

Classification by linear methods:

(Ch4 Hastie Tibshirani Freedman).

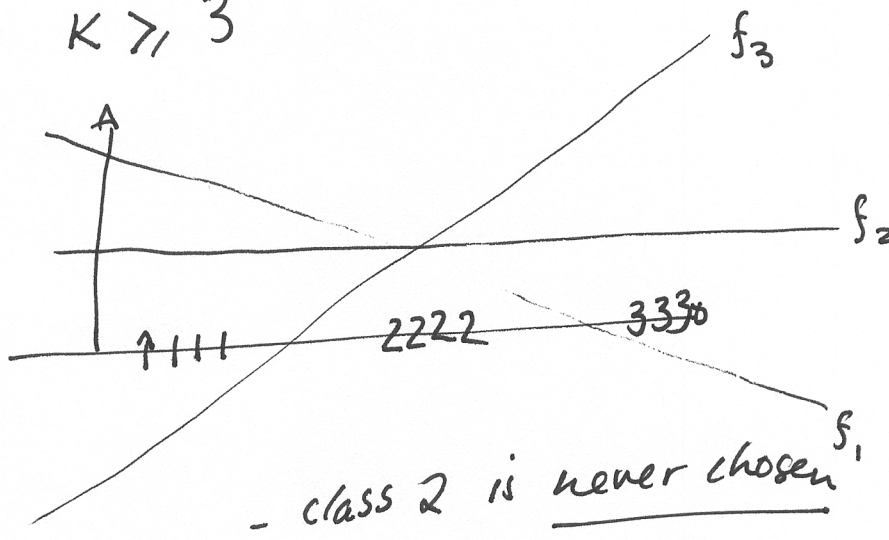
- k. classes
- generally, fit ~~for~~ k fns, and choose class w/ largest discriminant fns.

Simplest:

- | | | | |
|-------|--|---|----------------------------------|
| train | | - | use linear regression |
| | | | - fit k l.r.s to indicator vars. |
| test | | - | eval all k, |
| | | | - choose largest. |

• Serious problem for $k \geq 3$

• masking



• one might fix this w/ polynomial terms, but this is expensive

Linear Discriminant Analysis

1) Notice that decision theory gives us:
choose k for which $P(k|x)$ is largest

this is
$$\frac{P(x|k) \cdot P(k)}{\sum_{i \in \text{classes}} P(x|i) P(i)}$$

So that a model of $P(x|i)$ is very valuable

2) Assume: $P(x|i) = N(\mu_i; \Sigma)$
 \uparrow all classes have the same

so that $P(x|i)$

$$= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2}(x-\mu_i)^T \Sigma^{-1} (x-\mu_i)}$$

3) we can classify with $\log P(k|x)$ because it's log monotone

$$\log \left[\frac{P(k|x)}{P(l|x)} \right] = \log \left[\frac{P(x|k)}{P(x|l)} \right] + \log \left[\frac{\pi_k}{\pi_l} \right]$$

and in our case, this is

③ of

$$-\frac{1}{2} [(\mu_k + \mu_e)' \Sigma^{-1} (\mu_k - \mu_e)] + x^T \Sigma^{-1} (\mu_k - \mu_e) + \log\left(\frac{\pi_k}{\pi_e}\right)$$

↑ linear in x (Don't forget the priors!)

we can describe this decision rule w/
linear discriminant fun

$$S_k(x) = x' \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k' \Sigma^{-1} \mu_k + \log \pi_k$$

classification rule: choose class w/ largest

Note: Σ, μ_i unknown \therefore must estimate
notice we should have a quite good est for Σ .

Similar trick works for situation where covariance varies ("heteroskedastic")

• but now we have quadratic disc

Quadratic Discs

(4) of

$$g_k = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)' \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

- quadratic decision boundary
- more estimation \Rightarrow ~~less~~ more variance

Uncomfortable truths:

- linear decision boundaries are v. effective
- LDA is a good way to estimate them
- LDA / QDA are very competitive classifiers

Feature trick:

- a boundary that isn't linear in x might be linear in $\varphi(x)$.

Regularization:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}_k^{-1}$$

or even
$$\hat{\Sigma}_k(\gamma) = \gamma \hat{\Sigma}_k + (1-\gamma) \sigma^2 I$$

α, γ , est by cross-validation

Fisher linear discriminants
 (or) discriminant coordinates or canonical variates ⑤ of

Find $z = \alpha^T X$ st between class var of within class var

cf. Principal coordinates
 find $z = \alpha^T X$ st variance of z is maximised

between class var $\Sigma_B = \frac{1}{\# \text{ of classes} - 1} \left[\sum_k (\mu_k - \mu_{\text{data}})(\mu_k - \mu_{\text{data}})^T \right]$

within class var $\Sigma_w = \frac{1}{\# \text{ data} - 1} \left[\sum_{i \in \text{data}} [x_i - \mu_{c(i)}][x_i - \mu_{c(i)}]^T \right]$

for $z = \alpha^T x$

b.c.v. = $z^T \Sigma_B z$

w.c.v. = $z^T \Sigma_w z$

$\therefore \max \frac{z^T \Sigma_B z}{z^T \Sigma_w z}$ | \leftarrow hard

$\max z^T \Sigma_B z$ st. $z^T \Sigma_w z = 1$

Logistic regression:

⑥ of

$$\log \left[\frac{P(i|x)}{P(k|x)} \right] = \beta_i' x$$

there's a 1 in this.

$$\therefore P(i|x) = \frac{\exp[\beta_i' x]}{1 + \sum_{l=1}^{k-1} [\exp(\beta_l' x)]}$$

← k class

Note:

- this gives a P.L. decision boundary
- if $k = \# \text{ of classes} = 2$ linear

Fitting:

- by max likelihood

2 class case:

$$P(1|x) = P(x; \beta) = 1 - P(0|x)$$

$$P(x; \beta) = \frac{\exp(\beta' x)}{1 + \exp(\beta' x)}$$

$$l(\beta) = \sum_{i \in \text{data}} \left[\log \left[P(y_i | x_i; \beta) \right] \right]$$

$$= \sum_{i \in \text{data}} \left[y_i \log P(x_i; \beta) + (1 - y_i) \log [1 - P(x_i; \beta)] \right]$$

$$= \sum_i \left[y_i \beta' x_i - y_i \log (1 + \exp(\beta' x_i)) - (1 - y_i) \log (1 + \exp(\beta' x_i)) \right]$$

Max yields

$$\nabla l = 0 = \nabla \left[\sum_i \left\{ y_i (\beta' x_i - \log(1 + e^{\beta' x_i})) \right\} \right]$$

$$= \sum_i x_i (y_i - p(x_i; \beta))$$

some mild reorg

$$\text{and } H_l = - \sum_i x_i x_i' \left[p(x_i; \beta) (1 - p(x_i; \beta)) \right]$$

write:

$$y = \text{vec } y_i ; \quad X = \begin{bmatrix} x_1' \\ \vdots \\ x_N' \end{bmatrix} ;$$

$$\underline{P}^{(n)} = \text{vec } p(x_i; \beta^{(n)})$$

$$\underline{W}^{(n)} = \text{diag} \left[p(x_i; \beta^{(n)}) \cdot (1 - p(x_i; \beta^{(n)})) \right]$$

Newton's method becomes

$$\beta^{(n+1)} = \beta^{(n)} + \underbrace{\left(X^T W^{(n)} X \right)^{-1}}_{\text{Hessian}} \underbrace{X^T (y - P)}_{\text{gradient}}$$

$$= (X^T W^{(n)} X)^{-1} \cdot [X^T W z]$$

where

$$z = X \beta^{(n)} + W^{-1} (y - P)$$

But this is a form of iteratively reweight
least squares.

problem is convex but the is mis

- We could focus on search for decision $\textcircled{8}$ of boundary (rather than on prob model)
- e.g. choose best linear boundary

Perceptron

$$y_i \in \{1, -1\}$$

minimize

$$Err = - \sum_{i \in \text{Misclassified}} y_i (x_i' \beta + \beta_0)$$

$$\nabla E = \begin{bmatrix} - \sum_{i \in M} y_i x_i \\ - \sum_{i \in M} y_i \end{bmatrix}$$

minimize by Stochastic gradient descent.

- known to converge in finite # of steps for separable data

* but

- more than one soln
- can be many steps
- if not separable, mischief including long cycles.

Sideline:

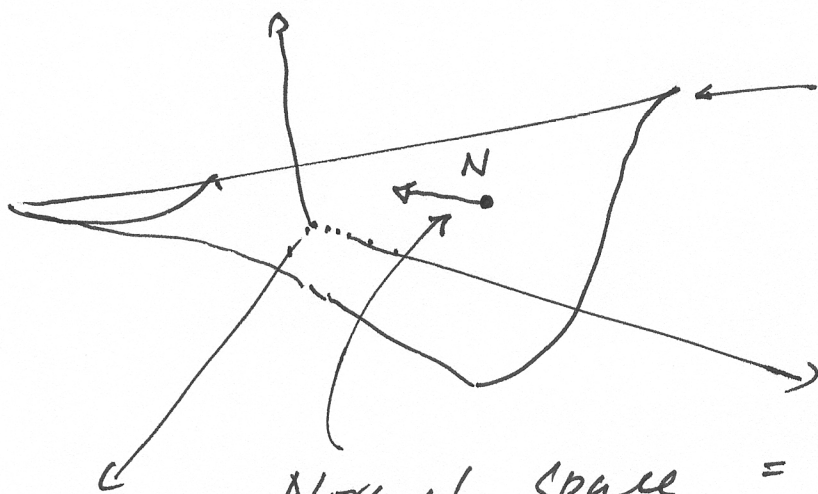
② 17

· constrained optimization

Consider

$$\max f_0(x) \quad \text{s.t.} \quad f_1(x) \cdots f_k(x) = 0$$

at an extremal value, ~~the~~ ~~the~~ ∇f will be normal to constraint surface



$$\begin{aligned} f_1 &= 0 \\ f_2 &= 0 \\ &\vdots \\ f_k &= 0 \end{aligned}$$

$$\text{Normal space} = \text{span} [\nabla f_1 \cdots \nabla f_k]$$

$$\therefore \nabla f_0 \in \text{span} [\nabla f_1 \cdots \nabla f_k]$$

$$\therefore \nabla f_0 + \lambda_1 \nabla f_1 + \lambda_2 \nabla f_2 \cdots \lambda_k \nabla f_k = 0$$

and

$$\begin{aligned} f_1 &= 0 \\ f_2 &= 0 \\ f_3 &= 0 \\ &\vdots \end{aligned}$$

This is what Lagrangian is all about

Inequalities:

$$\max f_0(x) \quad \text{st} \quad f_i(x) \cdot f_k = 0$$

$$g_i(x) \cdot g_{nr} \geq 0$$

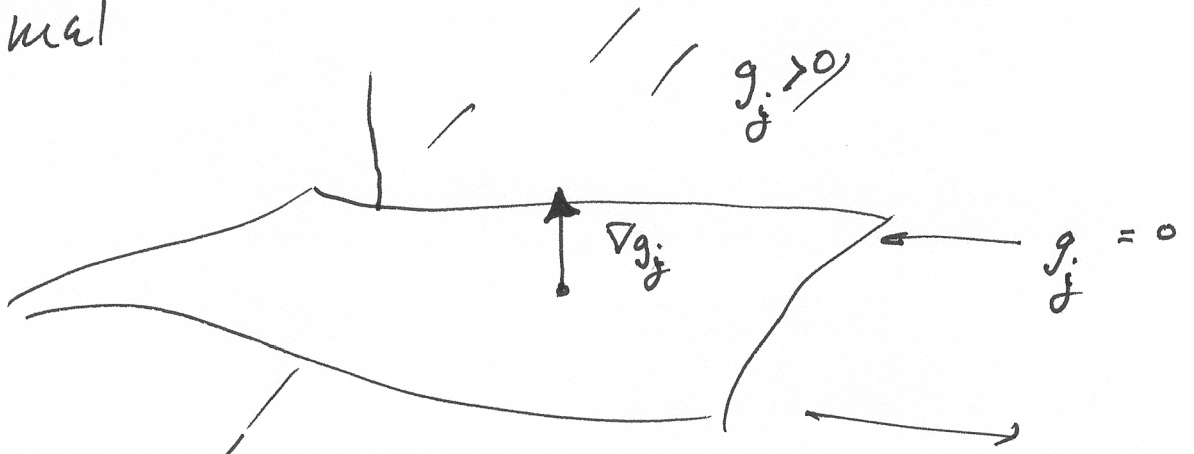
a constraint i is active $\equiv g_i = 0$
inactive $\equiv g_i > 0$

i active $\Rightarrow \nabla f_0$ is normal to g_i , too
i inactive $\Rightarrow g_i$ irrelevant.

$$\therefore \nabla f_0 + \sum_i \lambda_i f_i + \sum_{j \in \text{active}} \nu_j g_j = 0$$

notice that we care about ν_j of

Normal



so ν_j must be negative

Whence

(11)

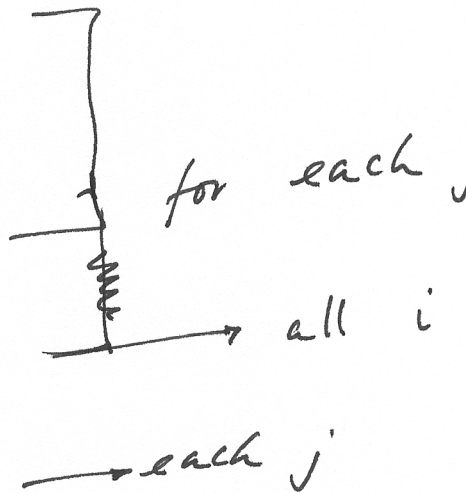
$$\nabla f_0 + \sum_i \lambda_i f_i + \sum_{j \in \text{all}} \nu_j g_j = 0$$

$$\nu_j \leq 0$$

$$\nu_j g_j = 0$$

$$f_i = 0$$

$$g_j \geq 0$$



Karush - Kuhn - Tucker conditions