beta+N)}{\Gamma(\alpha+h)\Gamma(\beta+N-h)} \theta^{(\alpha+h-1)} (1-\theta)^{(\beta+N-h-1)}.$$

Worked example 15.10 *More swearsy politicians*

Example 15.5 gives some data from a swearsy politician. Assume we have only the first 10 intervals of observations, and we wish to estimate the intensity using a Poisson model. Write θ for this parameter. Use a Gamma distribution as a prior, and write out the posterior.

Solution: We have that

$$\begin{aligned} p(\mathcal{D}|\theta) &= \left(\frac{\theta^0 e^{-\theta}}{0!}\right)^5 \left(\frac{\theta^1 e^{-\theta}}{1!}\right)^2 \times \\ &\quad \left(\frac{\theta^2 e^{-\theta}}{2!}\right)^2 \left(\frac{\theta^3 e^{-\theta}}{3!}\right)^1 \\ &= \frac{\theta^9 e^{-10\theta}}{12} \end{aligned}$$

and

$$p(\theta|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{(\alpha-1)} e^{-\beta\theta}$$

This means that

$$p(\theta|\mathcal{D}) \propto \theta^{(\alpha-1+9)} e^{-(\beta+10)\theta}.$$

Notice this has the form of another Gamma distribution, so we can write

$$p(\theta|\mathcal{D}) = \frac{(\beta+10)^{(\alpha+9)}}{\Gamma(\alpha+9)} \theta^{(\alpha-1+9)} e^{-(\beta+10)\theta}$$

15.2.2 Normal Prior and Normal Likelihood

When both $P(\mathcal{D}|\theta)$ and $P(\theta)$ are normal, some important simplifications occur. First, the prior is conjugate to the likelihood. Second, the posterior is also normal. And third, the mean and standard deviation of the posterior take a simple form.

We start with an example. Assume we drop a measuring device down a borehole. It is designed to stop falling and catch onto the side of the hole after it has fallen μ_0 meters. On board is a device to measure its depth. This device reports a known constant times the correct depth plus a zero mean normal random variable, which we call “noise”. The device reports depth every second.

The first question to ask is what depth do we believe the device is at *before* we receive any measurement? We designed the device to stop at μ_0 meters, so we are not completely ignorant about where it is. However, it may not have worked absolutely correctly. We choose to model the depth at which it stops as μ_0 meters plus a zero mean normal random variable. The second term could be caused by error in the braking system, etc. We could estimate the standard deviation of the second term (which we write σ_0) either by dropping devices down holes, then measuring

with tape measures, or by analysis of likely errors in our braking system. The depth of the object is the unknown parameter of the model; we write this depth θ . Now the model says that θ is a normal random variable with mean μ_0 and standard deviation σ_0 .

Notice that this model probably isn't exactly right — for example, there must be some probability in the model that the object falls beyond the bottom of the hole, which it can't do — but it captures some important properties of our system. The device should stop at or close to μ_0 meters most of the time, and it's unlikely to be too far away.

Now assume we receive a single measurement — what do we now know about the device's depth? The first thing to notice is that there is something to do here. Ignoring the prior and taking the measurement might not be wise. For example, imagine that the noise in the wireless system is large, so that the measurement is often corrupted — our original guess about the device's location might be better than the measurement. Write x_1 for the measurement. Notice that the scale of the measurement may not be the same as the scale of the depth, so the mean of the measurement is $c_1\theta$, where c_1 is a change of scale (for example, from inches to meters). We have that $p(x_1|\theta)$ is normal with mean $c_1\theta$ and standard deviation σ_{n1} . We would like to know $p(\theta|x_1)$.

We have that

$$\begin{aligned} \log p(\theta, x_1) &= \log p(x_1|\theta) + \log p(\theta) \\ &= -\frac{(x_1 - c_1\theta)^2}{2\sigma_{n1}^2} - \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \\ &\quad + \text{terms not depending on } \theta \text{ or } x. \end{aligned}$$

We have two estimates of the position, θ , and we wish to come up with a representation of what we know about θ . One is x_1 , which is a *measurement* — we know its value. The expected value of x_1 is $c_1\theta$, so we could infer θ from x_1 . But we have another estimate of the position, which is μ_0 . The posterior, $p(\theta|x_1)$, is a probability distribution on the variable θ ; it depends on the known values x_1 , μ_0 , σ_0 and σ_{n1} . We need to determine its form. We can do so by some rearrangement of the expression for $\log p(\theta, x_1)$.

Notice first that this expression is of degree 2 in θ (i.e. it has terms θ^2 , θ and things that don't depend on θ). This means that $p(\theta|x_1)$ must be a normal distribution, because we can rearrange its log into the form of the log of a normal distribution. This yields a fact of crucial importance.

Useful Fact: 15.1 *Normal distributions are conjugate*

A normal prior and a normal likelihood yield a normal posterior.

Write μ_1 for the mean of this distribution, and σ_1 for its standard deviation. The log of the distribution must be

$$-\frac{(\theta - \mu_1)^2}{2\sigma_1^2} + \text{terms not depending on } \theta.$$

The terms not depending on θ are not interesting, because if we know σ_1 those terms must add up to

$$\log\left(\frac{1}{\sqrt{2\pi}\sigma_1}\right)$$

so that the probability density function sums to one. Our goal is to rearrange terms into the form above. Notice that

$$-\frac{(\theta - \mu_1)^2}{2\sigma_p^2} = -\theta^2 \left(\frac{1}{2\sigma_1^2}\right) + 2\theta \frac{\mu_1}{2\sigma_p^2} + \text{term not depending on } \theta$$

We have

$$\begin{aligned} \log p(\theta|x_1) &= -\frac{(c_1\theta - x_1)^2}{2\sigma_{n1}^2} - \frac{(\theta - \mu_0)^2}{2\sigma_0^2} + \text{terms not depending on } \theta \\ &= -\theta^2 \left(\frac{1}{2 \left(\frac{\sigma_{n1}^2 \sigma_0^2}{\sigma_{n1}^2 + c_1^2 \sigma_0^2} \right)} \right) + 2\theta \left(c_1 \frac{x_1}{2\sigma_{n1}^2} + \frac{\mu_0}{2\sigma_0^2} \right) \\ &+ \text{terms not depending on } \theta \end{aligned}$$

which means that

$$\sigma_1^2 = \frac{\sigma_{n1}^2 \sigma_0^2}{\sigma_{n1}^2 + c_1^2 \sigma_0^2}$$

and

$$\begin{aligned} \mu_1 &= 2 \left(c_1 \frac{x_1}{2\sigma_{n1}^2} + \frac{\mu_0}{2\sigma_0^2} \right) \frac{\sigma_{n1}^2 \sigma_0^2}{\sigma_{n1}^2 + c_1^2 \sigma_0^2} \\ &= \left(\frac{c_1 x_1 \sigma_0^2 + \mu_0 \sigma_{n1}^2}{\sigma_{n1}^2 \sigma_0^2} \right) \frac{\sigma_{n1}^2 \sigma_0^2}{\sigma_{n1}^2 + c_1^2 \sigma_0^2} \\ &= \frac{c_1 x_1 \sigma_0^2 + \mu_0 \sigma_{n1}^2}{\sigma_{n1}^2 + c_1^2 \sigma_0^2}. \end{aligned}$$

These equations “make sense”. Imagine that σ_0 is very small, and σ_{n1} is very big; then our new expected value of θ — which is μ_1 — is about μ_0 . Equivalently, because our prior was very accurate, and the measurement was unreliable, our expected value is about the prior value. Similarly, if the measurement is reliable (i.e. σ_{n1} is small) and the prior has high variance (i.e. σ_0 is large), then our expected value of θ is about x_1/c_1 — i.e. the measurement, rescaled. I have put these equations, in a more general form, in a box below.

Useful Fact: 15.2 *Normal posteriors*

Assume we wish to estimate a parameter θ . The prior distribution for θ is normal, with known mean μ_π and known standard deviation σ_π . We receive a single data item x . The likelihood of this data item is normal with mean $c\theta$ and standard deviation σ_m , where c and σ_m are known. Then the posterior, $p(\theta|x, c, \sigma_m, \mu_\pi, \sigma_\pi)$, is normal, with mean

$$\frac{cx\sigma_\pi^2 + \mu_\pi\sigma_m^2}{\sigma_m^2 + c^2\sigma_\pi^2}$$

and standard deviation

$$\sqrt{\frac{\sigma_m^2\sigma_\pi^2}{\sigma_m^2 + c^2\sigma_\pi^2}}.$$

15.2.3 MAP Inference

Look at example 15.1, where we estimated the probability a coin would come up heads with maximum likelihood. We could not change our estimate just by knowing the coin was fair, but we could come up with a number for $\theta = p(H)$ (rather than, say, a posterior distribution). A natural way to produce a point estimate for θ that incorporates prior information is to choose $\hat{\theta}$ such that

$$\hat{\theta} = \operatorname{argmax}_{\theta} P(\theta|\mathcal{D}) = \operatorname{argmax}_{\theta} \frac{P(\theta, \mathcal{D})}{P(\mathcal{D})}$$

This is the MAP estimate. If we wish to perform MAP inference, $P(\mathcal{D})$ doesn't matter (it changes the value, but not the location, of the maximum). This means we can work with $P(\theta, \mathcal{D})$, often called the **joint** distribution.

Worked example 15.11 *Flipping a coin - II*

We have a coin with probability θ of coming up heads when flipped. We model the prior on θ with a Beta distribution, with parameters $\alpha > 0$, $\beta > 0$. We then flip the coin N times, and see h heads. What is the MAP estimate of θ ?

Solution: We have that

$$P(\theta|N, h, \alpha, \beta) = \frac{\Gamma(\alpha + \beta + N)}{\Gamma(\alpha + h)\Gamma(\beta + N - h)} \theta^{(\alpha+h-1)}(1-\theta)^{(\beta+N-h-1)}.$$

You can get the MAP estimate by differentiating and setting to 0, yielding

$$\hat{\theta} = \frac{\alpha - 1 + h}{\alpha + \beta - 2 + N}.$$

This has rather a nice interpretation. You can see α and β as extra counts of heads (resp. tails) that are added to the observed counts. So, for example, if you were fairly sure that the coin should be fair, you might make α and β large and equal. When $\alpha = 1$ and $\beta = 1$, we have a uniform prior as in the previous examples.

Worked example 15.12 *More swears politicians*

We observe our swearing politician for N intervals, seeing n_i swear words in the i 'th interval. We model the swearing with a Poisson model. We wish to estimate the intensity, which we write θ . We use a Gamma distribution for the prior on θ . What is the MAP estimate of θ ?

Solution: Write $T = \sum_{i=1}^N n_i$. We have that

$$p(\theta|\mathcal{D}) = \frac{(\beta + N)^{(\alpha+T)}}{\Gamma(\alpha + T)} \theta^{(\alpha-1+T)} e^{-(\beta+T)\theta}$$

and the MAP estimate is

$$\hat{\theta} = \frac{(\alpha - 1 + T)}{(\beta + N)}$$

(which you can get by differentiating with respect to θ , then setting to zero). Notice that if β is close to zero, you can interpret α as extra counts; if β is large, then it strongly discourages large values of $\hat{\theta}$, even if the counts are large.

15.2.4 Filtering

We can make online estimates of the maximum likelihood value of mean and standard deviation for a normal distribution. Assume, rather than seeing N elements

of a dataset in one go, you get to see each one once, and you cannot store them. Assume that this dataset is modelled as normal data. Write $\hat{\mu}_k$ for the maximum likelihood estimate of the mean based on data items $1 \dots k$ (and $\hat{\sigma}_k$ for the maximum likelihood estimate of the standard deviation, etc.). Notice that

$$\hat{\mu}_{k+1} = \frac{(k\hat{\mu}_k) + x_{k+1}}{(k+1)}$$

and that

$$\hat{\sigma}_{k+1} = \sqrt{\frac{(k\hat{\sigma}_k^2) + (x_{k+1} - \hat{\mu}_{k+1})^2}{(k+1)}}$$

This means that you can incorporate new data into your estimate as it arrives without keeping all the data. This process of updating a representation of a dataset as new data arrives is known as **filtering**.

This is particularly useful in the case of normal posteriors. Recall that if we have a normal prior and a normal likelihood, the posterior is normal. Now consider a stream of incoming measurements. Our initial representation of the parameters we are trying to estimate is the prior, which is normal. We allow ourselves to see one measurement, which has normal likelihood; then the posterior is normal. You can think of this posterior as a prior for the parameter estimate based on the *next* measurement. But we know what to do with a normal prior, a normal likelihood, and a measurement; so we can incorporate the measurement and go again. This means we can exploit our expression for the posterior mean and standard deviation in the case of normal likelihood and normal prior and a single measurement to deal with multiple measurements very easily.

Assume a second measurement, x_2 arrives. We know that $p(x_2|\theta, c_2, \sigma_{n2})$ is normal with mean $c_2\theta$ and standard deviation σ_{n2} . In the example, we have a new measurement of depth — perhaps in a new, known, scale — with new noise (which might have larger, or smaller, standard deviation than the old noise) added. Then we can use $p(\theta|x_1, c_1, \sigma_{n1})$ as a prior to get a posterior $p(\theta|x_1, x_2, c_1, c_2, \sigma_{n1}, \sigma_{n2})$. Each is normal, by useful fact 15.1. Not only that, but we can easily obtain the expressions for the mean μ_2 and the standard deviation σ_2 *recursively* as functions of μ_1 and σ_1 .

Applying useful fact 15.2, we have

$$\mu_2 = \frac{c_2x_2\sigma_1^2 + \mu_1\sigma_{n2}^2}{\sigma_{n2}^2 + c_2^2\sigma_1^2}$$

and

$$\sigma_2^2 = \frac{\sigma_{n2}^2\sigma_1^2}{\sigma_{n2}^2 + c_2^2\sigma_1^2}.$$

But what works for 2 and 1 will work for $k+1$ and k . We know the posterior after k measurements will be normal, with mean μ_k and standard deviation σ_k . The $k+1$ 'th measurement x_{k+1} arrives, and we have $p(x_{k+1}|\theta, c_{k+1}, \sigma_{n(k+1)})$ is normal. Then the posterior is normal, and we can write the mean μ_{k+1} and the standard deviation σ_{k+1} recursively as functions of μ_k and σ_k . The result is worth putting in a box.

Useful Fact: 15.3 *Online updating of normal posteriors*

Assume we wish to estimate a parameter θ . The prior distribution for θ is normal, with known mean μ_π and known standard deviation σ_π . All data is normal conditioned on θ . We have already received k data items. The posterior $p(\theta|x_1, \dots, x_k, c_1, \dots, c_k, \sigma_{n1}, \dots, \sigma_{nk}, \mu_\pi, \sigma_\pi)$ is normal, with mean μ_k and standard deviation σ_k . We receive a new data item x_{k+1} . The likelihood of this data item is normal with mean $c\theta$ and standard deviation $\sigma_{n(k+1)}$, where c_{k+1} and $\sigma_{n(k+1)}$ are known. Then the posterior, $p(\theta|x_1, \dots, x_{k+1}, c_1, \dots, c_k, c_{k+1}, \sigma_{n1}, \dots, \sigma_{n(k+1)}, \mu_\pi, \sigma_\pi)$, is normal, with mean

$$\mu_{k+1} = \frac{c_{k+1}x_{k+1}\sigma_k^2 + \mu_k\sigma_{n(k+1)}^2}{\sigma_{n(k+1)}^2 + c_{k+1}^2\sigma_k^2}$$

and

$$\sigma_{k+1}^2 = \frac{\sigma_{n(k+1)}^2\sigma_k^2}{\sigma_{n(k+1)}^2 + c_{k+1}^2\sigma_k^2}.$$

Again, notice the very useful fact that, if everything is normal, we can update our posterior representation when new data arrives using a very simple recursive form.

Worked example 15.13 *Normal data*

Assume you see N datapoints x_i which are modelled by a normal distribution with unknown mean θ and with known standard deviation σ . You model the prior on θ using a normal distribution with mean μ_0 and standard deviation σ_0 . What is the MAP estimate of the mean?

Solution: Recall that the maximum value of a normal distribution occurs at its mean. Now problem is covered by useful fact 15.2, but in this case we have $c_i = 1$ for each data point, and $\sigma_i = \sigma$. We can write

$$\mu_N = \frac{x_N \sigma_{N-1}^2 + \mu_{N-1} \sigma^2}{\sigma^2 + \sigma_{N-1}^2}$$

and

$$\sigma_N^2 = \frac{\sigma^2 \sigma_{N-1}^2}{\sigma^2 + \sigma_{N-1}^2}.$$

and evaluate the recursion down to μ_0, σ_0 .

15.2.5 Cautions about Bayesian Inference

Just like maximum likelihood inference, bayesian inference is not a recipe that can be applied without thought. It turns out that, when there is a lot of data, the prior has little inference on the outcome of the inference, and the MAP solution looks a lot like the maximum likelihood solution. So the difference between the two approaches is most interesting when there is little data, where the prior matters. The difficulty is that it might be hard to know what to use as a good prior. In the examples, I emphasized mathematical convenience, choosing priors that lead to clean posteriors. There is no reason to believe that nature uses conjugate priors (even though conjugacy is a neat property). How should one choose a prior for a real problem?

This isn't an easy point. If there is little data, then the choice could really affect the inference. Sometimes we're lucky, and the logic of the problem dictates a choice of prior. Mostly, we have to choose and live with the consequences of the choice. Often, doing so is succesful in applications.

The fact we can't necessarily justify a choice of prior seems to be one of life's inconveniences, but it represents a significant philosophical problem. It's been at the core of a long series of protracted, often quite intense, arguments about the philosophical basis of statistics. I haven't followed these arguments closely enough to summarize them; they seem to have largely died down without any particular consensus being reached.

15.3 SAMPLES, URNS AND POPULATIONS

Very often the data we see is a small part of the data we could have seen, if we'd been able to collect enough data. We need to know how the measurements we make

on the dataset relate to the measurements we could have made, if we had all the data. This situation occurs very often. For example, imagine we wish to know the average weight of a rat. This isn't random; you could weigh every rat on the planet, and then average the answers. But doing so would be absurd (among other things, you'd have to weigh them all at the same time, which would be tricky). Instead, we weigh a small set of rats, chosen rather carefully. If we have chosen sufficiently carefully, then the answer from the small set is quite a good representation of the answer from the whole set.

The data we could have observed, if we could have seen everything, is the **population**. The data we actually have is the **sample**. We would like to know the mean of the population, but can see only the sample; surprisingly, we can say a great deal from the sample alone, assuming that it is chosen appropriately.

15.3.1 Estimating the Population Mean from a Sample

Assume we have a population $\{x\}$, for $i = 1, \dots, N_p$. Notice the subscript here — this is the number of items in the population. The population could be unreasonably big: for example, it could consist of all the people in the world. We want to know the mean of this dataset, but we do not get to see the whole dataset. Instead, we see the sample.

How the sample is obtained is key to describing the population. We will focus on only one model (there are lots of others). In our model, the sample is obtained by choosing a fixed number of data items. Write k for the number of data items in the sample. We expect k is a lot smaller than N_p . Each item is chosen independently, and fairly. This means that each time we choose, we choose one from the entire set of N_p data items, and each has the same probability of being chosen. This is sometimes referred to as “sampling with replacement”.

One natural way to think about sampling with replacement is to imagine the data items as being written on tickets, which are placed in an urn (old-fashioned word for a jar, now used mainly by statisticians and morticians). You obtain the sample by repeating the following experiment k times: shake the urn; take a ticket from the urn and write down the data on the ticket; put it back in the urn. Notice that, in this case, each sample is drawn from the same urn. This is important, and makes the analysis easier. If we had not put the ticket back, the urn would change between samples.

We summarize the whole dataset with its mean, which we write $\text{popmean}(\{x\})$. This is known as the **population mean**. The notation is just to drive home the facts that it's the mean of the whole population, and that we don't, and can't, know it. The whole point of this exercise is to *estimate* this mean.

We would like to estimate the mean of the whole dataset from the items that we actually see. Imagine we draw k tickets from the urn as above, and average the values. The result is a random variable, because different draws of k tickets will give us different values. Write $X^{(k)}$ for this random variable, which is referred to as the **sample mean**. Because expectations are linear, we must have that

$$\mathbb{E}[X^{(k)}] = \frac{1}{k} \left(\mathbb{E}[X^{(1)}] + \dots + \mathbb{E}[X^{(1)}] \right) = \mathbb{E}[X^{(1)}]$$

(where $X^{(1)}$ is the random variable whose value is obtained by drawing one ticket

from the urn). Now

$$\begin{aligned}\mathbb{E}[X^{(1)}] &= \sum_{i \in 1, \dots, N_p} x_i p(i) \\ &= \sum_{i \in 1, \dots, N_p} x_i \frac{1}{N_p} && \text{because we draw fairly from the urn} \\ &= \frac{\sum_{i \in 1, \dots, N_p} x_i}{N_p} \\ &= \text{popmean}(\{x\})\end{aligned}$$

which is the mean value of the items in the urn. This means that

$$\mathbb{E}[X^{(k)}] = \text{popmean}(\{x_i\}).$$

Under our sampling model, the expected value of the sample mean is the population mean.

Useful Facts: 15.4 *Sample means and population means*

The sample mean is a random variable. It is random, because different samples from the population will have different values of the sample mean. The expected value of this random variable is the population mean.

We will not get the same value of $X^{(k)}$ each time we perform the experiment, because we see different data items in each sample. So $X^{(k)}$ has variance, and this variance is important. If it is large, then each estimate is quite different. If it is small, then the estimates cluster. Knowing the variance of $X^{(k)}$ would tell us how accurate our estimate of the population mean is.

15.3.2 The Variance of the Sample Mean

We write $\text{popstd}(\{x\})$ for the standard deviation of the whole population of $\{x\}$. Again, we write it like this to keep track of the facts that (a) it's for the whole population and (b) we don't — and usually can't — know it.

We can compute the variance of $X^{(k)}$ (the sample mean) easily. We have

$$\text{var}[X^{(k)}] = \mathbb{E}[(X^{(k)})^2] - \mathbb{E}[X^{(k)}]^2 = \mathbb{E}[(X^{(k)})^2] - (\text{popmean}(\{x\}))^2$$

so we need to know $\mathbb{E}[(X^{(k)})^2]$. We can compute this by writing

$$X^{(k)} = \frac{1}{k}(X_1 + X_2 + \dots + X_k)$$

where X_1 is the value of the first ticket drawn from the urn, etc. We then have

$$X^{(k)2} = \left(\frac{1}{k}\right)^2 \begin{pmatrix} X^2_1 + X^2_2 + \dots X^2_k + X_1X_2 + \dots \\ X_1X_k + X_2X_1 + \dots X_2X_k + \dots X_{k-1}X_k \end{pmatrix}.$$

Expectations are linear, so we have that

$$\mathbb{E}\left[(X^{(k)})^2\right] = \left(\frac{1}{k}\right)^2 \begin{pmatrix} \mathbb{E}[X^2_1] + \mathbb{E}[X^2_2] + \dots \mathbb{E}[X^2_k] + \mathbb{E}[X_1X_2] + \\ \dots \mathbb{E}[X_1X_k] + \mathbb{E}[X_2X_1] + \dots \mathbb{E}[X_2X_k] + \dots \mathbb{E}[X_{k-1}X_k] \end{pmatrix}.$$

The *order* in which the tickets are drawn from the urn doesn't matter, because each time we draw a ticket we draw from the same urn. This means that $\mathbb{E}[X^2_1] = \mathbb{E}[X^2_2] = \dots \mathbb{E}[X^2_k]$. You can think of this term as the expected value of the random variable generated by: drawing a single number out of the urn; squaring that number; and reporting the square. Notice that $\mathbb{E}[X^2_1] = \mathbb{E}[(X^{(1)})^2]$ (look at the definition of $X^{(1)}$).

Because the *order* doesn't matter, we also have that $\mathbb{E}[X_1X_2] = \mathbb{E}[X_1X_3] = \dots \mathbb{E}[X_{k-1}X_k]$. You can think of this term as the expected value of the random variable generated by: drawing a number out of the urn; writing it down; returning it to the urn; then drawing a second number from the urn; and reporting the product of these two numbers. So we can write

$$\mathbb{E}\left[X^{(k)2}\right] = \left(\frac{1}{k}\right)^2 \left(k\mathbb{E}\left[(X^{(1)})^2\right] + k(k-1)\mathbb{E}[X_1X_2]\right)$$

and these two terms are quite easy to evaluate.

Worked example 15.14 *Urn variances*

Show that

$$\mathbb{E}\left[(X^{(1)})^2\right] = \frac{\sum_{i=1}^{N_p} x_i^2}{N_p} = \text{popstd}(\{x\})^2 + \text{popmean}(\{x\})^2$$

Solution: First, we have $(X^{(1)})^2$ is the number obtained by taking a ticket out of the urn and squaring its data item. Now

$$\text{popstd}(\{x\})^2 = \mathbb{E}\left[(X^{(1)})^2\right] - \mathbb{E}\left[X^{(1)}\right]^2 = \mathbb{E}\left[(X^{(1)})^2\right] - \text{popmean}(\{x\})^2$$

so

$$\mathbb{E}\left[(X^{(1)})^2\right] = \text{popstd}(\{x\})^2 + \text{popmean}(\{x\})^2$$

Worked example 15.15 *Urn variances*

Show that

$$\mathbb{E}[X_1 X_2] = \text{popmean}(\{x\})^2$$

Solution: This looks hard, but isn't. Recall from the facts in chapter ?? (useful facts ??, page ??) that if X and Y are independent random variables, $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$. But X_1 and X_2 are independent — they are different random draws from the same urn. So

$$\mathbb{E}[X_1 X_2] = \mathbb{E}[X_1]\mathbb{E}[X_2]$$

but $\mathbb{E}[X_1] = \mathbb{E}[X_2]$ (they are draws from the same urn) and $\mathbb{E}[X] = \text{popmean}(\{x\})$. So

$$\mathbb{E}[X_1 X_2] = \text{popmean}(\{x\})^2.$$

Now

$$\begin{aligned} \mathbb{E}[(X^{(k)})^2] &= \frac{k\mathbb{E}[(X^{(1)})^2] + k(k-1)\mathbb{E}[X_1 X_2]}{k^2} \\ &= \frac{\mathbb{E}[(X^{(1)})^2] + (k-1)\mathbb{E}[X_1 X_2]}{k} \\ &= \frac{(\text{popsd}(\{x\})^2 + \text{popmean}(\{x\})^2) + (k-1)\text{popmean}(\{x\})^2}{k} \\ &= \frac{\text{popsd}(\{x\})^2}{k} + \text{popmean}(\{x\})^2 \end{aligned}$$

so we have

$$\begin{aligned} \text{var}[X^{(k)}] &= \mathbb{E}[(X^{(k)})^2] - \mathbb{E}[X^{(k)}]^2 \\ &= \frac{\text{popsd}(\{x\})^2}{k} + \text{popmean}(\{x\})^2 - \text{popmean}(\{x\})^2 \\ &= \frac{\text{popsd}(\{x\})^2}{k}. \end{aligned}$$

This is a very useful result which is well worth remembering together with our facts on the sample mean, so we'll put them in a box together.

Useful Fact: 15.5 *The sample mean*

The sample mean is a random variable. Write $X^{(k)}$ for the mean of k samples. We have that:

$$\begin{aligned}\mathbb{E}[X^{(k)}] &= \text{popmean}(\{x\}) \\ \text{var}[X^{(k)}] &= \frac{\text{popsd}(\{x\})^2}{k} \\ \text{std}(\{X^{(k)}\}) &= \frac{\text{popsd}(\{x\})}{\sqrt{k}}\end{aligned}$$

The consequence is this: If you draw k samples, the standard deviation of your estimate of the mean is

$$\frac{\text{popsd}(\{x\})}{\sqrt{k}}$$

which means that (a) the more samples you draw, the better your estimate becomes and (b) the estimate improves rather slowly — for example, to halve the standard deviation in your estimate, you need to draw four times as many samples. The standard deviation of the estimate of the mean is often known as the **standard error** of the mean. This allows us to draw a helpful distinction: the population has a standard deviation, and our estimate of its mean (or other things — but we won't go into this) has a standard error.

Notice we cannot state the standard error of our estimate exactly, because we do not know $\text{popsd}(\{x\})$. But we could make a good estimate of $\text{popsd}(\{x\})$, by computing the standard deviation of the examples that we have. It is now helpful to have some notation for the particular sample we have. I will write $\sum_{i \in \text{sample}}$ for a sum over the sample items, and we will use

$$\text{mean}(\{x\}) = \frac{\sum_{i \in \text{sample}} x_i}{k}$$

for the mean of the sample — that is, the mean of the data we actually see; this is consistent with our old notation, but there's a little reindexing to keep track of the fact we don't see all of the population. Similarly, I will write

$$\text{std}(\{x\}) = \sqrt{\frac{\sum_{i \in \text{sample}} (x_i - \text{mean}(\{x_i\}))^2}{k}}$$

for the sample standard deviation. Again, this is the standard deviation of the data we actually see; and again, this is consistent with our old notation, again with a

little reindexing to keep track of the fact we don't see all of the population. We could estimate

$$\text{popstd}(\{x\}) \approx \text{std}(\{x\})$$

and as long as we have enough examples, this estimate is good. If the number of samples k is small, it is better to use

$$\text{popstd}(\{x\}) \approx \sqrt{\frac{\sum_{i \in \text{sample}} (x_i - \text{mean}(\{x\}))^2}{k-1}}.$$

In fact, much more is known about the distribution of $X^{(k)}$.

15.3.3 The Probability Distribution of the Sample Mean

The sample mean is a random variable. We know an expression for its mean, and we can estimate its variance. In fact, we can determine its probability distribution, though I won't do this rigorously. In section 14.4.3, I mentioned that adding a number of independent random variables almost always got you a normal random variable, a fact sometimes known as the central limit theorem. I didn't prove it, and I'm not going to now. But when we form $X^{(k)}$, we're adding random variables. This means that $X^{(k)}$ is a normal random variable, for sufficiently big k (for some reason, $k > 30$ is usually seen as right).

This is important, because it has the following consequence. Draw a large number of different samples of k elements from the population. Each is a dataset of k items. Compute $\text{mean}(\{x\})$ for each, and regard the resulting numbers e_1, \dots, e_r as data items. Convert the e_i to standard coordinates s_i , where

$$s_i = \frac{(e_i - \text{mean}(\{e_i\}))}{\text{std}(e_i)}$$

(i.e. by subtracting the mean of the e_i , and dividing by their standard deviation). Now construct a histogram of the s . If r is sufficiently large, the histogram will be close to the standard normal curve.

15.3.4 When The Urn Model Works

In our model, there was a population of N_p data items x_i , and we saw k of them, chosen at random. In particular, each choice was fair (in the sense that each data item had the same probability of being chosen) and independent. These assumptions are very important for our analysis to apply. If our data does not have these properties, bad things can happen.

For example, assume we wish to estimate the percentage of the population that has beards. This is a mean (the data items take the value 1 for a person with a beard, and 0 without a beard). If we select people according to our model, then ask them whether they have a beard, then our estimate of the percentage of beards should behave as above.

The first thing that should strike you is that it isn't at all easy to select people according to this model. For example, we might select phone numbers at random, then call and ask the first person to answer the phone whether they have a beard;

but many children won't answer the phone because they are too small. The next important problem is that errors in selecting people can lead to massive errors in your estimate. For example, imagine you decide to survey all of the people at a kindergarten on a particular day; or all of the people in a women's clothing store; or everyone attending a beard growing competition (they do exist). In each case, you will get an answer that is a very poor estimate of the right answer, and the standard error might look very small. Of course, it is easy to tell that these cases are a bad choice.

It may not be easy to tell what a good choice is. You should notice the similarity between estimating the percentage of the population that wears a beard, and estimating the percentage that will vote for a particular candidate. There is a famous example of a survey that mispredicted the result of the Dewey-Truman presidential election in 1948; poll-takers phoned random phone numbers, and asked for an opinion. But at that time, telephones tended to be owned by a small percentage of rather comfortable households, who tended to prefer one candidate, and so the polls mispredicted the result rather badly.

Sometimes, we don't really have a choice of samples. For example, we might be presented with a small dataset of (say) human body temperatures. If we can be satisfied that the people were selected rather randomly, we might be able to use this dataset to predict expected body temperature. But if we knew that the subjects had their temperatures measured because they presented themselves at the doctor with a suspected fever, then we most likely cannot use it to predict expected body temperature.

One important and valuable case where this model works is in simulation. If you can guarantee that your simulations are independent (which isn't always easy), this model applies to estimates obtained from a simulation. Notice that it is usually straightforward to build a simulation so that the i 'th simulation reports an x_i where $\text{popmean}(\{x\})$ gives you the thing you want to measure. For example, imagine you wish to measure the probability of winning a game; then the simulation should report one when the game is won, and zero when it is lost. As another example, imagine you wish to measure the expected number of turns before a game is won; then your simulation should report the number of turns elapsed before the game was won.

15.4 YOU SHOULD

15.4.1 be able to:

- Write out the likelihood for a set of independent data items produced by models from chapter 14 (at least Normal, Binomial, Multinomial, Poisson, Beta, Gamma, Exponential).
- Write out the log likelihood for a set of independent data items produced by models from chapter 14 (at least Normal, Binomial, Multinomial, Poisson, Beta, Gamma, Exponential).
- Find maximum likelihood solutions for parameters of these models from a set of independent data items.
- Describe situations where maximum likelihood estimates might not be reliable.
- Describe the difference between maximum likelihood estimation and Bayesian inference.
- Write an expression for the posterior or log-posterior of model parameters given a set of independent data items.
- Compute the MAP estimate for the cases shown in the worked examples.
- Compute on-line estimates of the maximum likelihood estimate of the mean and standard deviation of a normal model.
- Compute on-line estimates of the MAP estimate of the mean and standard deviation in the case of a normal prior and a normal likelihood.
- Estimate the population mean from a sample mean.
- Estimate the standard error of the estimate of a population mean.

15.4.2 remember:

New term: Inference	347
New term: point estimate	347
New term: interval estimates	347
New term: hypothesis testing	347
New term: likelihood	349
Definition: Likelihood	349
New term: maximum likelihood principle	349
Definition: The maximum likelihood principle	350
New term: independent and identically distributed	350
New term: IID	350
Definition: The log-likelihood of a dataset under a model	354
New term: consistency	357
New term: prior probability distribution	359
New term: posterior	359

New term: Bayesian inference 359
 Definition: Bayesian inference 359
 New term: maximum a posteriori estimate 359
 New term: MAP estimate 359
 Definition: MAP estimate 359
 New term: conjugate 361
 Useful fact: Normal distributions are conjugate 364
 Useful fact: Normal posteriors 366
 New term: joint 366
 New term: filtering 368
 Useful fact: Online updating of normal posteriors 369
 New term: population 371
 New term: sample 371
 New term: population mean 371
 New term: sample mean 371
 Useful facts: Sample means and population means 372
 Useful fact: The sample mean 375
 New term: standard error 375

Markov Chains and Simulation

There are many situations where one must work with a sequence of random variables. For example, consider a bus queue; people arrive at random, and so do buses. What is the probability that the queue gets to a particular length? and what is the expected length of the queue? As another example, we could have a random model of purchases from a shop — under any particular rule for restoring inventory, what is the largest (resp. smallest) amount of stock on the shelf? what is the expected amount of stock on the shelf? It turns out that there is a simple and powerful model that applies to many sequences of random variables. One can often use this model to make closed-form predictions. In cases where closed-form predictions aren't available, one can use simulation methods to estimate probabilities and expectations.

16.1 MARKOV CHAINS

A Markov chain is a sequence of random variables which has important independence properties. We will describe these properties in some detail below; first, we give some examples. Markov chains are easily represented with the language of finite state machines. For some Markov chains, it is easy to determine probabilities and expectations of interest in closed form using simple methods. For others, it is tricky (straightforward, but unreasonable amounts of straightforward work). In these cases, we can estimate the relevant probabilities and expectations by simulating the finite state machine.

16.1.1 Motivating Example: Multiple Coin Flips

We start with three examples, each of which is easy to work. Each suggests a much harder question, however, which we need new machinery to handle.

Worked example 16.1 *Multiple Coin Flips - 1*

You choose to flip a fair coin until you see two heads in a row, and then stop. What is the probability that you flip the coin twice?

Solution: Because you stopped after two flips, you must have seen two heads. So $P(2 \text{ flips}) = P(\{HH\}) = P(\{H\})^2 = 1/4$.

Worked example 16.2 *Multiple Coin Flips - 2*

You choose to flip a fair coin until you see two heads in a row, and then stop. What is the probability that you flip the coin three times?

Solution: Because you stopped after three flips, you must have seen T , then H , then H ; any other sequence either doesn't stop, or stops too early. We write this with the last flip last, so THH . So $P(3 \text{ flips}) = P(\{THH\}) = P(\{T\})P(\{H\})^2 = 1/8$.

Worked example 16.3 *Multiple Coin Flips - 3*

You choose to flip a fair coin until you see two heads in a row, and then stop. What is the probability that you flip the coin four times?

Solution: This is more interesting. The last three flips must have been THH (otherwise you'd go on too long, or end too early). But, because the second flip must be a T , the first could be either H or T . This means there are two sequences that work: $HTHH$ and $TTHH$. So $P(4 \text{ flips}) = 2/8 = 1/4$.

The harder question here is to ask what is $P(N)$, for N some number (larger than 4, because we know the answers in the other cases). It is unattractive to work this case by case (you could try this, if you don't believe me). One very helpful way to think about the coin flipping experiment is as a finite state machine (Figure 16.1). If you think of this machine as a conventional finite state machine, it accepts any string of T , H , that (a) ends with HH , and (b) has no other HH in it. Alternatively, you could think of this machine as encoding a probabilistic process. At each state, an event occurs with some probability (in this case, a coin comes up H or T , with probability $(1/2)$). The event causes the machine to follow the appropriate edge. If we take this view, this machine stops when we have flipped two heads in succession. It encodes our problem of computing the probability of flipping a coin N times then stopping — this is just the probability that the machine hits the end state after N transitions.

It turns out to be straightforward to construct a recurrence relation for $P(N)$ (i.e. an equation for $P(N)$ in terms of $P(N-1)$, $P(N-2)$, and so on). This property is quite characteristic of repeated experiments. For this example, first, notice that $P(1) = 0$. Now imagine you have flipped the coin N times and stopped. We can assume that $N > 3$, because we know what happens for the other cases. The *last* three flips must have been THH . Write Σ_N for a sequence that is: (a) N flips long; (b) ends in HH ; and (c) contains no other subsequence HH . Equivalently, this is a string accepted by the finite state machine of figure 16.1. We must have that any

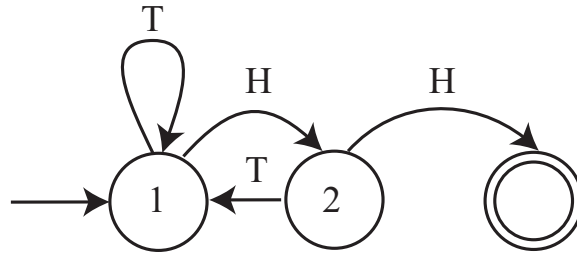


FIGURE 16.1: A finite state machine representing the coin flip example. By convention, the end state is a double circle, and the start state has an incoming arrow. I've labelled the arrows with the event that leads to the transition, but haven't bothered to put in the probabilities, because each is 0.5.

Σ_N has either the form $T\Sigma_{N-1}$ or the form $HT\Sigma_{N-2}$. But this means that

$$\begin{aligned} P(N) &= P(T)P(N-1) + P(HT)P(N-2) \\ &= (1/2)P(N-1) + (1/4)P(N-2) \end{aligned}$$

It is possible to solve this recurrence relation to get an explicit formula, but doing so would take us out of our way. You will check this recurrence relation in an exercise.

One really useful way to think of this recurrence relation is that represents an exercise in counting. We want to count all sequences that are of length N , and are accepted by the finite state machine of figure 16.1. This means they: (a) are N flips long; (b) end in HH ; and (c) contain no other subsequence HH . Now work backward along the FSM. The only way to arrive at the final state is to be in state 1, then see HH . So you can obtain an acceptable N element sequence by (a) prepending a T to an acceptable $N-1$ element sequence or (b) prepending TH (which takes you to 2, then back to 1) to an acceptable $N-2$ element sequence. This line of reasoning can be made much more elaborate. There are a few examples in the exercises.

16.1.2 Motivating Example: The Gambler's Ruin

Another useful example is known as the **gambler's ruin**. Assume you bet \$1 a tossed coin will come up heads. If you win, you get \$1 and your original stake back. If you lose, you lose your stake. But this coin has the property that $P(H) = p < 1/2$. We will study what happens when you bet repeatedly.

Assume you have \$ s when you start. You will keep betting until either (a) you have \$0 (you are ruined; you can't borrow money) or (b) the amount of money you have accumulated is \$ j , where $j > s$. The coin tosses are independent. We will compute $P(\text{ruined, starting with } s|p)$ the probability that you leave the table with nothing, when you start with \$ s . For brevity, we write $P(\text{ruined, starting with } s|p) = p_s$. You can represent the gambler's ruin problem with a finite state machine as well. Figure 16.2 shows a representation of the gambler's ruin problem.

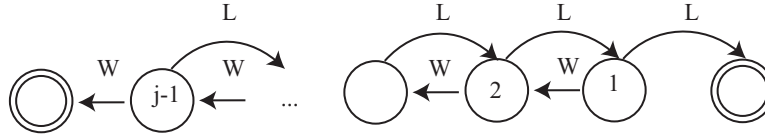


FIGURE 16.2: A finite state machine representing the gambler's ruin example. I have labelled each state with the amount of money the gambler has at that state. There are two end states, where the gambler has zero (is ruined), or has j and decides to leave the table. The problem we discuss is to compute the probability of being ruined, given the start state is s . This means that any state except the end states could be a start state. I have labelled the state transitions with "W" (for win) and "L" for lose, but have omitted the probabilities.

Worked example 16.4 *The gambler's ruin - 1*

Using the notation above, determine p_0 and p_j

Solution: We must have $p_0 = 1$, because if you have \$0, you leave the table. Similarly, if you have \$ j , you leave the table with \$ j , so you don't leave the table with nothing, so $p_j = 0$.

Worked example 16.5 *The gambler's ruin - 2*

Using the notation above, write a recurrence relation for p_s (the probability that you leave the table with nothing when you started with \$ s).

Solution: Assume that you win the first bet. Then you have \$ $s + 1$, so your probability of leaving the table with nothing now becomes p_{s+1} . If you lose the first bet, then you have \$ $s - 1$, so your probability of leaving the table with nothing now becomes p_{s-1} . The coin tosses are independent, so we can write

$$p_s = pp_{s+1} + (1-p)p_{s-1}.$$

Some fairly lively work with series, relegated to the end of the chapter as exercises, yields

$$p_s = \frac{\left(\frac{1-p}{p}\right)^j - \left(\frac{1-p}{p}\right)^s}{\left(\frac{1-p}{p}\right)^j - 1}.$$

This expression is quite informative. Notice that, if $p < 1/2$, then $(1-p)/p > 1$. This means that as $j \rightarrow \infty$, we have $p_s \rightarrow 1$. If you gamble repeatedly on an unfair

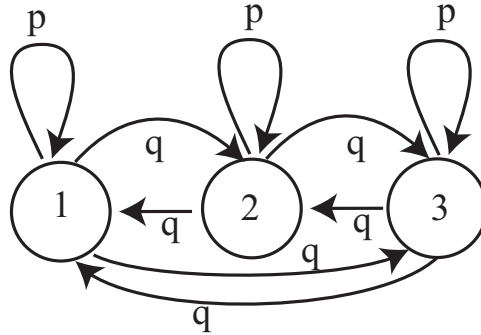


FIGURE 16.3: A virus can exist in one of 3 strains. At the end of each year, the virus mutates. With probability α , it chooses uniformly and at random from one of the 2 other strains, and turns into that; with probability $1 - \alpha$, it stays in the strain it is in. For this figure, we have transition probabilities $p = (1 - \alpha)$ and $q = (\alpha/2)$.

coin, the probability that you run out of money before you hit some threshold ($\$j$ in this case) tends to one.

16.1.3 Motivating Example: A Virus

Problems represented with a finite state machine don't need to have an end state. As one example, (which I got from ACC Coolen's lecture notes), we have a virus that can exist in one of k strains. At the end of each year, the virus mutates. With probability α , it chooses uniformly and at random from one of the $k - 1$ other strains, and turns into that; with probability $1 - \alpha$, it stays in the strain it is in (Figure 16.3 shows an example with three strains). For that figure, we have $p = (1 - \alpha)$ and $q = (\alpha/2)$. The virus just keeps on changing, and there is no end state. But there are a variety of very interesting questions we can try to answer. We would like to know, for example, the expected time to see a strain a second time, i.e. if the virus is in strain 1 at time 1, what is the expected time before it is in strain 1 again? If the virus has mutated many times, what is the probability that it is in each strain? and do these probabilities depend on the start strain?

16.1.4 Markov Chains

In Figure 16.1 and Figure 16.2, I showed the event that caused the state transitions. It is more usual to write a probability on the figure, as I did in Figure 16.3, because the probability of a state transition (rather than what caused it) is what really matters. The underlying object is now a weighted directed graph, because we have removed the events and replaced them with probabilities. These probabilities are known as **transition probabilities**; notice that the sum of transition probabilities over *outgoing* arrows must be 1.

We can now think of our process as a **biased random walk** on a weighted directed graph. A bug (or any other small object you prefer) sits on one of the

graph's nodes. At each time step, the bug chooses one of the outgoing edges at random. The probability of choosing an edge is given by the probabilities on the drawing of the graph (equivalently, the transition probabilities). The bug then follows that edge. The bug keeps doing this until it hits an end state.

This bug produces a sequence of random variables. If there are k states in the finite state machine, we can label each state with a number, $1 \dots k$. At the n 'th time step, the state of the process — the node that the bug is sitting on — is a random variable, which we write X_n . These random variables have an important property. The probability that X_n takes some particular value depends only on X_{n-1} , and not on any other previous state. If we know where the bug is at step $n - 1$ in our model, we know where it could go, and the probability of each transition. Where it was at previous times does not affect this, *as long as we know its state at step $n - 1$* .

A sequence of random variables X_n is a **Markov chain** if it has the property that, $P(X_n = j | \text{values of all previous states}) = P(X_n = j | X_{n-1})$, or, equivalently, only the last state matters in determining the probability of the current state. The probabilities $P(X_n = j | X_{n-1} = i)$ are the **transition probabilities**. Any model built by taking the transitions of a finite state machine and labelling them with probabilities must be a Markov chain. However, this is not the only way to build or represent a Markov chain.

One representation of a Markov chain uses a matrix of transition probabilities. We define the matrix \mathcal{P} with $p_{ij} = P(X_n = j | X_{n-1} = i)$. Notice that this matrix has the properties that $p_{ij} \geq 0$ and

$$\sum_j p_{ij} = 1$$

because at the end of each time step the model must be in some state. Equivalently, the sum of transition probabilities for outgoing arrows is one. Non-negative matrices with this property are **stochastic matrices**. By the way, you should look very carefully at the i 's and j 's here — Markov chains are usually written in terms of *row* vectors, and this choice makes sense in that context.

Worked example 16.6 Viruses

Write out the transition probability matrix for the virus of Figure 16.3, assuming that $\alpha = 0.2$.

Solution: We have $P(X_n = 1 | X_{n-1} = 1) = (1 - \alpha) = 0.8$, and $P(X_n = 2 | X_{n-1} = 1) = \alpha/2 = P(X_n = 3 | X_{n-1} = 1)$; so we get

$$\begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{pmatrix}$$

Now imagine we do not know the initial state of the chain, but instead have a probability distribution. This gives $P(X_0 = i)$ for each state i . It is usual to take these k probabilities and place them in a k -dimensional row vector, which is usually written π . For example, we might not know what the initial strain of the virus is, but just that each strain is equally likely. So for a 3-strain virus, we would have $\pi = [1/3, 1/3, 1/3]$. From this information, we can compute the probability distribution over the states at time 1 by

$$\begin{aligned} P(X_1 = j) &= \sum_i P(X_1 = j, X_0 = i) \\ &= \sum_i P(X_1 = j | X_0 = i) P(X_0 = i) \\ &= \sum_i p_{ij} \pi_i. \end{aligned}$$

If we write $\mathbf{p}^{(n)}$ for the row vector representing the probability distribution of the state at step n , we can write this expression as

$$\mathbf{p}^{(1)} = \pi \mathcal{P}.$$

Now notice that

$$\begin{aligned} P(X_2 = j) &= \sum_i P(X_2 = j, X_1 = i) \\ &= \sum_i P(X_2 = j | X_1 = i) P(X_1 = i) \\ &= \sum_i p_{ij} \left(\sum_k p_{ki} \pi_k \right). \end{aligned}$$

so that

$$\mathbf{p}^{(n)} = \pi \mathcal{P}^n.$$

This expression is useful for simulation, and also allows us to deduce a variety of interesting properties of Markov chains.

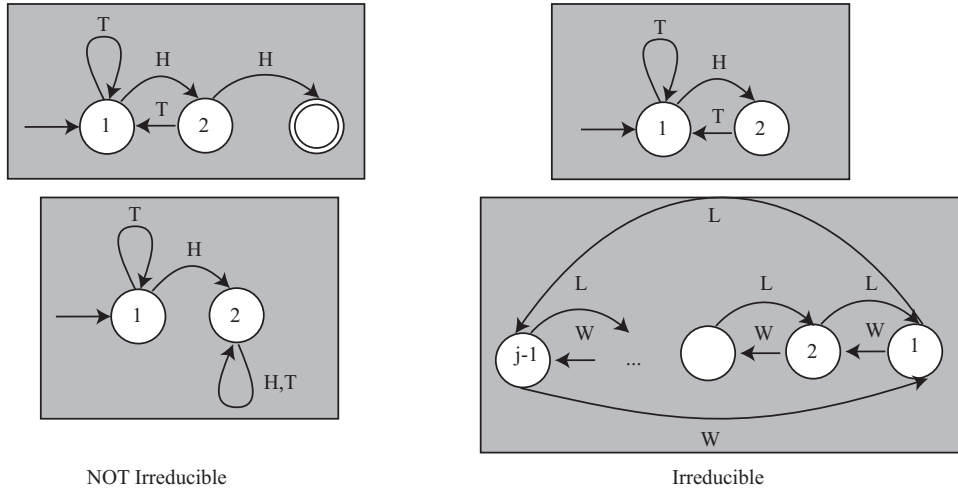


FIGURE 16.4: Examples of finite state machines that give rise to Markov chains that are NOT irreducible (left) and irreducible (right). To obtain Markov chains from these drawings, we would have to give the probabilities of the events that lead to state transitions. We'll assume that none of the probabilities are zero; after that, the values don't matter for irreducibility analysis. The **top left** FSM is not irreducible, because it has an end state; once the bug reaches this point, it can't go anywhere else. The **bottom left** FSM is not irreducible, because the bug can get stuck in state 2.

Worked example 16.7 Viruses

We know that the virus of Figure 16.3 started in strain 1. After two state transitions, what is the distribution of states when $\alpha = 0.2$? when $\alpha = 0.9$? What happens after 20 state transitions? If the virus starts in strain 2, what happens after 20 state transitions?

Solution: If the virus started in strain 1, then $\pi = [1, 0, 0]$. We must compute $\pi(\mathcal{P}(\alpha))^2$. This yields $[0.66, 0.17, 0.17]$ for the case $\alpha = 0.2$ and $[0.4150, 0.2925, 0.2925]$ for the case $\alpha = 0.9$. Notice that, because the virus with small α tends to stay in whatever state it is in, the distribution of states after two years is still quite peaked; when α is large, the distribution of states is quite uniform. After 20 transitions, we have $[0.3339, 0.3331, 0.3331]$ for the case $\alpha = 0.2$ and $[0.3333, 0.3333, 0.3333]$ for the case $\alpha = 0.9$; you will get similar numbers even if the virus starts in strain 2. After 20 transitions, the virus has largely “forgotten” what the initial state was.

In example 16.7, the distribution of virus strains after a long interval appeared not to depend much on the initial strain. This property is true of many Markov chains. Assume that any state can be reached from any other state, by some

sequence of transitions. Such chains are called **irreducible**; notice this means there is no end state, like the virus example. Irreducibility also means that the chain cannot get “stuck” in a state or a collection of states (Figure 16.4). Then there is a unique vector \mathbf{s} , usually referred to as the **stationary distribution**, such that for *any* initial state distribution π ,

$$\lim_{n \rightarrow \infty} \pi \mathcal{P}^{(n)} = \mathbf{s}.$$

Equivalently, if the chain has run through many steps, it no longer matters what the initial distribution is. You expect that the probability distribution over states is \mathbf{s} .

16.1.5 Example: Particle Motion as a Markov Chain

One can find Markov chains in quite unexpected places, often with useful consequences. In this example, I will obtain a Markov chain without reasoning about graphs or finite state machines. We will investigate a particle moving under gravity, in 3 dimensions. Write the position of the particle as $\mathbf{p}(t)$, its velocity as $\mathbf{v}(t)$, its acceleration as $\mathbf{a}(t)$, its mass as m , and the gravitational force as \mathbf{g} . Then we know that

$$\begin{aligned} \mathbf{v}(t) &= \frac{d\mathbf{p}}{dt} \\ \mathbf{a}(t) &= \frac{d\mathbf{v}}{dt} \\ &= \mathbf{g}. \end{aligned}$$

Now stack the position and the acceleration into a single vector $X(t) = (\mathbf{p}(t), \mathbf{v}(t))^T$. We could write these equations as

$$\frac{dX}{dt} = \mathcal{A}X + \mathbf{b}$$

where

$$\mathcal{A} = \begin{pmatrix} 0 & \mathcal{I} \\ 0 & 0 \end{pmatrix}$$

and

$$\mathbf{b} = \begin{pmatrix} 0 \\ \mathbf{g} \end{pmatrix}.$$

Now imagine that we look at the position, velocity and acceleration of the particle at fixed time instants, so we are really interested in $X_i = X(t_0 + i\Delta t)$. In this case, we can approximate $\frac{dX}{dt}$ by $(X_{i+1} - X_i)/\Delta t$. We then have

$$X_{i+1} = X_i + (\Delta t)(\mathcal{A}X_i + \mathbf{b}).$$

This is beginning to look like a Markov chain, because X_{i+1} depends only on X_i but not on any previous X . But there isn't any randomness here. We can fix that by assuming that the particle is moving in, say, a light turbulent wind. Then the acceleration at any time consists of (a) the acceleration of gravity and (b) a small,

random force produced by the wind. Write \mathbf{w}_i for this small, random force at time i , and $P(\mathbf{w}_i|i)$ for its probability distribution, which could reasonably depend on time. We can then rewrite our equation by writing

$$X_{i+1} = X_i + (\Delta t)(\mathcal{A}X_i + \mathbf{b}_r).$$

where

$$\mathbf{b}_r = \begin{pmatrix} 0 \\ \mathbf{g} + \mathbf{w}_i \end{pmatrix}.$$

Now the X_i are clearly random variables. We could get $P(X_{i+1}|X_i)$ by rearranging terms. If I know X_{i+1} and X_i , I know the value of \mathbf{w}_i ; I could plug this into $P(\mathbf{w}_i|i)$ to yield $P(X_{i+1}|X_i)$.

Using a finite state machine in this example would be a bit unnatural. This example turns out to be surprisingly useful in applications, because (with other algorithmic machinery we can't go into here) it offers the basis for algorithms that can track moving objects by predicting where they will go next. Applications are widespread. Military applications include tracking aircraft, UFO's, missiles, etc. Civilian applications include surveillance in public places; games console interfaces that track moving people; and methods that can follow experimental mice as they move around their cages (useful for testing medicines).

16.2 SIMULATION

Many problems in probability can be worked out in closed form if one knows enough combinatorial mathematics, or can come up with the right trick. Textbooks are full of these, and we've seen some. Explicit formulas for probabilities are often extremely useful. But it isn't always easy or possible to find a formula for the probability of an event in a model. An alternative strategy is to build a simulation, run it many times, and count the fraction of outcomes where the event occurs. This is a simulation experiment.

16.2.1 Obtaining Random Numbers

Simulation is at its most useful when we try to estimate probabilities that are hard to calculate. Usually we have processes with some random component, and we want to estimate the probability of a particular outcome. To do so, we will need a supply of random numbers to simulate the random component.

Uniform Random Numbers

I will describe features in Matlab, because I'm used to Matlab. It's a good programming environment for numerical work, particularly for working with matrices and vectors. You can expect pretty much any programming environment to provide a random number generator that returns a number uniformly distributed in the range $[0 - 1]$. If you know anything about representing numbers, you'll already have spotted something funny. The number that comes out of the random number generator must be represented by a small set of bits, but almost all numbers in the interval $[0 - 1]$ require an infinite number of bits to represent. We'll sweep this issue under the general carpet of floating point representations; it doesn't matter

for anything we need to do. The Matlab function to do this is called `rand`. It can also return matrices; for example, `rand(10, 20)` will get you a 10×20 table of independent uniformly distributed random numbers in the range $[0 - 1]$.

It is useful to get a uniformly distributed random integer in a particular range (for example, you might want to choose a random element in an array). You can do so with a random number generator and a function (like Matlab's `floor`) that returns the largest integer smaller than its argument. If I want a discrete random variable with uniform distribution, maximum value 100 and minimum value 7, I habitually choose a very tiny number (for this example, say $1e - 7$) and do `floor((100-7-1e-7)*rand()+7)`. I use the $1e - 7$ because I can never remember whether `rand` produces a number no larger than one, or one that is guaranteed to be smaller than one, and I never need to care about the very tiny differences in probability caused by the $1e-7$. You might be able to do something cleaner if you bothered checking this point.

Normal Random Numbers

In pretty much any programming environment, you would also expect to find a random number generator that returns a normal random variable, with mean zero and standard deviation one. In Matlab, this function is called `randn`. Conveniently, `randn(3, 4)` will give you a 3×4 table of such numbers, which are independent of each other. As you would expect from section 16.20, to change the mean of this random number, you add a constant; to change the variance, you multiply by a constant. So in Matlab, if you want a normal random variable with mean 3 and standard deviation 4, you use `4*randn()+3`.

16.2.2 Computing Expectations with Simulations

Simulation is also a very good way to estimate expectations. Imagine we have a random variable X with probability distribution $P(X)$ that takes values in some domain D . Assume that we can easily produce independent simulations, and that we wish to know $\mathbb{E}[f]$, the expected value of the function f under the distribution $P(X)$.

The weak law of large numbers tells us how to proceed. Define a new random variable $F = f(X)$. This has a probability distribution $P(F)$, which might be difficult to know. We want to estimate $\mathbb{E}[f]$, the expected value of the function f under the distribution $P(X)$. This is the same as $\mathbb{E}[F]$. Now if we have a set of IID samples of X , which we write x_i , then we can form a set of IID samples of F by forming $f(x_i) = f_i$. Write

$$F_N = \frac{\sum_{i=1}^N f_i}{N}.$$

This is a random variable, and the weak law of large numbers gives that, for any positive number ϵ

$$\lim_{N \rightarrow \infty} P(\{\|F_N - \mathbb{E}[F]\| > \epsilon\}) = 0.$$

You can interpret this as saying that, that for a set of IID random samples x_i , the

probability that

$$\frac{\sum_{i=1}^N f(x_i)}{N}$$

is very close to $\mathbb{E}[f]$ is high for large N

Worked example 16.8 *Computing an Expectation*

Assume the random variable X is uniformly distributed in the range $[0 - 1]$, and the random variable Y is uniformly distributed in the range $[0 - 10]$. X and Z are independent. Write $Z = (Y - 5X)^3 - X^2$. What is $\text{var}(\{Z\})$?

Solution: With enough work, one could probably work this out in closed form. An easy program will get a good estimate. We have that $\text{var}(\{Z\}) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$. My program computed 1000 values of Z (by drawing X and Y from the appropriate random number generator, then evaluating the function). I then computed $\mathbb{E}[Z]$ by averaging those values, and $\mathbb{E}[Z]^2$ by averaging their squares. For a run of my program, I got $\text{var}(\{Z\}) = 2.76 \times 10^4$.

16.2.3 Computing Probabilities with Simulations

You can compute a probability using a simulation, too, because a probability can be computed by taking an expectation. Recall the property of indicator functions that

$$\mathbb{E}[\mathbb{I}_{\{\mathcal{E}\}}] = P(\mathcal{E})$$

(Section ??). This means that computing the probability of an event \mathcal{E} involves writing a function that is 1 when the event occurs, and 0 otherwise; we then estimate the expected value of that function.

The weak law of large numbers justifies this procedure. An experiment involves drawing a sample from the relevant probability distribution $P(X)$, then determining whether event \mathcal{E} occurred or not. We define a new random variable E , which takes the value 1 when the event \mathcal{E} occurs during our experiment (i.e. with probability $P(\mathcal{E})$) and 0 when it doesn't (i.e. with probability $P(\mathcal{E}^c)$). Our experiment yields a sample of this random variable. We perform N experiments yielding a set of IID samples of this distribution e_i and compute $E_N = \frac{\sum_{i=1}^N e_i}{N}$. From the weak law of large numbers, we have that for any $\epsilon > 0$,

$$\lim_{N \rightarrow \infty} P(\{\|E_N - \mathbb{E}[E]\| \geq \epsilon\}) = 0$$

meaning that for large enough N , E_N will be very close to $\mathbb{E}[E] = P(\mathcal{E})$.

Worked example 16.9 *Computing a Probability for Multiple Coin Flips*

You flip a fair coin three times. Use a simulation to estimate the probability that you see three H 's.

Solution: You really should be able to work this out in closed form. But it's amusing to check with a simulation. I wrote a simple program that obtained a 1000x3 table of uniformly distributed random numbers in the range $[0 - 1]$. For each number, if it was greater than 0.5 I recorded an H and if it was smaller, I recorded a T . Then I counted the number of rows that had 3 H 's (i.e. the expected value of the relevant indicator function). This yielded the estimate 0.127, which compares well to the right answer.

Worked example 16.10 *Computing a Probability*

Assume the random variable X is uniformly distributed in the range $[0 - 1]$, and the random variable Y is uniformly distributed in the range $[0 - 10]$. Write $Z = (Y - 5X)^3 - X^2$. What is $P(\{Z > 3\})$?

Solution: With enough work, one could probably work this out in closed form. An easy program will get a good estimate. My program computed 1000 values of Z (by drawing X and Y from the appropriate random number generator, then evaluating the function) and counted the fraction of Z values that was greater than 3 (which is the relevant indicator function). For a run of my program, I got $P(\{Z > 3\}) \approx 0.619$

16.2.4 Simulation Results as Random Variables

The estimate of a probability or of an expectation that comes out of a simulation experiment is a random variable, because it is a function of random numbers. If you run the simulation again, you'll get a different value (unless you did something silly with the random number generator). Generally, you should expect this random variable to behave like a normal random variable. You can check this by constructing a histogram over a large number of runs. The mean of this random variable is the parameter you are trying to estimate. It is useful to know that this random variable tends to be normal, because it means the standard deviation of the random variable tells you a lot about the likely values you will observe.

Another helpful rule of thumb, which is almost always right, is that the standard deviation of this random variable behaves like

$$\frac{C}{\sqrt{N}}$$

where C is a constant that depends on the problem and can be very hard to evaluate, and N is the number of runs of the simulation. What this means is that if you want

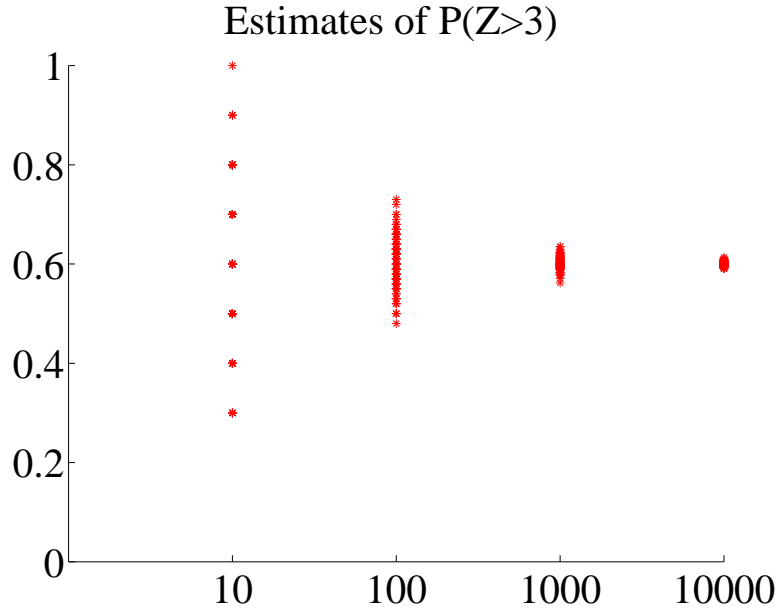


FIGURE 16.5: Estimates of the probability from example 16.10, obtained from different runs of my simulator using different numbers of samples. In each case, I used 100 runs; the number of samples is shown on the horizontal axis. You should notice that the estimate varies pretty widely when there are only 10 samples, but the variance (equivalently, the size of the spread) goes down sharply as the number of samples increases to 1000. Because we expect these estimates to be roughly normally distributed, the variance gives a good idea of how accurate the original probability estimate is.

to (say) double the accuracy of your estimate of the probability or the expectation, you have to run four times as many simulations. Very accurate estimates are tough to get, because they require immense numbers of simulation runs.

Figure 16.5 shows how the result of a simulation behaves when the number of runs changes. I used the simulation of example 16.10, and ran multiple experiments for each of a number of different samples (i.e. 100 experiments using 10 samples; 100 using 100 samples; and so on).

Small probabilities can be rather hard to estimate, as we would expect. In the case of example 16.10, let us estimate $P(\{Z > 950\})$. A few moments with a computer will show that this probability is of the order of $1e-3$ to $1e-4$. I obtained a million different simulated values of Z from my program, and saw 310 where $Z > 950$. This means that to know this probability to, say, three digits of numerical accuracy might involve a daunting number of samples. Notice that this does not contradict the rule of thumb that the standard deviation of the random variable defined by a simulation estimate behaves like $\frac{C}{\sqrt{N}}$; it's just that in this case, C is very large indeed.

16.3 SIMULATION EXAMPLES

Computing probabilities and expectations from simulations should be so natural to a computer science student that it can be hard to see the magic. You estimate the probability of an event \mathcal{E} by writing a program that runs N independent simulations of an experiment, counts how many times the event occurs (which we write $\#(\mathcal{E})$), and reports

$$\frac{\#(\mathcal{E})}{N}$$

as an estimate of $P(\mathcal{E})$. This estimate will not be exact, and may be different for different runs of the program. It's often a good, simple estimate.

Choosing N depends on a lot of considerations. Generally, a larger N means a more accurate answer, and also a slower program. If N is too small, you may find that you report 1 or 0 for the probability (think of what would happen if you measured $P(H)$ for a coin with one flip). One strategy is to run several simulations, report the mean, and use the standard deviation as some guide to the accuracy.

16.3.1 Simulating Experiments

Worked example 16.11 *Getting 14's with 20-sided dice*

You throw 3 fair 20-sided dice. Estimate the probability that the sum of the faces is 14 using a simulation. Use $N = [1e1, 1e2, 1e3, 1e4, 1e5, 1e6]$. Which estimate is likely to be more accurate, and why?

Solution: You need a fairly fast computer, or this will take a long time. I ran ten versions of each experiment for $N = [1e1, 1e2, 1e3, 1e4, 1e5, 1e6]$, yielding ten probability estimates for each N . These were different for each version of the experiment, because the simulations are random. I got means of $[0, 0.0030, 0.0096, 0.0100, 0.0096, 0.0098]$, and standard deviations of $[0.00670, 0.00330, 0.00090, 0.00020, 0.0001]$. This suggests the true value is around 0.0098, and the estimate from $N = 1e6$ is best. The reason that the estimate with $N = 1e1$ is 0 is that the probability is very small, so you don't usually observe this case at all in only ten trials.

Worked example 16.12 *Comparing simulation with computation*

You throw 3 fair six-sided dice. You wish to know the probability the sum is 3. Compare the true value of this probability with estimates from six runs of a simulation using $N = 10000$. What conclusions do you draw?

Solution: I ran six simulations with $N = 10000$, and got [0.0038, 0.0038, 0.0053, 0.0041, 0.0056, 0.0049]. The mean is 0.00458, and the standard deviation is 0.0007, which suggests the estimate isn't that great, but the right answer should be in the range [0.00388, 0.00528] with high probability. The true value is $1/216 \approx 0.00463$. The estimate is tolerable, but not super accurate.

16.3.2 Simulating Markov Chains

Worked example 16.13 *Coin Flips with End Conditions*

I flip a coin repeatedly until I encounter a sequence HTHT, at which point I stop. What is the probability that I flip the coin nine times?

Solution: You might well be able to construct a closed form solution to this if you follow the details of example 16.20 and do quite a lot of extra work. A simulation is really straightforward to write; notice you can save time by not continuing to simulate coin flips once you've flipped past nine times. I got 0.0411 as the mean probability over 10 runs of a simulation of 1000 experiments each, with a standard deviation of 0.0056.

Worked example 16.14 *A Queue*

A bus is supposed to arrive at a bus stop every hour for 10 hours each day. The number of people who arrive to queue at the bus stop each hour has a Poisson distribution, with intensity 4. If the bus stops, everyone gets on the bus and the number of people in the queue becomes zero. However, with probability 0.1 the bus driver decides not to stop, in which case people decide to wait. If the queue is ever longer than 15, the waiting passengers will riot (and then immediately get dragged off by the police, so the queue length goes down to zero). What is the expected time between riots?

Solution: I'm not sure whether one could come up with a closed form solution to this problem. A simulation is completely straightforward to write. I get a mean time of 441 hours between riots, with a standard deviation of 391. It's interesting to play around with the parameters of this problem; a less conscientious bus driver, or a higher intensity arrival distribution, lead to much more regular riots.

Worked example 16.15 *Inventory*

A store needs to control its stock of an item. It can order stocks on Friday evenings, which will be delivered on Monday mornings. The store is old-fashioned, and open only on weekdays. On each weekday, a random number of customers comes in to buy the item. This number has a Poisson distribution, with intensity 4. If the item is present, the customer buys it, and the store makes \$100; otherwise, the customer leaves. Each evening at closing, the store loses \$10 for each unsold item on its shelves. The store's supplier insists that it order a fixed number k of items (i.e. the store must order k items each week). The store opens on a Monday with 20 items on the shelf. What k should the store use to maximise profits?

Solution: I'm not sure whether one could come up with a closed form solution to this problem, either. A simulation is completely straightforward to write. To choose k , you run the simulation with different k values to see what happens. I computed accumulated profits over 100 weeks for different k values, then ran the simulation five times to see which k was predicted. Results were 21, 19, 23, 20, 21. I'd choose 21 based on this information.

For example 16.15, you should plot accumulated profits. If k is small, the store doesn't lose money by storing items, but it doesn't sell as much stuff as it could; if k is large, then it can fill any order but it loses money by having stock on the shelves. A little thought will convince you that k should be near 20, because that is the expected number of customers each week, so $k = 20$ means the store can expect to sell all its new stock. It may not be exactly 20, because it must depend

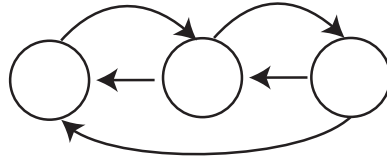


FIGURE 16.6: A very small fraction of the web, drawn to suggest a finite state machine; each state represents a page, and each directed edge represents an outgoing link. A random web surfer could either (a) follow an outgoing link, chosen at random or (b) type in the URL of any page, chosen at random. Such a surfer would see lots of pages that have many incoming links from pages that have lots of incoming links, and so on. Pages like this are likely important, so that finding important pages is analogous to simulating a random web surfer.

a little on the balance between the profit in selling an item and the cost of storing it. For example, if the cost of storing items is very small compared to the profit, an very large k might be a good choice. If the cost of storage is sufficiently high, it might be better to never have anything on the shelves; this point explains the absence of small stores selling PC's.

16.3.3 Example: Ranking the Web by Simulating a Markov Chain

Perhaps the most valuable technical question of the last thirty years has been: Which web pages are interesting? Some idea of the importance of this question is that it was only really asked about 20 years ago, and at least one gigantic technology company has been spawned by a partial answer. This answer, due to Larry Page and Sergey Brin, and widely known as PageRank, starts with a Markov chain.

Figure 16.6 shows a picture of (a very small fraction of) the world wide web. I have drawn the web using a visual metaphor that should strongly suggest a finite state machine, and so a Markov chain. Each page is a state. Directed edges from page to page represent links. I count only the first link from a page to another page. Some pages are linked, others are not. I want to know how important each page is.

One way to think about importance is to think about what a random web surfer would do. The surfer can either (a) choose one of the outgoing links on a page at random, and follow it or (b) type in the URL of a new page, and go to that instead. As Figure 16.6 suggests, it is useful to think of this as a random walk on a finite state machine. We expect that this random surfer should see a lot of pages that have lots of incoming links from other pages that have lots of incoming links that (and so on). These pages are important, because lots of pages have linked to them.

For the moment, ignore the surfer's option to type in a URL. Write $r(i)$ for the importance of the i 'th page. We model importance as leaking from page to page across outgoing links (the same way the surfer jumps). Page i receives importance down each incoming link. The amount of importance is proportional to the amount of importance at the other end of the link, and inversely proportional to the number of links leaving that page. So a page with only one outgoing link transfers all its

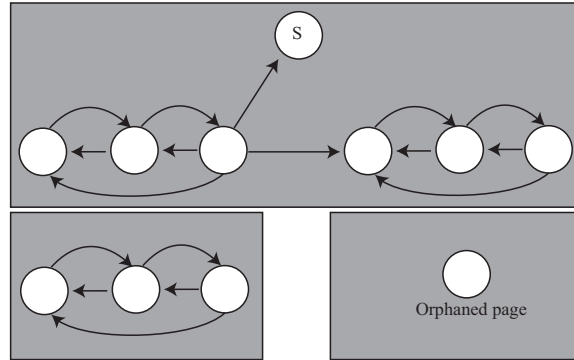


FIGURE 16.7: The web isn't really like Figure 16.6. It's more like this figure (only bigger). In this figure, we see sections of the web that can't be reached by following links from other sections (the gray boxes); orphaned pages; and pages where a random walker would get stuck (S). Any algorithm for assigning importance must be able to deal with these effects.

importance down that link; and the way for a page to receive a lot of importance is for it to have a lot of important pages link to it alone. We write

$$r(j) = \sum_{i \rightarrow j} \frac{r(i)}{|i|}$$

where $|i|$ means the total number of links pointing *out* of page i . We can stack the $r(j)$ values into a *row* vector \mathbf{r} , and construct a matrix \mathcal{P} , where

$$p_{ij} = \begin{cases} \frac{1}{|i|} & \text{if } i \text{ points to } j \\ 0 & \text{otherwise} \end{cases}$$

With this notation, the importance vector has the property

$$\mathbf{r} = \mathbf{r}\mathcal{P}$$

and should look a bit like the stationary distribution of a random walk to you, except that \mathcal{P} isn't stochastic — there may be some rows where the row sum of \mathcal{P} is zero, because there are *no* outgoing links from that page. We can fix this easily by replacing each row that sums to zero with $(1/n)\mathbf{1}$, where n is the total number of pages. Call the resulting matrix \mathcal{G} (it's quite often called the **raw Google matrix**). Notice that doing this doesn't really change anything significant; pages with no outgoing links leak a tiny fraction of their importance to every other page.

Figure 16.6 isn't a particularly good model of the web (and not because it's tiny, either). Figure 16.7 shows some of the problems that occur. There are pages with no outgoing links (which we've dealt with), pages with no incoming links, and even pages with no links at all. Worse, a random walk can get trapped (in one of the gray boxes). One way to fix all this would be to construct a process that wandered around the web inserting links that clean up the structure of the

graph. This strategy is completely infeasible, because the real web is much too big. Allowing the surfer to randomly enter a URL sorts out all of these problems, because it inserts an edge of small weight from every node to every other node. Now the random walk cannot get trapped.

There are a variety of possible choices for the weight of these inserted edges. The original choice was to make each inserted edge have the same weight. Write $\mathbf{1}$ for the n dimensional column vector containing a 1 in each component, and let $0 < \alpha < 1$. We can write the matrix of transition probabilities as

$$\mathcal{G}(\alpha) = \alpha \frac{(\mathbf{1}\mathbf{1}^T)}{n} + (1 - \alpha)\mathcal{G}$$

where \mathcal{G} is the original Google matrix. An alternative choice is to choose a weight for each web page, using anything from advertising revenues to page visit statistics to thaumaturgy (Google keeps quiet about the details). Write this weight vector \mathbf{v} , and require that $\mathbf{1}^T \mathbf{v} = 1$ (i.e. the coefficients sum to one). Then we could have

$$\mathcal{G}(\alpha, \mathbf{v}) = \alpha \frac{(\mathbf{1}\mathbf{v}^T)}{n} + (1 - \alpha)\mathcal{G}.$$

Now the importance vector \mathbf{r} is the (unique, though I won't prove this) *row* vector \mathbf{r} such that

$$\mathbf{r} = \mathbf{r}\mathcal{G}(\alpha, \mathbf{v}).$$

How do we compute this vector? One natural algorithm is to start with some initial estimate, and propagate it. We write $\mathbf{r}^{(k)}$ for the estimated importance after the k 'th step. We define updates by

$$(\mathbf{r}^{(k)}) = (\mathbf{r}^{(k-1)})\mathcal{G}(\alpha, \mathbf{v}).$$

We can't compute this directly, either, because $\mathcal{G}(\alpha, \mathbf{v})$ is unreasonably big so (a) we can't form or store $\mathcal{G}(\alpha, \mathbf{v})$ and (b) we can't multiply by it either. But we could estimate \mathbf{r} with a random walk, because \mathbf{r} is the stationary distribution of a Markov chain. If we simulate this walk for many steps, the probability that the simulation is in state j should be $r(j)$, at least approximately.

This simulation is easy to build. Imagine our random walking bug sits on a web page. At each time step, it transitions to a new page by either (a) picking from all existing pages at random, using \mathbf{v} as a probability distribution on the pages (which it does with probability α); or (b) chooses one of the outgoing links uniformly and at random, and follows it (which it does with probability $1 - \alpha$). The stationary distribution of this random walk is \mathbf{r} . Another fact that I shall not prove is that, when α is sufficiently large, this random walk very quickly "forgets" its initial distribution. As a result, you can estimate the importance of web pages by starting this random walk in a random location; letting it run for a bit; then stopping it, and collecting the page you stopped on. The pages you see like this are independent, identically distributed samples from \mathbf{r} ; so the ones you see more often are more important, and the ones you see less often are less important.

16.3.4 Example: Simulating a Complicated Game

I will build several examples around a highly simplified version of a real card game. This game is Magic: The Gathering, and is protected by a variety of trademarks, etc. My version — MTGDFAF — isn't very interesting as a game, but is simple enough to be studied, and interesting enough it casts some light on the real game. The game is played with decks of 60 cards. There are two types of card: Lands, and Spells. Lands can be placed on the play table and stay there permanently; Spells are played and then disappear. A Land on the table can be “tapped” or “untapped”. Players take turns. Each player draws a hand of seven cards from a shuffled deck. In each turn, a player first untaps any Lands on the table, then draws a card, then plays a land onto the table (if the player has one in hand to play), then finally can play one or more spells. Each spell has a fixed cost (of $1, \dots, 10$), and this cost is played by “tapping” a land (which is not untapped until the start of the next turn). This means that the player can cast only cheap spells in the early turns of the game, and expensive spells in the later turns.

Worked example 16.16 *MTGDFAF — The number of lands*

Assume a deck of 60 cards has 24 Lands. It is properly shuffled, and you draw seven cards. You could draw $0, \dots, 7$ Lands. Estimate the probability for each, using a simulation. Furthermore, estimate the error in your estimates.

Solution: The matlab function `randperm` produces a random permutation of given length. This means you can use it to simulate a shuffle of a deck, as in listing 16.1. I then drew 10, 000 random hands of seven cards, and counted how many times I got each number. Finally, to get an estimate of the error, I repeated this experiment 10 times and computed the standard deviation of each estimate of probability. This produced

0.0218 0.1215 0.2706 0.3082 0.1956 0.0686 0.0125 0.0012

for the probabilities (for 0 to 7, increasing number of lands to the right) and

0.0015 0.0037 0.0039 0.0058 0.0027 0.0032 0.0005 0.0004

for the standard deviations of these estimates.

Listing 16.1: Matlab code used to simulate the number of lands

```

simcards=[ones(24, 1); zeros(36, 1)]
% 1 if land, 0 otherwise
ninsim=10000;
nsims=10;
counts=zeros(nsims, 8);
for i=1:10
    for j=1:10000
        shuffle=randperm(60);
        hand=simcards(shuffle(1:7));
        %useful matlab trick here
        nlands=sum(hand);
        %ie number of lands
        counts(i, 1+nlands)=...
            counts(i, 1+nlands)+1;
        % number of lands could be zero
    end
end
probs=counts/ninsim;
mean(probs)
std(probs)
%%

```

Worked example 16.17 *MTGDAF — The number of lands*

What happens to the probability of getting different numbers of lands if you put only 15 Lands in a deck of 60? It is properly shuffled, and you draw seven cards. You could draw 0, . . . , 7 Lands. Estimate the probability for each, using a simulation. Furthermore, estimate the error in your estimates.

Solution: You can change one line in the listing to get

```
0.1159 0.3215 0.3308 0.1749 0.0489 0.0075 0.0006 0.0000
```

for the probabilities (for 0 to 7, increasing number of lands to the right) and

```
0.0034 0.0050 0.0054 0.0047 0.0019 0.0006 0.0003 0.0000
```

for the standard deviations of these estimates.

Worked example 16.18 *MTGDFA — Playing spells*

Assume you have a deck of 24 Lands, 10 Spells of cost 1, 10 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 2 Spells of cost 6. Assume you always only play the cheapest spell in your hand (i.e. you never play two spells). What is the probability you will be able to play at least one spell on each of the first four turns?

Solution: This simulation requires just a little more care. You draw the hand, then simulate the first four turns. In each turn, you can only play a spell whose cost you can pay, and only if you have it. I used the matlab of listing 16.2 and listing 16.3; I found the probability to be 0.64 with standard deviation 0.01. Of course, my code might be wrong....

Worked example 16.19 *MTGDFA — Playing spells*

Now we use a different distribution of cards. Assume you have a deck of 20 Lands, 9 Spells of cost 1, 5 Spells of cost 2, 5 Spells of cost 3, 5 Spells of cost 4, 5 Spells of cost 5, and 11 Spells of cost 6. Assume you always only play the cheapest spell in your hand (i.e. you never play two spells). What is the probability you will be able to play at least one spell on each of the first four turns?

Solution: This simulation requires just a little more care. You draw the hand, then simulate the first four turns. In each turn, you can only play a spell whose cost you can pay, and only if you have it. I found the probability to be 0.33 with standard deviation 0.05. Of course, my code might be wrong....

One engaging feature of the real game that is revealed by these very simple simulations is the tension between a players goals. The player would like to have few lands — so as to have lots of spells — but doing so means that there's a bigger chance of not being able to play a spell. Similarly, a player would like to have lots of powerful (=expensive) spells, but doing so means there's a bigger chance of not being able to play a spell. Players of the real game spend baffling amounts of time arguing with one another about the appropriate choices for a good set of cards.

Listing 16.2: Matlab code used to simulate the four turns

```

simcards=[zeros(24, 1); ones(10, 1);...
          2*ones(10, 1);3*ones(10, 1); ...
          4*ones(2, 1); 5*ones(2, 1); 6*ones(2, 1)];
nsims=10;
ninsim=1000;
counts=zeros(nsims, 1);
for i=1:nsims
    for j=1:ninsim
        % draw a hand
        shuffle=randperm(60);
        hand=simcards(shuffle(1:7));
        %reorganize the hand
        cleanhand=zeros(7, 1);
        for k=1:7
            cleanhand(hand(k)+1)=cleanhand(hand(k)+1)+1;
            % ie count of lands, spells, by cost
        end
        landsontable=0;
        [playedspell1, landsontable, cleanhand]=...
            playround(landsontable, cleanhand, shuffle, ...
                simcards, 1);
        [playedspell2, landsontable, cleanhand]=...
            playround(landsontable, cleanhand, shuffle, ...
                simcards, 2);
        [playedspell3, landsontable, cleanhand]=...
            playround(landsontable, cleanhand, shuffle, ...
                simcards, 3);
        [playedspell4, landsontable, cleanhand]=...
            playround(landsontable, cleanhand, shuffle, ...
                simcards, 4);
        counts(i)=counts(i)+...
            playedspell1*playedspell2*...
            playedspell3*playedspell4;
    end
end
end

```

Worked example 16.20 *MTGDAF* — How long before you can play a spell of cost 3?

Assume you have a deck of 15 Lands, 15 Spells of cost 1, 14 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 2 Spells of cost 6. What is the expected number of turns before you can play a spell of cost 3? Assume you always play a land if you can.

Solution: I get 6.3, with a standard deviation of 0.1. The problem is it can take quite a large number of turns to get three lands out. I used the code of listings 16.4 and 16.5

Recall that you can reasonably expect that the probability you compute from a simulation behaves like normal data. Run a simulation experiment a large num-

Listing 16.3: Matlab code used to simulate playing a turn

```

function [playedspell, landsontable, cleanhand]=...
    playground(landsontable, cleanhand, shuffle, simcards, ...
        turn)
    % draw
    ncard=simcards(shuffle(7+turn));
    cleanhand(ncard+1)=cleanhand(ncard+1)+1;
    % play land
    if cleanhand(1)>0
        landsontable=landsontable+1;
        cleanhand(1)=cleanhand(1)-1;
    end
    playedspell=0;
    if landsontable>0
        i=1; done=0;
        while done==0
            if cleanhand(i)>0
                cleanhand(i)=cleanhand(i)-1;
                playedspell=1;
                done=1;
            else
                i=i+1;
                if i>landsontable
                    done=1;
                end
            end
        end
    end
end

```

ber of times and construct a data set whose entries are the probability from each experiment. This data should be normal. This means that, if you subtract the mean and divide by the standard deviation, you should get a histogram that looks like the standard normal curve. Figure 16.8 shows some examples of this effect. This effect is important, because it means that the right answer should be very few standard deviations away from the mean you compute — the standard deviation gives you quite a good idea of the accuracy of your estimate.

16.4 YOU SHOULD

16.4.1 remember

- A Markov chain is a sequence of random variables, with particular independence properties. It is valuable to think of a Markov chain as being produced by a biased random walk on a directed graph. In this view, the values of the random variables are given by the labels at each node of the graph.
- In some Markov chains, there is a special state that the chain cannot leave (an absorbing state). Not all chains have such states.
- Markov chains that satisfy reasonable structural assumptions have a stationary distribution. The random variables emitted by the chain will be samples

Listing 16.4: Matlab code used to estimate number of turns before you can play a spell of cost 3

```

simcards=[zeros(15, 1); ones(15, 1);...
          2*ones(14, 1);3*ones(10, 1); ...
          4*ones(2, 1); 5*ones(2, 1); 6*ones(2, 1)];
nsims=10;
ninsim=1000;
counts=zeros(nsims, 1);
for i=1:nsims
    for j=1:ninsim
        % draw a hand
        shuffle=randperm(60);
        hand=simcards(shuffle(1:7));
        %reorganize the hand
        cleanhand=zeros(7, 1);
        for k=1:7
            cleanhand(hand(k)+1)=cleanhand(hand(k)+1)+1;
            % ie count of lands, spells, by cost
        end
        landsontable=0;
        k=0; played3spell=0;
        while played3spell==0;
            [played3spell, landsontable, cleanhand]=...
                play3round(landsontable, cleanhand, shuffle, ...
                    simcards, k+1);
            k=k+1;
        end
        counts(i)=counts(i)+k;
    end
end
counts(i)=counts(i)/ninsim;
end

```

(but not necessarily independent samples) of that stationary distribution once the chain has changed state often enough.

- One can often write expressions for various probabilities associated with chains; these are usually recurrence relations.
- Simulating Markov chains is usually straightforward, and often extremely useful.
- Simulation is a powerful method to estimate the probabilities of particular events. The method can be used to estimate expectations as well.
- In general, one simulates a random process with a program, then counts the different outputs encountered in a large number of runs of the simulation. In most cases, for a sufficiently large number of runs, the number produced by a simulation is the value of a normal random variable whose expected value is the number sought and whose variance is k/N where k is an unknown constant and N is the number of simulation runs. This means that the accuracy of simulation based estimates improves increasingly slowly as N increases.

Listing 16.5: Matlab code used to simulate a turn to estimate the number of turns before you can play a spell of cost 3

```

function [played3spell, landsontable, cleanhand]=...
    play3round(landsontable, cleanhand, shuffle, simcards, ...
        turn)
    % draw
ncard=simcards(shuffle(7+turn));
cleanhand(ncard+1)=cleanhand(ncard+1)+1;
% play land
if cleanhand(1)>0
    landsontable=landsontable+1;
    cleanhand(1)=cleanhand(1)-1;
end
played3spell=0;
if (landsontable>=3)&&(cleanhand(4)>0)
    played3spell=1;
end

```

- Simulation methods can be used to estimate a wide range of useful properties, including the importance of web documents and the value of various strategies in games.

PROBLEMS

16.1. Multiple coin flips:

- For example ??, show that $P(5) = (3/32)$ by directly writing out the sequences of flips, and summing their probabilities.
- Now check the recurrence relation $P(N) = (1/2)P(N-1) + (1/4)P(N-2)$ for the case $N = 5$.
- What is $P(7)$?

16.2. Multiple die rolls:

You roll a fair die until you see a 5, then a 6; after that, you stop. Write $P(N)$ for the probability that you roll the die N times.

- What is $P(1)$?
- Show that $P(2) = (1/36)$.
- Draw a finite state machine encoding all the sequences of die rolls that you could encounter. Don't write the events on the edges; instead, write their probabilities. There are 5 ways not to get a 5, but only one probability, so this simplifies the drawing.
- Show that $P(3) = (1/36)$.
- Now use your finite state machine to argue that $P(N) = (5/6)P(N-1) + (25/36)P(N-2)$.

16.3. More complicated multiple coin flips:

You flip a fair coin until you see either HTH or THT , and then you stop. We will compute a recurrence relation for $P(N)$.

- Figure ?? shows a finite state machine. Check that this finite state machine represents all coin sequences that you will encounter.
- Write Σ_N for some string of length N accepted by this finite state machine. Use this finite state machine to argue that Sigma_N has one of four forms:

- $TT\Sigma_{N-2}$

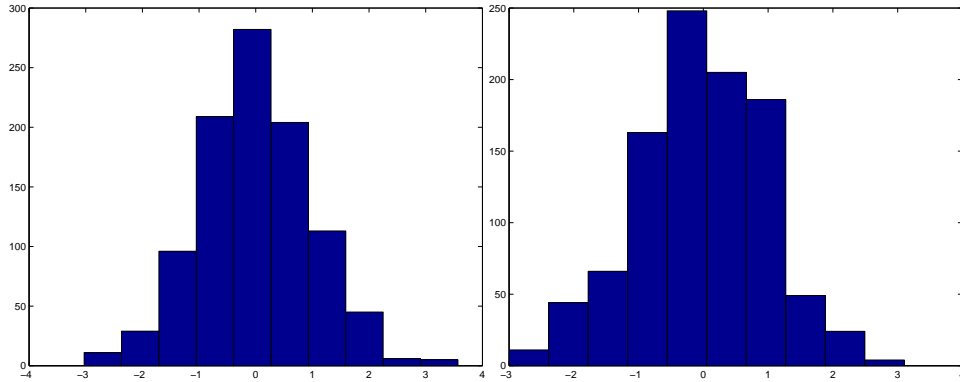


FIGURE 16.8: Estimates of probabilities produced by simulation typically behave like normal data. On the left, I show a histogram of probabilities of having a hand of 3 Lands in the simulation of example 16.16; these are plotted in standard coordinates. On the right, I show a histogram of probability of playing a spell in each of the first four turns (example 16.19), from 1000 simulation experiments; again, these are plotted in standard coordinates. Compare these to the standard normal histograms of the previous chapter.

- 2. $HH\Sigma_{N-3}$
- 3. $THH\Sigma_{N-2}$
- 4. $HTT\Sigma_{N-3}$

(c) Now use this argument to show that $P(N) = (1/2)P(N-2) + (1/4)P(N-3)$.

16.4. The gambler’s ruin:

(a) Show that you can rearrange the recurrence relation of example 16.20 to get

$$p_{s+1} - p_s = \frac{(1-p)}{p} (p_s - p_{s-1}).$$

Now show that this means that

$$p_{s+1} - p_s = \left(\frac{(1-p)}{p}\right)^2 (p_{s-1} - p_{s-2})$$

so that

$$\begin{aligned} p_{s+1} - p_s &= \left(\frac{(1-p)}{p}\right)^s (p_1 - p_0) \\ &= \left(\frac{(1-p)}{p}\right)^s (p_1 - 1). \end{aligned}$$

(b) Now we need a simple result about series. Assume I have a series u_k , $k \geq 0$, with the property that

$$u_k - u_{k-1} = cr^{k-1}.$$

Show that

$$u_k - u_0 = c \left(\frac{r^k - 1}{r - 1} \right).$$

(c) Use the results of the last two steps to show that

$$p_s - 1 = (p_1 - 1) \left(\frac{\left(\frac{1-p}{p} \right)^s - 1}{\left(\frac{1-p}{p} \right) - 1} \right)$$

(d) Use the fact that $p_j = 0$ and the result of the last exercise to show

$$(p_1 - 1) = \frac{-1}{\left(\frac{\left(\frac{1-p}{p} \right)^j - 1}{\left(\frac{1-p}{p} \right) - 1} \right)}.$$

(e) Use the results of the previous exercises to show that

$$p_s = \frac{\left(\frac{1-p}{p} \right)^j - \left(\frac{1-p}{p} \right)^s}{\left(\frac{1-p}{p} \right)^j - 1}.$$

PROGRAMMING EXERCISES

- 16.5.** Assume you have an MTG-DAF deck of 24 Lands, 10 Spells of cost 1, 10 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 2 Spells of cost 6. It is properly shuffled, and you draw seven cards. At each turn, you play a land if it is in your hand, and you always only play the most expensive spell in your hand that you are able to play and you never play two spells.
- Write a program to estimate the probability that the first time you play a spell is on turn 4.
 - Write a program to estimate the probability that the first spell you play has cost 4.
 - Extend your program to prepare a 6 x 6 table of turns by spells. In the t , s 'th cell of the table, your program should place the probability that the first spell you play is played on turn t , and has cost s . Notice that many cells easily contain a zero.
- 16.6.** Assume you have an MTG-DAF deck of 10 Lands, 10 Spells of cost 1, 10 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 16 Spells of cost 6. It is properly shuffled, and you draw seven cards. At each turn, you play a land if it is in your hand, and you always only play the cheapest spell in your hand that you are able to play and you never play two spells.
- Write a program to estimate the probability that the first time you play a spell is on turn 4.
 - Write a program to estimate the probability that the first spell you play has cost 4.
 - Extend your program to prepare a 6 x 6 table of turns by spells. In the t , s 'th cell of the table, your program should place the probability that the first spell you play is played on turn t , and has cost s . Notice that many cells easily contain a zero.

- 16.7.** You have three MTG-DAF decks. Deck 1 has 24 Lands, 10 Spells of cost 1, 10 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 2 Spells of cost 6. Deck 2 has 10 Lands, 10 Spells of cost 1, 10 Spells of cost 2, 10 Spells of cost 3, 2 Spells of cost 4, 2 Spells of cost 5, and 16 Spells of cost 6. Deck 3 has 44 Lands, and 16 Spells of cost 6. These decks will be played by players who rigorously implement the following rules: At each turn, a player plays a land if it is in hand, and always plays only the most expensive spell in hand that can be played, and never plays two spells. The first player to play a spell of cost greater than 4 wins. Write a program to estimate the probability of winning for each deck in each possible matchup. In the matchups, you should consider (for example) Deck 1 vs Deck 2, 1 plays first, AND Deck 1 vs Deck 2, 2 plays first.

CHAPTER 17

Math Resources

17.1 USEFUL MATERIAL ABOUT MATRICES

Terminology:

- A matrix \mathcal{M} is **symmetric** if $\mathcal{M} = \mathcal{M}^T$. A symmetric matrix is necessarily square.
- We write \mathcal{I} for the identity matrix.
- A matrix is **diagonal** if the only non-zero elements appear on the diagonal. A diagonal matrix is necessarily symmetric.
- A symmetric matrix is **positive semidefinite** if, for any \mathbf{x} such that $\mathbf{x}^T \mathbf{x} > 0$ (i.e. this vector has at least one non-zero component), we have $\mathbf{x}^T \mathcal{M} \mathbf{x} \geq 0$.
- A symmetric matrix is **positive definite** if, for any \mathbf{x} such that $\mathbf{x}^T \mathbf{x} > 0$, we have $\mathbf{x}^T \mathcal{M} \mathbf{x} > 0$.
- A matrix \mathcal{R} is **orthonormal** if $\mathcal{R}^T \mathcal{R} = \mathcal{I} = \mathcal{I}^T = \mathcal{R} \mathcal{R}^T$. Orthonormal matrices are necessarily square.

Orthonormal matrices: You should think of orthonormal matrices as rotations, because they do not change lengths or angles. For \mathbf{x} a vector, \mathcal{R} an orthonormal matrix, and $\mathbf{u} = \mathcal{R}\mathbf{x}$, we have $\mathbf{u}^T \mathbf{u} = \mathbf{x}^T \mathcal{R}^T \mathcal{R} \mathbf{x} = \mathbf{x}^T \mathcal{I} \mathbf{x} = \mathbf{x}^T \mathbf{x}$. This means that \mathcal{R} doesn't change lengths. For \mathbf{y}, \mathbf{z} both unit vectors, we have that the cosine of the angle between them is $\mathbf{y}^T \mathbf{x}$; but, by the same argument as above, the inner product of $\mathcal{R}\mathbf{y}$ and $\mathcal{R}\mathbf{x}$ is the same as $\mathbf{y}^T \mathbf{x}$. This means that \mathcal{R} doesn't change angles, either.

Eigenvectors and Eigenvalues: Assume \mathcal{S} is a $d \times d$ symmetric matrix, \mathbf{v} is a $d \times 1$ vector, and λ is a scalar. If we have

$$\mathcal{S}\mathbf{v} = \lambda\mathbf{v}$$

then \mathbf{v} is referred to as an **eigenvector** of \mathcal{S} and λ is the corresponding **eigenvalue**. Matrices don't have to be symmetric to have eigenvectors and eigenvalues, but the symmetric case is the only one of interest to us.

In the case of a symmetric matrix, the eigenvalues are real numbers, and there are d distinct eigenvectors that are normal to one another, and can be scaled to have unit length. They can be stacked into a matrix $\mathcal{U} = [\mathbf{v}_1, \dots, \mathbf{v}_d]$. This matrix is orthonormal, meaning that $\mathcal{U}^T \mathcal{U} = \mathcal{I}$. This means that there is a diagonal matrix Λ such that

$$\mathcal{S}\mathcal{U} = \mathcal{U}\Lambda.$$

In fact, there is a large number of such matrices, because we can reorder the eigenvectors in the matrix \mathcal{U} , and the equation still holds with a new Λ , obtained by reordering the diagonal elements of the original Λ . There is no reason to keep track of this complexity. Instead, we adopt the convention that the elements of \mathcal{U} are always ordered so that the elements of Λ are sorted along the diagonal, with the largest value coming first.

Diagonalizing a symmetric matrix: This gives us a particularly important procedure. We can convert any symmetric matrix \mathcal{S} to a diagonal form by computing

$$\mathcal{U}^T \mathcal{S} \mathcal{U} = \Lambda.$$

This procedure is referred to as **diagonalizing** a matrix. Again, we assume that the elements of \mathcal{U} are always ordered so that the elements of Λ are sorted along the diagonal, with the largest value coming first. Diagonalization allows us to show that positive definiteness is equivalent to having all positive eigenvalues, and positive semidefiniteness is equivalent to having all non-negative eigenvalues.

Factoring a matrix: Assume that \mathcal{S} is symmetric and positive semidefinite. We have that

$$\mathcal{S} = \mathcal{U} \Lambda \mathcal{U}^T$$

and all the diagonal elements of Λ are non-negative. Now construct a diagonal matrix whose diagonal entries are the positive square roots of the diagonal elements of Λ ; call this matrix $\Lambda^{(1/2)}$. We have $\Lambda^{(1/2)} \Lambda^{(1/2)} = \Lambda$ and $(\Lambda^{(1/2)})^T = \Lambda^{(1/2)}$. Then we have that

$$\mathcal{S} = (\mathcal{U} \Lambda^{(1/2)}) (\Lambda^{(1/2)} \mathcal{U}^T) = (\mathcal{U} \Lambda^{(1/2)}) (\mathcal{U} \Lambda^{(1/2)})^T$$

so we can factor \mathcal{S} into the form $\mathcal{X} \mathcal{X}^T$ by computing the eigenvectors and eigenvalues.

17.1.1 The Singular Value Decomposition

For any $m \times p$ matrix \mathcal{X} , it is possible to obtain a decomposition

$$\mathcal{X} = \mathcal{U} \Sigma \mathcal{V}^T$$

where \mathcal{U} is $m \times m$, \mathcal{V} is $p \times p$, and Σ is $m \times p$ and is diagonal. If you don't recall what a diagonal matrix looks like when the matrix *isn't* square, it's simple. All entries are zero, except the i, i entries for i in the range 1 to $\min(m, p)$. So if Σ is tall and thin, the top square is diagonal and everything else is zero; if Σ is short and wide, the left square is diagonal and everything else is zero. Both \mathcal{U} and \mathcal{V} are orthonormal (i.e. $\mathcal{U} \mathcal{U}^T = \mathcal{I}$ and $\mathcal{V} \mathcal{V}^T = \mathcal{I}$).

Notice that there is a relationship between forming an SVD and diagonalizing a matrix. In particular, $\mathcal{X}^T \mathcal{X}$ is symmetric, and it can be diagonalized as

$$\mathcal{X}^T \mathcal{X} = \mathcal{V} \Sigma^T \Sigma \mathcal{V}^T.$$

Similarly, $\mathcal{X} \mathcal{X}^T$ is symmetric, and it can be diagonalized as

$$\mathcal{X} \mathcal{X}^T = \mathcal{U} \Sigma \Sigma^T \mathcal{U}.$$

17.1.2 Approximating A Symmetric Matrix

Assume we have a $k \times k$ symmetric matrix \mathcal{T} , and we wish to construct a matrix \mathcal{A} that approximates it. We require that (a) the rank of \mathcal{A} is precisely $r < k$ and (b) the approximation should minimize the **Frobenius norm**, that is,

$$\|(\mathcal{T} - \mathcal{A})\|_F^2 = \sum_{ij} (T_{ij} - A_{ij})^2.$$

It turns out that there is a straightforward construction that yields \mathcal{A} .

The first step is to notice that if \mathcal{U} is orthonormal and \mathcal{M} is any matrix, then

$$\|\mathcal{U}\mathcal{M}\|_F = \|\mathcal{M}\mathcal{U}\|_F = \|\mathcal{M}\|_F.$$

This is true because \mathcal{U} is a rotation (as is $\mathcal{U}^T = \mathcal{U}^{-1}$), and rotations do not change the length of vectors. So, for example, if we write \mathcal{M} as a table of row vectors $\mathcal{M} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k]$, then $\mathcal{U}\mathcal{M} = [\mathcal{U}\mathbf{m}_1, \mathcal{U}\mathbf{m}_2, \dots, \mathcal{U}\mathbf{m}_k]$. Now $\|\mathcal{M}\|_F^2 = \sum_{j=1}^k \|\mathbf{m}_j\|^2$, so $\|\mathcal{U}\mathcal{M}\|_F^2 = \sum_{i=1}^k \|\mathcal{U}\mathbf{m}_k\|^2$. But rotations do not change lengths, so $\|\mathcal{U}\mathbf{m}_k\|^2 = \|\mathbf{m}_k\|^2$, and so $\|\mathcal{U}\mathcal{M}\|_F = \|\mathcal{M}\|_F$. To see the result for the case of $\mathcal{M}\mathcal{U}$, just think of \mathcal{M} as a table of row vectors.

Notice that, if \mathcal{U} is the orthonormal matrix whose columns are eigenvectors of \mathcal{T} , then we have

$$\|(\mathcal{T} - \mathcal{A})\|_F^2 = \|\mathcal{U}^T(\mathcal{T} - \mathcal{A})\mathcal{U}\|_F^2.$$

Now write Λ_r for $\mathcal{U}^T\mathcal{A}\mathcal{U}$, and Λ for the diagonal matrix of eigenvalues of \mathcal{T} . Then we have

$$\|(\mathcal{T} - \mathcal{A})\|_F^2 = \|\Lambda - \Lambda_A\|_F^2,$$

an expression that is easy to solve for Λ_A . We know that Λ is diagonal, so the best Λ_A is diagonal, too. The rank of \mathcal{A} must be r , so the rank of Λ_A must be r as well. To get the best Λ_A , we keep the r largest diagonal values of Λ , and set the rest to zero; Λ_A has rank r because it has only r non-zero entries on the diagonal, and every other entry is zero.

Now to recover \mathcal{A} from Λ_A , we know that $\mathcal{U}^T\mathcal{U} = \mathcal{U}\mathcal{U}^T = \mathcal{I}$ (remember, \mathcal{I} is the identity). We have $\Lambda_A = \mathcal{U}^T\mathcal{A}\mathcal{U}$, so

$$\mathcal{A} = \mathcal{U}\Lambda_A\mathcal{U}^T.$$

We can clean up this representation in a useful way. Notice that only the first r columns of \mathcal{U} (and the corresponding rows of \mathcal{U}^T) contribute to \mathcal{A} . The remaining $k - r$ are each multiplied by one of the zeros on the diagonal of Λ_A . Remember that, by convention, Λ was sorted so that the diagonal values are in descending order (i.e. the largest value is in the top left corner). We now keep only the top left $r \times r$ block of Λ_A , which we write Λ_r . We then write \mathcal{U}_r for the $k \times r$ matrix consisting of the first r columns of \mathcal{U} . Then

$$\mathcal{A} = \mathcal{U}_r\Lambda_r\mathcal{U}_r^T$$

This is so useful a result, I have displayed it in a box; you should remember it.

Procedure: 17.1 *Approximating a symmetric matrix with a low rank matrix*

Assume we have a symmetric $k \times k$ matrix \mathcal{T} . We wish to approximate \mathcal{T} with a matrix \mathcal{A} that has rank $r < k$. Write \mathcal{U} for the matrix whose columns are eigenvectors of \mathcal{T} , and Λ for the diagonal matrix of eigenvalues of \mathcal{A} (so $\mathcal{A}\mathcal{U} = \mathcal{U}\Lambda$). Remember that, by convention, Λ was sorted so that the diagonal values are in descending order (i.e. the largest value is in the top left corner).

Now construct Λ_r from Λ by setting the $k - r$ smallest values of Λ to zero, and keeping only the top left $r \times r$ block. Construct \mathcal{U}_r , the $k \times r$ matrix consisting of the first r columns of \mathcal{U} . Then

$$\mathcal{A} = \mathcal{U}_r \Lambda_r \mathcal{U}_r^T$$

is the best possible rank r approximation to \mathcal{T} in the Frobenius norm.

Now if \mathcal{A} is positive semidefinite (i.e. if at least the r largest eigenvalues of \mathcal{T} are non-negative), then we can factor \mathcal{A} as in the previous section. This yields a procedure to approximate a symmetric matrix by factors. This is so useful a result, I have displayed it in a box; you should remember it.

Procedure: 17.2 *Approximating a symmetric matrix with low dimensional factors*

Assume we have a symmetric $k \times k$ matrix \mathcal{T} . We wish to approximate \mathcal{T} with a matrix \mathcal{A} that has rank $r < k$. We assume that at least the r largest eigenvalues of \mathcal{T} are non-negative. Write \mathcal{U} for the matrix whose columns are eigenvectors of \mathcal{T} , and Λ for the diagonal matrix of eigenvalues of \mathcal{A} (so $\mathcal{A}\mathcal{U} = \mathcal{U}\Lambda$). Remember that, by convention, Λ was sorted so that the diagonal values are in descending order (i.e. the largest value is in the top left corner).

Now construct Λ_r from Λ by setting the $k - r$ smallest values of Λ to zero and keeping only the top left $r \times r$ block. Construct $\Lambda_r^{(1/2)}$ by replacing each diagonal element of Λ with its positive square root. Construct \mathcal{U}_r , the $k \times r$ matrix consisting of the first r columns of \mathcal{U} . Then write $\mathcal{V} = (\mathcal{U}_r \Lambda_r^{(1/2)})$

$$\mathcal{A} = \mathcal{V}\mathcal{V}^T$$

is the best possible rank r approximation to \mathcal{T} in the Frobenius norm.

Error functions and Gaussians: The **error function** is defined by

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

and programming environments can typically evaluate the error function. This fact is made useful to us by a simple change of variables. We get

$$\frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{u^2}{2}} du = \frac{1}{\sqrt{\pi}} \int_0^{\frac{x}{\sqrt{2}}} e^{-t^2} dt = \frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right).$$

A particularly useful manifestation of this fact comes by noticing that

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^0 e^{-\frac{t^2}{2}} dt = 1/2$$

(because $\frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$ is a probability density function, and is symmetric about 0). As a result, we get

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt = 1/2(1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right)).$$

Inverse error functions: We sometimes wish to know the value of x such that

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt = p$$

for some given p . The relevant function of p is known as the **probit function** or the **normal quantile function**. We write

$$x = \Phi(p).$$

The probit function Φ can be expressed in terms of the **inverse error function**. Most programming environments can evaluate the inverse error function (which is the inverse of the error function). We have that

$$\Phi(p) = \sqrt{2} \operatorname{erf}^{-1}(2p - 1).$$

One problem we solve with some regularity is: choose u such that

$$\int_{-u}^u \frac{1}{\sqrt{2\pi}} \exp(-x^2/2) dx = p.$$

Notice that

$$\begin{aligned} \frac{p}{2} &= \frac{1}{\sqrt{2\pi}} \int_0^u e^{-\frac{t^2}{2}} dt \\ &= \frac{1}{2} \operatorname{erf}\left(\frac{u}{\sqrt{2}}\right) \end{aligned}$$

so that

$$u = \sqrt{2} \operatorname{erf}^{-1}(p).$$