

- Book:
- Hastie, Tibshirani + Friedman
 - Material from Jordan's course notes
 - Papers, etc. + Semi-supervised + Bishop

- Themes:
- Classification
 - Regression
 - Inference

Regression: • predict some value y from vars x

- Choose some loss function

$$L(y, f(x))$$

and obtain f by minimizing

$$E[L(y, f(x))]$$

with respect to θ

important case: • Square error loss

$$L(Y, f(x)) = (y - f(x))^2$$

hence

$$\min E((y - f(x))^2) \quad \text{EPE} = \text{expected prediction error}$$

$$= E_x \left[E_{Y|X} \left[(y - f(x))^2 \right] \right]$$

and so we can minimize

$$E_{Y|X} \left[(y - f(x))^2 \right]$$

$$= \int (y - f(x))^2 P(y|x) dy$$

at each x .

but the solution is

$$f(x) = E(Y|X=x)$$

$$= \int y P(Y|X=x) dy.$$

Example cases:

linear regression

$$f(x) = \underline{x}^T \underline{\beta}$$

we usually ensure that one column of x is a 1 to simplify notation

$$\begin{aligned} EPE &= E[(Y - \underline{x}^T \underline{\beta})^2] \\ &= E[Y^2] - 2E[\underline{x}^T Y] + \underline{\beta}' E[\underline{x} \underline{x}'] \underline{\beta} \end{aligned}$$

so

$$\underline{\beta} = E[\underline{x} \underline{x}']^{-1} E[\underline{x} Y]$$

alternatively

what we're accustomed to is

$$Y = \begin{bmatrix} y_1 \\ \vdots \end{bmatrix}; \quad X = \begin{bmatrix} x_1 \\ \vdots \end{bmatrix}$$

$$\min \| Y - X \underline{\beta} \|^2 \quad \underline{\beta} = (X' X)^{-1} (X' Y)$$

→ this replaces $E_{x,y}$ with Σ_{data}

k-nearest neighbours

$$f(x) = \frac{1}{K} \sum_{\substack{\text{K nearest examples} \\ \text{closest to } x}} y_i$$

now this approximates $E(Y|X=x)$

by assuming that x changes
"slowly" and then replacing an expect
with a sum.

Notice:

for most $P(X, Y)$ as $N, k \rightarrow \infty$
and $\frac{k}{N} \rightarrow 0$

$$f(x) \rightarrow E(Y|X=x).$$

(5)

How does this relate to classification?

- regression where y is a categorical var
- for now, $y \in \{0, 1\}$

Zero-One loss

$$L(y, f(x)) = \begin{cases} 1 & \text{don't agree} \\ 0 & \text{otherwise} \end{cases}$$

$$EPE = E_x \left[E_{y|x} [L(y, f(x))] \right]$$

$$= E_x \left[L(1, f(x)) P(1|x=x) + L(0, f(x)) P(0|x=x) \right]$$

→ we can minimize pointwise.

by choosing the ~~class~~ ^{$f(x)$} such that

$$L(\text{class}, L(1, f(x)) P(1|x=x) + L(0, f) P(0|x=x))$$

is min.

(6)

for 0-1 loss, we get

1 if $P(1/x) < P(0/x)$
0 if $P(1/x) > P(0/x)$
(doesn't matter) if $P(1/x) = P(0/x)$

We could do this with k-NN easily

- but there are problems in high dimensions

The curse of Dimension

1) Where is volume in HD?

- on the skin

2) What % of data is nearby?

- very little

e.g. →

(7)

Unit cube, p dim, Uniform data.

- we want fraction r of data.

- This is r of volume
 $\Rightarrow \sim^{1/p}$ of edge length

10 dim, to capture 10% of data

we need edge length of 0.8

- hardly local

Bias and Variance:

• consider $y = f(x) = e^{-\|x\|^2}$

(Deterministic f_n)

- we will draw a bunch of (x, y) samples T
and approx f_n using 1-NN

- All error is due to choice of samples.

- consider $\hat{y}_0 =$ 1-NN est of $y(x_0)$

8

and we can write

$$MSE(x_0) = E_T \left[(f(x_0) - \hat{y}_0)^2 \right]$$

$$= f(x_0)^2 - 2 E_T[y_0] f(x_0) + E_T[y_0^2]$$

$$= f(x_0)^2 - 2 E_T[y_0] f(x_0) + (E_T[y_0])^2$$

$$+ (E_T[y_0])^2 - 2 (E_T[y_0])^2 + E_T[y_0^2]$$

$$= (f(x_0) - E_T[y_0])^2 \quad \leftarrow \text{BIAS}^2$$

$$+ E_T \left[(y_0 - E_T[y_0])^2 \right] \quad \leftarrow \text{VARIANCE}$$

- this sort of decomposition is universal.
- In our example, as inclusion goes up, $E_T[y_0]$ goes down fast (why?)
so bias goes up fast

~~linear~~ Stabilizing linear regression

- wish to reduce prediction error
 - by ditching variables
 - Disadvantages:
 - search
 - variance could go up.
- penalizing large coeffs
(in the hope they shrink to zero)

$$\text{ridge } R(\beta) = \{ \|Y - \beta_0 \mathbf{1} - X\beta\|^2 + \lambda \beta' \beta \}$$

- we don't want to penalize the constant because that would result in fits that aren't invariant under translation in y — doesn't make sense.

- for this discussion, center y 's

$$(i.e. \sum_i y_i = 0) \Rightarrow \beta_0 = 0$$

Then $(Y - X\beta)'(Y - X\beta) + \lambda\beta'\beta$

or $\beta^{ridge} = \underbrace{(X'X + \lambda I)^{-1}}_{\uparrow}$ $X'Y$.

Deals with possible rank problems.

Notice that

$$\min (Y - X\beta)'(Y - X\beta) + \lambda\beta'\beta$$

$$\lambda \geq 0$$

|||

$$\min (Y - X\beta)'(Y - X\beta)$$

$$st \quad \beta'\beta \leq S$$

Notice that:

- ridge regression isn't covar under scaling of inputs (why?)
- qualitative arg:
 - β can be poorly determined in presence of correlated vars

Ridge regression + SVD ↙ diagonal.

write $X = U D V^T$

$N \times P$ $P \times P$

↖ orthogonal

$$\beta^{LS} = (X^T X)^{-1} X^T y$$

$$\begin{aligned} \text{so } X \beta^{LS} &= X (X^T X)^{-1} X^T y \\ &= U \underbrace{U^T y} \end{aligned}$$

↳ these are coords of y in U basis.

$$X\beta^{\text{ridge}} = X(X^T X + \lambda I)^{-1} X^T y$$

$$= U D (D^2 + \lambda I)^{-1} D U^T y$$

$$= \sum_j u_j \left[\frac{d_j^2}{d_j^2 + \lambda} \right] u_j^T y$$

col of u

- we shrink coords by $\frac{d_j^2}{d_j^2 + \lambda}$

- shrinkage is most pronounced for dims where d_j^2 is small

\Rightarrow components of x with low variance.

$\hookrightarrow \Rightarrow$ poor estimates of gradient of y in this dir.

Lasso :

$$\beta = \underset{\text{min}}{\text{arg}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2$$

st

$$\sum_j |\beta_j| \leq t$$

• can no longer use linear alg; this is a quadratic programming problem

• sufficiently small t forces some

$$\beta_j = 0$$