

ol notes ← James et al notes

SL P1

01, 1

1) Statistical learning could be defined as the statistical analysis of multivariate data. Machine learning, data mining, big data and predictive analytics are synonymous.

2) 01, 30 A case or observation consists of  $k$  random variables measured on one person or thing. The  $i$ th case  $\underline{z}_i = (z_{i1}, \dots, z_{ik})^T$ .

The training data consists of  $\underline{z}_1, \dots, \underline{z}_n$ .

The statistical model is fit (trained) on the training data. The test data is  $\underline{z}_{n+1}, \dots, \underline{z}_{n+m}$  and is used to evaluate the quality of the fitted model.

3) 02 The focus of supervised learning is predicting a future value of a response variable  $Y$  given  $\underline{x}$  and training data  $(Y_1, \underline{z}_1), \dots, (Y_n, \underline{z}_n)$ . The focus of

unsupervised learning is to group

$x_1, \dots, x_n$  into clusters. Data mining  
is looking for relationships in large data sets.

ex)  $y =$  college GPA  $x_1 =$  high school GPA

$x_2 =$  ACT or SAT score  $x_3 =$  gender

marry data

ex) Oex 3.2 p 96  $y =$  # of women married  
to civilians in  $n = 26$  districts in  
Prussia in 1843.

$x_1 =$  constant = trivial predictor

$x_2 =$  pop of district in 1843

$x_3 =$  mmen = # married civilian men in the  
district

$x_4 =$  mmilmen = # married men in the military in the district

$x_5 =$  milwmn = # women married to husband in the military

$n = 26 > 5p$ ,  $p = 5$ , sometimes the surveyor

would not count the spouse if the spouse was

not at home but  $y = x_3 + e$  is a good model,

ex) p 4  $n = 64$  cancer cell lines  $p = 6830$ ,  $x_1, \dots, x_p$   
from 14 cancer types

or 6830 gene expression measurements

#### 4) Classical statistics

i)  $n \geq 10p$

ii) all  $p$  variables  $X_1, \dots, X_p$  are used in the model

iii) 1 model see ii)  
eg  $Y = \beta_1 + \beta_2 X_3 + e$   
for many data

iv) inference focus is more on hypothesis testing than prediction

#### Statistical Learning SL P2

$n \geq p$  or  $n < p$ , often  $n \geq 10d$   
where  $d =$  model degrees of freedom

there is variable selection  
 $X_{i_1}, \dots, X_{i_d}$  are used

model selection  
choose 1 model from several

eg  $Y = \beta_1 + \beta_2 X_3 + e$

$Y = \beta_1 + \beta_2 X_3 + \beta_3 X_2 + e$

$Y = \beta_1 + \beta_2 X_3 + \beta_3 X_2 + \beta_4 X_4 + e$

$Y = \beta_1 + \beta_2 X_3 + \beta_3 X_2 + \beta_4 X_4 + \beta_5 X_5 + e$

inference focus is more on prediction

#### 5) 0.4 Statistical learning principles:

i) there is more interest in prediction or classification eg producing  $\hat{Y}$

than other types of inference such as parameter estimation, hypothesis testing, confidence intervals or which model fits best.

ii) often the focus is on extracting information

when  $N/p$  is not large eg  $p > N$ ,

If  $d$  is a crude estimator of fitted model degrees of freedom (df), we want  $N/d$  large. A sparse model

has few nonzero coefficients and often  $d = \#$  nonzero coefficients

( $\hat{y} = \hat{\beta}_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$  ;  $p$  is large but most  $\hat{\beta}_i = 0$ .)

iii) Interest is in performance on test data.

Performance on training data, is overly optimistic for estimating performance on training data.

iv) Some methods are flexible while

others are unflexible. ( $sp = \beta_1 + \beta_2 x_2$  is a line  
 $sp = \beta_1 + \beta_2 x_2 + \beta_3 x_3$  is a plane.

$sp = \sum \beta_j x_j$  is a hyper plane and is unflexible

Flexibility tends to increase with  $d$ . Unflexible regression methods are often used when the mean function  $E(y|x) = M(x^T \beta)$  is known up to the  $p$  unknown parameters  $\beta$ . Flexible methods lead to be

6] The response variable is the variable of interest; the variable you want to predict. SL P3  
 The predictors or features  $x_1, \dots, x_p$  are used to predict  $Y$ .

7] GLM Regression investigates how  $Y$  changes with the  $p \times 1$  vector  $\underline{x}$  of predictors. Often this conditional distribution  $Y | \underline{x}$  is described by a 1D regression model where  $Y$  is conditionally independent of  $\underline{x}$  given the sufficient predictor  $SP = h(\underline{x})$ , written  $Y \perp\!\!\!\perp \underline{x} | SP$  or  $Y \perp\!\!\!\perp \underline{x} | h(\underline{x})$  where  $h: \mathbb{R}^p \rightarrow \mathbb{R}$ . The estimated sufficient predictor  $ESP = \hat{h}(\underline{x})$ .

[ex]  $Y = m(\underline{x}) + e \rightarrow h(\underline{x}) = m(\underline{x}), \hat{h}(\underline{x}) = \hat{m}(\underline{x})$

$Y \sim \text{Poisson}[\exp(\underline{\beta}'\underline{x})] \rightarrow h(\underline{x}) = \underline{\beta}'\underline{x}, \hat{h}(\underline{x}) = \hat{\underline{\beta}}'\underline{x}$   
 $Y \sim \text{binomial}\left(1, \frac{\alpha + \sum_{j=1}^p s_j(x_j)}{1 + \alpha + \sum_{j=1}^p s_j(x_j)}\right),$   $h(\underline{x}) = \alpha + \sum_{j=1}^p s_j(x_j)$   
 SP for a generalized linear model GLM  
 SP for a generalized additive model GAM

$\hat{h}(\underline{x}) = \hat{\alpha} + \sum_{j=1}^p \hat{s}_j(x_j)$

8) <sup>05</sup> know

A response plot is a plot of the ESP  $\hat{m}(x_i)$  vs.  $Y_i$ . A residual plot is a plot of the ESP vs. the residuals  $\hat{\epsilon}_i$ .

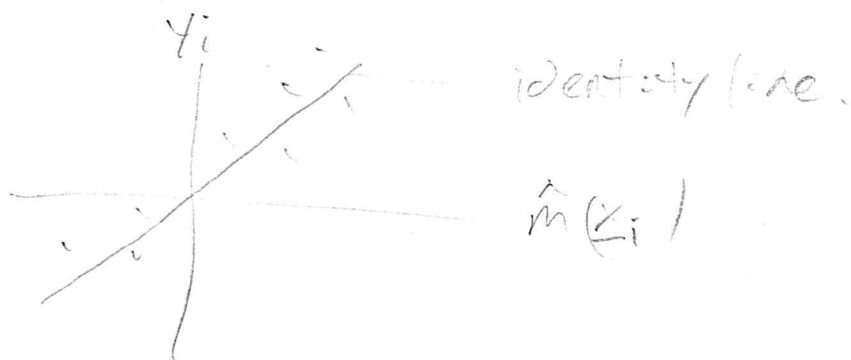
ex] If  $Y_i = m(x_i) + \epsilon_i$  then the  $i$ th residual

$$\hat{\epsilon}_i = Y_i - \hat{m}(x_i) \quad i=1, \dots, n.$$

Often add the <sup>estimated</sup> model mean function and a scatterplot smoother to the response plot.

ex] For  $Y = m(x) + \epsilon$  the identity line is the <sup>unit slope & intercept</sup>  $Y = m(x)$   
estimated mean function  $\hat{m}(x)$

$Y = \hat{m}(x)$  is the identity line in a plot of  $\hat{m}(x)$  vs  $Y$



9) 03.2.2 PB response transformations

Let  $Z$  be the variable of interest and let

$Y = Z(\beta) = Z'\beta + \epsilon$  follow the multiple linear regression (MLR) model.

10) 09,10 Assume all values of the "response"  $z$  are positive. The

ladder of powers  $\Lambda_L = \{-1, -\frac{1}{2}, -\frac{1}{3}, 0, \frac{1}{3}, \frac{1}{2}, 1\}$ .

A power transformation has the form

$$Y = t_\lambda(z) = z^\lambda \quad \text{for } \lambda \neq 0$$

$$Y = t_0(z) = \log|Y| \quad \text{for } \lambda = 0.$$

A modified power transformation has the form

$$t_\lambda(z) = \frac{z^\lambda - 1}{\lambda} \quad \lambda \neq 0$$

$$t_0(z) = \log(z) \quad \lambda = 0.$$

For the additive error regression model  $Y = \mu + \epsilon$

11) 10M for MMR  $Y = X\beta + \epsilon$ , a graphical

method for response transformations computes the "fitted values"  $\hat{w}_i$  using  $w_i = t_\lambda(z_i)$  as

the response. A transformation plot is a plot

of  $\hat{w}_i$  vs  $w_i$  and is made for each of the 7

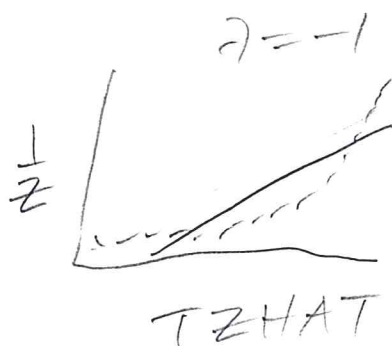
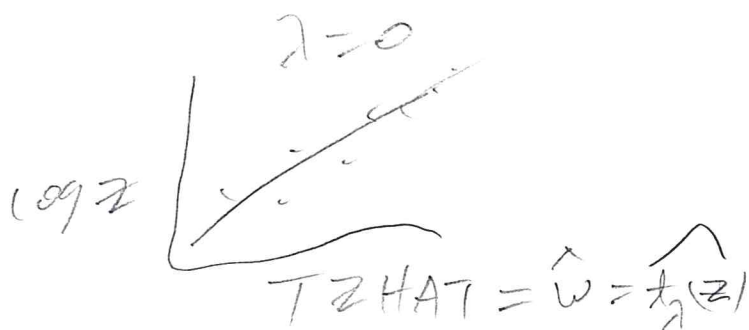
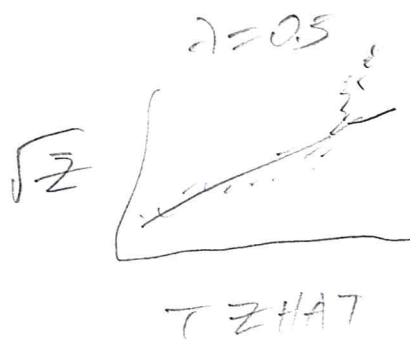
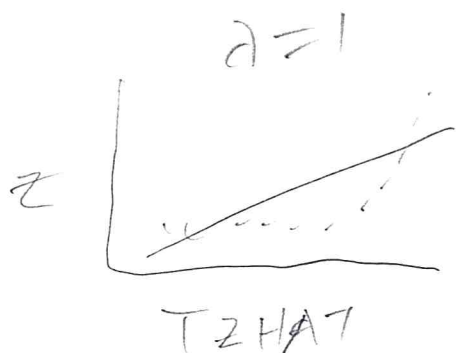
values of  $\lambda \in \Lambda_L$  with the identity line added

as a visual aid. Use the transformation

10)  $Y = f_{\lambda^*}(z)$  if the plotted points in the transformation plot for  $\lambda^*$  follow the identity line in a roughly evenly populated band.

See Oex 1.3 and HW1 C.

12) know for exam 1 how to pick the transformation from a few plots



pick  $Y = \log(z)$  since this transformation has a linear transformation plot.

13) If more than 1 transformation plot is linear take the simplest or most reasonable transformation that



that makes the most sense to subject SL5 matter experts. Also check which of the competing transformations has the best "residual plot" of  $\hat{w}$  vs  $w - \hat{w}$ .

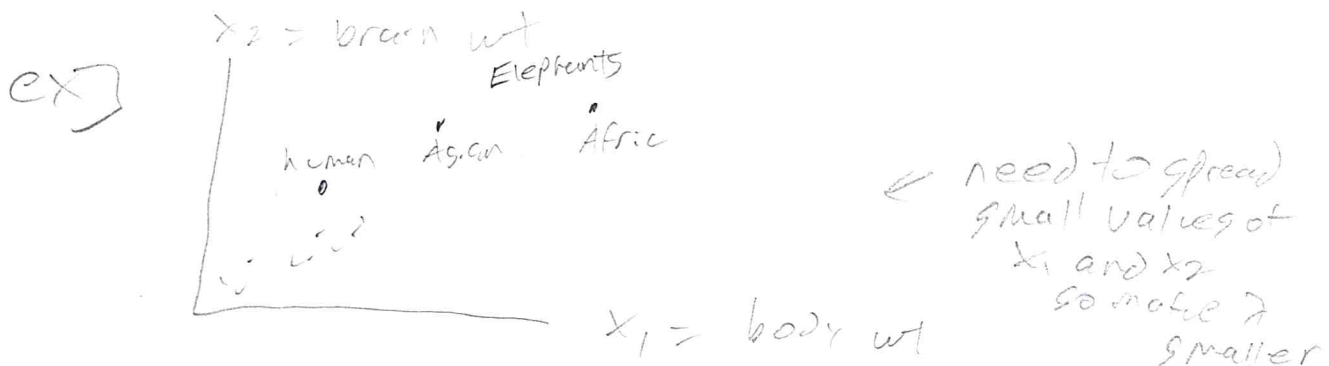
The values of  $\lambda$  in decreasing order of importance are  $1, 0, \frac{1}{2}, -1$  and  $\frac{1}{3}$  so a log transformation is preferred over a cube root transformation if both transformation plots look equally good.

14} know for exam 1. The log rule says that a positive variable  $w > 0$  with  $\frac{\max(w)}{\min(w)} > 10$  suggests using

$\log(w)$  instead of  $w$ . This rule removes skew from predictors and can greatly improve the model.

15} know for exam 1 The following rule is used for plots of 2 variables, usually ESP vs  $Y$  in this class for response variables.

Ladder rule In a plot of  $x_1$  vs  $x_2$ ,  $\begin{matrix} x_2 \\ \swarrow \\ x_1 \end{matrix}$   
 to spread small values of a variable,  
 make  $\lambda$  smaller. To spread large values  
 of a variable make  $\lambda$  larger.



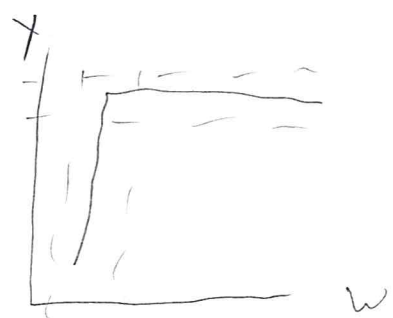
Both variables are right skewed so there are lots of small values and a few large values

$$\frac{\max x_2}{\min x_2} \approx \frac{6654.2}{0.005} > 10$$

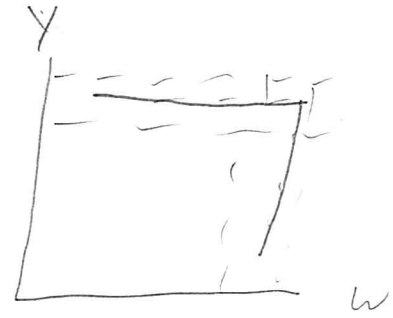
$$\frac{\max x_1}{\min x_1} \approx \frac{5711.9}{0.14} > 10$$

} log rule

ex Fig 1.2 09



small values of  $w$  and large values of  $y$  need spreading



large values of both variables need spreading



small values of both variables



small values of  $y$  and large values of  $w$



ex) the median is the middle ordered value or average of the 2 middle values.

1, 2, 3, 4, 5, 6, 7, 8, 9  
 $\uparrow$   
 $MED(n) = 5$

$|Y_i - MED(n)|^0: +4, +3, +2, +1, 0, 1, 2, 3, 4$

ordered  $\rightarrow$  0, 1, 1, 2, 2, 3, 3, 4, 4  
 $\uparrow$

$MAD(n) = 2$

22) Collect the data:  $W = X = \begin{pmatrix} x_{11} \\ \vdots \\ x_{n1} \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$   
 Ath case  $\rightarrow$   $x_{n1}$   $\downarrow$  Pth variable

$= \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$   
 $\uparrow$   
 2nd variable

23) <sup>014</sup> The coordinatewise median  $MED(W) =$

$(MED(x_1), \dots, MED(x_n))^T$  where  $MED(x_i)$  is the gaultle median of the  $i$ th column

247 know that  $\dots$  for details coordinatewise

Median for a small data set.

see O'ex 1.5 problems 1.1 and 1.2

HW 2

SL 7

see notes 6.5

2.5) 015 The sample covariance  $S_{ij}$  estimates the pop covariance  $\sigma_{ij} = E[(X_i - E(X_i))(X_j - E(X_j))]$  and

$$S_{ij} = \frac{1}{n-1} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j),$$

The sample variance  $S_{ii} = S_i^2$ ,  $\sigma_{ii} = \sigma_i^2 = \text{pop var.}$

The sample correlation  $r_{ij}$  estimates the pop

corr  $\rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$  and  $r_{ij} = \frac{S_{ij}}{S_i S_j} = \frac{S_{ij}}{\sqrt{S_{ii} S_{jj}}}$

2.6) <sup>015-16</sup> Let  $\underline{x}_1, \dots, \underline{x}_n$  be the data.  
 $p \times 1$

The sample covariance matrix  $S = (S_{ij})_{p \times p}$

The classical estimator of multivariate location and dispersion (MLD) is

$(T, S) = (\bar{\underline{x}}, S)$ . The sample correlation matrix  $R = (r_{ij})$ .

2.7) 017 know the  $i$ th Mahalanobis distance  $D_i = \sqrt{(T_i - T)^T W (T_i - T)}$  where  $D_i^2 = D_i^2(T(W), C(W)) =$

$$(\underline{x}_i - T(w))^\top C^{-1}(w) (\underline{x}_i - T(w))$$

$$= (\underline{x}_i - T)^\top C^{-1} (\underline{x}_i - T) = D_{\underline{x}_i}^2(T, C)$$

$$D_{\underline{x}}^2(T, C) = (\underbrace{(\underline{x} - T)^\top}_{\substack{\uparrow \\ \text{center of hyperellipsoid}}} C^{-1} (\underline{x} - T))$$

eigenvectors determine axes of hyperellipsoids

28) Note that  $D_T^2(\underline{x}, C) = D_{\underline{x}}^2(T, C)$ .

29) The pop squared Mahalanobis distance

$$D_{\underline{x}}^2(\underline{\mu}, \Sigma) = (\underline{x} - \underline{\mu})^\top \Sigma^{-1} (\underline{x} - \underline{\mu}) \text{ where}$$

$\underline{\mu}$  = pop location vector and  $\Sigma$  = pop dispersion matrix, often the pop cov. matrix,

30) The squared Euclidean distance of  $\underline{x}$  from  $T$

$$\text{is } (\underline{x} - T)^\top (\underline{x} - T) = D_{\underline{x}}^2(T, I_p) \text{ where } I_p = \text{diag}(1, \dots, 1)$$

is the  $p \times p$  identity matrix.

31) § 13.3 Outlier detection if  $p > n$ ; more than  $\frac{n}{2}$  cases are in the belt of the data.

a) Use Euclidean distances from the coordinatewise median  $D_i(\text{MED}(w), I_p)$ .

ex 4 measurements on 5 trees

SL 6.5

N	E	S	W
72	66	76	77
60	53	66	63
56	57	64	58
41	29	36	38
32	32	35	36

		Ordered	
32	41	56	60
29	32	53	57
35	36	64	66
36	38	58	63
		72	77

MED(W) = 58  
 MED(E) = 64  
 MED(S) = 66  
 MED(N) = 72

$\bar{z}$  261 237 277 272  
 ← show work →

$$\frac{1}{5} \begin{pmatrix} 261 \\ 237 \\ 277 \\ 272 \end{pmatrix} = \begin{pmatrix} 52.2 \\ 47.4 \\ 55.4 \\ 54.4 \end{pmatrix} = \bar{x}$$



Options

Rem lin trend

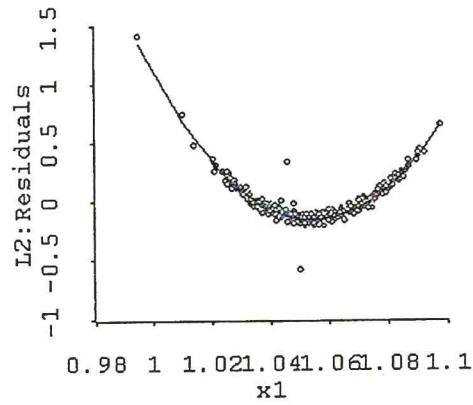
Zero line

Join points

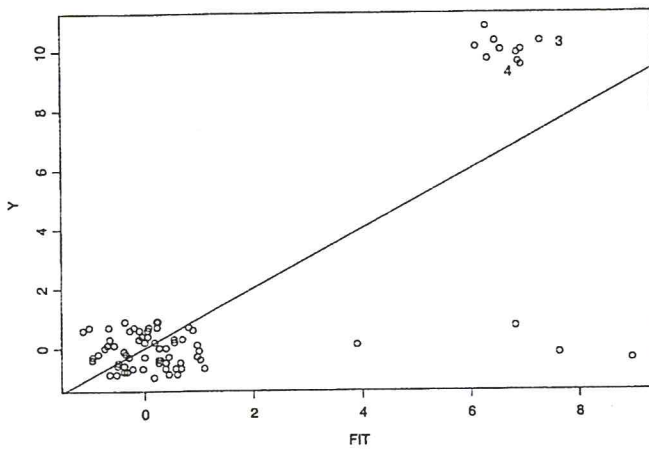
OLS  2

lowess  NIL

Case deletions



a)



b)



b) Let  $MED_j$  be the coordinatewise median SL 8 computed from cases  $\underline{x}_i$  with  $D_i^2 \leq$

$MED(D_i^2(MED_{j-1}, IP))$  where  $MED_0 = MED(w)$ .

Often use  $j=0$  or  $j=9$ . Let  $D_i = D_i(MED_j, IP)$ .

Let  $w_i = \begin{cases} 1 & \text{if } D_i < MED(D_1, \dots, D_n) + 5 \text{MAD}(D_1, \dots, D_n) \\ 0 & \text{else.} \end{cases}$

The convex set  $B$  consists of the  $m \geq \frac{n}{2}$  cases with weight  $w_i = 1$ . The convex estimator

$$(T, C) \text{ has } T = \frac{\sum_{i=1}^m w_i \underline{x}_i}{\sum_{i=1}^m w_i}, \quad C = \frac{\sum_{i=1}^m w_i (\underline{x}_i - T)(\underline{x}_i - T)^T}{\sum_{i=1}^m w_i - 1}$$

which is the sample mean and sample covariance matrix computed from the cases in set  $B$ .

32] Let  $\underline{w}_i = (y_i, \underline{x}_i)$  and let the continuous predictors from  $\underline{x}_i$  be  $\underline{v}_i$  (the predictors that take on many values, so not gender).

Apply the regression method to the  $m$  cases  $\underline{w}_i$  corresponding to the convex set  $B = \{i_1, \dots, i_m\}$  applied to  $\underline{u}_1, \dots, \underline{u}_n$ , MLR, GLMs, GAMs, LDA, QDA, k-NN etc

33) 019 } The function `ddplot5` plots

$$D_i: (\text{MED}(w), IP) \quad \text{VS} \quad D_i: \left( \frac{F_{\text{count}_2}}{\text{count}_2}, IP \right).$$

The plotted points tend to cluster about the identity line with outliers in the upper right corner of the plot with a gap between the bulk of the data and the outliers.



34) To detect outliers in one group and the bulk of the data in another group, the distance of the outliers from the bulk of the data increases roughly with  $\sqrt{p}$ .

### 05/4 Large Sample Theory Skim

35) Know 024 Central Limit Theorem (CLT):

Let  $Y_1, \dots, Y_n$  be iid with  $E Y_i = \mu$  and  $V(Y_i) = \sigma^2$ .

Then  $\sqrt{n} (\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2)$ .

36) If  $W_n \xrightarrow{D} X$ , then  $X$  is the asymptotic distribution or limiting distribution of  $W_n$ .  $X$

does not depend on  $n$ . The approximate

$\bar{Y} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$ , which does depend on  $n$ . SL 9

37) 027 Let  $\{Z_n, n=1,2,\dots\}$  be a sequence of RVs with CDFs  $F_n$  and let  $X$  be a RV with CDF  $F$ . Then  $Z_n$  converges in distribution to  $X$ , written  $Z_n \xrightarrow{D} X$ , if  $\lim_{n \rightarrow \infty} F_n(x) = F(x)$  at each continuity point of  $F$ .

38) If  $X_n \xrightarrow{D} X$ , then  $P(a < X_n \leq b) = F_n(b) - F_n(a) \rightarrow F(b) - F(a) = P(a < X \leq b)$  if  $F$  is continuous at  $a$  and  $b$ . So  $F$  can be used to approximate probabilities and percentiles.

See ex's 1.13 and 1.14.

39) 029  $X_n$  converges in probability to  $X$ ,

$X_n \xrightarrow{P} X$ , if  $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0 \quad \forall \varepsilon > 0$

$\Leftrightarrow \lim_{n \rightarrow \infty} P(|X_n - X| < \varepsilon) = 1 \quad \forall \varepsilon > 0.$

40) A sequence of estimators  $T_n$  is consistent for  $\gamma(\theta)$  if  $T_n \xrightarrow{P} \gamma(\theta) \quad \forall \theta \in \Theta$ . Then  $T_n$  is a consistent estimator of  $\gamma(\theta)$ ,

41) i) If  $\lim_{n \rightarrow \infty} \text{Var}(T_n) = 0$  and  $\lim_{n \rightarrow \infty} E_\theta(T_n) = \tau(\theta)$

both  $\forall \theta \in \Theta$ , then  $T_n$  is a consistent estimator of  $\tau(\theta)$ .

ii) Let  $0 < \delta \leq 1$ . If  $n^\delta (T_n - \tau(\theta)) \xrightarrow{D} N(0, \sigma_\theta)$ ,

$\forall \theta \in \Theta$ , then  $T_n$  is a consistent estimator of  $\tau(\theta)$ .

42)  $X_n \xrightarrow{P} \tau(\theta)$  iff  $X_n \xrightarrow{D} \tau(\theta)$ .

43) If  $n^\delta (W_n - \tau) \xrightarrow{D} X$  for some nondegenerate RV  $X$ , then  $W_n$  has rate  $n^\delta$ , and

$n^{1/2}$  rate = square root  $n$  consistency is good.

Multivariate Limit Th's

44) 038 Let  $\underline{X}_n$  have joint CDF  $F_n(\underline{x})$  and  $\underline{X}$  have joint CDF  $F(\underline{x})$ .

a)  $\underline{X}_n \xrightarrow{D} \underline{X}$  if  $F_n(\underline{x}) \rightarrow F(\underline{x})$  at all continuity points  $\underline{x}$  of  $F$ .

b)  $\underline{X}_n \xrightarrow{P} \underline{X}$  if  $\forall \varepsilon > 0, P(\|\underline{X}_n - \underline{X}\| > \varepsilon) \rightarrow 0$

45) know | Multivariate CLT (MCLT)

SL 10

If  $\underline{X}_1, \dots, \underline{X}_n$  are iid  $k \times 1$  random vectors with  $E(\underline{X}) = \underline{\mu}$  and  $\text{cov}(\underline{X}) = \underline{\Sigma}$ , then

$$\sqrt{n} (\underline{\bar{X}}_n - \underline{\mu}) \xrightarrow{D} N_k(\underline{0}, \underline{\Sigma}).$$

46) a) If estimator  $\underline{g}(\underline{T}_n) \xrightarrow{P} \underline{g}(\underline{\theta}) \quad \forall \underline{\theta} \in \Theta$ ,

then  $\underline{g}(\underline{T}_n)$  is a consistent estimator of  $\underline{g}(\underline{\theta})$ .

b) If  $n^b (\underline{g}(\underline{T}_n) - \underline{g}(\underline{\theta})) \xrightarrow{D} \underline{X}$ , then  $\underline{g}(\underline{T}_n) \xrightarrow{P} \underline{g}(\underline{\theta})$ .

c) If  $\underline{x}_n \xrightarrow{P} \underline{x}$  then  $\underline{x}_n \xrightarrow{D} \underline{x}$ .

d)  $\underline{x}_n \xrightarrow{P} \underline{g}(\underline{\theta})$  iff  $\underline{x}_n \xrightarrow{D} \underline{g}(\underline{\theta})$ .

47) 040 Continuous Mapping Th. Let  $g$  be a continuous function  $g: \mathbb{R}^k \rightarrow \mathbb{R}^d$ . If

$$\underline{x}_n \xrightarrow{D} \underline{x}, \text{ then } g(\underline{x}_n) \xrightarrow{D} g(\underline{x}).$$

48) <sup>041</sup> know | a) If  $\sqrt{n} (\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_p(\underline{0}, \underline{\Sigma})$  and

$A$  is a  $q \times p$  constant matrix, then

$$A \sqrt{n} (\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_q(A \underline{\theta}, A \underline{\Sigma} A^T)$$

b) Let  $\Sigma > 0$  be positive definite.

If  $\sqrt{n}(\bar{T}_n - \mu) \xrightarrow{D} N_p(0, \Sigma)$  and if  $C$  is a consistent estimator of  $\Sigma$ , then

$$\sqrt{n}(\bar{T}_n - \mu)^T C^{-1}(\bar{T}_n - \mu) \xrightarrow{D} \chi_p^2$$

c) If  $\Sigma > 0$ ,  $\bar{T}_n \xrightarrow{P} \mu$  and  $C \xrightarrow{P} \Sigma$  then

$$\begin{aligned} D_x^2(\bar{T}_n, C) &= (\bar{x} - \bar{T}_n)^T C^{-1}(\bar{x} - \bar{T}_n) \xrightarrow{D} D_x^2(\mu, \Sigma) \\ &= (\bar{x} - \mu)^T \Sigma^{-1}(\bar{x} - \mu). \end{aligned}$$

So  $D_x^2(\bar{T}_n, C) \xrightarrow{P} D_x^2(\mu, \Sigma)$  where  $P(D_x^2(\bar{x}, C) \leq D_x^2(\bar{x}, \Sigma)) = 1 - \alpha$   
 often  $\chi_{p, 1-\alpha}^2$

## Och 2 James et al §5.2 Prediction Regions and Bootstrap

047 Consider predicting a future test value  $Y_F$  given a  $p \times 1$  vector  $\underline{x}_F$  and training data  $(x_1, y_1), \dots, (x_n, y_n)$ . A large sample  $(100(1-\alpha))\%$  prediction interval PI for  $Y_F$  has the form

$[\hat{L}_n, \hat{U}_n]$  where  $P\{\hat{L}_n \leq \theta \leq \hat{U}_n\} \rightarrow 1-\delta$  as  $\frac{(\text{SL})}{n} \rightarrow \infty$ .

A large sample  $100(1-\delta)\%$  confidence interval for an unknown parameter  $\theta$  is  $[\hat{L}_n, \hat{U}_n]$  where  $P\{\hat{L}_n \leq \theta \leq \hat{U}_n\} \rightarrow 1-\delta$  as  $n \rightarrow \infty$ .

2) 05/55 Consider predicting a  $p \times 1$  future test vector  $\underline{x}$  or estimating an  $r \times 1$  unknown parameter vector  $\underline{\mu}$  given training data  $\underline{x}_1, \dots, \underline{x}_n$ .

A large sample  $100(1-\delta)\%$  prediction region is a set  $A_n$  such that  $P(\underline{x} \in A_n) \rightarrow 1-\delta$  as  $n \rightarrow \infty$ .

A large sample  $100(1-\delta)\%$  confidence region for  $\underline{\mu}$  is a set  $A_n$  such that  $P(\underline{\mu} \in A_n) \rightarrow 1-\delta$  as  $n \rightarrow \infty$ .

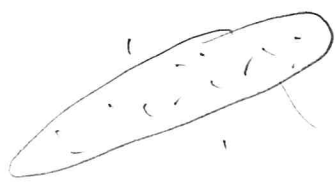
A CI and PI are special cases where  $p=r=1$ .

3) A large sample hypothesis test for  $H_0: \underline{\mu} = \underline{\mu}_0$  vs.  $H_A: \underline{\mu} \neq \underline{\mu}_0$  (with asymptotic level  $\delta$ ) rejects  $H_0$  if  $\underline{\mu}_0 \notin A_n$  and fails to reject  $H_0$  if  $\underline{\mu}_0 \in A_n$ .  
100(1-\delta)% confidence region for  $\underline{\mu}$ .

4) Interpretation Collect data, form  $A_{n,i}$  and see if  $\underline{w}_i \in A_{n,i}$  for  $i=1, \dots, k$  where the  $k$

trials are independent from the same pop  
 and  $\underline{w}_i = \underline{x}_i$  or  $\underline{w}_i = \underline{\theta}$ . Then the  
 number of times  $\underline{w}_i \in A_n$  - bin( $k, 1-\delta_n$ )  
 where  $\delta_n \rightarrow \delta$  as  $n \rightarrow \infty$ . So for a 95% region  
 and  $k=100$ , expect about 95 of the  $\underline{w}_i \in A_n$ .

5) The volume of a confidence region  $\rightarrow 0$  as  $n \rightarrow \infty$   
 while the volume of a prediction region goes  
 to that of a population region that would  
 contain a new  $\underline{x}_f$  with prob  $1-\delta$ .



95% pop prediction region

ex) Consider predicting the height of the next person  
 to come through the door.  $[4ft, 7ft] \approx 100\% PI$ .

$[5ft, 6ft]$  might be an 80% PI since the class  
 has some short ladies.

6) <sup>084</sup>  $100(1-\delta)\%$  coverage is the nominal coverage.

If the actual coverage  $100(1-\delta_n) > 100(1-\delta)$ ;

the region is conservative if  $100(1-\delta_n) < 100(1-\delta)$

is liberal if  $100(1-\delta_n) > 100(1-\delta)$



Being 5% conservative is "much better" than SLI2  
being 5% liberal. Want small volume  
with coverage near or higher than the nominal  
coverage.

7) <sup>060</sup> In simulations with 5000 runs, simulated  
coverage (prop times  $\underline{w} \in A_{n_i}$ )  $\in [0.94, 0.98]$   
suggests actual coverage is close to the  
nominal 95% coverage.

8) The bootstrap generates pseudodata  
 $T_1^*, \dots, T_B^*$  for a statistic  $T_n$  that  
estimates a parameter  $\mu$ , such that,  
under regularity conditions, applying certain  
large sample  $100(1-\delta)\%$  prediction regions  
to the bootstrap sample results in  
large sample  $100(1-\delta)\%$  confidence regions.

9) ~~known EI~~  
Let  $\underbrace{z_1, \dots, z_n}_{\substack{\text{training} \\ z_i}}, \underbrace{z_{n+1}}_{\text{test}}$  be iid. The Shortest CI  
interval is the shortest closed interval containing  
at least  $C$  of the  $z_i$ . For a small data  
set order the data and compute the lengths  
of intervals containing  $\geq C$  cases!

$$\{(\bar{z}_1, z_{c1}), (\bar{z}_2, z_{c+1}), \dots, (\bar{z}_{n-c+1}, z_n)\}$$

$$z_{c1} - z_{c1} \quad z_{c+1} - z_{c2} \quad z_n - z_{n-c+1}$$

Shortest  $c$  is  $[\bar{z}_c, z_{c+c-1}]$  is the interval with shortest length.

ex) Find shortest 4

0, 1, 3, 6, 9, 10, 11

$$\begin{array}{r} 6-0=6 \\ \hline 9-1=8 \\ \hline 10-3=7 \\ \hline 11-6=5 \end{array}$$

$$[6, 11] = \text{Shortest}(4)$$

Let  $\lceil x \rceil$  be the smallest integer  $\geq x$ . So  $\lceil 7.7 \rceil = 8$ ,  $\lceil 7 \rceil = 7$ .

10) If  $\frac{c_n}{n} \rightarrow 1-\delta$ , the Shortest  $(c_n)$  interval is a large sample  $100(1-\delta)\%$  PI,  $c_n \geq k_n = \lceil n(1-\delta) \rceil$  contains

$\approx 100(1-\delta)\%$  of the training data so has

$k_n < 1-\delta$  for test data. Frey (2013)

showed the max undercoverage  $\approx 1.12\sqrt{\frac{\delta}{n}}$

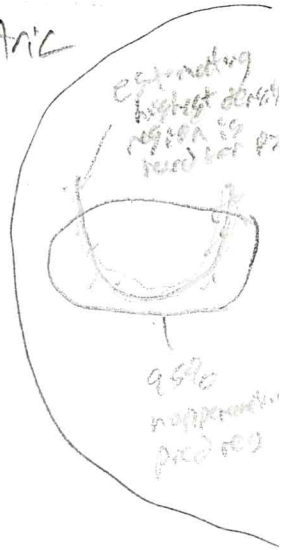
and used  $c_n = \min\left[n, \left\lceil n\left[1-\delta + 1.12\sqrt{\frac{\delta}{n}}\right]\right\rceil\right]$

11) The Shortest  $(c_n)$  PI estimates the  $100(1-\delta)\%$  highest density region if the  $x_i$  are iid from a unimodal bimodal  $1-\delta$  shortest

12} <sup>052</sup> Let  $\underbrace{x_1, \dots, x_n}_{\text{training}}$   $\stackrel{p \times 1}{x}$  be iid from SL 13

a distribution with nonsingular cov. matrix  $\Sigma$ .

A large sample 100(1- $\delta$ )% nonparametric prediction region for  $x$  is



$$A_n = \{x : D_x^2(\bar{x}, S) \leq D(c)\} =$$

$$\{x : D_x^2(\bar{x}, S) \leq D(c)\} \text{ where}$$

$$c = c_n \downarrow \lceil n(1-\delta) \rceil \text{ as } n \rightarrow \infty.$$

↙ increase training data coverage by 5% if  $n \leq 20p$

$$\text{Let } \delta_n = \begin{cases} \min(1-\delta + 0.05, 1-\delta + \frac{p}{n}) & \delta \geq 0.1 \\ \min(1-\frac{\delta}{2}, 1-\delta - \frac{10\delta p}{n}) & \delta \leq 0.1 \end{cases}$$

$$\text{↖ 95\% PR uses 97.5\% coverage of training data if } n \leq 20p$$

If  $1-\delta < 0.999$  and  $\delta_n < 1-\delta + 0.001$  set  $\delta_n = 1-\delta$ .

Then  $D(c)$   <sup>$= D(c_n)$</sup>  is the  $100\delta_n$ th quantile of  $D_1, \dots, D_n$

$$\text{eg use } c = \lceil n\delta_n \rceil. \quad (*)$$

13} 049 A large sample PI for  $Y_i$  when  $Y_i = m(x_i) + \epsilon_i$  finds the shortest CI of the residuals  $\{\hat{r}(s), \{\hat{s}(c)\}$ .

$$\text{Then } \hat{r}(s) = \hat{m}(x(s)) + b_n \{r(s), \hat{m}(x) + b_n \{r(s), \hat{m}(x)\}\} \text{ where}$$

$$b_n = \begin{cases} \left(1 + \frac{15}{n}\right) \sqrt{\frac{n+2d}{n-d}} & d < \frac{8n}{9} \\ 5 \left(1 + \frac{15}{n}\right) & d \geq \frac{8n}{9} \end{cases}$$

and  $d =$  crude estimator of model  $d^*$ .

Here  $c = \lceil n b_n \rceil$  and  $g_n$  replaces  $p$  by  $d$  in (\*).  
 This PI roughly uses the strength of the pseudodata  $\{c + r_i, i=1, \dots, n\}$ .  
 0.3.3 Bootstrapping James et al §5.2

14) Nonparametric bootstrap: Let data  $z_1, \dots, z_n$  be iid and  $T_n = T(z_1, \dots, z_n)$ . Draw a sample of size  $n$   $z_1^*, \dots, z_n^*$  with replacement from  $z_1, \dots, z_n$ , compute  $T_{n,1}^* = T(z_1^*, \dots, z_n^*)$  and repeat  $B$  times to generate  $T_{n,1}^*, \dots, T_{n,B}^*$ .

15) Know for  $E|$  For any  $B$  given  $B$  bootstrap samples compute  $T_1^*, \dots, T_B^*$ . Offer  $T_n = \bar{X}_n$  or  $MED_n$ .

ex) range =  $x_{(n)} - x_{(1)} = T$ .  $\{x_1, x_2\} = \{0, 4\}$   
 $T_n = 4 - 0 = 4$

bootstrap sample $x_i^*$	$T_i^*$
(0, 0)	0
(0, 4)	4
(4, 0)	4

Problem 02.4

(SL 14)

ex}  $X_1, \dots, X_n = 1, 2, 5, 10, 50$

$MED(n) = 5 = T_n$

bootstrap sample	ordered	
2, 10, 1, 2, 2	1, 2, 2, 2, 10	$T_1^* = 2$
50, 10, 50, 2, 2	2, 2, 10, 50, 50	$T_2^* = 10$
10, 50, 2, 1, 1	1, 1, 2, 10, 50	$T_3^* = 2$

ex] If  $(Y_i, X_i)$  are iid with  $Y_i = X_i^T \beta + \epsilon_i$

draw bootstrap sample, compute  $\hat{\beta}_j^*$   $j=1, \dots, B$ ,  
for the nonparametric bootstrap.

16) <sup>063</sup> An alternative to the above ex is the residual bootstrap. The  $i$ th residual

$$r_i = Y_i - \hat{Y}_i = Y_i - X_i^T \hat{\beta}$$

In matrix form  $\underline{Y} = \underline{X} \underline{\beta} + \underline{e}$ . Regress  $\underline{Y}$  on  $\underline{X}$

to obtain  $\hat{\beta}$   $\begin{matrix} \hat{\beta} \\ \sim \\ p \times 1 \end{matrix}$ ,  $\begin{matrix} \hat{Y} \\ \sim \\ n \times 1 \end{matrix}$ . Draw a sample

of size  $n$  with replacement from the residuals

$r_1^*, \dots, r_n^*$  to form vector  $\underline{Y}_j^*$  with  $i$ th  
element  $Y_{ij}^* = \hat{Y}_i + r_i^*$ . For  $j=1, \dots, B$

regress  $\underline{Y}_j^*$  on  $\underline{X}$  to form  $\hat{\beta}_1^*, \dots, \hat{\beta}_B^*$

The residual bootstrap works well with least squares

17)  $0 \leq \delta \leq 1$  It  $D_{\underline{\mu}}(T_n, \hat{\underline{\mu}}_T) \xrightarrow{D} D^2$  and  $D_{\underline{\mu}}(T_n, \hat{\underline{\mu}}_T) / B \xrightarrow{D} D^2$

and  $\hat{D}_{1-\delta}^2 \xrightarrow{P} D_{1-\delta}^2$  then  $R_C =$

$\{ \underline{w} : D_{\underline{w}}^2(T_n, \hat{\underline{\mu}}_T) \leq \hat{D}_{1-\delta}^2 \}$  is a large sample

100(1- $\delta$ )% cont. reg. for  $\underline{\mu}$ , and if  $\underline{\mu}$  is known

$R_p = \{ \underline{w} : D_{\underline{w}}(\underline{\mu}, \hat{\underline{\mu}}_T) \leq \hat{D}_{1-\delta}^2 \}$  is a large

sample 100(1- $\delta$ )% prediction region for a future value of the statistic  $T_n$ . Region  $R_C$  contains  $\underline{\mu}$

iff region  $R_p$  contains  $T_n$ .

Problem usually  $B \neq 1$  is not large

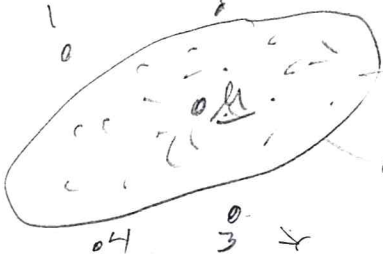
18) Idea: suppose there is an iid sample  $T_{n1}, \dots, T_{nB}$

of size  $B$  of the statistic. A large sample

100(1- $\delta$ )% prediction region for  $T_n$  is

$\{ \underline{w} : D_{\underline{w}}^2(\underline{\mu}, \hat{\underline{\mu}}_T) \leq D_C^2 \}$  if  $\underline{\mu}$  is known.

$$D_{T_i}^2(\underline{\mu}, \hat{\underline{\mu}}_T) = D_{\underline{\mu}}^2(T_i, \hat{\underline{\mu}}_T)$$

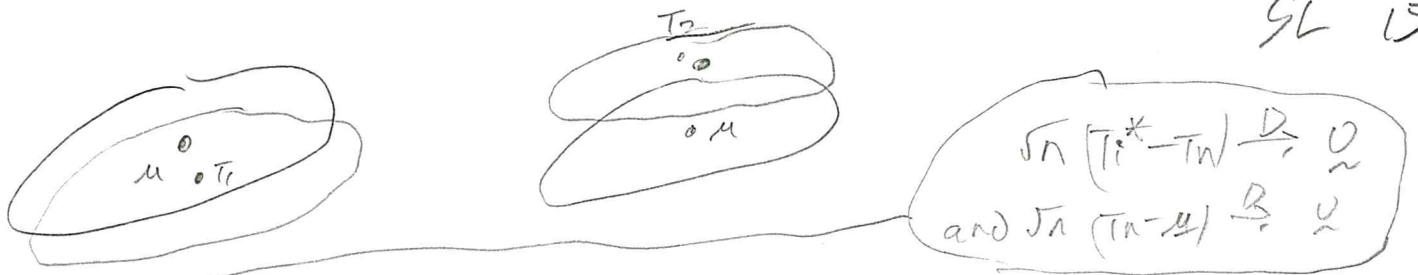


95 values of  $T_i$  in the pred region

Contains  $\approx 95\%$  of the  $T_i$  if  $n$  and  $B$  are large

95 values are in 5 are out

so  $\underline{\mu} \in \{ \underline{w} : D_{\underline{w}}^2(T_i, \hat{\underline{\mu}}_T) \leq D_C^2 \}$  iff  $T_i \in \{ \underline{w} : D_{\underline{w}}^2(\underline{\mu}, \hat{\underline{\mu}}_T) \leq D_C^2 \}$



19] The bootstrap sample basically takes the cloud of  $T_i$  centered at  $\mu$  and shifts the cloud to be centered at  $T_n$ .

20] Need  $B \geq 500$  if  $\mu$  is  $p \times 1$  and need  $n$  large, sometimes extremely large. If  $n$  is not large enough, undercoverage ( $<$  nominal  $100(1-\alpha)\%$  coverage) often occurs.

21] <sup>obj</sup> Prediction region method to test

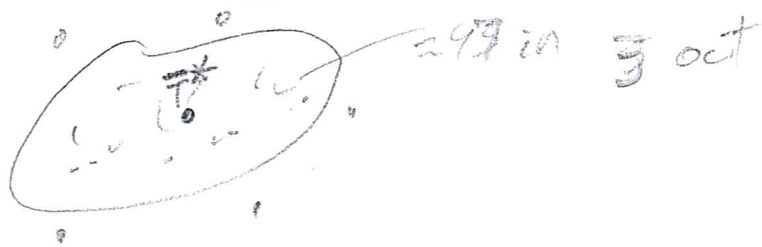
$$H_0: \underline{\mu} = \underline{c} \quad \text{vs} \quad H_A: \underline{\mu} \neq \underline{c}$$

$\uparrow$   
 $r \times 1$

Make a bootstrap sample  $\underline{w}_i = \underline{\hat{\mu}}_i^* - \underline{c}$

for  $i = 1, \dots, B$ . Apply the nonparametric prediction region to get a large sample  $100(1-\alpha)\%$  confidence region for  $\underline{\mu} - \underline{c}$ . Reject  $H_0$  if  $\underline{0}$  is not in the region. Fail to reject  $H_0$  if  $\underline{0}$  is in the region. If  $r = 1$  apply the short  $n(n)$  PI to get a large sample  $100(1-\alpha)\%$  CI. For  $n = 1000$  the 95% conf. region covers  $\sim 0.5\%$ .

of the  $w_i$  and  $\approx 5\%$  are not covered



22) 068 The large sample 100(1- $\alpha$ )% confidence region

for  $\mu$  is  $\left\{ \underline{w} : D^2(\bar{T}^*, S_T^*) \leq D_{(1-\alpha)}^2 \right\} = R$  where

$D_{(1-\alpha)}^2$  is computed from  $D_i^2 = (\underline{T}_i^* - \bar{T}^*)^T [S_T^*]^{-1} (\underline{T}_i^* - \bar{T}^*)$

for  $i=1, \dots, B$  and  $D_{(1-\alpha)}^2$  is the 100(1- $\alpha$ )th sample

quantile of the  $D_i^2$ . Reject  $H_0$  if  $\underline{w} \notin R$ .

23) Sufficient conditions:

$\underline{x}_1, \dots, \underline{x}_n$  iid  $T_n = t(\underline{x}_1, \dots, \underline{x}_n)$  for nonparametric bootstrap

$\underline{T}_n$  from a good estimator  $\hat{\underline{\beta}}$  for the residual bootstrap where  $\hat{\underline{\mu}} = A \hat{\underline{\beta}} = \underline{T}_n$ ,  $\underline{\mu} = A \underline{\beta}$ .

Solar only proved it:  $\sqrt{n}(\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_r(0, \Sigma)$

$\Sigma \succ 0$ . Conjectured to work for  $\sqrt{n}(\underline{T}_n - \underline{\mu}) \xrightarrow{D} \underline{X}$

for many other distributions for  $\underline{X}$ .

Sto. mod 2.3 We will come back to O § 2.3.3

later.

and optimal material



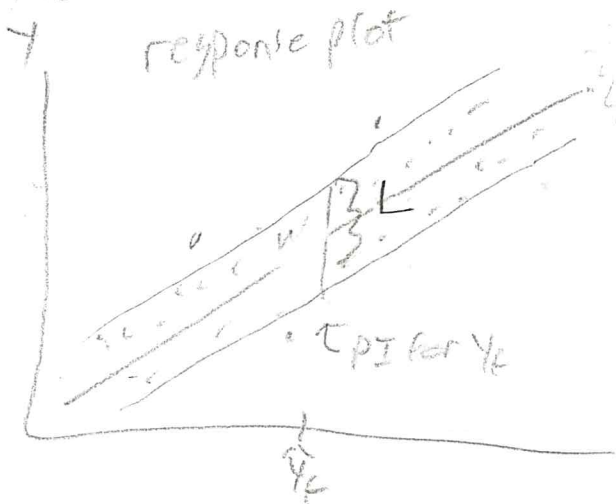
Multiple Linear Regression MLR

Och 3 5 ch 3, 6 part of 5

1} Also see Olive (2010, ch 2) with URL given on O p. 268 (M484 material)

2} 049 <sup>pointwise</sup> Prediction interval bands for  $\frac{n}{d}$  large

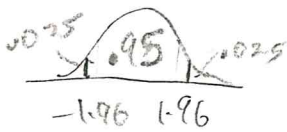
$\approx$  2 lines parallel to the identity line covering 100(1- $\alpha$ )% of the data with smallest vertical length  $L$



$\approx 9.7\sigma$   
 3  $\sigma$  rule for 95% PI  
 for  $\alpha = 100$   
 a bit more than 95% of the training data are within the bands

$\hat{y} = \underline{x}^T \underline{\beta}$

$e_i \sim N(0, 1) \Rightarrow L = 2(1.96) = 3.92$  as  $\frac{n}{d} \rightarrow \infty$



text uses  $\Sigma_i$

3} <sup>p. 71</sup> MLR  $y_i = \underline{x}_i^T \underline{\beta} + e_i$

$y_1 = x_{11} \beta_1 + \dots + x_{1p} \beta_p + e_1$

$\vdots$

$y_n = x_{n1} \beta_1 + \dots + x_{np} \beta_p + e_n$

or  $\underline{y} = \underline{X} \underline{\beta} + \underline{e}$   
 $n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1$

4) <sup>p69</sup> The ordinary least squares OLS full model  
 dummy variable  $\beta$  from  $y = X\beta + e$

$\hat{\beta}_{OLS}$  minimizes  $Q_{OLS}(\beta) = \sum_{i=1}^n r_i^2(\beta) = RSS(\beta)$

$= (y - X\beta)' (y - X\beta)$  estimating equations

$\hat{\beta}_{OLS}$  estimates  $\beta$  from  $y = X\beta + e$

Good criterion if  $n \geq 10p$ .

Analogy Maximum likelihood estimator

$y_i$  iid with pdf  $f(y_i|\theta)$   
 $\theta$  parameter

$L(\theta) = \prod_{i=1}^n f(y_i|\theta)$

dummy variable

but  $\hat{\theta} = \arg \max L(\theta)$  estimates parameter  $\theta$

5) OLS LS CLT:  $y_i = x_i' \beta + e_i$ , the zero

mean errors  $e_i$  are iid with  $V(e_i) = \sigma^2$ .

Assume  $p$  is fixed and  $n \rightarrow \infty$ ,  $\max(h_1, \dots, h_n) \rightarrow 0$

$h_i = h_{ii}$  where  $H = X(X'X)^{-1}X' = (h_{ij})$ , and

$\frac{X'X}{n} \rightarrow V^{-1}$  as  $n \rightarrow \infty$  Then

$\sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{D} N_p(0, \sigma^2 V)$  and

$\sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{D} N_p(0, \sigma^2 I_p)$

$\frac{\sum X^2}{n}$  estimates  $V^{-1}$  so

SL 17

$$\left(\frac{\sum X^2}{n}\right)^{-1} = n(\sum X^2)^{-1} \text{ estimates } V$$

$$\hat{\beta}_{OLS} \sim AN_p(\beta, \text{MSE}(\sum X^2)^{-1})$$

where  $\text{MSE} = \frac{1}{n-p} \sum_{i=1}^n r_i^2$

$$\left( \begin{aligned} \sqrt{n}(\hat{\beta}_{OLS} - \beta) &\sim AN_p\left[\begin{matrix} 0 \\ \sigma^2 n(\sum X^2)^{-1} \end{matrix}\right] \text{ multiply by } \frac{1}{\sqrt{n}} \\ \hat{\beta}_{OLS} - \beta &\sim AN_p\left[\begin{matrix} 0 \\ (\sum X^2)^{-1} \end{matrix}\right] \text{ so} \end{aligned} \right)$$

6) Let the nontrivial predictors  $\underline{u}_i = (x_{i2}, \dots, x_{ip})^T$ .

Let the  $n \times (p-1)$  matrix of standardized nontrivial predictors be  $\underline{W} = (w_{ij})$  where

$$\sum_{i=1}^n w_{ij} = 0 \quad \text{and} \quad \sum_{i=1}^n w_{ij}^2 = n$$

$\underbrace{\sum_{i=1}^n w_{ij}}_{\text{jth standardized predictor has sample mean } = 0 = \bar{w}_j}$

$$\text{so } \sigma_j^2 = \frac{1}{n} \sum_{i=1}^n w_{ij}^2 = 1$$

$\uparrow$   
biased sample variance.

Then the sample correlation matrix of the nontrivial predictors is  $R_{\underline{u}} = \frac{\underline{W}^T \underline{W}}{n}$ .

Many MLR methods fit

$$\underline{z} = \underline{W} \underline{\eta} + \underline{\epsilon} \quad \text{and then find } \underline{\beta} \text{ such}$$

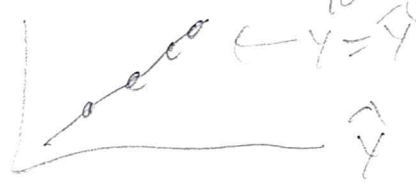
$$\text{use } \underline{\hat{y}} = \underline{\hat{z}} + \bar{y}$$

7} Problems i) when  $n > p$   $X$  is not invertible

$$\text{but if } n = p \quad \underline{\hat{y}} = H \underline{y} = X (X^T X)^{-1} X^T \underline{y} = I_n \underline{y} = \underline{y}$$

regardless of how bad the predictors are.

"Overfitting"  
"fitting noise"



PI for  $y_i$  is  
 $[\hat{y}_i - \hat{\sigma}_i, \hat{y}_i + \hat{\sigma}_i]$  95% coverage.

ii) If  $n < p$   $\underline{\hat{y}} = \underline{y}$  or program fails,

iii) Need  $n \geq 5p$  with  $J \geq 5$  and preferably

$J \geq 10$ . If  $n < 5p$  the model is overfitting

(not enough data to estimate  $\hat{\beta}$  well).

variable	coef.	SE	short 95% CI for $\beta_i$
$x_1 = \text{intercept} = \text{constant}$	$\hat{\beta}_1$	$SE(\hat{\beta}_1)$	$[\hat{L}_1, \hat{U}_1]$
$x_2$	$\hat{\beta}_2$	$SE(\hat{\beta}_2)$	$[\hat{L}_2, \hat{U}_2]$
$\vdots$			
$x_p$	$\hat{\beta}_p$	$SE(\hat{\beta}_p)$	$[\hat{L}_p, \hat{U}_p]$

→ ... all large sample 95% CI

For  $\beta_i$  is  $\hat{\beta}_i \pm 1.96 SE(\hat{\beta}_i)$ .

SL 18

Consider  $H_0: \beta_i = 0$   $H_A: \beta_i \neq 0$

If  $0 \in CI$  for  $\beta_i$ , then fail to reject  $H_0$  and conclude  $X_i$  is not needed in the MLR model given the other predictors are in the model.

If  $0 \notin CI$  for  $\beta_i$ , <sup>reject  $H_0$  and</sup> conclude  $X_i$  is needed in the MLR model. (see HW 4 B)

0  $\phi$  2,3,5,3,2 J  $\phi$  3,1 Forward Selection

9] Variable selection is a search for a subset of predictors that can be deleted without important loss of information if  $n \geq 10p$  so that the model  $I$  is good for prediction  $\uparrow$  if  $n < 5p$ .  
 $\uparrow$  remaining predictors

10]  $I_{min}$  is the model, among the  $p$  candidates, that  $\hat{MSE}$  minimizes  $C_p$  if  $n \geq 10p$ . EBIC if  $n < 10p$ .

$I_j$ Model	$X_2$	$X_3$	$X_4$	$X_5$	$\hat{\beta}_j$
$I_1$		*			$(\hat{\beta}_1, 0, \hat{\beta}_3, 0, 0)^T$
$I_2$		*	*		$(\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4, 0)^T$
$I_3$		*	*	*	$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, 0)^T$
$I_4$	*	*	*	*	$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, \hat{\beta}_5)^T = \hat{\beta}_{OLS}$
$I_5$	*	*	*	*	$\hat{\beta}_{OLS}$
out of $C_p$	$C_p(I_2)$	$C_p(I_3)$	$C_p(I_4)$	$C_p(I_5)$	$\uparrow$ only if model

The models also contain  $X_1 = \text{constant} = 1$

11) For  $n \geq 100$  find the OLS full model residuals for the residual bootstrap. Then bootstrap (see HW 4 A).

$p \times 1$	explanatory covar	missing if $\hat{\beta}_{I_{min}, 0} = 0$	95% CI for $\beta_j$
$X_1 = \text{int} = \text{const}$	$\hat{\beta}_1$	$SE(\hat{\beta}_1)$	$[\hat{L}_1, \hat{U}_1]$
$X_2$	0		$[\hat{L}_2, \hat{U}_2]$
$X_3$	0		$[\hat{L}_3, \hat{U}_3]$
$X_4$	$\hat{\beta}_4$	$SE(\hat{\beta}_4)$	$[\hat{L}_4, \hat{U}_4]$
$X_5$	$\hat{\beta}_5$	$SE(\hat{\beta}_5)$	$[\hat{L}_5, \hat{U}_5]$

OLS  $T_n = \hat{\beta}_{I_{min}, 0}$  OLS SE not accurate unless for  $\hat{\beta}_{I_{min}, 0}$

if  $I_{min}$  is chosen before seeing the data

or if  $\frac{X^T X}{n} \rightarrow V^{-1} = \text{diag}(d_1, \dots, d_p)$ , all  $d_i > 0$

Suppose the nontrivial predictors are i.i.d. 0 mean finite variance, or the predictors are orthogonal. Then

$$W = (w_{ij}) \quad \text{with } w_{ij} = \frac{\underline{v}_i^T \underline{v}_j}{n} = \frac{\sum_{k=1}^n x_{ki} x_{kj}}{n} \stackrel{P.L.}{\rightarrow} E(x_i x_j)$$

where  $\underline{V} = [\underline{v}_1 \dots \underline{v}_p]$ . If  $\underline{v}_i \perp \underline{v}_j$  then  $\underline{v}_i^T \underline{v}_j = 0$  for  $i \neq j$

If  $x_i \perp x_j$  for  $i, j > 1$ , i.i.d. then  $w_{ij} \rightarrow E(x_i x_j) = 0$  (0)

if  $x_i \in \{1\}$  then  $w_{ij} \rightarrow \bar{x}_j \cdot P_j E(x_j) = 0$ . So ii) holds.

Note that if  $\hat{\beta}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_4, \hat{\beta}_5)^T$ , then

$\hat{\beta}_{I_{min}, 0} = (\hat{\beta}_1, 0, 0, \hat{\beta}_4, \hat{\beta}_5)^T$ : pad  $\hat{\beta}_{I_{min}}$  with 0's to and  $\hat{\beta}_1$

12) <sup>074</sup> A model for variable selection is

SL 19

$$Y = \underline{X}^T \underline{B} + e = \underline{X}_S^T \underline{B}_S + \underline{X}_E^T \underline{B}_E + e = \underline{X}_S^T \underline{B}_S + e \quad (*)$$

where  $\underline{X} = (\underline{X}_S^T, \underline{X}_E^T)^T$  is  $p \times 1$ ,  $\underline{X}_S$  is  $q_S \times 1$

and  $\underline{X}_E$  is  $(p - q_S) \times 1$ , Given  $\underline{X}_S$  is in the model,

$\underline{B}_E = \underline{0}$  and  $E$  denotes the subset of terms that can be eliminated given subset  $S$  is in the model.

13)  $S$  is unknown. Let  $\underline{X}_I$  be the vector of  $k$  terms from a candidate subset  $I$  and let  $\underline{X}_0$  be the vector of remaining predictors (not of the candidate subset).

$$\text{Then } Y = \underline{X}_I^T \underline{B}_I + \underline{X}_0^T \underline{B}_0 + e, \quad \forall S \subseteq I$$

and (\*) holds, then

$$\underline{X}_I^T \underline{B} = \underline{X}_S^T \underline{B}_S = \underline{X}_S^T \underline{B}_S + \underline{X}_{I \setminus S}^T \underline{B}_{I \setminus S} + \underline{X}_0^T \underline{B}_0 = \underline{X}_I^T \underline{B}_I$$

where  $\underline{X}_{I \setminus S}$  denotes predictors in  $I$  that are not in  $S$ .

$$14) \quad Y = \underline{X}^T \underline{B} + e \quad = \text{full model}$$

$$Y = \underline{X}_I^T \underline{B}_I + e \quad = \text{submodel}$$

... is a submodel.

15) Underfitting occurs if  $S \not\subseteq I$  so  $X_I$  is missing important predictors, underfitting occurs if  $\frac{X_I}{\hat{f}(X_I)}$  with  $d = k < n$ . Overfitting occurs

if  $n < 5k$  or if  $S \subseteq I$  but  $S \neq I$ .

not enough data to estimate the  $k$  parameters well,

Overfitting is serious if  $n < 5k$  but not

"much of a problem" if  $n > 10p$  or  $n > 20p$ .

Underfitting is a serious problem.

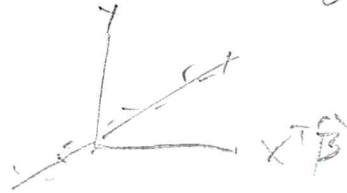
$$Y = X_I^T \beta_I + e_I$$

$\text{var}(e_I) > \text{var}(e) = \sigma^2$   
 may not be constant  
 could depend on case  $i$ .

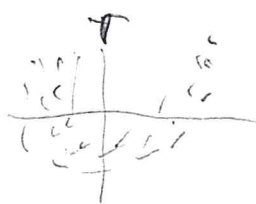
and the model may no longer be linear. (check with response and residual plots.)

ex]  $Y = \beta_1 + \beta_2 X + \beta_3 X^2 + e$

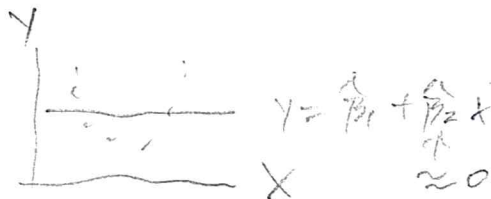
$SP = X^T \beta$  is a hyperplane that is a quadratic in  $X$



If  $SP(I) = \beta_1 + \beta_2 X = \text{linear in } X$



$$X^T \hat{\beta}_I$$





15) Forward Selection is a shrinkage method:

sp models are produced and except for the full model, some  $|\hat{\beta}_i|$  are shrunk to 0. Lasso and elastic net are also shrinkage methods. Ridge regression is a shrinkage method but  $|\hat{\beta}_i|$  is not shrunk to 0. Shrinkage methods that shrink  $\hat{\beta}_i$  to 0 are also variable selection methods.

$B_0$  is  $q \times 1$

17) The POP MLR model is sparse if  $a_0$  is small. The pop model is dense if  $\frac{n}{a_0} < J$  where  $J=5$  or  $10$ , say.

The fitted model  $\hat{\beta} = \hat{\beta}_{\tilde{I}_{min,d}}$  is sparse

if  $d = \#$  nonzero coefficients is small.

The fitted model is dense if  $\frac{n}{d} < J$  where  $J=5$  or  $10$ .

18)	<sup>097</sup> $I_j$		# models fitted
	$I_1$	$x_1^*$	0 (or 1; often do not fit this model)
	$I_2$	$x_1^*, x_2^*$	$p-1$
	$I_3$	$x_1^*, x_2^*, x_3^*$	$p-2$
	$\vdots$		
	$I_p$	$x_1^*, \dots, x_p^*$	1
			$p(p-1)$

Forward selection forms a sequence of submodels

$$I_1, \dots, I_m \text{ eg } m = \min(p, \lfloor \frac{n}{5} \rfloor)$$

$I_1$  uses  $x_i^* = x_1$  a constant but no nontrivial predictors.

To form  $I_2$  consider all models  $I$  with 2 predictors including  $x_1^*$ . Compute

$$Q_2(I) = SSE(I) = RSS(I) = \mathbf{r}^T(I) \mathbf{r}(I) = \sum_{i=1}^n r_i^2(I) \\ = \sum_{i=1}^n (y_i - \hat{y}_i(I))^2 \quad (\text{SSE stands for sum of squared errors; RSS is more accurate.})$$

Let  $I_2$  minimize  $Q_2$  and use  $x_1^*, x_2^*$ .

In general to form  $I_j$ , consider all models  $I$  with  $j$  predictors including  $x_1^*, \dots, x_{j-1}^*$ .

Let  $I_j$  minimize  $Q_j(I) = RSS(I)$  and use predictors  $x_1^*, \dots, x_j^*$ .

$$[9] \text{ If } m = \lfloor \frac{n}{5} \rfloor \approx \frac{n}{5} \left( \frac{2p - \frac{n}{5}}{2} \right) \approx \frac{n(2p - m)}{2} \text{ models are}$$

test. There are fast updating formulas for OLS (for adding 1 predictor and branch and bound algorithms shorten time) but forward selection can be done at least  $n \dots$

10) <sup>098</sup> Need to choose the final model  $S_L \subseteq I$  from the sequence of  $M$  models  $I_1, \dots, I_M$ .

Let  $X_I$  and  $\hat{\beta}_I$  be  $a \times 1$ . For a given data set,  $p, n$ , and  $\hat{\sigma}^2$  act as constants.

A criterion below may add a constant or be divided by a <sup>positive</sup> constant without changing the subset  $I_{\min}$  that minimizes the criterion.

Let criteria  $C_S(I)$  have the form

$$C_S(I) = SSE(I) + a k_n \hat{\sigma}^2$$

$C_P(I) = AIC_S(I)$  uses  $k_n = 2$  while  $BIC_S(I)$

uses  $k_n = \log(n)$ , if  $n \geq 5p$ ,  $J \geq 5$

preferably  $J \geq 10$ ,  $\hat{\sigma}^2 = MSE = \frac{1}{n} \sum_{i=1}^n \frac{r_{i,OLS}^2}{n-p}$

The following criterion still need  $\frac{n}{p}$  large

$$AIC(I) = n \log \left( \frac{SSE(I)}{n} \right) + 2a$$

$$BIC(I) = n \log \left( \frac{SSE(I)}{n} \right) + a \log(n),$$

The EBIC criterion may work when  $\frac{n}{p}$  is not large.

$$EBIC(I) = BIC(I) + 2 \log \left[ \binom{p}{a} \right]$$

$$\binom{p}{a} = \frac{p!}{(p-a)! a!} \quad \text{so } \log \binom{p}{a} = \log(p!) - \log[(p-a)!] - \log[a!]$$

2.11} Fix  $p$ , and let  $I_{min}$  minimize  $C_p$  or AIC.

Let  $m = p$ , The probability that  $I_{min}$  underfits  $\rightarrow 0$  as  $n \rightarrow \infty$ .

If  $S \subseteq I$  then  $\hat{\beta}_{I,0}$  is a  $\sqrt{n}$  consistent estimator of  $\beta$  if  $y = X\beta + e = \underline{X}_S^T \beta_S + e$ .

Since there are at most  $2^p$  regression models  $I$  that contain  $S$  and the prob that  $I_{min}$  picks one of these models goes

to 0.  $\hat{\beta}_{I_{min},0}$  is a  $\sqrt{n}$  consistent estimator of  $\beta$ .

Hence the large sample 100(1- $\alpha$ )% PI (2.7) works for OLS (forward selection if  $p$  is fixed)

$\hat{\beta}_{I_{min}}$  is  $d \times 1$  and  $\hat{m}(x) = x^T \hat{\beta}_{I_{min},0} = x^T \hat{\beta}_{I_{min},0} = \hat{m}(x_{I_{min}})$

2.12}  $\hat{\beta}_{I_{min},0}$  is not asymptotically normal, so residual bootstrap inference has not yet been proven to work, but in simulations, bootstrap inference was superior to that of the OLS full model if  $a_S < p$   $\frac{n}{p}$  is large and  $B \geq 50p$  is large.

Exception  $\hat{\beta}_{I_{min},S}$  is asymptotically normal if

$\frac{X^T X}{n} \rightarrow \text{diag}(d_1, \dots, d_p)$  with all  $d_i > 0$ .

ex 0.91-83) <sup>133</sup>  $n=100, p=4, \beta=(1,1,0,0)^T$  SL 22

so  $\underline{\beta}_2 = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \underline{x} = (1 \ 0^T)^T$ . Let  $m=A$ .

$\underline{w}_i \sim N_{p-1}(0, I), \underline{u}_i = A \underline{w}_i \sim N_{p-1}(0, A^2)$

where  $A=A^T=(a_{ij}), a_{ij} = \begin{cases} \psi & i \neq j \\ 1 & i = j \end{cases}, 0 \leq \psi < 1$

so  $\Sigma_0 = (\sigma_{ij})$  with  $\sigma_{ij} = \begin{cases} 1 + (m-1)\psi^2 & i=j \\ 2\psi + (m-2)\psi^2 & i \neq j \end{cases}$

so  $\rho = \text{cor}(x_i, x_j) = \frac{2\psi + (m-2)\psi^2}{1 + (m-1)\psi^2}$  for  $i \neq j, i, j > 1$ .

If  $\psi=0, \rho=0$ , if  $\psi = \frac{1}{\sqrt{cp}}$   $\rho \rightarrow \frac{1}{c+1}$  as  $p \rightarrow \infty$ .

If  $\psi$  is close to 1 or  $\psi \rightarrow 0$  fixed and  $p$  large, then  $\rho$  gets close to 1 and the non-trivial predictor vectors cluster about a line in the direction of  $\underline{1}$ .

The prediction region method was used to test  $H_0(\beta_2) = (0)$ ,  $e_i \sim N(0,1) \rightarrow \text{pop cor len} = 0.392$  if  $\psi=0$

model	$\psi$	corr/len	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	test	len of $(0, D_{c-1})$
reg	0	corr	.9496	.9474	.9496	.9474	.9442	$\sqrt{\lambda^2} = .95 \approx 2.4$
		len	.3961	.3997	.3988	.3992	2.4503	
vg	0	corr	.9472	.9470	.9980	.9980	.9936	different digits so these are not comparable not $\lambda^2$
		len	.3964	.3991	.3246	.3233	2.6936	
reg	0.9	corr	.9432	.9512	.9500	.9498	.9442	not $\lambda^2$
		len	.3963	3.2621	3.2613	3.2611	2.4505	

Testing  $\beta_0 = 0$  has smaller volume (length) for bootstraps because  $\beta_i^* = 0$  often for  $\beta_i$  a component of  $\beta_0$ .

23 } 083 To my knowledge, this is the 1st simulation where bootstrap inference is as good or better than OLS full model inference if  $\frac{n}{p}$  is large. When  $\frac{n}{p}$  is small OLS full model inference is bad and should not be used (eg interpolates the data if  $n=p$ ).

24 } <sup>0B2</sup> Suppose  $y = X^T \beta + e$ ,  $n \geq 10p$ ,  $B \geq 50p$  and  $\beta_{2i}^* = \beta_{3i}^* = \beta_{6i}^* = 0$  for  $i=1, \dots, B$

$$\begin{matrix} \beta_1^* & \beta_2^* & \beta_3^* & \beta_4^* & \beta_5^* & \beta_6^* \\ & 0 & 0 & & & 0 \\ & \vdots & \vdots & & & \vdots \\ & 0 & 0 & & & 0 \end{matrix}$$

Consider testing  $H_0: \beta_0 = \begin{pmatrix} \beta_2 \\ \beta_3 \\ \beta_6 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$  vs  $H_A: \beta_0 \neq 0$

Conjecture: fail to reject  $H_0$  and method can be used after looking at the data. (after the bootstrap).

Idea: prediction region method covers  $\geq 100(1-\alpha)\%$  of  $\beta_{0i}^*$  so prediction region confidence region is  $\sum \alpha_i$  for

25) 0133 Another variable selection model SL 23

$$i \in \underline{X^T B} = \underline{X_{S_i}^T} \underline{\beta_{S_i}} \quad (\text{for } i=1, \dots, J)$$

So that there are  $J \geq 2$  non-nested "true" models where  $\underline{\beta_{S_i}}$  is  $a_{S_i} \times 1$ .

If  $x_1, x_2, x_3$  are in  $S_1$  and  $x_1, x_2, x_4$  in  $S_2$  then  $x_3$  and  $x_4$  are (exactly) both often selected and omitted by forward selection in the  $B$  bootstrap samples. So omitting all predictors with  $\beta_{i0}^* = 0$  would result in underfitting (using model  $x_1, x_2$ ).

26) Criteria (like  $AIC_{(S)}$ ,  $BIC_{(S)}$  and EBIC

$$\text{attempt to minimize } \frac{1}{n} \sum_{i=1}^n E \left( Y_{i,\text{new}} - \hat{Y}_{i,\text{new}}(F) \right)^2 \\ = \frac{1}{n} E \left[ \| Y_{\text{new}} - \hat{Y}_{\text{new}}(F) \|^2 \right]$$

(based on training data)

where  $Y_{i,\text{new}}$  uses  $\underline{x}_i = x_i$  for  $i=1, \dots, n$ .

Using the PI length and coverage with cross validation may be better than EBIC when  $\frac{n}{p}$  is small. (covered later)

0.3.3 J 6.3.1, 6.7.1 PCR

27) 0102 Suppose  $A$   $p \times p$  is symmetric so  $A = A^T$

Then  $A$  has eigenvalue  $\lambda$  with eigenvector  $\underline{x} \neq \underline{0}$  if  $A\underline{x} = \lambda\underline{x}$ . Let  $\underline{e}$  be an eigenvector of  $A$  with unit length  $\|\underline{e}\|_2 =$

$\sqrt{\underline{e}^T \underline{e}} = 1$ . Then  $\underline{e}$  and  $-\underline{e}$  are eigenvectors with unit length, then

$A$  has  $p$  eigenvalue eigenvector pairs  $(\lambda_1, \underline{e}_1), \dots, (\lambda_p, \underline{e}_p)$  such that  $\underline{e}_i^T \underline{e}_j = 0$ .

Symmetric  $A \succ 0$  is positive definite if

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$  and  $A \succeq 0$  is

positive semidefinite if  $\lambda_p \geq 0$ .

28) PCR Principal components Regression

for  $\underline{y} = \underline{X}\underline{\beta} + \underline{e}$  uses  $\underline{z} = \frac{\underline{W}}{\sqrt{\lambda_i}} \underline{y} + \underline{e}$

and finds  $\hat{\underline{\beta}}$  and  $\hat{\underline{y}}$  using  $\hat{\underline{y}} = \underline{R}_0^{-1} \underline{z}$ .  $\underline{R}_0 = \frac{\underline{W}^T \underline{W}}{n}$

has eigenvalue eigenvector pairs  $(\hat{\lambda}_1, \hat{\underline{e}}_1), \dots, (\hat{\lambda}_k, \hat{\underline{e}}_k)$

where  $k = \min(n, p-1)$ . Let  $\underline{w}_1, \dots, \underline{w}_k$  be the standardized vectors of the non trivial predictors

with  $\underline{W} = \begin{pmatrix} \underline{w}_1^T \\ \vdots \end{pmatrix}$  Then the  $k$  principal



components corresponding to the  $j$ th case  $\underline{w}_j$

are  $P_{j1} = \hat{e}_1^T \underline{w}_j, \dots, P_{jk} = \hat{e}_k^T \underline{w}_j$

The  $j$ th principal component consists of  $(\hat{e}_1^T \underline{w}_j, \dots, \hat{e}_k^T \underline{w}_j)^T = \underline{P}_j$

PCR does OLS regression of  $\underline{z}$  on  $\underline{P}_1$ , then  $\underline{P}_2, \dots, \underline{P}_k$ , resulting in  $k$  models.

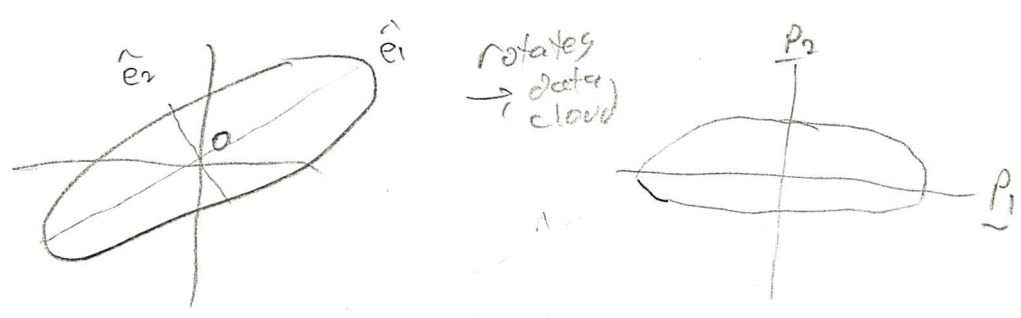
10 fold CV (0.3.11) is used to select the models.

29) Suppose  $n \geq 10p$  so  $k=p-1$ . The OLS full model = PCR<sub>k</sub>. Note that  $\bar{\underline{w}} = 0$  and the

hyperellipsoid  $\{\underline{w} \mid D_{\underline{w}}^2 @ R^p \leq h^2\} = \{\underline{w} \mid \underline{w}^T R_0^{-1} \underline{w} \leq h^2\}$

is centered at  $\underline{0}$  with axes given by the eigenvectors  $\hat{e}_i$  where the half length in the direction  $\hat{e}_i$  is  $h\sqrt{\lambda_i}$ .  $\underline{P}_1$  is obtained by projecting the  $\underline{w}_i$  on

the longest axis,  $\underline{P}_2$  on the next longest axis, ...,  $\underline{P}_k$  on the shortest axis.



30) Problems i) In general Bonferroni is an inconsistent

estimator of  $\beta$  unless  $P(p \rightarrow \infty) =$

$$P(\hat{\beta}_{PCR_{j_i}} = \hat{\beta}_{OLS}) \rightarrow 1.$$

ii) There is no reason why  $\beta_1, \beta_2, \dots, \beta_k$  should decrease in importance for predicting  $y$ .

0.3.4 J 6.3.2 6.7.2 PLS

31) Partial least squares uses PLS components  $\beta_1, \dots, \beta_m$ . Let  $V_i = X_i^T U$  and model

J: contain  $V_1, \dots, V_m$ . For  $i=1, \dots, m$ .

If  $m=p-1$ , PLS = OLS full model.

PLS uses  $Z = WU + e$  then gets  $\hat{\beta}$  and  $\hat{Y}$  from  $\hat{M}, \hat{Z}$  and  $\hat{Y}$ . Unlike PCR, PLS

uses  $Y$  in forming the PLS components: want components highly correlated with  $Y$ .

32) The 1st PLS component does OLS on  $Z$  and  $W_1$  without a constant to produce  $\hat{\beta}_{i1}$  for  $i=1, \dots, p-1$ . Then  $\underline{\gamma}_1 = (\gamma_{11}, \dots, \gamma_{(p-1)1})^T$ .

Important point  
PLS  
PLS

33)  $\hat{\beta}_{PLS}$  is not a consistent estimator of  $\beta$  unless  $\frac{p}{n} \rightarrow 0$ .

0.3.5 J 6.2.1, 6.6.1 RR

34}  $\underline{y} = \underline{X}\underline{\beta} + \underline{e}$  and

9L 25

$\underline{z} = \underline{W}\underline{\eta} + \underline{e}$  is used to fit ridge regression.

Then  $\hat{\underline{\eta}}$  and  $\bar{y}$  are used to find  $\hat{\underline{\beta}}$  and  $\hat{y}$ .

Let  $\lambda \geq 0$ . Then the ridge regression estimator  $\hat{\underline{\eta}}_R$  minimizes the ridge regression

criterion  $Q_R(\underline{\eta}) = \frac{1}{a} (\underline{z} - \underline{W}\underline{\eta})^T (\underline{z} - \underline{W}\underline{\eta}) + \frac{\lambda n}{a} \sum_{i=1}^{p-1} \eta_i^2$   
 $\underbrace{\hspace{10em}}_{\underline{\eta}^T \underline{\eta}}$

$= \frac{1}{a} \text{RSS}(\underline{\eta}) + \frac{\lambda n}{a} \|\underline{\eta}\|_2^2$  with  $a = 1, 2, n, 2n$

common. If  $\lambda n = 0$ ,  $\hat{\underline{\eta}}_{RR} = \hat{\underline{\eta}}_{OLS}$ . As

$\lambda n \rightarrow \infty$ ,  $\hat{\underline{\eta}}_{RR} \rightarrow \underline{0}$  and  $\hat{y} \rightarrow \bar{y}$ . (RR is a shrinkage method.)

$\hat{\underline{\eta}}_R = (\underline{W}^T \underline{W} + \lambda n \underline{I}_{p-1})^{-1} \underline{W}^T \underline{z}$ .

$\hat{\underline{\eta}}_{OLS} = (\underline{W}^T \underline{W})^{-1} \underline{W}^T \underline{z}$   $\left\{ \begin{array}{l} \text{OLS full model} \\ \text{if the inverse exists.} \end{array} \right.$

Warning! the literature typically uses  $\lambda = \frac{\lambda n}{a}$ .

35}  $\underline{W}^T \underline{W}$  is symmetric and square so  $\underline{W}^T \underline{W} \geq 0$  with

eigenvalues  $\psi_1 \geq \dots \geq \psi_p \geq 0$ . If  $\psi_p = 0$  then  $(\underline{W}^T \underline{W})^{-1}$  does not exist.

Let  $(\psi, \underline{g})$  be an eigenvalue eigenvector pair of  $W^T W = nR\underline{v}$ . Then  $(W^T W + \lambda_{in} I_{p-1}) \underline{g}$

$$= W^T W \underline{g} + \lambda_{in} \underline{g} = \psi \underline{g} + \lambda_{in} \underline{g} = (\psi + \lambda_{in}) \underline{g}$$

So  $(\underbrace{\psi + \lambda_{in}}_{> 0 \text{ if } \lambda_{in} > 0}, \underline{g})$  is an eigenvalue eigenvector pair of  $W^T W + \lambda_{in} I_{p-1}$  positive definite  
 $> 0$  if  $\lambda_{in} > 0$ .

Hence  $(W^T W + \lambda_{in} I_{p-1})^{-1}$  exists  $\forall \lambda_{in} > 0$ , even if  $W^T W$  is singular or ill conditioned.

36) Usually, a grid of  $M \approx 100 \lambda_{in}$  values

$$0 \leq \lambda_1 < \lambda_2 < \lambda_3 < \dots < \lambda_M \text{ is used where } \lambda_i = \lambda_{in} i.$$

to find CV is used to select  $\lambda_s = \hat{\lambda}_{in}$

for the final model. Model selected from  $M$  models, not variable selection.

37) 0108 Let augmented matrices

$$W_A = \begin{pmatrix} W \\ \sqrt{\lambda_{in}} I_{p-1} \end{pmatrix}, \quad \underline{z}_A = \begin{pmatrix} \underline{z} \\ \underline{0} \end{pmatrix} \text{ where } \underline{0} \text{ is } (p-1) \times 1.$$

For  $\lambda_n > 0$ , the OLS estimator from regressing  $\underline{z}_A$  on  $\underline{w}_A$  is

$$\hat{\underline{y}}_A = (\underline{w}_A^T \underline{w}_A)^{-1} \underline{w}_A^T \underline{z}_A = \hat{\underline{y}}_R. \text{ So we}$$

$$\underline{w}_A^T \underline{z}_A = \underline{w}^T \underline{z} \quad \text{and} \quad \underline{w}_A^T \underline{w}_A = \underline{w}^T \underline{w} + \lambda_n \mathbb{I}_{p-1}.$$

$$38) \text{ OIII) } \hat{\underline{y}}_R = \underbrace{\underline{w}^T (\underline{w} \underline{w}^T + \lambda_n \mathbb{I}_n)^{-1}}_{n \times n} \underline{z}$$

$$= \underbrace{(\underline{w}^T \underline{w} + \lambda_n \mathbb{I}_{p-1})^{-1}}_{(p-1) \times (p-1)} \underline{w}^T \underline{z}$$

$A^T$  has  $O(n^3)$  complexity so use 1st formula  $n \times n$   
if  $n < p-1$  and 2nd formula otherwise.

39) Suppose  $n > p$  and  $(\underline{w}^T \underline{w})^{-1}$  exists.

$$\text{Then } \hat{\underline{y}}_R = \underbrace{A_n}_{\substack{\text{An} \\ (\underline{w}^T \underline{w} + \lambda_n \mathbb{I}_{p-1})^{-1} \underline{w}^T \underline{w} (\underline{w}^T \underline{w})^{-1} \underline{w}^T \underline{z}}} \hat{\underline{y}}_{OLS}$$

So  $\hat{\underline{y}}_R = A_n \hat{\underline{y}}_{OLS}$ . HW 4 showed (problem 3.3)

$$A_n = B_n = \mathbb{I}_{p-1} - \lambda_n (\underline{w}^T \underline{w} + \lambda_n \mathbb{I}_{p-1})^{-1}$$

40) suppose  $P$  is fixed  $\frac{W^T W}{n} \xrightarrow{P} V^T (= \Sigma_0)$

and the OLS CLT holds so that

$$\sqrt{n} (\hat{\underline{\eta}}_{OLS} - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V)$$

If  $\frac{\lambda_n}{n} \rightarrow 0$  then  $\frac{W^T W + \lambda_n I_{p-1}}{n} \xrightarrow{P} V^T$

and  $n (W^T W + \lambda_n I_{p-1})^{-1} \xrightarrow{P} V$ . Also

$$A_n = \frac{n}{n} A_n = \left( \frac{W^T W + \lambda_n I_{p-1}}{n} \right)^{-1} \frac{W^T W}{n} \xrightarrow{P} V V^T = I_{p-1}$$

41) 0112) RD CLT: Assume  $p$  is fixed

and LS CLT holds for  $\underline{z} = W\underline{\eta} + \underline{e}$  as in 39.

a) If  $\frac{\lambda_n}{\sqrt{n}} \xrightarrow{P} 0$  then  $\sqrt{n} (\hat{\underline{\eta}}_R - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V)$

so OLS full model and RR are asymptotically equivalent

b) If  $\frac{\lambda_n}{\sqrt{n}} \xrightarrow{P} \gamma \geq 0$  then

$$\sqrt{n} (\hat{\underline{\eta}}_R - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underbrace{-\gamma V \underline{\eta}}_{\text{inferior to OLS unless this term} = 0}, \sigma^2 V)$$

Proof)  $\hat{\underline{\eta}}_0 = B_n \hat{\underline{\eta}}_{OLS} = I_{p-1} \underline{\eta} - \lambda_n (W^T W + \lambda_n I_{p-1})^{-1} \hat{\underline{\eta}}_{OLS}$

Hence  $\sqrt{n} (\hat{\underline{\mu}}_R - \underline{\mu}) = \sqrt{n} \left[ \hat{\underline{\mu}}_R - \hat{\underline{\mu}}_{OLS} + \hat{\underline{\mu}}_{OLS} - \underline{\mu} \right]$  SL 27

$$= \sqrt{n} (\hat{\underline{\mu}}_{OLS} - \underline{\mu}) - \sqrt{n} \frac{\hat{\lambda}_{in}}{\lambda} \underbrace{\left( \underbrace{\mathbf{W}^T \mathbf{W}}_{P \rightarrow V} + \underbrace{\hat{\lambda}_{in} \mathbf{I}_{p-1}}_{P \rightarrow V} \right)^{-1}}_{P \rightarrow V} \underbrace{\hat{\underline{\mu}}_{OLS}}_{P \rightarrow \underline{\mu}}$$

$\frac{n}{n} = 1$

$$\xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V) - \gamma V \underline{\mu} \sim N_{p-1}(\underbrace{-\gamma V \underline{\mu}}_{0 \text{ if } \gamma = 0}, \sigma^2 V).$$

□

Knight and Fu (2000) have a harder proof.

42)  $\hat{\underline{\mu}}_R$  is  $\sqrt{n}$  consistent if  $\frac{\hat{\lambda}_{in}}{\lambda} \xrightarrow{P} \gamma$

and  $\hat{\underline{\mu}}_R$  is a consistent estimator of  $\underline{\mu}$

if  $\frac{\hat{\lambda}_{in}}{\lambda} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ .

43) Problems i) RR large sample theory is worse or as good as that of OLS full model if  $\frac{n}{p}$  is large.

ii) 10 fold CV does not appear to guarantee

that  $\frac{\hat{\lambda}_{in}}{\lambda} \xrightarrow{P} 0$  or  $\frac{\hat{\lambda}_{in}}{\lambda} \xrightarrow{P} 0$ .

iii) RR tends to undertit if  $a_s$  is more than about 20 and the predictors are highly correlated

44) RR is a lot better than OLS full model

if a)  $X^T X$  is singular or ill conditioned or b)  $\frac{n}{p}$  is small,

0.63.6  $\rightarrow$  0.62.2 0.6.62 lasso

$$45) \underline{y} = \underline{X}\underline{\beta} + \underline{e} \quad \text{fit } \underline{z} = \underline{W}\underline{\eta} + \underline{e}$$

The lasso estimator  $\hat{\underline{\eta}}_L$  minimizes the criterion

$$Q_L(\underline{\eta}) = \frac{1}{n} (\underline{z} - \underline{W}\underline{\eta})^T (\underline{z} - \underline{W}\underline{\eta}) + \frac{\lambda_1 n}{a} \sum_{i=1}^{p-1} |\eta_i|$$

$$= \frac{1}{n} \text{RSS}_{\underline{W}}(\underline{\eta}) + \frac{\lambda_1 n}{a} \|\underline{\eta}\|_1, \quad \leftarrow L_1 \text{ norm}$$

a shrinkage method and a variable selection method some  $\hat{\eta}_i = 0$ . Like ridge regression, the lasso criterion is convex so fast algorithms to compute  $\hat{\underline{\eta}}_L$  exist.

46) A grid of  $M \approx 100$   $\lambda_1 n$  values is used

$$0 \leq \lambda_1 < \lambda_2 < \dots < \lambda_M \quad \text{and } \lambda_i = \lambda_1 n_i.$$

10 fold CV is used to select  $\lambda_S = \hat{\lambda}_1 n$ .

$\lambda_M$  is the smallest value of  $\lambda$  such that

$$\hat{\underline{\eta}}_{\lambda_M} = \underline{0}. \quad \text{Hence } \hat{\underline{\eta}}_{\lambda_i} \neq \underline{0} \text{ for } i < M.$$

47) By Karush Kuhn Tucker (KKT) conditions for convex optimality (math <sup>see</sup> 471),

$$-\underline{W}^T (\underline{z} - \underline{W}\hat{\underline{\eta}}_L) + \frac{\lambda_1 n}{2} \underline{s}_n = \underline{0} \quad \text{where } s_{in} \in [-1, 1].$$

$$\text{Thus } \hat{\underline{\eta}}_L = (\underline{W}^T \underline{W})^{-1} \underline{W}^T \underline{z} - n (\underline{W}^T \underline{W})^{-1} \frac{\lambda_1 n}{2} \underline{s}_n$$



48) Lasso CLT: Assume  $p$  is fixed and 9L-28  
 the OLS CLT holds for  $\underline{z} = \underline{w}\eta + \underline{\varepsilon}$ .

a) If  $\hat{\lambda}_n / \sigma_n \xrightarrow{P} 0$ , then  $\sqrt{n}(\hat{\underline{\mu}}_L - \underline{\mu}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V)$ .

b) If  $\hat{\lambda}_n / \sigma_n \xrightarrow{P} \tau \geq 0$  and  $\underline{s}_n \xrightarrow{P} \underline{s} = \underline{s}_n$ , then

$$\sqrt{n}(\hat{\underline{\mu}}_L - \underline{\mu}) \xrightarrow{D} N_{p-1}\left(-\frac{\tau}{2} V \underline{s}, \sigma^2 V\right).$$

Proof  $\sqrt{n}(\hat{\underline{\mu}}_L - \underline{\mu}) = \sqrt{n}(\hat{\underline{\mu}}_L - \hat{\underline{\mu}}_{OLS} + \hat{\underline{\mu}}_{OLS} - \underline{\mu})$

$$= \sqrt{n}(\hat{\underline{\mu}}_{OLS} - \underline{\mu}) - \sqrt{n} \frac{\hat{\lambda}_n}{2n} (\underline{w}^T \underline{w})^{-1} \underline{s}_n \xrightarrow{D}$$

$$N_{p-1}(\underline{0}, \sigma^2 V) - \frac{\tau}{2} V \underline{s} \sim N_{p-1}\left(-\frac{\tau}{2} V \underline{s}, \sigma^2 V\right) \square$$

49)  $\hat{\underline{\mu}}_L$  is  $\sqrt{n}$  consistent if  $\frac{\hat{\lambda}_n}{\sigma_n} \xrightarrow{P} \tau$  and

$\hat{\underline{\mu}}_L$  is a consistent estimator of  $\underline{\mu}$  if  $\frac{\hat{\lambda}_n}{\sigma_n} \xrightarrow{P} 0$

as  $n \rightarrow \infty$ ,

50) At most  $n$  <sup>coeffs</sup>  $\hat{\underline{\mu}}_{OLS} \neq 0$  even if  $p > n$ .

51) Problems i) Lasso. Large sample theory is worse or as good as that of the OLS full model if  $\frac{n}{p}$  is large. ii) CV does not appear to guarantee that  $\frac{\hat{\lambda}_n}{n} \xrightarrow{P} 0$  or  $\frac{\hat{\lambda}_n}{\sigma_n} \xrightarrow{P} 0$ . iii) Lasso tends to under fit if  $\sigma_g \geq 0$  and the predictors are highly correlated.

iv) RR can be better than lasso if  $a_2 > n$ .

52}  $p$  large and iii) is bad, program should include  $\lambda_1 = 0$  if  $n \geq 5p$ .

Also include a value like  $\lambda_2 \approx \frac{\sqrt{n}}{\log(n)}$ .

53} Lasso can be a lot better than the OLS full model if

a)  $X^T X$  is singular or ill conditioned

b)  $\frac{n}{p}$  is small.

With  $n=100$  lasso can be much faster than forward selection if  $n$  and  $p$  are large.

54} The  $l_1$  relaxed lasso estimator

$\hat{\beta}_{RL}$  is OLS fit to the  $j$  variables,

including a constant, that have  $\hat{\eta}_{iL} \neq 0$ .

So  $\hat{\beta}_{RL}$  is an alternative to forward selection. Let  $X_A$  denote the matrix corresponding to the  $j$  "active" lasso variables.

$$\text{Then } \hat{\beta}_{RL} = (X_A^T X_A)^{-1} X_A^T Y$$

55} Relaxed lasso should be  $J_n$  consistent when lasso is  $J_n$  consistent.

$$56) \quad \underline{Y} = \underline{X} \underline{\beta} + \underline{e} \quad \underline{Z} = \underline{W} \underline{\eta} + \underline{e}$$

The elastic net estimator  $\hat{\underline{\eta}}_{EN}$  minimizes

$$\text{the criterion } RSS_w(\underline{\eta}) + \lambda_1 \|\underline{\eta}\|_2^2 + \lambda_2 \|\underline{\eta}\|_1 \quad (*)$$

where  $\lambda_1 = (1-\alpha)\lambda_n$ ,  $\lambda_2 = 2\alpha\lambda_n$  and  $0 \leq \alpha \leq 1$ .

Note that  $\alpha=1$  corresponds to lasso using  $2\lambda_n$  and  $\alpha=0$  corresponds to RR.

$$57) \quad \text{Let } \underline{W}_A = \begin{pmatrix} \underline{W} \\ \sqrt{\lambda_1} \underline{I}_{p-1} \end{pmatrix}, \quad \underline{Z}_A = \begin{pmatrix} \underline{Z} \\ \underline{0} \end{pmatrix}.$$

Then  $\hat{\underline{\eta}}_{EN}$  can be obtained from the lasso of  $\underline{Z}_A$  on  $\underline{W}_A$ . See proof on 0121.

58) By KKT conditions for convex optimality of (\*),  $2\underline{W}^T \underline{W} \hat{\underline{\eta}}_{EN} - 2\underline{W}^T \underline{Z} + 2\lambda_1 \hat{\underline{\eta}}_{EN} + \lambda_2 \underline{s}_n = \underline{0}$

$$\text{or } (\underline{W}^T \underline{W} + \lambda_1 \underline{I}_{p-1}) \hat{\underline{\eta}}_{EN} = \underline{W}^T \underline{Z} - \frac{\lambda_2}{2} \underline{s}_n \quad \text{or}$$

$$\hat{\underline{\eta}}_{EN} = \hat{\underline{\eta}}_{RR} - n (\underline{W}^T \underline{W} + \lambda_1 \underline{I}_{p-1})^{-1} \frac{\lambda_2}{2n} \underline{s}_n. \quad s_i \in \{-1, 1\}.$$

If  $\hat{\lambda}_n/n \xrightarrow{P} \tau$  and  $\hat{\alpha} \xrightarrow{P} \psi$ , then  $\frac{\hat{\lambda}_1}{\hat{\lambda}_n} \xrightarrow{P} (1-\psi)/\tau$  and  $\frac{\hat{\lambda}_2}{\hat{\lambda}_n} \xrightarrow{P} 2\psi\tau$ .

59} EN CLT Assume  $p$  is fixed and the LSCLT holds. a) If  $\hat{\sigma}_n / \sigma_n \xrightarrow{P} 0$ , then

$$\sqrt{n} (\hat{\mu}_{EN} - \underline{\mu}) \xrightarrow{D} N_{p+1}(\underline{0}, \sigma^2 V) \quad (\text{asy equiv to OLS}).$$

b) If  $\frac{\hat{\sigma}_n}{\sqrt{n}} \xrightarrow{P} \tau$ ,  $\hat{\alpha} \xrightarrow{P} \psi \in \mathbb{R}^d$  and  $\underline{s}_n \xrightarrow{P} \underline{s} = \underline{s}_n$ , then

$$\sqrt{n} (\hat{\mu}_{EN} - \underline{\mu}) \xrightarrow{D} N_{p+1}(-V[(1-\psi)\tau\underline{\eta} + \psi\tau\underline{s}], \sigma^2 V).$$

see proof on OLS use  $\sqrt{n}(\hat{\mu}_{EN} - \underline{\mu}) =$

$$\sqrt{n} (\hat{\mu}_{EN} - \hat{\mu}_{RR} + \hat{\mu}_{RR} - \underline{\mu}) \text{ and RR CLT.}$$

60} Theenet function uses 10 fold CV and

a grid of  $\alpha$  values  $\left\{ 0, \frac{1}{a_m}, \frac{2}{a_m}, \dots, \frac{a_m}{a_m} = 1 \right\}$   
 $a_m \geq 1$ ,  $a_m + 1$  values

so enet takes about  $(a_m + 1)$  times as long as

lasso as lasso or RR. The default is  $a_m = 10$ .

Both RR and lasso have problems with highly correlated predictors, so EN likely does, too.

OBS 3.9 61} PI for  $Y_e$  is  $\left[ \hat{Y}_e + b_n \hat{\Sigma}_{s_1}, \hat{Y}_e + b_n \hat{\Sigma}_{s_2} \right]$

asymptotically adds  $\frac{Y_e}{n}$  to the shorth applied to the residuals.

$$b_n = \left(1 + \frac{15}{n}\right) \sqrt{\frac{n+2d}{n-d}} \quad \text{if } d \leq \frac{8n}{9}, \quad b_n = 5 \left(1 + \frac{15}{n}\right) \text{ otherwise.}$$

where  $d =$  crude estimate of  $d_f$ .

SL 30

If  $n \geq 10p$   $d = p$  works ok for OLS, FS, PCR, PLS, RL, L, RL, EN. FS, PCR, PLS, L, RL use variables

$\underline{x}_1^*, \dots, \underline{x}_d^*$  and a better value for  $d$  is

$d = \#$  variables used, including a constant  $= \# \hat{\beta}_i \neq 0$ .

EN and RL have a better  $d = \text{tr}(H)$ , but it is hard to get  $d$  from the software.

FS, L, RL have  $\underline{x}_i^* = \underline{x}_j = (0, \dots, 0, \underset{\substack{\uparrow \\ \text{in position } j}}{1}, 0, \dots, 0) \underline{x}$

PCR and PLS have  $\underline{x}_i^* = \underline{\gamma}_i^T \underline{x}$ , some linear combination of the predictors  $\underline{x}$ .

62] If  $\frac{n}{p}$  is small, FS with EBIC, lasso and RL can work well under strong regularity conditions for sparse pop models.

63] PIs work badly if  $\underline{x}_f$  is not like the training data  $\underline{x}_1, \dots, \underline{x}_n$ , eg if  $\underline{x}_f$  is an outlier or not in a covering hyperellipsoid (or hypersphere) for  $\underline{x}_1, \dots, \underline{x}_n$ .



$\underline{x}_f$   
outlier

0  $\$ 3.10$  5  $\$ 5.1$  Choosing from  $M$  models

64) If  $\frac{n}{p} \geq 20$  could always pick OLS full model,  
Then inference is easy.

65) Program uses  $C_p$  if  $n \geq 10p$  and  $EBIC$  if  $n < 10p$   
for FS.

66) For  $k$ -fold cross validation ( $k$ -fold CV),  
randomly divide the training data into  $k$  groups  
or folds of approx equal size  $n_j \approx \frac{n}{k}$  for  
 $j=1, \dots, k$ . Leave out the 1st fold, fit the  
method to the  $k-1$  remaining folds, then compute  
some criterion for the 1st fold. Repeat for  
folds  $2, \dots, k$ .

67) For  $\underline{y} = \underline{X}\underline{\beta} + \underline{e}$  MLR, compute  
 $\hat{y}_i(j)$  for each  $y_i$  in the fold  $j$  left out.

Then  $MSE_j = \frac{1}{n_j} \sum_{i=1}^{n_j} (y_i - \hat{y}_i(j))^2$  and

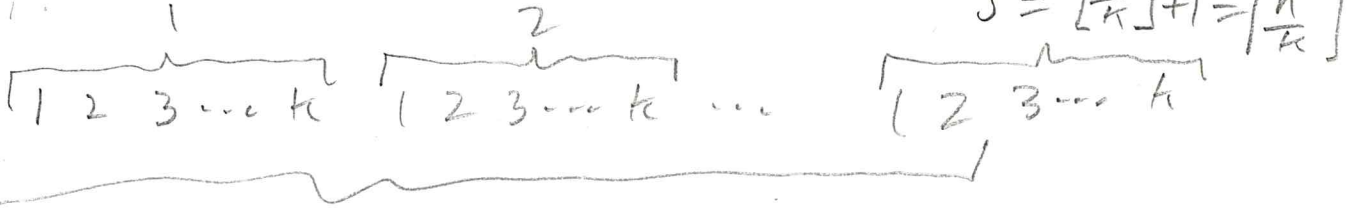
the overall criterion is  $CV_{(k)} = \frac{1}{k} \sum_{j=1}^k MSE_j$ .

If each  $n_j = \frac{n}{k}$ , then  $CV_{(k)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i(j))^2$ .

Pick model  $\underline{\beta}$  that minimizes  $CV_{(k)}(\underline{\beta})$ . Usually  $k=10$  or

68) J p 250 code is not very good. SL 31

It is better to generate



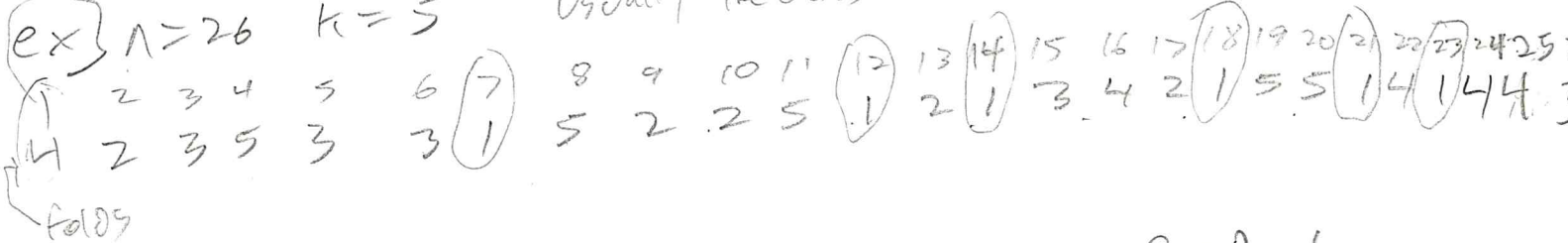
get vector  $J$  folds  
 $n \times 1$

Then select a random permutation of folds.

cases

ex  $n=26$   $k=5$

usually the cases are not given, see E2 rev 91).



cases 7, 12, 14, 18, 21 and 23 are in fold 1  
the other folds have 5 cases.

69) If  $n \geq 5p$ , lasso, RR, EN should

use  $\lambda_1 = 0$  and  $\lambda_2 \approx \frac{\sqrt{n}}{\log(n)}$ . They may

use  $\lambda_i = \frac{\lambda_m i}{100}$ ,  $i=1, \dots, 100$  where  $\lambda_m \stackrel{\text{lasso}}{=} 2 \max_j \left( \frac{s_j^T z}{\|z\|} \right)$ .

So  $\lambda_i$  may not be small enough for

RR, L or EN CLT's to hold.

$j$ 'th column of  $W$

$\lambda_m \propto n$  is possible if  $s_j$  is highly correlated with  $z$ .

70} For FS, PCR, and PLS,

$\bar{I}_1, \bar{I}_2, \dots, \bar{I}_p$   
 $x_1^*, x_2^*, \dots, x_p^*$  are well defined

For RR, lasso and EN,  $\bar{I}_i$  uses  $\frac{\bar{I}_i}{2n} = \frac{\bar{I}_{ni}}{2n}$ , but  
10 fold CV uses  $K = n - n_j \approx .9n$  cases when fold  $j$   
is left out. Maybe use  $\frac{\bar{I}_{in}}{2n} \approx \frac{\bar{I}_{ik}}{2K} \approx c$  are

comparable, so  $\bar{I}_{ik} \approx \frac{K}{n} \bar{I}_{in}$ . It is not clear  
what the code does.

71} Modify 10 fold CV to compute PI coverage  
and ave PI length. So  $n$  PIs are made using  
 $\bar{I}_i = \bar{I}_i$  for  $i=1, \dots, n$ . The coverage is the prop  
of times the  $n$  PIs contained  $y_i$ . Choose model  
 $I_d$  with the shortest ave PI length such  
that the nominal 100(1- $\delta$ )% PI had coverage

$$\geq c_n = \max\left(1-\delta - \frac{1}{35n}, 1-\delta - 0.02\right).$$

If no model  $I_i$  has coverage  $\geq c_n$ , choose  
the model with the largest coverage. See 0.119, Fig 3.1.

Read 0.3.11: all  $\beta_{0i}^* = 0 \Rightarrow \bar{I}_i$  is the cont reg. 0.01

3.12 simulations 72}  $\hat{\beta}_1, \dots, \hat{\beta}_p$  test for FS, lasso, RR



73) Compare methods with PIS:

want PI len near pop shorth

$e_i \sim N(0,1)$   $\bar{x}_3$   $E\|e\|_1 = 1.90$   $1.9N(0,1) + 1N(0,100)$   
 pop len 3.92 6.365 2.996 1.90 13.490

10 For RR used same data as lasso.

74) For  $\frac{n}{p}$  large OLS forward sel PIS work,

PLS and PCR simulated well. Need  $\frac{n}{p} \approx 100$  to get close to pop len. Lasso, RR, (RL) had

trouble if  $k \geq 19$   $\psi \geq 0.5$ .  $k = \#$  nontrivial active predictors =  $q_3 - 1$ .

$$B = \underbrace{(1, \dots, 1)}_{k+1} \underbrace{(0, \dots, 0)}_{p-k-1}^T$$

n	p	$\psi$	k	cov/len	FS	L	RL	RR	PLS	PCR	
200	20	0	19	cov	.979	.977	.979	.979	.955	.979	PLS best given cov $\geq .94$
				len	4.961	4.964	4.961	5.046	4.321	4.961	
400	20	0.9	19	cov	.967	.975	.960	.973	.955	.954	PCR best
				len	4.512	10.609	4.562	10.663	4.002	3.977	

75) For  $\frac{n}{p} < 5$  FS and lasso worked ok if k is small (sparse model) severe undercoverage

n	p	$\psi$	k	cov/len	FS	L	RL	RR	PLS	PCR	
100	200	0	1	cov	.965	.976	.927	.958	.662	.992	FS best given cov $\geq .90$
				len	4.427	4.976	4.225	6.161	2.770	12.412	
400	400	0.9	19	cov	.935	.964	.956	.963	.946	.948	PLS best
				len	4.369	47.761	4.853	48.021	4.291	4.476	

$\psi = .9$  is favorable for PCR, PLS: 1 dominant principal component

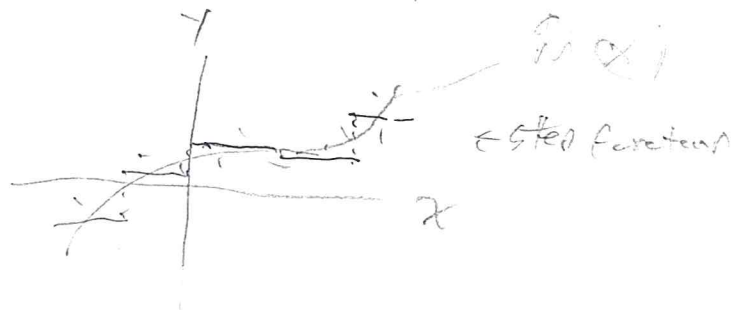
76) For a dense model  $\frac{n}{q_3}$  or  $\frac{n}{k}$  is not small  $\neq$

$(y_i, u_i) \sim N_p(\mu, \Sigma)$  every subset  $J_i$  follows MLR model index 2

0 ch4 5 ch7

507.1-7.3

Suppose  $Y = m(x) + e$ .



There are several flexible scatterplot smoothers  $\hat{m}(x)$  with  $df > 2$

including loess, lowess, smoothing splines

and the additive error GAM  $Y = \sum \beta_j x_j + e$ , or normal,

polynomial regression uses  $Y = \sum_{i=1}^J \beta_i x^i + e$

$$= \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_J x^{J-1} + e.$$

We could use step functions: divide  $x$  into  $J$  intervals with  $\approx \frac{n}{J}$  cases compute  $\bar{Y}$  in each interval

Basis functions use  $Y = \sum_{i=1}^J \beta_i b_i(x) + e$  where  $b_i(x) \equiv 1$ .

PI 0(2.7) often works well.

2) 507.4 We could divide  $x$  into  $k+1$  intervals

$$[a_0, a_1] [a_1, a_2], \dots, [a_k, a_{k+1}] [a_k, a_{k+1}]$$

where  $a_1, a_2, \dots, a_k$  are called knots.

Fit a cubic polynomial to each interval.

3) we get a cubic spline by having  $\hat{m}(x)$  continuous with 2nd derivatives (so the curve  $\hat{m}(x)$  is smooth). A cubic spline with  $k$  knots

$$\text{can be modeled as } Y = \beta_1 + \beta_2 b_2(x) + \dots + \beta_{k+4} b_{k+4}(x) + e$$

$1, x, x^2, x^3, \dots, h(x, \xi_1), \dots, h(x, \xi_m)$  where 9L 33

$\xi_1, \dots, \xi_m$  are the knots and

$$h(x, \xi) = \begin{cases} (x - \xi)^3, & x > \xi \\ 0, & \text{else.} \end{cases}$$

Programs use general sets of knots at percentiles and CV.

§7.5 4) A smoothing spline  $\hat{g}$  minimizes   
 use CV or BEV and various values of  $\lambda$

$$\underbrace{\sum_{i=1}^n (y_i - g(x_i))^2}_{RSS(g)} + \lambda \int [g''(t)]^2 dt$$

$\hat{Y} = H_2 Y$  and  $d = \text{tr}(H_2)$ . Then  $\hat{g}(x) = \hat{m}(x)$ .

§7.6 describes loess and loess

§7.7 OCH4

6) 4 important regression models are

1)  $Y = SP + e$ ,  $\begin{cases} SP = m \times 1 \\ \text{response plot and PIs like MLR} \end{cases}$  AER

2)  $Y|X \sim \text{Poisson}(e^{SP})$  poisson reg

3)  $Y|X \sim \text{bin}(1, \frac{e^{SP}}{1 + e^{SP}})$  binary reg

4)  $Y_i | x_i \sim \text{bin}(m_i, \frac{e^{SP_i}}{1 + e^{SP_i}})$  binomial reg

3) is a special case of 4)   
 when  $m_i = 1$

If  $SP = \underline{x}^T \underline{\beta}$  the model tends to be a generalized linear model GLM.

If  $SP = \beta_1 + \sum_{j=2}^p g_j(x_j) = AP =$  additive predictor, the model tends to be a generalized additive model GAM. The GLM is a special case with  $g_j(x_j) = \beta_j x_j$ . A GAM is flexible while a GLM is inflexible.

7) For  $n \geq 5p$ ,  $J \geq 5$  and preferably,  $J \geq 10$  AIC is used for forward selection for a GLM.

Also plot  $\underline{x}_i^T \underline{\beta}_{OLS}$  vs  $\underline{x}_i^T \underline{\beta}_{GLM}$ . If the correlation is high do MLR forward selection and fit GLM to  $J_{min}$  variables using Cp if  $n \geq 10p$ .

8) <sup>0.54.6.2</sup> LASSO-EN, type criterion for binary, binomial?, and poisson regression (add family = binomial or family = poisson to the cv.glmnet arguments).

A relaxed lasso GLM fits the GLM to the predictors with nonzero lasso coefficients.

Or do MLR lasso and fit GLM to  $J_{min}$  variables.

Using MLR forward selection with EBIC and applying  $J_{min}$  variables led to severe underfitting.

$$EAP = ESP = \hat{\beta}_1 + \sum_{j=2}^p \hat{\gamma}_j(x_j)$$

SL 34

Often smoothing splines are used iteratively to find  $\hat{\gamma}_j$  for  $j=2, \dots, p$ . If a plot of  $x_j$  vs  $\hat{\gamma}_j$  is linear,  $\hat{\beta}_j x_j$  can be used. If a plot of  $x_j$  vs  $\hat{\gamma}_j$  is a quadratic add  $x_j^2$ , cubic add  $x_j^3$  and  $x_j^4$ . So the GAM can be useful for checking the GLM and for suggesting predictor transformations for the GLM.

10) Poisson regression response plot

$$E(Y|x) = e^{ESP}$$



Counts should track  $e^{ESP}$  except possibly for large ESP.

11) Binary regression response plot

(For binomial regression replace  $Y_i$  by  $\frac{Y_i}{m_i}$  for a similar plot)

$$E(Y|x) = \frac{e^{ESP}}{1 + e^{ESP}}$$



0 if  $Y=0$   
1 if  $Y=1$

Add a step function with step height =  $\bar{y}$  for cases in interval. Want step function to track logistic curve closely.

1) In supervised classification, there are  $G \geq 2$  known groups and  $m$  test cases to be classified. Each case is predicted to be in exactly one group based on its measurements

$w_i$ .

ex) patient with heart attack symptoms:

3 tests are done  $\underline{w} = (w_1, w_2, w_3)^T$  and the patient is classified as i) had a heart attack or ii) did not have a heart attack.

ex) Person applies for credit card based on  $(\text{salary}, \text{credit rating})^T$  and is classified as acceptable or not acceptable.

2) Suppose there are  $G \geq 2$  groups or populations

with pdf  $f_j$  ( $\geq 1$  for  $j=1, \dots, G$ ) and  $\underline{x}$  is  $p \times 1$ .

So if  $\underline{x}$  comes from pop  $j$ , then  $\underline{x}$  has pdf  $f_j$  ( $\geq 1$ ).

Assume there is training data consisting of a random sample of  $n_j$  cases  $\underline{x}_{1j}, \dots, \underline{x}_{n_j j}$  for each group. Let  $(\bar{\underline{x}}_j, s_j)$  be the sample

mean and covariance matrix for each group. SL 35

Let  $\underline{w}_i$  be a new  $p \times 1$  <sup>test</sup> random vector from 1 of the  $G$  groups, but the group is unknown. Usually there are  $\underline{w}_1, \dots, \underline{w}_m$  = test data, and discriminant analysis or classification attempts to allocate the  $\underline{w}_i$  to the correct groups.

3) Let  $\underline{\tilde{w}}_i$  be a random vector and  $\underline{w}_i$  the observed random vector. Let  $Y_i = j$  if  $\underline{w}_i$  comes from the  $j$ th group for  $j=1, \dots, G$ .

Then  $\pi_j = P(Y=j)$  and the posterior probability that  $Y=k$  or  $\underline{w}_i$  belongs to group  $k$

$$\text{is } P_k(\underline{w}_i) = P(Y=k | \underline{\tilde{w}}_i = \underline{w}_i) = \frac{\pi_k f_k(\underline{w}_i)}{\sum_{j=1}^G \pi_j f_j(\underline{w}_i)}$$

4) i) The maximum likelihood discriminant rule

allocates  $\underline{w}_i$  to group  $a$  (predicts  $Y=a$ ) if

$\hat{f}_a(\underline{w}_i)$  maximizes the  $\hat{f}_j(\underline{w}_i)$  for  $j=1, \dots, G$ .

ii) The Bayesian discriminant rule allocates

$\underline{w}_i$  to group  $a$  if  $\hat{P}_a(\underline{w}_i)$  maximizes

$$\hat{P}_k(\underline{w}_i) = \frac{\hat{\pi}_k \hat{f}_k(\underline{w}_i)}{\sum_{j=1}^G \hat{\pi}_j \hat{f}_j(\underline{w}_i)} \quad \text{for } k=1, \dots, G.$$

iii) The pop Bayes Classifier allocates

$\underline{w}_i$  to group  $a$  if  $P_a(\underline{w}_i)$  maximizes

$P_k(\underline{w}_i)$  for  $k=1, \dots, G$ .

5 } p38-39, 139 The Bayes classifier has the lowest possible expected test error rate out of all classifiers using the same  $p$  predictors  $\underline{w}$ .  
Problem: usually the  $\pi_j$  and  $f_j$  are unknown.

6) The maximum likelihood rule and Bayesian discriminant rule are equivalent if  $\hat{\pi}_j \equiv \frac{1}{G}$  for  $j=1, \dots, G$ .

7) General discriminant rules can be modified to incorporate  $\pi_j$  and costs of incorrect and incorrect classification. We will assume that the costs of correct allocation are unknown or 0 and that the costs of incorrect allocation are unknown



the probabilities  $\pi_j$  are unknown or GL 36

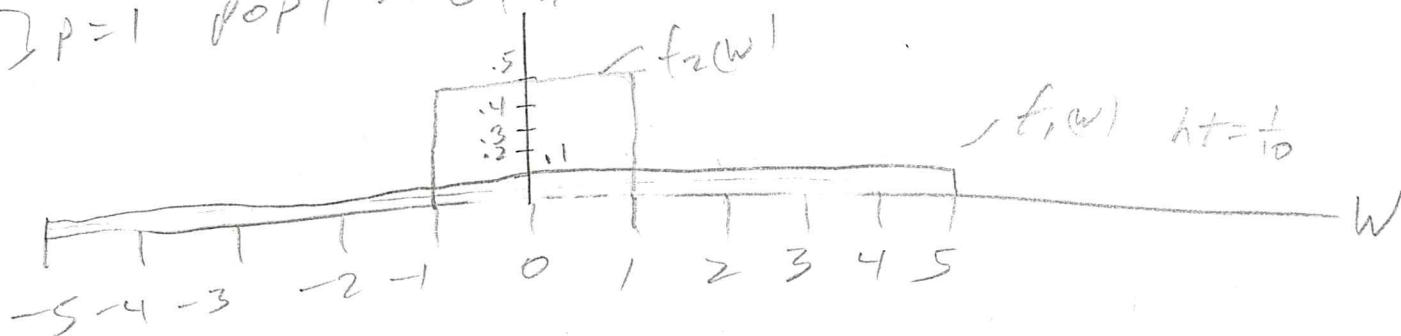
$$\pi_j \equiv \frac{1}{G} \text{ for } j=1, \dots, G,$$

8) A regularized estimator attempts to use  $d$  such that  $\frac{n}{d}$  is large.

(Lasso and RR with  $\lambda > 0$  were regularized MLE estimators.)

9) Given  $p=1$  and a graph of  $f_1, \dots, f_G$ , give the maximum likelihood rule.

ex)  $p=1$  pop 1  $\sim U(-5, 5)$  pop 2  $\sim U(-1, 1)$



ML rule: allocate  $w$  to pop 2 if  $-1 < w < 1$   
 allocate  $w$  to pop 1 if  $-5 < w < -1$  or  $1 < w < 5$ .

10) The pooled covariance matrix is

$$S_{\text{pool}} = \frac{1}{n-G} \sum_{j=1}^G (n_j - 1) S_j.$$

This estimator is good if the  $G$  groups have the same

covariance matrix  $\Sigma_x$ .

$\hat{\beta}_j$  and  $S_j$  have  $p + (p-1) + \dots + 1 = \frac{p(p+1)}{2}$  unknown parameters

So the  $G$   $\hat{\beta}_j$  have  $G \frac{p(p+1)}{2}$  unknown parameters.

The  $d_j$  is often = # estimated parameters.

So  $S_{pool}$  is regularized compared to using

$S_1, \dots, S_G$  separately.

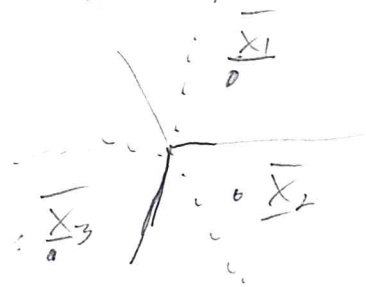
11) \* The linear discriminant analysis (LDA) rule

allocates  $\underline{w}$  to the group with the largest value of

$$d_j(\underline{w}) = \underline{\bar{x}}_j^T S_{pool}^{-1} \underline{w} - \frac{1}{2} \underline{\bar{x}}_j^T S_{pool}^{-1} \underline{\bar{x}}_j = \hat{\alpha}_j + \hat{\beta}_j^T \underline{w}.$$

LDA is the most used rule, and basically separates the

$G$  groups with  $G-1$  hyperplanes



12) \* The quadratic discriminant analysis (QDA) rule allocates  $\underline{w}$  to the

group with the largest value of

$$Q_j(\underline{w}) = -\frac{1}{2} \log |S_j| - \frac{1}{2} (\underline{w} - \underline{\bar{x}}_j)^T S_j^{-1} (\underline{w} - \underline{\bar{x}}_j)$$

$$= -\frac{1}{2} \log |S_j| + D_w^2(\underline{\bar{x}}_j, S_j)$$

QDA  $\Delta \gg$  LDA  $\Delta$

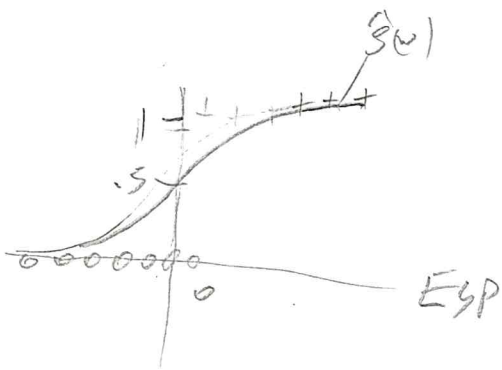
13)\* Let  $G=2$  with groups 0 and 1.

SL 37

Let  $S(\underline{w}) = P(\underline{w} \in \text{group } 1) = P(Y=1|\underline{w})$ . Let  $\hat{S}(\underline{w})$  be the binary logistic regression estimator of  $S(\underline{w})$ . Then  $ESP = \hat{\alpha} + \hat{\beta}^T \underline{w}$  and  $\hat{S}(\underline{w}) = \frac{e^{ESP}}{1 + e^{ESP}} =$

$\frac{\exp(\hat{\alpha} + \hat{\beta}^T \underline{w})}{1 + \exp(\hat{\alpha} + \hat{\beta}^T \underline{w})}$ . The LR discriminant rule

allocates  $\underline{w}$  to group  $\begin{cases} 1 & \text{if } \hat{S}(\underline{w}) \geq 0.5 \text{ (ESP} \geq 0) \\ 0 & \text{if } \hat{S}(\underline{w}) < 0.5 \text{ (ESP} < 0) \end{cases}$



training data  
The misclassification rate  $\approx \min(\hat{S}, 1 - \hat{S})$ .  
ESP  $> 2$  or ESP  $< -2$  will have low error (misclassification) rate if the LR model is good.

14)\* Training data:  $x_{1j}, \dots, x_{n_j j}$  group  $j$ .

Used discriminant analysis method to classify the training data (optimistic compared to test data! estimated error rate is too low). If  $m_j$  of the  $n_j$  group  $j$  cases are correctly classified, the apparent error rate for group  $j$  is  $1 - \frac{m_j}{n_j}$ .

Let  $m_A = \sum_{j=1}^G m_j$ . Then the  
apparent error rate  $AER = 1 - \frac{m_A}{n}$ .

15) leave one out cross validation CV: For  $i=1, \dots, n$   
leave case  $i$  out, compute the discriminant rule  
and see if case  $i$  is correctly classified. Let  
 $m_C$  be the number of cases correctly classified.  
The CV error rate is  $1 - \frac{m_C}{n}$ .

16) Leave out a subset with enough cases.  $N_V$   
(10% to 50%) so that a good error estimate  
can be obtained. Compute the discriminant rule  
from the cases not left out. Let  $m_L$  be the  
number of cases correctly classified. The

"leave a subset out" error rate is  $1 - \frac{m_L}{N_V}$ .

17)  $k$ -fold CV. Let  $m_K$  be the number of  
cases correctly classified. Then the

$k$ -fold CV error rate is  $1 - \frac{m_K}{n}$ .

18) \* <sup>3.2.2</sup> The  $k$ -nearest neighbors <sup>kNN</sup> method SL 38

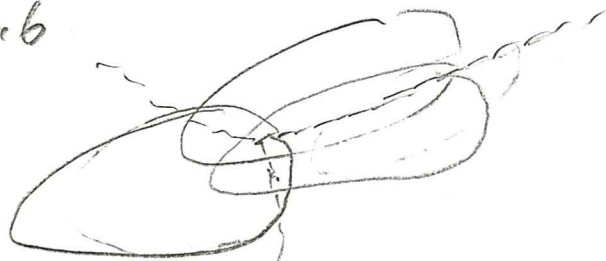
identifies the  $k$  cases in the training data closest to the test case  $w$ . Suppose  $m_j$  of the  $k$  cases are from group  $j$ .

The  $k$  nearest neighbor kNN discriminant rule allocates  $w$  to group  $a$  if  $m_a$  maximizes  $m_j$  for  $j=1, \dots, G$ .

(If there is a tie, we might use the group, among the tied groups, that has the smallest sum of distances.)

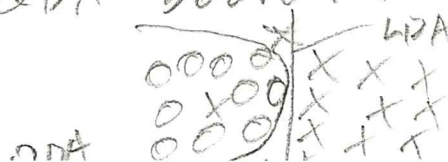
19) LDA roughly makes  $G$  hyperellipsoids  $\{D_j^2(\bar{x}_j, S_p) \leq h^2\}$  of the same volume and shape that cover most of the training data. The hyperplanes minimize overlap of the hyperellipsoids.

See 3 Fig 4.6



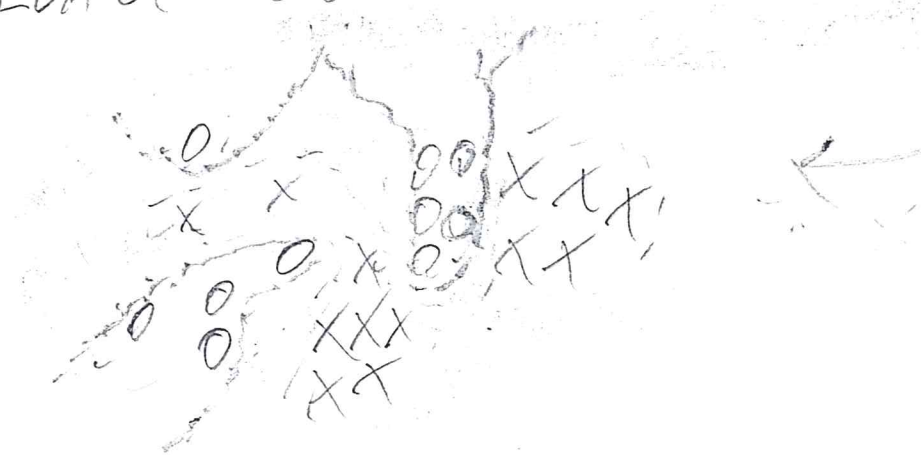
20) For  $p=2$  LDA boundaries are lines  
QDA boundaries are quadratics.

See 3 F 4.9



21) KNN is flexible see § 7.15 7.16

22) Adding terms like  $x_i^2$  &  $x_i x_j$  to LDA increases the LDA df and makes LDA more like QDA.



23) For Fisher's discriminant analysis (FDA), let  $\hat{W}$  be a  $p \times p$  dispersion matrix and let  $\hat{B}$  be a  $p \times p$  symmetric matrix used to measure variability between classes. Let the eigenvalue eigenvector pairs of  $\hat{W}^{-1} \hat{B}$  be  $(\hat{\lambda}_1, \hat{e}_1), \dots, (\hat{\lambda}_p, \hat{e}_p)$  with  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$ . Then  $\max_{\underline{a} \neq 0} \frac{\underline{a}^T \hat{B} \underline{a}}{\underline{a}^T \hat{W} \underline{a}} = \hat{\lambda}_1$ , with  $\underline{a} = \hat{e}_1$ . Similarly  $(\hat{\lambda}_i, \hat{e}_i)$  achieves the max among all unit vectors orthogonal to  $\hat{e}_1, \dots, \hat{e}_{i-1}$ .

24) Let  $(T_i, C_i)$  be an estimator of multivariate location and dispersion for the  $i$ th group, eg

$$(T_i, C_i) = (\bar{x}_i, S_i). \text{ Let } \bar{T} = \frac{1}{G} \sum_{i=1}^G T_i \text{ and}$$

$$\hat{B}_T = \sum_{i=1}^G (T_i - \bar{T})(T_i - \bar{T})^T \text{ So } \frac{\hat{B}_T}{G-1} \text{ is the sample}$$

covariance matrix of  $T_1, \dots, T_G$ . Let  $\hat{W}_T = \sum_{i=1}^G C_i$ .

Let  $\hat{W}_L = G \hat{\Sigma}_{\text{pool}}$ . Other choices are possible.

SL 39

25) The FDA discriminant rule allocates  $\underline{w}$  to group  $a$  that minimizes  $|\hat{\underline{e}}_i^T \underline{w} - \hat{\underline{e}}_i^T \underline{T}_i|$  for  $i=1, \dots, G$ .

26) Let the  $i$ th group have pop mean and covariance matrix  $(\underline{\mu}_i, \underline{\Sigma}_i)$ . Let  $\bar{\underline{\mu}} = \frac{1}{G} \sum_{i=1}^G \underline{\mu}_i$ ,  $\underline{B} = \sum_{i=1}^G (\underline{\mu}_i - \bar{\underline{\mu}})(\underline{\mu}_i - \bar{\underline{\mu}})^T$  and  $\underline{W} = \sum_{i=1}^G \underline{\Sigma}_i$ . Then the between group variability

$$b(\underline{a}) = \underline{a}^T \underline{B} \underline{a} = \sum_{i=1}^G |\underline{a}^T (\underline{\mu}_i - \bar{\underline{\mu}})|^2, \text{ the within group}$$

$$\text{variability } w(\underline{a}) = \underline{a}^T \underline{W} \underline{a} = \sum_{i=1}^G \underline{a}^T \underline{\Sigma}_i \underline{a} = \sum_{i=1}^G \text{Var}(\underline{a}^T \underline{x}_i)$$

where  $\underline{x}_i$  is a RV from the pop corresponding to group  $i$ .

$$\text{Then } \max_{\underline{a} \neq \underline{0}} \frac{b(\underline{a})}{w(\underline{a})} = \max_{\underline{a} \neq \underline{0}} \frac{\underline{a}^T \underline{B} \underline{a}}{\underline{a}^T \underline{W} \underline{a}} = \lambda_1 \text{ with } \underline{a} = \underline{e}_1$$

where  $\lambda_1$  is the largest eigenvalue of  $\underline{W}^{-1} \underline{B}$ .

→ FDA attempts to approximate 26), and will likely work well with  $\hat{\underline{W}}_A$  and  $\hat{\underline{B}}_T$  if  $\hat{\underline{a}}^T \hat{\underline{W}}_A \underline{a}$  can be made small.

28) If  $G=2$ ,  $(\tau_i, \underline{c}_i) = (\bar{\underline{x}}_i, \underline{s}_i)$ ,  $\hat{\underline{B}} = \hat{\underline{B}}_T$  and

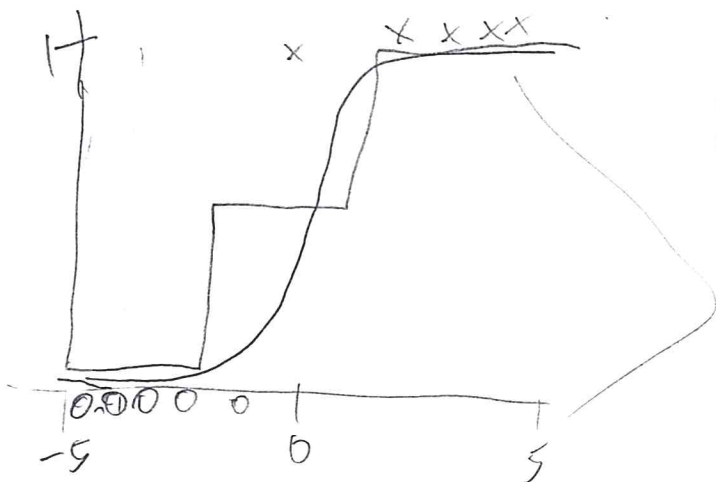
$\hat{\underline{W}} = 2 \hat{\Sigma}_{\text{pool}}$ , then LDA and FDA are equivalent.

The OLS program is actually an FDA program.

29) For kNN if  $k=1$   $AER=1$  but if  $\frac{n_j}{n} \rightarrow \pi_j$  then the test error rate  $L_n$  of  $k=1$  kNN converges in probability to  $L$  where  $L_B \leq L \leq 2L_B$  and  $L_B$  is the test error rate of the Bayes classifier. If  $k=k_n$  satisfies  $k_n \rightarrow \infty$  and  $\frac{k_n}{n} \rightarrow 0$  as  $n \rightarrow \infty$  then the kNN test error rate  $L_n \xrightarrow{P} L_B$ . The leave one out CV error rate  $\hat{L}_n$  is a good estimator of  $L_n$ .

30) If  $\underline{w}$  is  $p \times 1$  where  $p$  is large often  $\underline{w}_{I_k}$   $k \times 1$  will do better at classification than  $\underline{w}$  for some subset  $I_k$  of the  $p$  predictors.

31) LR response plots



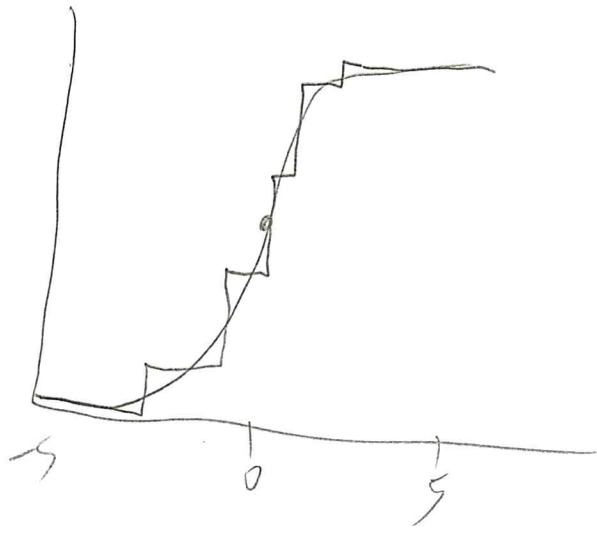
good response plot

group 1 and 0 cases are well separated

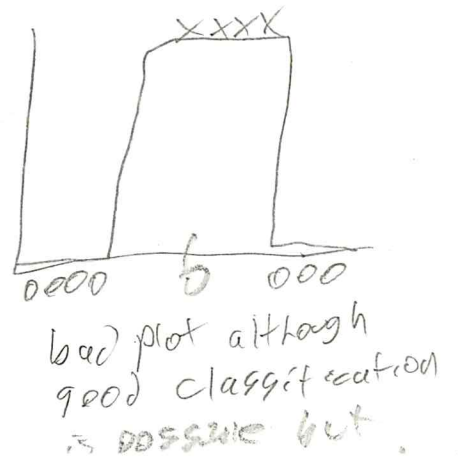
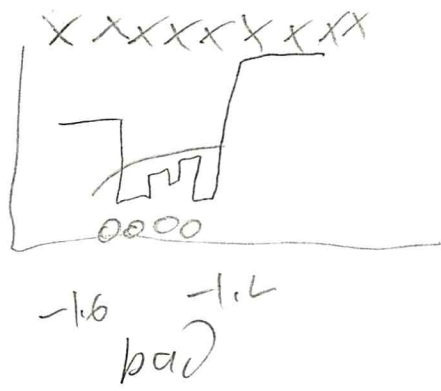
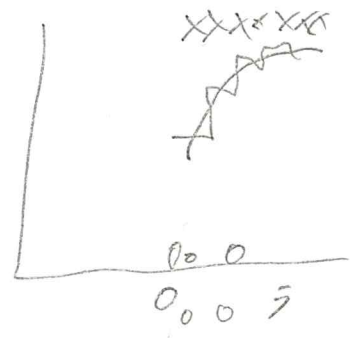
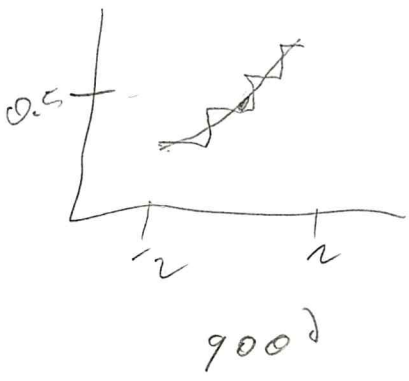




good response plot  
 perfect classification  
 of training data



good response plot



32} LR output

label	coef estimate	std err	Est/SE	Pr >  z
constant	$\hat{\alpha}$			
$x_1$	$\hat{\beta}_1$			
$\vdots$	$\vdots$			
$x_p$	$\hat{\beta}_p$			

ex label Estimate

	Estimate	0 = F	1 = M	Gender
constant	-19.7762			
head measurements	circum	0.0244688		
	length	0.0371472		

Let circum =  $x_1 = 550$ , length =  $x_2 = 200$

a) Find ESP for  $x$ .

$$ESP = \hat{\alpha} + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 = -19.7762 + 0.0244688(550) + 0.0371472(200) = \underline{1.1105}$$

b) Is  $x$  classified in group 0 or group 1?

group 1 since  $ESP > 0$

c) Find  $P(x) = P(ESP)$

$$= \frac{e^{ESP}}{1 + e^{ESP}} = \frac{3.0376}{4.0376} = 0.7523$$

33] variable selection: find a subset (GL 41)  
of variables that does a good job of  
classification

Crude forwards selection  $x_1, \dots, x_p$  are potential  
predictors

Step 1) Choose  $w_1 = x_j$  that minimizes AER

2) Keep  $w_1$  in the model and add  $w_2 = x_j$   
that minimizes the AER. So  $w_1$  and  $w_2$  are in  
the model.

⋮

k)  $w_1, \dots, w_{k-1}$  are in the model.

Add  $w_k = x_j$  that minimizes the AER.

⋮

p)  $w_1, \dots, w_p = x_1, \dots, x_p$

Final model: might be the one with the smallest AER  
if  $n \gg p$  so there is no overfitting.

Crude backward elimination

Step 1)  $w_1, \dots, w_p = x_1, \dots, x_p$  are in the model

2) Delete  $w_p = x_{j_1}$  such that the model with  
 $p-1$  variables minimizes the AER.  $w_1, \dots, w_{p-1}$  are  
in the model.

3) Delete  $w_{p-1} = x_{j_2}$  such that the model with  
 $p-2$  variables minimizes the AER.  $w_1, \dots, w_{p-2}$   
are in the model

k)  $w_1, \dots, w_{p-k+2}$  are in the model.

Delete  $w_{p-k+2} = x_{jk}$  such that the model with  $p-k+1$  variables minimizes the AER.

$w_1, \dots, \underbrace{w_{p-k+1}}_{p-(k-1)}$  are in the model.

⋮

p)  $w_1$  and  $w_2$  are in the model, Delete  $w_2$  such that the model with  $w_1$  minimizes the AER.

34) Other Criteria can be used.

Proc stepdisc in SAS does variable selection, see OBS.8 and HW 8 for crude variable selection with LDA and QDA.

0 ch 6 Regularizing a correlation matrix  $R$ .

1) Classification is a multivariate analysis technique. Many multivariate analysis techniques are based on the covariance matrix  $S = \hat{\Sigma}_X$  or correlation matrix  $R = \hat{\rho}_X$ . If  $X$  is  $p \times 1$   $\hat{\Sigma}_X$  has  $p + (p-1) + \dots + 1 = \underline{p(p+1)}$  parameters and  $\hat{\rho}_X$

has  $(p-1) + (p-2) + \dots + 1 = p(p+1)/2$  parameters. SL42

We want  $n \geq 10p$  to use R or S. If

$n \leq 5p$ , these matrices are being overfit:

the degrees of freedom is too large for the sample size  $n$ , and the matrices may be ill conditioned.

2) Also many multivariate techniques need

$\hat{\Sigma}_x^{-1}$  or  $\hat{S}_x^{-1}$ . S and R need

$n > p$  to be invertible.

3) Plug in estimators are often used.

ex "LDA" uses  $\hat{\mu}_j^T \hat{\Sigma}_{pool}^{-1} \underline{w} = -\frac{1}{2} \hat{\mu}_j^T \hat{\Sigma}_{pool}^{-1} \hat{\mu}_j$

where the usual LDA has  $(\hat{\mu}_j, \hat{\Sigma}_{pool}) = (\bar{X}_j, S_{pool})$ ,

4) A regularized estimator  ~~$\hat{\Sigma}_x$~~  attempts to

use well conditioned  ~~$\hat{\Sigma}_x$~~  or

~~$\hat{\Sigma}_x$~~  or  $\hat{S}_x$ .

5) A common technique is to use

$\hat{\Sigma}_x = S_d = \text{diag}(S)$  or  $\hat{S}_x = R_d = \text{diag}(R) = I_p$ .

(like using  $\hat{\beta}_M = \underline{0}$  for lasso, RR or EN).

6) Let  $S = (s_{ij})$  and  $D = \text{diag}(\sqrt{s_{11}}, \dots, \sqrt{s_{pp}}) = \text{diag}(\sqrt{s_{ii}})$

Then  $S = D R D$  and  $R = D^{-1} S D^{-1}$

7) Let  $S^{-1} = (s^{ij})$  and  $R^{-1} = (r^{ij})$

Let  $\underline{x}_{(ij)}$  be the vector  $\underline{x}$  with  $x_i$  and  $x_j$  deleted  $i \neq j$ . Then partial correlation

$r_{ij \underline{x}_{(ij)}}$  between  $x_i$  and  $x_j$  is found by

regressing  $x_i$  and  $x_j$  on  $\underline{x}_{(ij)}$  get the two sets of residuals then find the correlation

$r_{ij \underline{x}_{(ij)}}$  of the two sets of residuals.

$$r_{ij \underline{x}_{(ij)}} = \frac{-s^{ii}}{\sqrt{s^{ii} s^{jj}}} = \frac{-r^{ii}}{\sqrt{r^{ii} r^{jj}}}$$

and  $R^{-1} = D S^{-1} D$ . Also  $r^{ii} = \text{VIF}_i = \frac{1}{1 - R_i^2}$

where  $R_i^2$  is the squared multiple correlation from regressing  $x_i$  on  $\underline{x}_{(i)}$  that omits predictor  $i$ .

8) Graphical lasso (Glasso) takes  $\underline{S} \succeq \underline{0}$  as an input and returns  $\hat{\beta}_G \succeq \underline{0}$ . Many elements  $\hat{\beta}_G$  are zero.  $\underline{S}$  is positive definite,  $\hat{\beta}_G$  is positive semi-definite.

9] Could use a robust estimator as an input.

SL43

10] There are a couple of books written on regularized versions of  $\hat{\beta}$ ,  $\hat{\beta}^T$ ,  $R$ ,  $R^T$ .

11]  $RR^T$  regularized  $WTW = nR$ .

12] Let  $\delta \geq 0$  and  $R_{p \times p} = (r_{ij})$ . A simple way to regularize  $R$  is to use

$$R_\delta = \frac{1}{1+\delta} (R + \delta I_p) = (t_{ij}) \quad \text{with } t_{ii} = 1$$

and  $t_{ij} = \frac{r_{ij}}{1+\delta}$ , for  $i \neq j$ .  $R_\delta = kR + (1-k)I_p$

where  $k = \frac{1}{1+\delta} \in (0, 1]$ . Note that each correlation  $r_{ij}$  is divided by the same factor  $1+\delta$ .

If  $\lambda_i$  is the  $i$ th eigenvalue of  $R$ , then  $\frac{\lambda_i + \delta}{1+\delta}$  is the  $i$ th eigenvalue of  $R_\delta$ . The eigenvectors of  $R$  and  $R_\delta$  are the same since if  $R\mathbf{x} = \lambda_i \mathbf{x}$  then

$$R_\delta \mathbf{x} = \frac{1}{1+\delta} (R + \delta I_p) \mathbf{x} = \frac{1}{1+\delta} (\lambda_i + \delta) \mathbf{x}.$$

13]  $R$  can be further regularized by setting

$t_{ij} \doteq 0$  if  $|t_{ij}| \leq \tau$  where  $\tau \in [0, 1]$  should have  $\tau < 0.5$ . Denote the resulting matrix by

$R(\delta, \tau)$ , Then  $R_\delta = R(\tau, 0)$ . The default is  $\tau = 0.05$ . Using  $\tau$  is known as thresholding.

14) A regularized covariance matrix

$$S(\delta, \tau) = D R(\delta, \tau) D \quad \text{with } D = \text{diag}(s_{ii})$$

↑  
See 6)

15) The condition number of a symmetric  $A > 0$  is  $\text{cond}(A) = \frac{\lambda_1(A)}{\lambda_p(A)} \geq 1$  where  $\lambda_1(A) \geq \dots \geq \lambda_p(A) > 0$

are the eigenvalues of  $A$ , A well conditioned matrix has  $\text{cond}(A) \leq C$  for some number

$C$  such as 50 or 500,  $R_\delta > 0$  for  $\delta > 0$

and well conditioned if  $\text{cond}(R_\delta) = \frac{\lambda_1 + \delta}{\lambda_p + \delta} \leq C$

or  $\delta = \max\left(0, \frac{\lambda_1 - c\lambda_p}{c-1}\right)$  if  $1 < C \leq 500$ ,

so  $C = 50$  suggests  $\delta = \max\left(0, \frac{\lambda_1 - 50\lambda_p}{49}\right)$ .

16) For  $R(\delta, \tau)$  compute  $I_p$  and

$R(\delta, 0)$ ,  $R(\delta, 0.05)$  for  $C = 50, 100, 200, 300, 400, \text{ and } 500$ .  
compute  $R$  if  $R > 0$ .

As  $\delta \rightarrow \infty$ ,  $R_\delta \rightarrow I_p$  and  $I_p$  corresponds to  $C = 1$ .



So  $S(\infty, 0) = \text{diag}(s_{ii}) = S_d$ . See 5.

SL 44

17) To get a regularized analog of  $R^{-1}$ , let

$$R_d = \begin{cases} R & \text{cond}(R) \leq 500 \\ R(s=0.0) & \text{else} \end{cases}$$

Let  $A = R_d^{-1}$  the analog of  $R^{-1}$  to be regularized.

$A$  acts like a covariance matrix. Let

$D_A = \text{diag}(\sqrt{A_{ii}})$ . Then a generalized correlation matrix  $R_I = D_A^{-1} A D_A$  is regularized with

$$R_{I,d} = R_I(b, 0) \text{ and } R_{I,t} = R_I(s, \tau). \text{ Then}$$

the regularized analogs of the inverse correlation matrix are  $R_{INV,d} = D_A R_{I,d} D_A$  and

$$R_{INV,t} = D_A R_{I,t} D_A.$$

$$\text{ex) } R = \begin{bmatrix} 1 & .4 \\ .4 & 1 \end{bmatrix}, \quad R_{s=1} = \begin{bmatrix} 1 & .2 \\ .2 & 1 \end{bmatrix}$$

$1+s=2$

$$R[1, .2] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad R[1, .3] = R_{s=1}$$

correlation  $\Delta$

0 ch 7 5 10.3, 10.4 cluster analysis

1) Cluster Analysis is an example of unsupervised learning; there are no labelled responses (groups for classification).

2) Discriminant analysis is called supervised classification while clustering is called unsupervised classification; we try to find groups but the groups are unknown both for the training data and the test data. Sometimes there is no test data.

3) K means clustering;

i) partition the  $n$  training data cases into  $k$  initial groups and find the mean of each group. Alternatively choose  $k$  initial seed points. These are groups of size 1 so the mean is equal to the seed point.

ii) Compute distances between each training data case and each mean. Assign each case to the cluster whose mean is nearest.

iii) Recalculate the mean of each cluster.

iv) Go to ii) and repeat until no changes occur.

4) 2 problems i) There could be more or less than  $k$  clusters. SL 45

ii) 2 initial seed points could actually belong to the same cluster.

(Could use several random sets of  $k$  seeds for several values of  $k$ , but need a method to choose the final group of clusters and  $k$ .)

5) Hierarchical clustering needs a distance.

Single linkage is the minimum distance between cases in cluster  $i$  and cases in cluster  $j$ .

Complete linkage is the maximum distance

between cases in cluster  $i$  and cases in cluster  $j$ .  
Sometimes the average distance is used.

i) Start with  $m = n$  clusters so each case forms a cluster. Compute the distance matrix for the  $n$  clusters. Let  $d_{u,v}$  be the smallest distance. Combine clusters  $u$  and  $v$  into a single cluster and set  $m = n - 1$ .  
Repeat i) with the new  $m$ . Continue

until there is a single cluster.

9ii) Plot the results in a dendrogram. Use the dendrogram to select  $k$  reasonable clusters of cases.

ex) (SW p 558-561) Consider the numbers 1-10 in 11 languages.

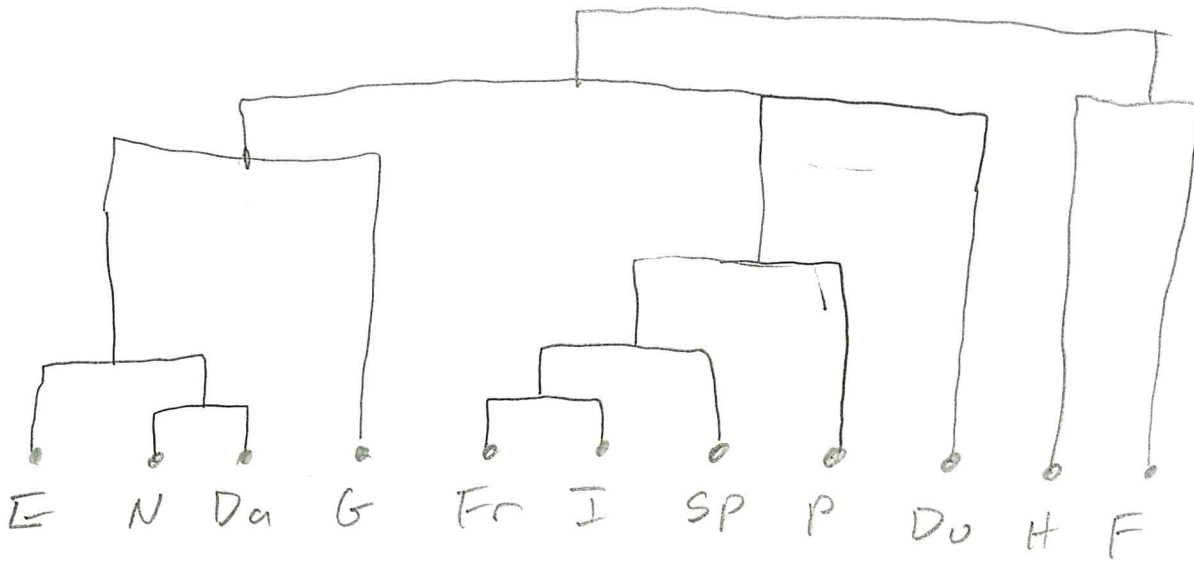
one en en een ein ün ünö ünö jeden egy ytsi e

Look at 1st letter of each word. Two different languages are concordant if they have the same 1st letter.

concordant 1st letters

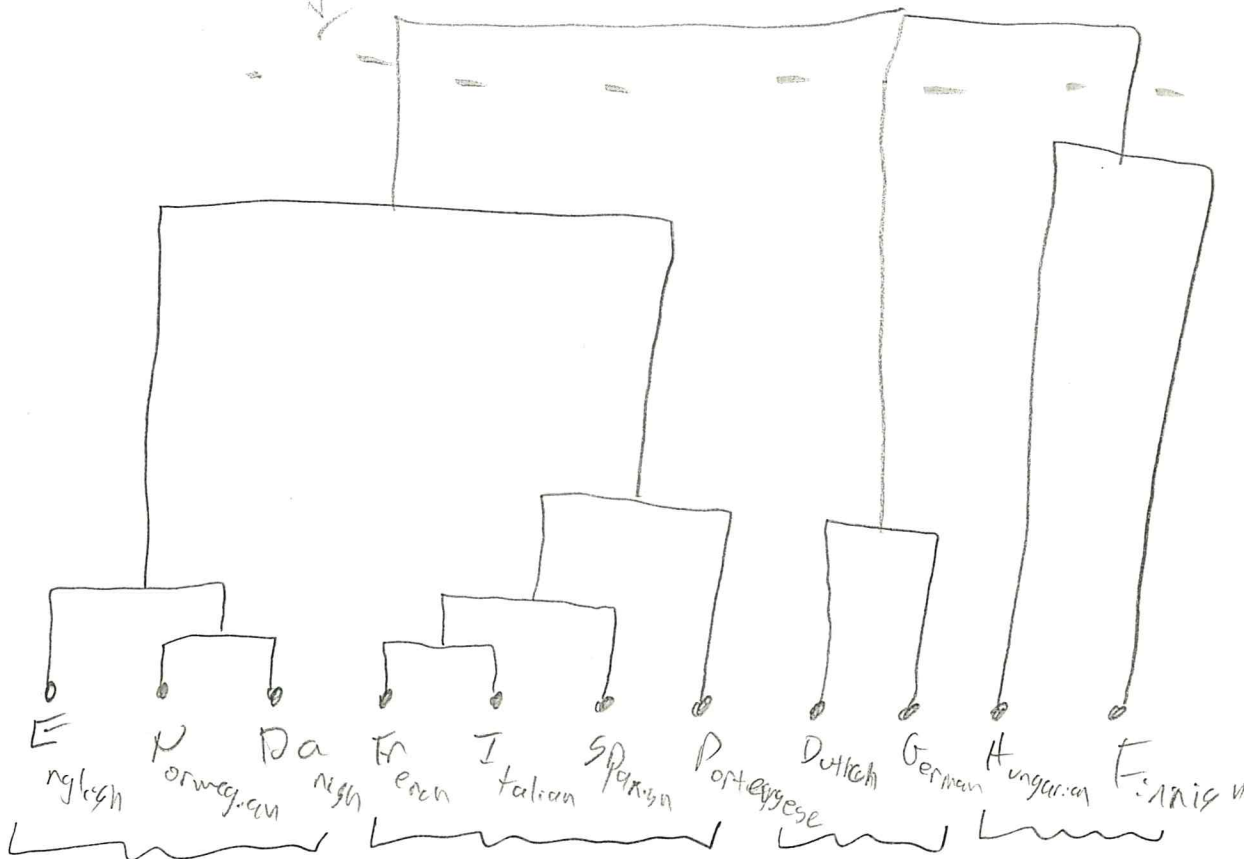
	E	N	Da	Dv	G	F	Sp	I	P	H	Fi
E	10										
N	8	10									
Da	8	9	10								
Dv											
F											
Sp											
I											
P											
H											
Fi	1	1	1					1	2	10	

10 - concordant matrix  
= distance matrix



complete linkage

3 clusters



single linkage

vertical height represents closeness of clusters

could cluster with a horizontal line

could use

4 clusters

could use 3 clusters.

6) Dendrograms are hard to use with large data sets. Need an automated method to choose the final  $k$  clusters for large data sets.

→ often take training data with known groups to check the cluster analysis

J Ch 8

Regression and Classification Trees

1) A regression tree is a flexible method for  $Y = m(X)$

or for  $Y_i = m(x_i) + \sigma_i e_i$ . A classification tree is

a flexible method for classification. Both methods produce graphs called trees

that look like dendrograms. Each

branch has a label like  $X_i > 7.56$

$X_i$  quantitative.

or  $X_j = a, c, d$

$X_j$  a factor taking on levels  $a, b, c, d$ , etc say

unless told otherwise,

If the condition is true, go to the left

of the branch. Otherwise go to the

right of the branch. <sup>Some software suggests this. could group problem.</sup> The bottom of the

tree has leaves that give  $\hat{Y}_i = \hat{Y}(x_i)$ .

$\hat{Y}$  is a number for a regression tree and

$\hat{Y}$  is a label for the classification 96 47 group for a classification tree. For

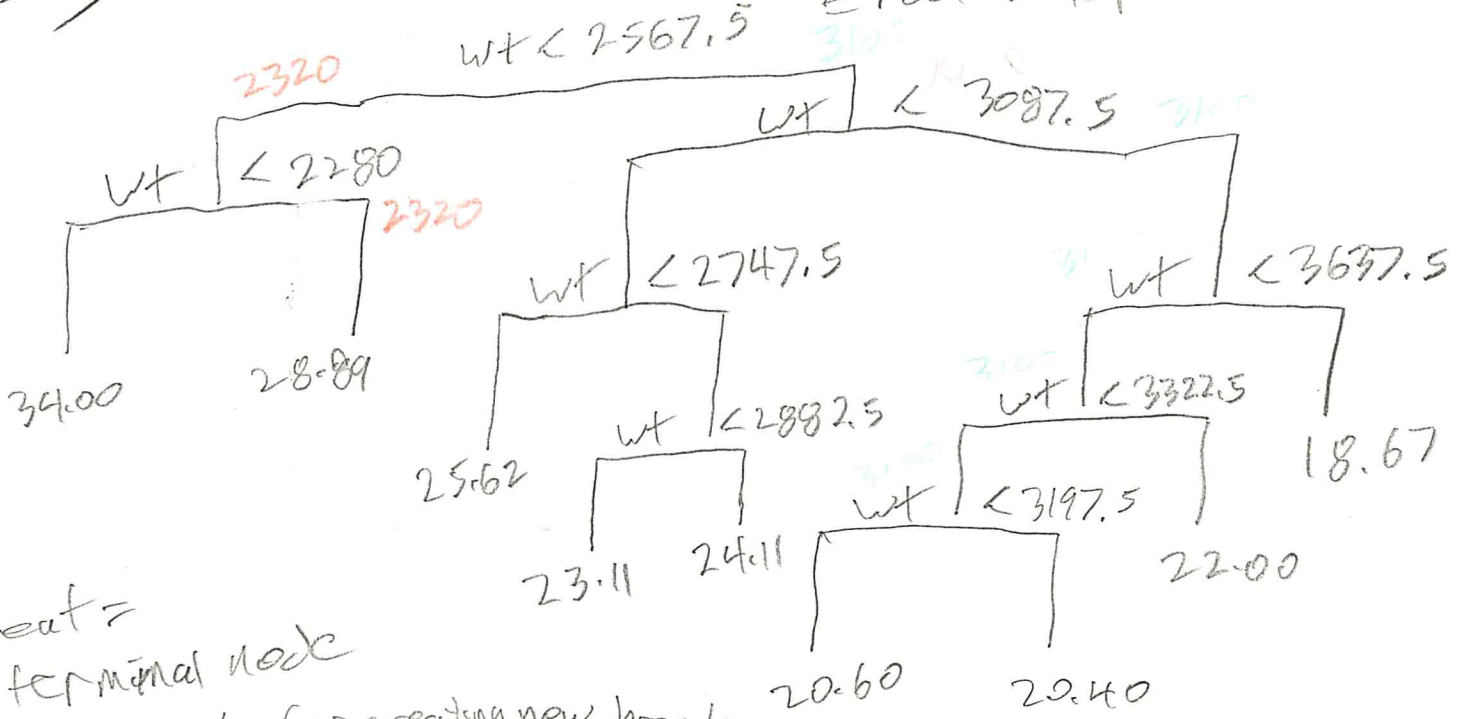
example if  $Y = \begin{cases} F = 0 \\ M = 1 \end{cases}$ , we could have

$\hat{Y} \in \{0, 1\}$  or  $\hat{Y} \in \{F, M\}$ .

2) know (Spivey 2000, Guide to Statistics, Vol I 372-374)

Given a tree and  $X$  values, find  $\hat{Y}$ . see HW10 c), d).

ex)  $Y = \text{gas mileage}$        $X = \text{weight of car}$   
 $\hookrightarrow \text{root} = \text{top node}$



leaf = terminal node

split = rule for creating new branches

a) Predict mileage if  $wt = 2320$   
 $\hat{Y} = 28.89$

b) Predict mileage if  $wt = 3100$        $\hat{Y} = 20.60$

ex] (Bert 2008 p145-147) Prison data

$Y =$  no misconduct, minor misconduct, or major misconduct  
 0 - 1 - 2

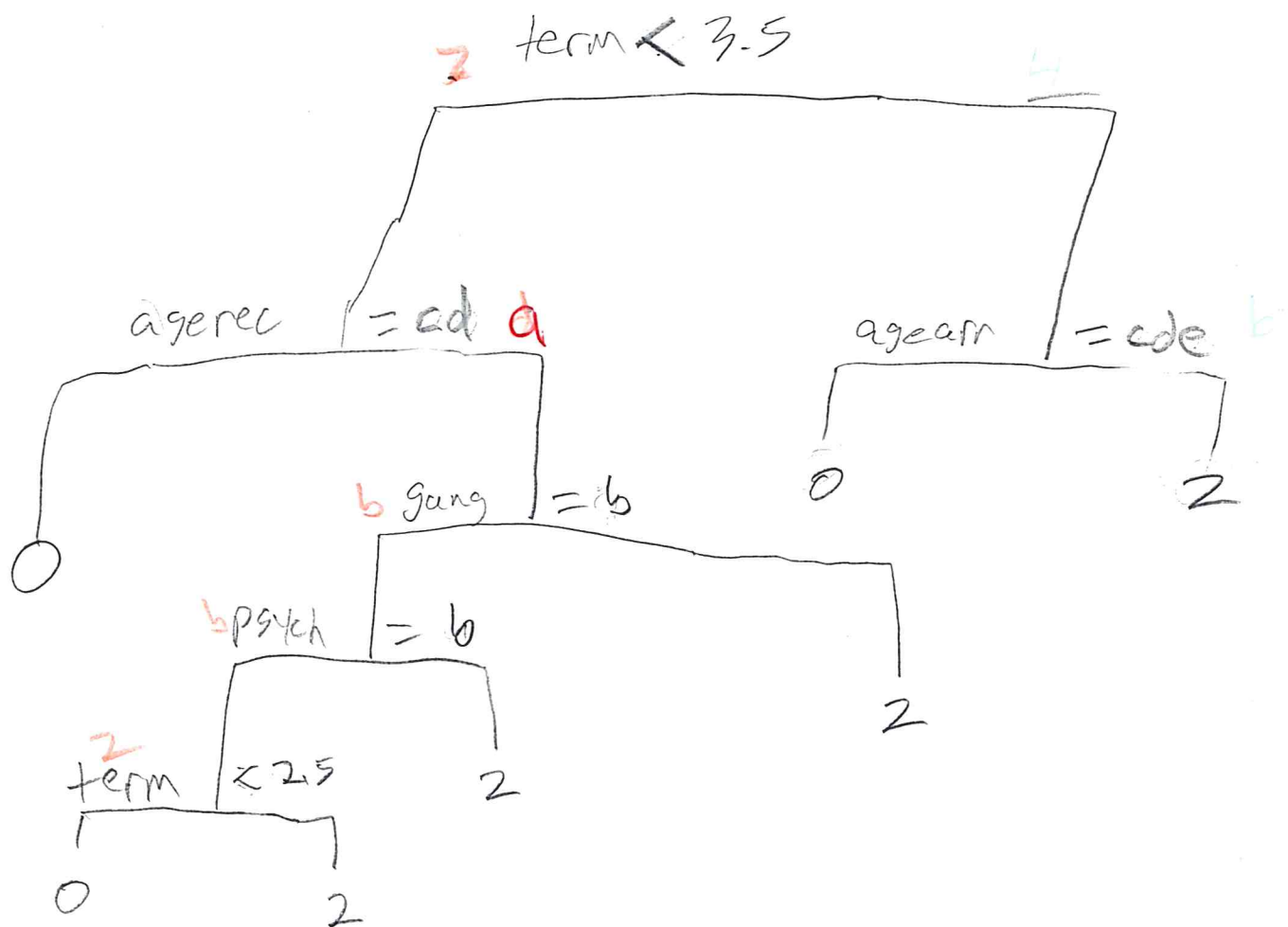
$X_1 =$  term  $\hat{=}$  sentence length in years

$X_2 =$  <sup>(enter prison)</sup> agerec: a: 16-20, b: 21-26, c: 27-35, d:  $\geq 36$

$X_3 =$  <sup>(1st arrest)</sup> agearr: a: 0-17, b: 18-21, c: 22-29, d: 30-35, e:  $\geq 36$

$X_4 =$  gang: a = gang activity, b = no gang activity

$X_5 =$  psych: a = mental illness, b = no mental illness



a) classify if  $X_1 = \text{term} = 2$   $X_2 = \text{agerec} = d$

$X_4 = \text{gang} = b$ ,  $X_5 = \text{psych} = b$

$Y = 0$

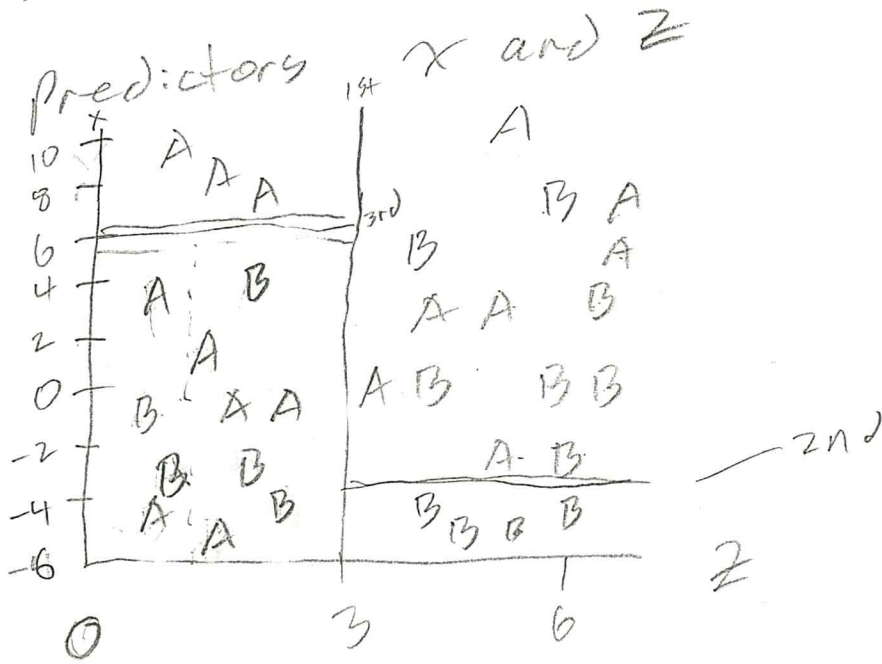
red

b) classify if  $X_1 = \text{term} = 4$ ,  $X_3 = \text{agearr} = b$



3) Trees that use recursive partitioning  
 for classification and regression trees use  
 the CART algorithm

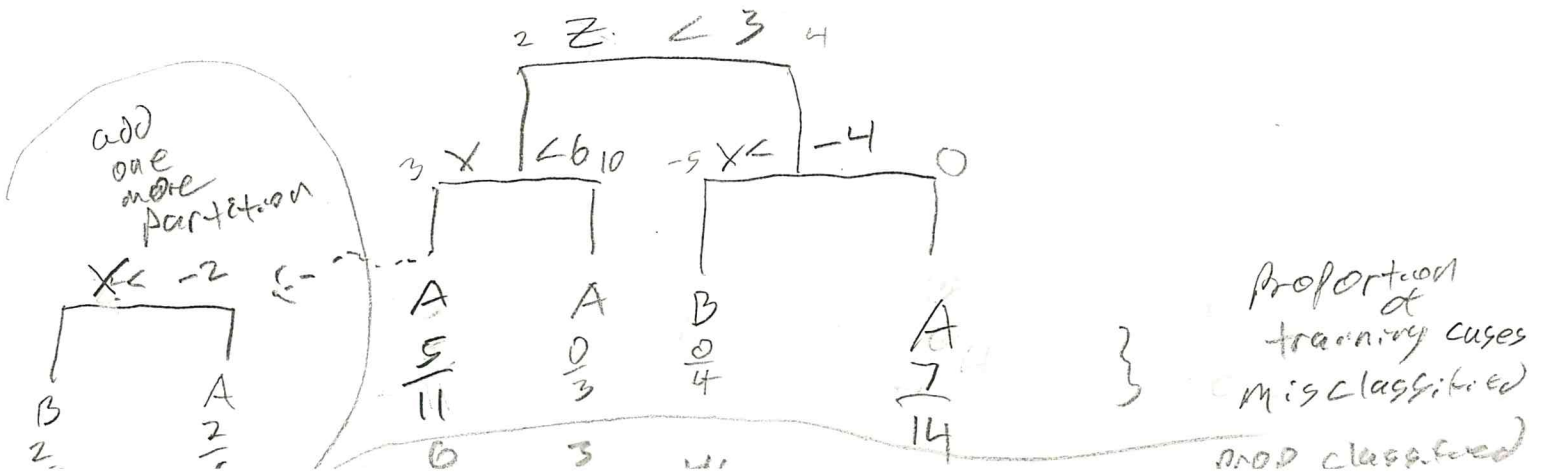
ex)  $Y = A \text{ or } B$  (classification) with



A single vertical line at  $Z=3$  is the 1st partition.

The horizontal line at  $X=6$  is the 2nd partition.

"  $X=-4$  3rd "



4} Trees i) give prediction rules that can be rapidly and repeatedly evaluated,

ii) are useful for screening predictors (variable selection, interactions)

iii) can be used to assess the adequacy of linear models

iv) can be used to summarize large multivariate data sets

5} If  $Y$  is a factor, <sup>has levels 1, ..., K</sup> classification rules are of the form if  $x_1 \leq 2.3$  and  $x_3 \in \{A, B\}$

then  $Y$  is most likely to be in level 5.

If  $Y$  is numerical, regression rules are of the form if  $x_2 \leq 2.3$  and  $x_9 \in \{C, D, F\}$  and  $x_5 \leq 3.5$

then  $\hat{Y} = 4.75$ .

6} Trees can be easier to interpret when some predictors are numerical and some categorical. Trees are invariant to monotone (increasing or decreasing) transformations of the predictor variable  $X_i$ .


Trees handle missing values better SL 49

than MLR. Trees beat MLR if there is nonadditive behavior. Trees can handle complex unknown interactions. For classification  $Y$  can have more than 2 levels.

7) In growing a tree, the binary partitioning algorithm recursively splits the data in each node until either the node is homogeneous (0 <sup>roughly</sup> training data misclassifications for a classification tree,  $Y \approx \text{constant}$  for a regression tree) or the node contains too few observations (default  $\leq 5$ ).

8) The deviance is a measure of node homogeneity and deviance = 0 for a perfectly homogeneous node.

9) Often use the mean of the region for  $\hat{Y}$ , regression tree  
mode classification tree  
mode = B



10) 5 p 306-7 a) Divide the predictor space  
 (= set of possible values for  $x_1, \dots, x_p$ ) into  
 training data

$J$  distinct and nonoverlapping regions  $R_1, \dots, R_J$ .

b) For every observation that falls in  
 region  $R_j$ , make the same prediction

$\hat{y}_{R_j} \equiv \begin{cases} \bar{y}_j & = \text{sample mean of training data } Y \text{ in } R_j, \text{ reg. tree} \\ \text{mode}_j & = \text{class tree} \end{cases}$

11) The  $R_j$  are high dimensional boxes.

Choose  $R_j$  so  $RSS = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$

is small.

$I(y_i \neq \hat{y}_{R_j})$  class

12) Let  $\{\bar{x} \mid x_j \leq s\}$  be the region in the predictor

space such that  $x_j \leq s$  where  $\bar{x} = (x_1, \dots, x_p)^T$ .

Define  $2$  ~~half hyperplanes~~ regions

$R_1(s) = \{\bar{x} \mid x_j \leq s\}$ ,  $R_2(s) = \{\bar{x} \mid x_j \geq s\}$  and seek

... "cutpoint"  $s$  and  $j$  to minimize

$$\sum_{i: \bar{x} \in R_1(s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: \bar{x} \in R_2(s)} (y_i - \hat{y}_{R_2})^2$$

This can be done "quickly" if  $p$  is small, SL 50 (could use order statistics).

Then repeat the process looking for the best predictor and best cutpoint in order to split the data further so as to minimize the RSS within each of the resulting regions.

Only split one of the regions. So now there are regions  $R_1, R_2$ , and  $R_3$ . Continue this process until a stopping criterion is reached such as no region contains more than 5 obs's (and stop splitting if the region is homogeneous eg all  $Y$  in region belong to class  $k$ ).

13) (For classification) such trees are usually not competitive with earlier techniques (Bagging, random forests and boosting makes trees more competitive).

14) Trees use regions  $R_1, \dots, R_J$ . If  $J$  is too large, the tree overfits. One strategy is to grow a large tree  $T_0$  with  $J_0$  regions

15) J p 308 Cost complexity pruning = weakest link pruning

Let  $T \subseteq T_0$ ,  $\alpha \geq 0$ , and  $|T| = \#$  terminal nodes of tree  $T$ . Each terminal node corresponds to a region (hyper rectangle)  $R_i$ . Let  $R_m$  be the region corresponding to the  $m$ th terminal node and  $\hat{y}_{R_m}$  be the predicted response for  $R_m$ .

For each value of  $\alpha > 0$ , there corresponds a subtree  $T \subseteq T_0$  such that  $\sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$  is as

\* 
$$\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$
 is as

small as possible. Note that  $\alpha = 0$  has  $T = T_0$

and (\*) =  $RSS(T_0)$  = training data  $RSS$  for  $T_0$ .

16) Much like lasso, as  $\alpha$  increases there is a sequence of nested subtrees

\*\*  $T_0 \supseteq T_{\alpha_1} \supseteq T_{\alpha_2} \supseteq \dots \supseteq T_{\alpha_m}$ . Branches get "pruned" from  $T_0$  in a nested and predictable fashion,

17) a) Build tree  $T_0$ , stopping when each terminal node has  $\leq 5$  obs.

c) Use  $K$ -fold CV to choose  $\alpha = \alpha_{j^*}$ . SL 51

For each  $i \in \{1, \dots, K\}$

i) Repeat steps a) and b) on all but the  $i$ th fold.

ii) Evaluate the mean square prediction error on data in the left out fold  $i$  as a function of  $\alpha$ .

Average the results for each value of  $\alpha$  and pick  $\alpha_{j^*}$  to minimize the average error.

d) Use tree  $T_{\alpha_{j^*}}$  from b),

For a regression tree

$$MSE_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_j - \hat{y}_j(i))^2 \quad \text{for data } y_j \text{ in}$$

left out fold  $i$ .

$$CV_{(K)} = \frac{1}{K} \sum_{i=1}^K MSE_i$$

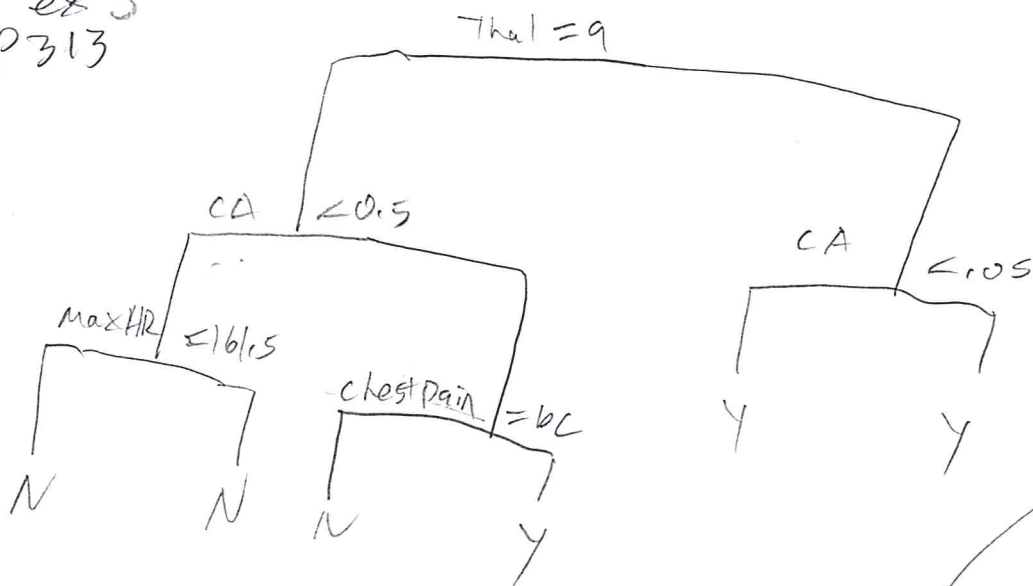
For a classification tree use

$$MSE_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbb{I}(y_j \neq \hat{y}_j(i)) = \text{proportion}$$

misclassified in the  $i$ th fold.

$$\text{If } n_i \equiv \frac{n}{K} \quad \text{then } CV_{(K)} = \begin{cases} \frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j(i))^2 & \text{reg tree} \\ \frac{1}{n} \sum_{j=1}^n \mathbb{I}(y_j \neq \hat{y}_j(i)) & \text{class tree} \end{cases}$$

ex 5  
p 313



303 patients with chest pain

Y = heart disease

N = no heart disease

← subtree  
6 terminal nodes  
Big tree to added many  
more branches and had  
18 terminal nodes.

Thal = Thallium stress test a = normal

Chest pain a = typical agonal b = atypical agonal

c = non-agonal pain d = asymptomatic

maxHR maximum heart rate? CA calcium in blood?

18) For a tree  $T_\alpha$ ,  $y_i = m(x_i) + \sigma_i e_i$  and

$$\hat{y} = \hat{m}(x_i) = \sum_{m=1}^{J_\alpha} c_m I(x_i \in R_m) \quad \text{where } T_\alpha$$

uses regions  $R_1, \dots, R_{J_\alpha}$ .

19) Trees can handle categorical variables (factors, qualitative variables) without creating indicators = dummy variables.



§ 8.2

9652

20] Bagging was used before.

compute  $T_1^*, \dots, T_B^*$  with the bootstrap and the sample mean  $\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^*$  is the bagging estimator.

20] For a regression tree, draw a sample of

size  $n$  with replacement from  $x_1, \dots, x_n$ .

Fit the tree and find  $\hat{f}_1^*(x)$ . Repeat to get

$\hat{f}_1^*(x), \dots, \hat{f}_B^*(x)$ . Then the bagging estimator

$$\hat{f}_{\text{bag}}^*(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i^*(x).$$

The trees are

not pruned so terminate when each terminal node has 5 or fewer obs's.

22] For classification draw a sample of size

$n_i$  from each group with replacement.

Let  $\hat{f}_i^*(x) = \hat{j}_i(x) \in \{1, \dots, G\}$  where  $\gamma$  takes on

levels  $1, \dots, G$ . Compute  $\hat{f}_1^*(x), \dots, \hat{f}_B^*(x)$  and

let  $M_k = \# \hat{j}_i(x) = k$  for  $k=1, \dots, G$ . Take

23) <sup>5 p197</sup> The prob that  $x_j$  is not in the bootstrap sample  $\rightarrow e^{-1} \approx .3679 \approx \frac{1}{3}$  as  $n \rightarrow \infty$ .

$\frac{1-\frac{1}{n}}{1} \dots \frac{1-\frac{1}{n}}{n}$  prob  $x_j$  is not the  $k$ th obs  $= 1-\frac{1}{n}$  since each obs has  $\frac{1}{n}$  chance of being selected for the  $k$ th position.

$(1-\frac{1}{n})^n \rightarrow e^{-1}$

24) For each bootstrap sample  $b$

let  $x_{i1}, \dots, x_{ik_b}$  be the  $k_b$  obs's not used. These are the out of bag (OOB) obs's, predict  $\hat{Y}$  for each OOB obs. Doing this for all  $B$  bootstraps produces about  $\frac{B}{3}$  predictions for each  $x_i$ . Let  $\hat{Y}_{io} = \begin{cases} \text{ave } \hat{Y}_i & \text{reg tree} \\ \text{mode level} & \text{class tree} \end{cases}$

OOB MSE =  $\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_{io})^2$  reg

$\frac{1}{n} \sum_{i=1}^n I(Y_i \neq \hat{Y}_{io})$  class

The OOB MSE is "virtually equivalent" to the leave one out CV estimate, for sufficiently large  $B$ .

25) Bagging typically gives better accuracy

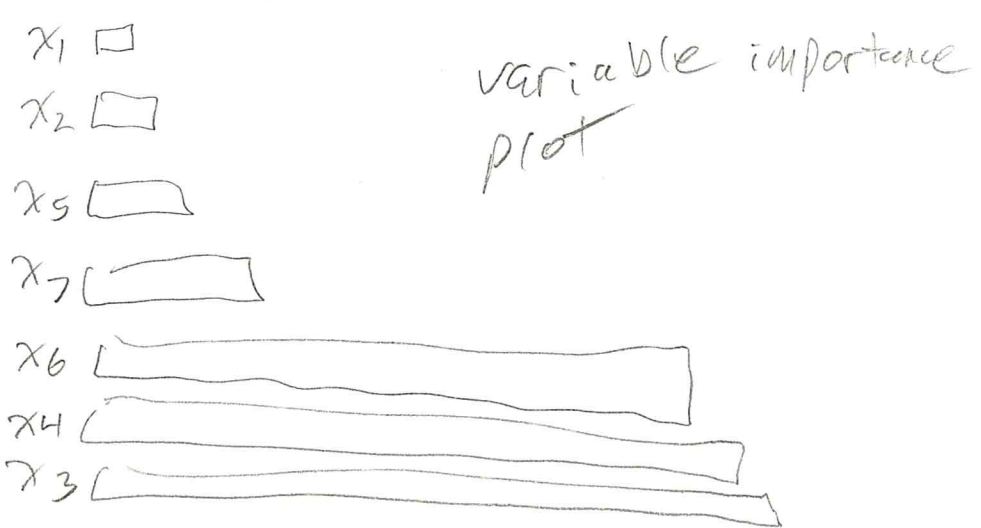
26) For classification trees, let

$\hat{P}_{mk}$  = proportion of training obs's in  $R_m$  from the  $k$ th class.

Gini's index =  $\sum_{k=1}^C \hat{P}_{mk} (1 - \hat{P}_{mk})$

is small if all  $\hat{P}_{mk}$  are close to 0 or 1.

27) <sup>p319</sup> For bagging trees with B trees a measure of variable importance can be computed for each variable using splits for each variable.



$x_3, x_4, x_6$  are the most important predictors for the tree

28) If  $Y = \alpha + \sum_{j=1}^p \beta_j(x_j) + \epsilon$  or  $Y = m(x) + \epsilon$  with

$m(x) = g(\alpha + \beta^T x)$  then slicing the  $\epsilon$

$\epsilon = \alpha + \beta^T x + \epsilon$

is more effective than partitioning the predictor space with hyper boxes  $R_k$ .

$$Y = g(\beta^T X) + \epsilon$$

Curse of Dimensionality  
IRP



If  $Y = g(w)$   
and you plot  
 $w$  vs  $Y$ , then  
you'll see  $g$ !  
So if  $Y = g(\beta^T X) + \epsilon$   
you'll see  $g$  up to  
noise of  $ESP \approx SP = w$   
are  $\epsilon_i$ .

$$ESP = \hat{\sigma} + \frac{\hat{\sigma}^2}{\beta^T X}$$

MR idea, Predict

29) If  $Y = m(x) + \epsilon$ , we can make prediction intervals for  $Y$  with the regression tree using  $\hat{Y} = ESP = \hat{m}(x)$  and  $\Gamma = Y - \hat{Y}$  as before.

### §8.2.2 Random Forests

30) For random forests the bootstrap is used, but each time a split is considered, a random sample of  $m \approx \sqrt{p}$  predictors is chosen as split candidates. Random forests produces bootstrap trees that are less correlated than the bagged trees (that use  $m=p$ ) and the random forests estimator tends to have a better test error and MSE error.

than the bagging estimator,

9L 54

31) B around a few hundred seeds to work,

32) If there is a single strong predictor bagged trees tend to use that predictor in the 1st split. For random forests, the strong predictor is not considered for  $\frac{p-m}{p}$  splits, on average. Trees from random forests also tend to be less correlated than bagged trees if there are many correlated predictors,

JP321

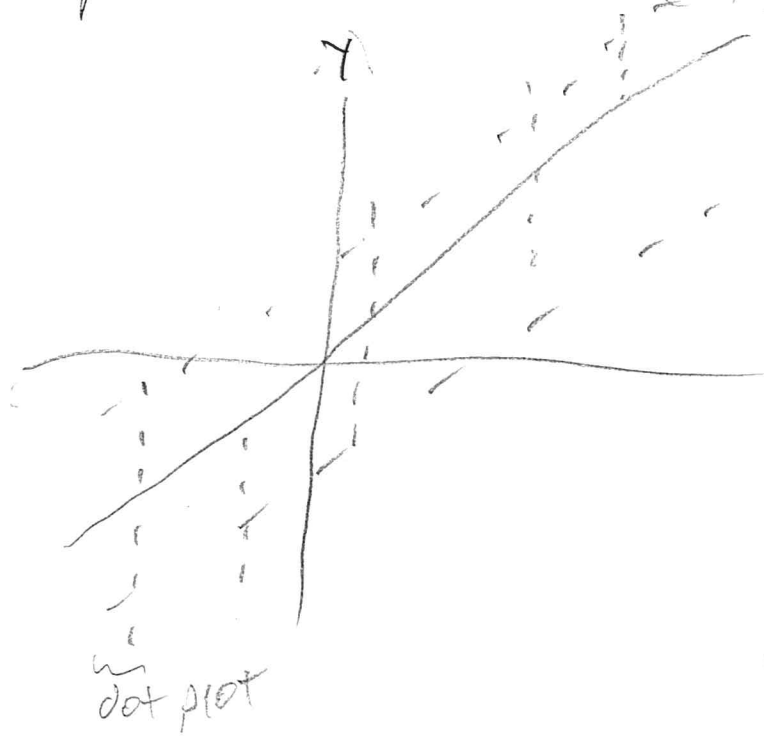
33) For classification, the null classifier has  $\hat{y} \equiv d$  where  $d$  is the dominant class.

so if  $k\%$  of obs's belong to the dominant class, then the test error =  $\frac{100-k}{100} \leq 1 - \frac{1}{G}$

where there are  $G$  groups. Since  $k \geq \frac{100}{G}$ .

Classifiers that do not beat the null classifier are very bad.

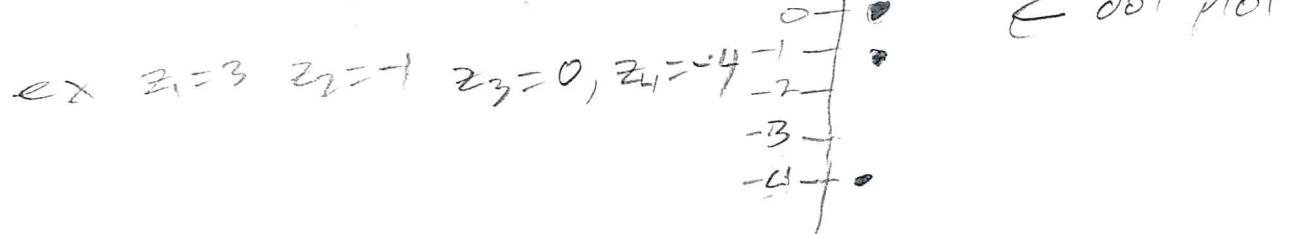
34) Since a tree uses  $J_d$  regions, the response plot of  $ESP = \hat{Y} = \hat{m}(x)$  vs  $Y$  looks like



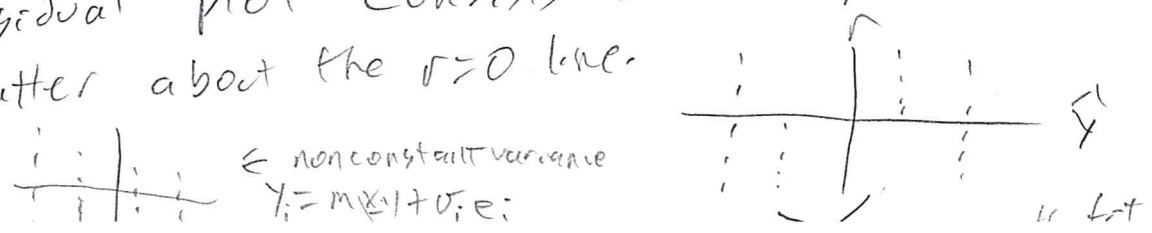
Pointwise  
PI Bands  
The plotted points scatter about the identity line, but there is a dot plot of  $n_m$  cases with  $\hat{Y} = \hat{Y}_{RM}$  for each of the  $J_d$  regions

(One way Anova models have a similar dot plot, but each dot plot crosses the identity line at  $\hat{Y}_i = \bar{Y}_i$ .)

35) A dot plot of  $Z_1, \dots, Z_m$  consists of an axis and  $m$  points corresponding to the values of  $Z_i$



36) The residual plot consists of dot plots that scatter about the  $r=0$  line.



37) 308.2.3 The third technique for improving trees is boosting. Like bagging, boosting can be applied to many statistical methods, not just trees. SL 55

38) Boosting for regression trees

first tree  $\hat{f}_1$  is fit on the response

i) set  $\hat{f}(x) = 0$  and  $r_i = y_i$ ,  $i = 1, \dots, N$   
training data

ii) For  $b = 1, \dots, B$

a) fit tree  $\hat{f}_b$  with  $d$  splits ( $d+1$  terminal nodes) to the training data  $(\mathcal{X}, r)$

b) update  $\hat{f}(x)$  by adding a shrinked version of the new tree:  $\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}_b(x)$   
and update the residuals  $r_i \leftarrow r_i - \lambda \hat{f}_b(x_i)$ .

iii) The boosted model

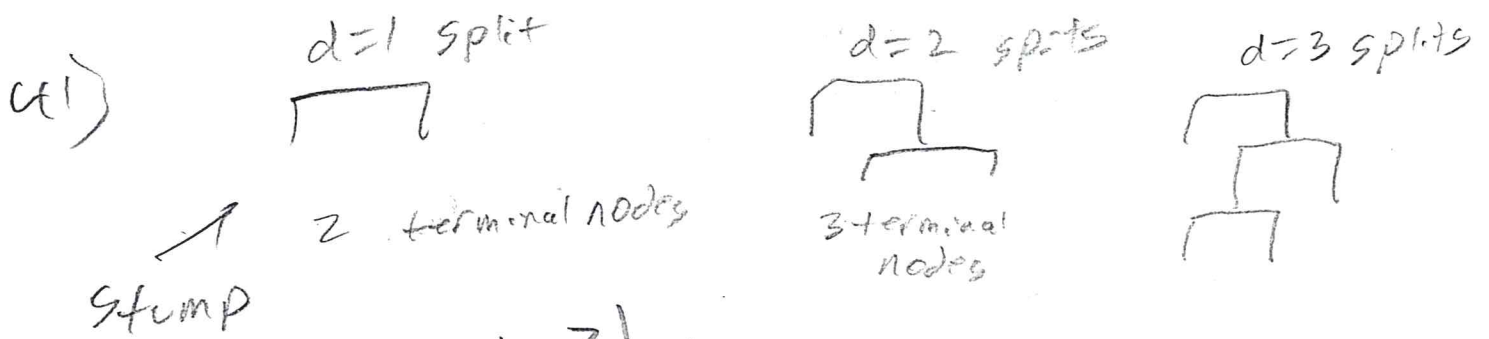
$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}_b(x).$$

39) The tree is fit to updated residuals rather than  $y$ . This slowly improves  $\hat{f}$  in areas where it does not perform well, and  $\lambda$  slows the learning process further.

As a rule of thumb, iterative techniques that learn slowly tend to perform well.

40) Boosting can overfit if  $B$  is too large.

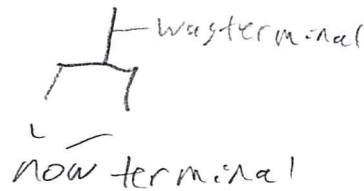
$k$  fold CV is used to select  $B$ .



A tree with  $d-1$  splits has  $d$  terminal nodes

Then a tree with  $d$  splits adds the split to one of the nodes

replacing the terminal node with 2 terminal nodes.



Hence a tree with  $d$  splits has  $d+1$  terminal nodes.

42) often  $d=1$  or  $2$  is used. A  $d=1$  tree is called a stump. The value  $d$  is called the interaction depth.

43) The value  $\lambda$  tends to be  $\infty$  or  $0.01$ .

very small  $\lambda$  tends to need very large  $B$  for good performance.



44] 5 p 332 Using the  $d=1$  stumps leads to an additive model

SL 56

$$\hat{f}(x) = \sum_{j=1}^P \hat{f}_j(x_j)$$

I can't prove this yet.

Each stump splits one variable

Let  $\underline{r}^{(0)} = \underline{y}$ ,  $\underline{r}^{(b)} = \underline{r}^{(b-1)} + \lambda \begin{pmatrix} \hat{f}_b(x_1) \\ \vdots \\ \hat{f}_b(x_n) \end{pmatrix}$  with  $\underline{r}^{(0)} = \underline{0}$

$$\underline{r}^{(b)} = \underline{r}^{(b-1)} - \lambda \begin{pmatrix} \hat{f}_b(x_1) \\ \vdots \\ \hat{f}_b(x_n) \end{pmatrix}$$

45] 5 p 332 For a binary classification tree with  $Y = 0$  or  $1$ , for a fixed value of  $x$  (a test data point), the bootstrap produces  $B$  estimates of  $P(Y=1|x)$ .

Two common ways to get  $\hat{Y}|x$  are

a)  $\hat{Y}|x =$  mode class of 0 or 1

b) average the prob  $\hat{Y}|x = \begin{cases} 0 & \text{if } \text{ave} \hat{P}(Y=1|x) < 0.5 \\ 1 & > \end{cases}$

ex  $\hat{P}(Y=1|x) \stackrel{B=10}{=} \begin{matrix} 0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, 0.75 \\ \text{class} \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \end{matrix}$

a) gives  $\hat{Y}|x = 1$   $\dots (11) \dots 151 - 4.5 - .45 \text{ so } \hat{Y}|x = 1$

# A note on the residual bootstrap for regression

Suppose  $\underline{Y} = X\underline{\beta} + \underline{e}$

$e_i$  iid with  $V(e_i) = \sigma^2$

$\text{cov}(\underline{e}) = \sigma^2 \underline{I}_n = \text{cov}(\underline{Y})$

$\sqrt{n} (\underline{\hat{\beta}} - \underline{\beta}) \xrightarrow{D} N_p(\underline{0}, \sigma^2 \underline{W})$

(Efron 1982 p36)

The residual bootstrap has

$\underline{Y}^* = X\underline{\hat{\beta}} + \underline{r}^w$

is a general linear model

with respect to the bootstrap sample since

the  $r_i^w$  are iid from the distribution of  $\underline{e}$

which is the empirical distribution of the

residuals  $\frac{r_1}{n} \dots \frac{r_n}{n}$

so  $E(r_i^w) = \frac{1}{n} \sum_{i=1}^n r_i = \bar{r} = 0$  since  $\sum_{i=1}^n r_i = 0$  for OLS residuals.

so  $V(r_i^w) = E(r_i^w)^2 = \frac{1}{n} \sum_{i=1}^n r_i^2 = \frac{n-p}{n} \frac{1}{n-p} \sum_{i=1}^n r_i^2 = \frac{n-p}{n} \text{MSE}$

and  $E(\underline{r}^w) = \underline{0}$ ,  $\text{cov}(\underline{Y}^*) = \text{cov}(\underline{r}^w) = \frac{n-p}{n} \text{MSE} \underline{I}_n$

$E(\underline{Y}^*) = X\underline{\hat{\beta}} = H\underline{y}$

so by standard linear model theory with  $\frac{X^T X}{n} \rightarrow \underline{W}^{-1}$

$\sqrt{n} (\underline{\hat{\beta}}^* - \underline{\hat{\beta}}) \xrightarrow{D} N_p(\underline{0}, \sigma^2 \underline{W})$ ,  $\text{cov}(\underline{\hat{\beta}}^*) = \frac{n-p}{n} \text{MSE} (X^T X)^{-1}$

Here  $\underline{\hat{\beta}}^* = (X^T X)^{-1} X^T \underline{y}^*$ ,  $\underline{\hat{\beta}}_{I_j}^* = (X_{I_j}^T X_{I_j})^{-1} X_{I_j}^T \underline{y}^*$


# JCH9 Support Vector Machines (SVMs) SL 57

1) Logistic regression is used a lot in biostatistics and epidemiology where the focus is statistical inference. Support vector machines are used in machine learning where the goal is classification accuracy.

2) When  $p \gg n$  there is often a hyperplane that perfectly separates 2 groups. The learning point for SVMs was finding the optimal separating hyperplane.

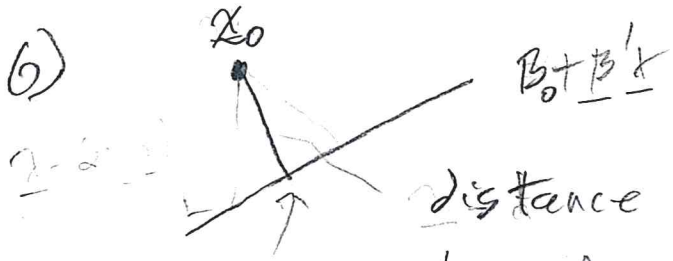
3) For 2 groups let  $SP = \beta_0 + \beta'x$   
 classify  $x$  in group 1 if  $ESP > 0$  and  
 group -1 if  $ESP < 0$  (just like LR, but code group 0 as -1).

so the classifier  $\hat{C}(x) = \text{sign}(ESP)$

4) 
 The estimated "optimal separating hyperplane"  $ESP$  has the largest margin on the training data.

5) Let  $ESP = \beta_0 + \beta'x$  and  $f(x) = SP = \beta_0 + \beta'x$

so  $\hat{f}(x) = ESP$ ,  $y_i \in \{-1, 1\}$ .



distance of  $x_0$  from the decision boundary

$$= \frac{\beta_0 + \beta'x_0}{\|\beta\|_2} = \frac{f(x_0)}{\|\beta\|_2}$$

Projected of  $x_0$  onto hyperplane must be  $\beta_0 + \beta'x_0$

7) Think of the hyperplane  $\beta_0 + \beta'x_i =$

$\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$  dividing  $\mathbb{R}^p$  into 2

halves. A separating hyperplane has  $SP > 0$  if  $x \in \text{group}$   
 $< 0$  if  $x \in \text{group}-1$

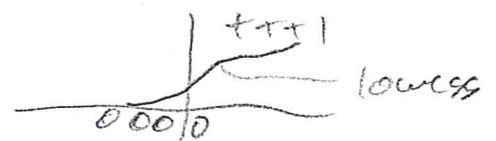
So  $y_i SP_i = y_i (\beta_0 + \beta'x_i) > 0$  for  $i=1, \dots, n$ .

8) Think of a binary classifier that uses  $ESP$  as a binary regression

$y|x \sim \text{bin}(n=1, s(x))$  where  $s(x) = s(ESP) = P(Y=1|x)$

where  $Y=-1$  has been recoded as  $Y=0$ , but

$s(ESP)$  is unknown, Response plot



\* The bootstrap with  $n_i$  cases selected with replacement from each group is likely useful if  $n \gg p$ . SL 58

9) CART splits  $\mathbb{R}^p$  with regions  $R_m \in \mathbb{R}^p$   
SVM splits  $\mathbb{R}^p$  into 2 regions using ESP  $E_R$  so there is dimension reduction. The SVM split tries to make the 2 "halves" or partitions as homogeneous as possible.

10) The hyperplanes parallel to ESP that form the boundaries of the margin are called fences. The fences pass through at least 2 training data cases. These cases form the support set  $S$  of support vectors. It turns out that

$$\hat{\beta}_M = \sum_{i \in S} \hat{\alpha}_i \underline{x}_i$$

↖ optimal margin classifier

11) Wide data has  $p \gg n$ . If  $n \leq |P|$

there is a separating hyperplane unless there are exact predictor trees across the class barrier (whatever that means).

12) Let  $M$  be the margin. The optimal margin classifier  $\hat{\beta}_{\text{opt}}, \hat{\beta}_M$  maximizes  $M$  subject

$$\text{to } \forall_i \quad y_i \cdot w_i = y_i (B_0 + B_1 x_{i1} + \dots + B_p x_{ip}) \geq M \quad \forall i=1, \dots, n. \quad (*)$$

This is called a hard margin classifier since no cases from either group can pass the fences of the classifier. Equivalently  $\min_{B_0, B} \|B\|_2$  subject to (\*).

13) A soft margin classifier allows cases from either group to pass the fences or to be misclassified. This classifier solves

$$\text{minimize } \|B\|_2 \quad \text{subject to}$$

$$B_0, B$$

$$y_i (B_0 + x_i^T B) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0 \quad \text{for } i=1, \dots, n$$

$$\text{and } \sum_{i=1}^n \epsilon_i \leq B.$$

slack variables are used in linear programming

14) This minimization is equivalent to

$$\text{minimize } \sum_{i=1}^n \left[ 1 - y_i (B_0 + x_i^T B) \right]_+ + \lambda \|B\|_2^2$$

$$B_0, B$$

$$\text{where } [w]_+ = \begin{cases} w & w \geq 0 \\ 0 & w \leq 0. \end{cases}$$

hinge loss =  $\begin{cases} 0 & \text{if } x_i \text{ is on the correct side} \\ \text{cost of } x_i \text{ being on the} \\ & \text{wrong side of the margin} \end{cases}$

This is similar to ridge regression.

15) A support vector machine that uses  $x_i$  minimizes the above loss criterion. For separable data

$$(\hat{\beta}_{\text{OSVM}}, \hat{\beta}_{\text{SVM}}) \rightarrow (\hat{\beta}_{\text{ORM}}, \hat{\beta}_{\text{M}}) \text{ as } \lambda \rightarrow 0.$$

16) It turns out that  $\hat{\beta}_{\text{SVM}} = \sum_{i \in S} \hat{\gamma}_i \underline{x}_i$

and  $\hat{\beta}_{\text{OSVM}} + \underline{x}^T \hat{\beta}_{\text{SVM}} = \hat{\beta}_{\text{OSVM}} + \sum_{i \in S} \hat{\gamma}_i \langle \underline{x}, \underline{x}_i \rangle$

where  $\langle \underline{x}, \underline{x}_i \rangle = \underline{x}^T \underline{x}_i$ . This quantity can be computed using the  $n \times n$  Gram matrix  $\sum_{n \times n} \underline{x} \underline{x}^T$  with  $O(n^2 p)$  complexity. Using  $\underline{x}^T \underline{x}$  has  $O(n p^2)$  complexity. Ridge regression could also be computed this way.

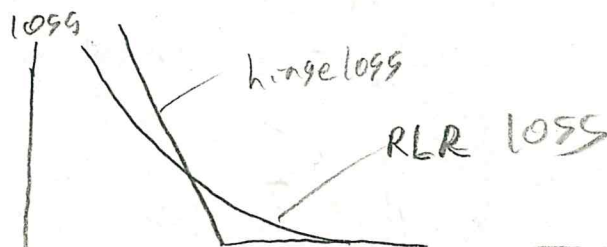
17) A lasso-SVM solves:

$$\text{Minimize}_{\beta_0, \beta} \sum_{i=1}^n \left[ 1 - \gamma_i (\beta_0 + \underline{x}_i^T \beta) \right]_+ + \lambda \|\beta\|_1$$

and does variable selection. (like lasso)

"Ridged" logistic regression" with  $\gamma_i \in \{-1, 1\}$

$$\text{Minimize}_{\beta_0, \beta} \sum_{i=1}^n \log \left[ 1 + e^{-\gamma_i (\beta_0 + \underline{x}_i^T \beta)} \right] + \lambda \|\beta\|_2^2$$



19) \* Truth table = confusion matrix

		truth	
predict	-1	7	1
	1	18	3

misclassified

$$\text{error rate} = 1 - \frac{18+7}{18+7+3}$$

$$= 1 - \frac{25}{28} = \frac{3}{28} = 0.1071$$

diagonal: correctly classified  
 off diagonal: incorrectly classified

		truth			
predict	a	b	c	d	
	a	10	0	0	12
b	0	100	11	0	
c	0	5	50	0	
d	6	0	0	30	

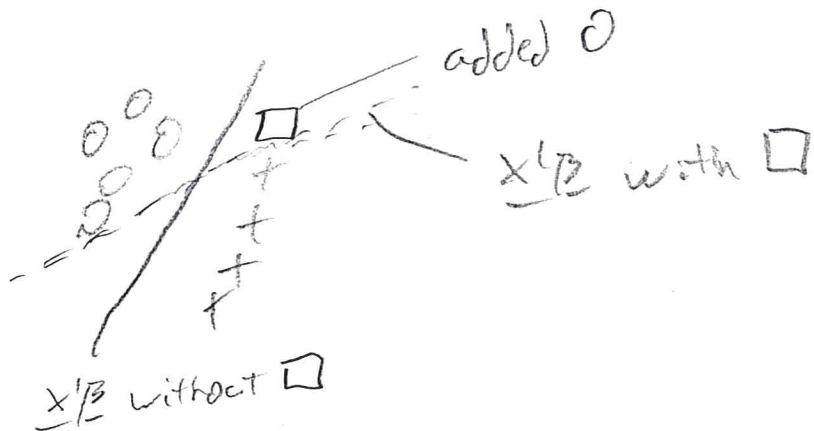
correct

$$\text{error rate} = \frac{6+5+11+12}{6+5+11+12+10+100+150+30}$$

$$= \frac{34}{224} = 0.1518$$

see HW 11

20) Sometimes 1 or a few obs's shift the maximal margin hyperplane. The SVM classifier is a soft margin classifier and can do better.



OM from 50

21) SVM maximizes pos. sep. given  $M = \text{width of margin}$  subj

$$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq B, \quad y_i (\beta_0 + \beta^T x_i) \geq M (1 - \epsilon_i) \quad \text{L.M.H.}$$



$$\hat{C}(x) = \text{sign}(\beta_0 + \beta'x)$$

SL 60

(I use  $G$  for # of classes)

21) A slack variable  $\xi_i = 0$  if  $\underline{x}_i$  is on the correct side of the margin. If  $\xi_i > 0$  then  $\underline{x}_i$  is on the wrong side of the hyperplane.

$y_i (\beta_0 + \beta'x_i) \geq M$  has  $\xi_i = 0$  and is necessary for  $\underline{x}_i$  to be on the correct side of the margin. If  $y_i (\beta_0 + \beta'x_i) \geq M(1 - \xi_i)$

with  $\xi_i > 0$  (but not if  $\xi_i = 0$ ), then

$\underline{x}_i$  is on the wrong side of the hyperplane. see 7

22) Let the kernel function be

$k(\underline{x}_i, \underline{x}_j)$ . A linear kernel

is  $k(\underline{x}_i, \underline{x}_j) = \underline{x}_i^T \underline{x}_j$ . A polynomial

kernel of degree  $d$  is  $k(\underline{x}_i, \underline{x}_j) = (1 + \underline{x}_i^T \underline{x}_j)^d$ .

A radial kernel is  $k(\underline{x}_i, \underline{x}_j) =$

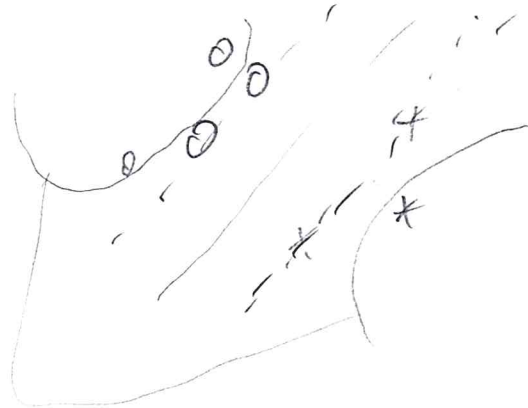
$$\exp\left[-\gamma \sum_{k=1}^p (x_{ik} - x_{jk})^2\right] = \exp\left[-\gamma \|\underline{x}_i - \underline{x}_j\|_2^2\right]$$

23) A SVM uses  $f(x) = \beta_0 + \sum_{i=1}^n \alpha_i k(x, \underline{x}_i)$

$$= \beta_0 + \sum_{i \in S} \alpha_i k(\underline{x}, \underline{x}_i) = ESP = ESP(\underline{x})$$

where  $S$  is the index of support vectors.

Note: Support vectors determine the hyperplane and margin: if they are moved the hyperplane moves too.



these points are not support vectors

24} Using  $k(\underline{x}, \underline{x}_i)$  leads to non-linear decision boundaries if  $k$  is nonlinear. The kernel is a bivariate transformation. There

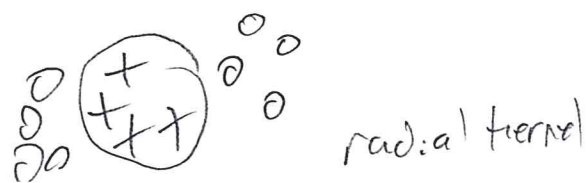
are  $\binom{n}{2} = \frac{n(n-1)}{2}$  distinct pairs  $(\underline{x}_i, \underline{x}_j)$ .

that are needed to estimate  $\beta_0$  and the  $\alpha_i$ .

'SVM with'  $f(\underline{x}) = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i k(\underline{x}, \underline{x}_i) = ESP = ESP(\underline{x})$

is a competitor for QDA while the

SVM with  $f(\underline{x}) = \hat{\beta}_0 + \hat{\beta}'\underline{x}$  is a competitor for LDA.



29) P393 If  $\underline{x}$  is far from  $\underline{x}_i$ , then

SL 6)

$\|\underline{x} - \underline{x}_i\|_2^2$  is large so

$K(\underline{x}, \underline{x}_i) = \exp\left[-\gamma \|\underline{x} - \underline{x}_i\|_2^2\right]$  is tiny

and  $\underline{x}_i$  has almost no contribution in  $f(\underline{x})$ .

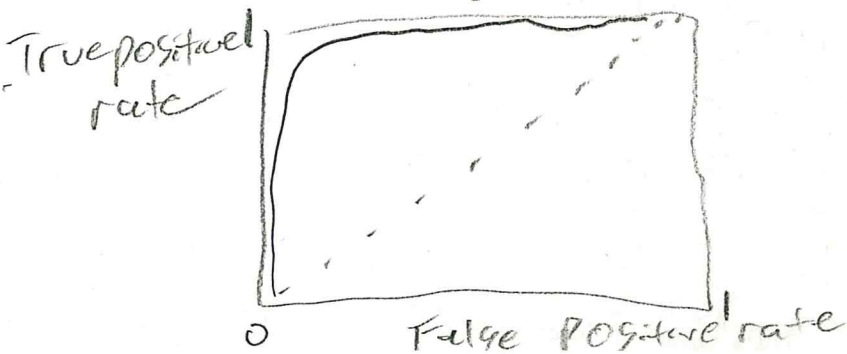
Analogy KNN.

(receiver operating characteristic from quality control)

253 p 147, 354

A ROC curve is used to evaluate binary classifiers. The overall performance is summarized by the area under the ROC curve (AUC). An ideal ROC curve is close to the top left corner of the plot, so the larger the AUC the better the classifier.

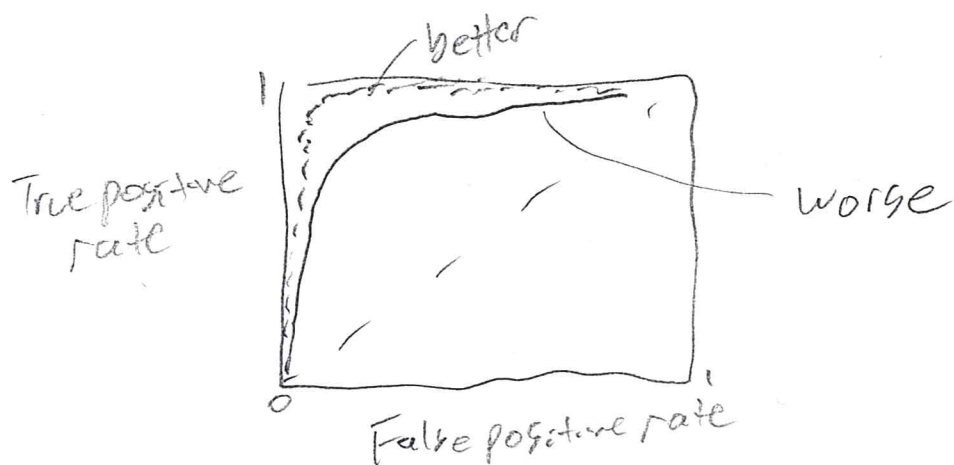
	truth -	FN False negative	total $N^*$
Predict	-   true negative TN	TP True positive	$P^*$
	FP False Positive		
total $N$	ROC Curve		



$$0 \leq AUC \leq 1$$

A classifier with  $AUC = 0.5$  does no better than chance

27) p 354 varies  $\gamma$  for the radial kernel and selects  $\hat{\gamma}$  with the best ROC curve



28) ROC from test data or validation data is better than ROC from training data.

29) The true positive rate is called the sensitivity and the false positive rate is  $1 - \text{specificity}$ .

$$\text{False positive rate} = \frac{FP}{N} \approx \text{type I error, } 1 - \text{specificity}$$

$$\text{True positive rate} = \frac{TP}{P} \approx 1 - \text{type II error, power, sensitivity, recall}$$

$$\text{Positive predicted value} \frac{TP}{P^*} \approx \text{precision, } 1 - \text{false discovery proportion}$$

$$\text{negative predicted value} \frac{TN}{N}$$

$$\text{error rate} = \frac{(FP + FN)}{(FP + FN + TN + TP)}$$

30) Sometimes one error is much more important than the other. ex loan default do not default misclassifying default should be small compared to misclassifying do not default

§ 9.4 SVMs with  $G > 2$  classes

9L 62

31) The one versus one or all pairs R package uses this

Classifier constructs  $\binom{G}{2}$  binary classifiers, one for each distinct pair of groups.

Classify  $\underline{x}$  with  $f_{ij}(\underline{x})$  and let

$m_i = \#$  times  $\underline{x}$  is predicted to be in class  $i$ .

Then  $\hat{Y}(\underline{x}) = d$  where  $m_d = \max(m_1, \dots, m_G)$ .

32) The one versus all classifier

fits  $G$  binary SVMs: group  $i = +1$  versus the  $G-1$  other classes with  $-1$  with  $f_i(\underline{x})$ .

Then  $\hat{Y}(\underline{x}) = d$  where  $f_d(\underline{x}) = \max(f_1(\underline{x}), \dots, f_G(\underline{x}))$ .

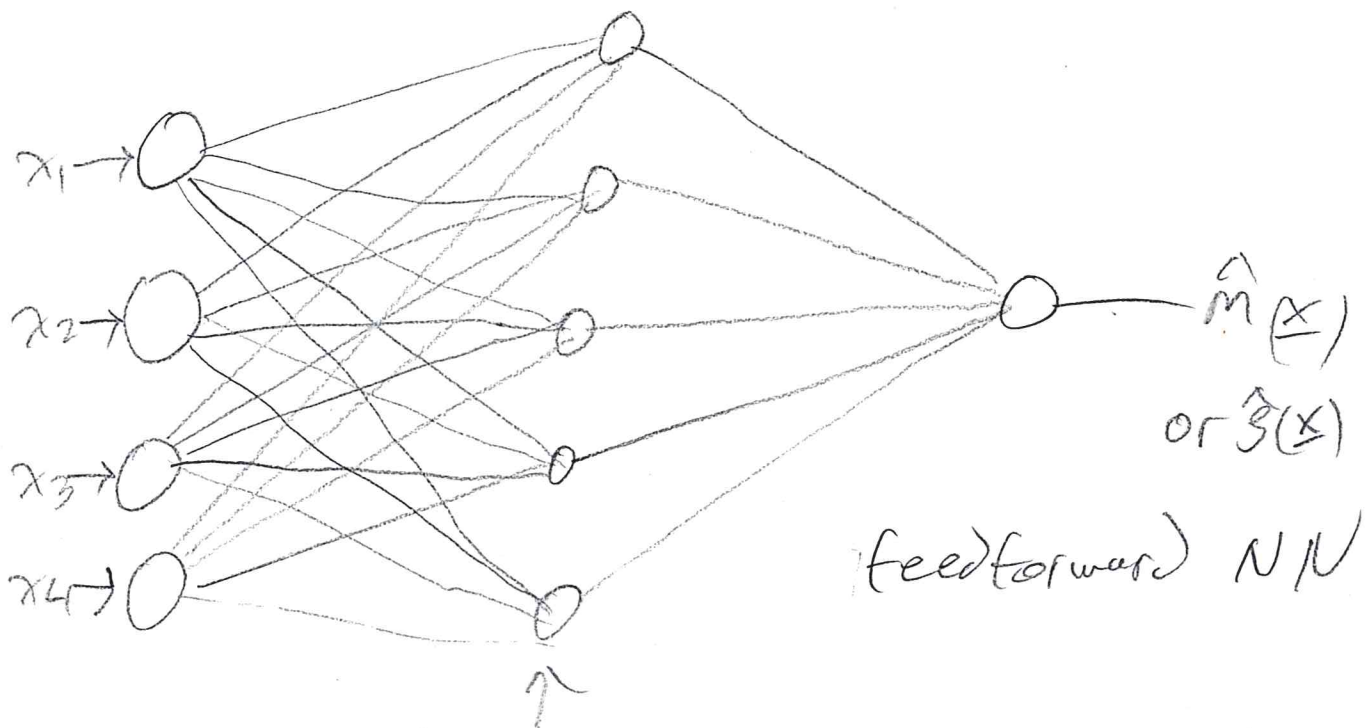
33) Rules 31) and 32) can be applied to other binary classifiers such as logistic regression,

Neural Networks not on final

Suppose  $Y = m(\underline{x}) + e$ , or Neural networks  $Y | \underline{x} \sim \text{bin}(1, \sigma(\underline{x}))$ .

(NN) will produce  $\hat{m}(x)$  or  $\hat{g}(x)$ .

2) Input layer  $L_1$  hidden layer  $L_2$  output layer  $L_3$



The 5 hidden units  $a_l(x) =$

$$a_l = g \left( \beta_{l0} + \underline{x}^T \underset{4 \times 1}{\beta_l^{(l)}} \right), \quad l=1, \dots, 5 \text{ and } g$$

single output unit  $o = h \left( \eta_0 + \underset{5 \times 1}{\underline{a}}^T \underline{\eta} \right)$ .

Often use a sigmoid function

$$g(x) = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-x}} \frac{e^x}{e^x} = \frac{e^x}{1+e^x}$$

For binary regression often  $h(x) = \frac{1}{1+e^{-x}}$  (sigmoid)

For  $y = m(x) + e$ , often  $h(x) = x$

SL 63  
(identity)

$$\text{So } \hat{m}(x) = m_0 + \underline{a}^T \underline{m}$$

Often write  $\underbrace{m_0}_{\text{"bias"}} = w_{00}^{(1)}$ ,  $\underline{a} = \underline{w}_p^{(1)}$  for weights

$$m_0 = w_0^{(2)}, \quad \underline{m} = \underline{w}^{(2)} \quad \text{where}$$

(1) = 1st layer and (2) = 2nd layer.

The NN attempts to learn  $m(x)$  or  $g(x)$ .

3) ANNs can have  $L_{m+2}$  layers

$L_1 = \text{input}$ ,  $L_{m+2} = \text{output}$ , and  $m$  hidden layers  $L_2, \dots, L_{m+1}$ . Each hidden layer

has  $K_i$  hidden units,  $i = 2, \dots, m+1$ .

4) NNs were popular 1985-1995, then boosting and SVM's became popular. Around 2010 NNs made a comeback: deep learning.

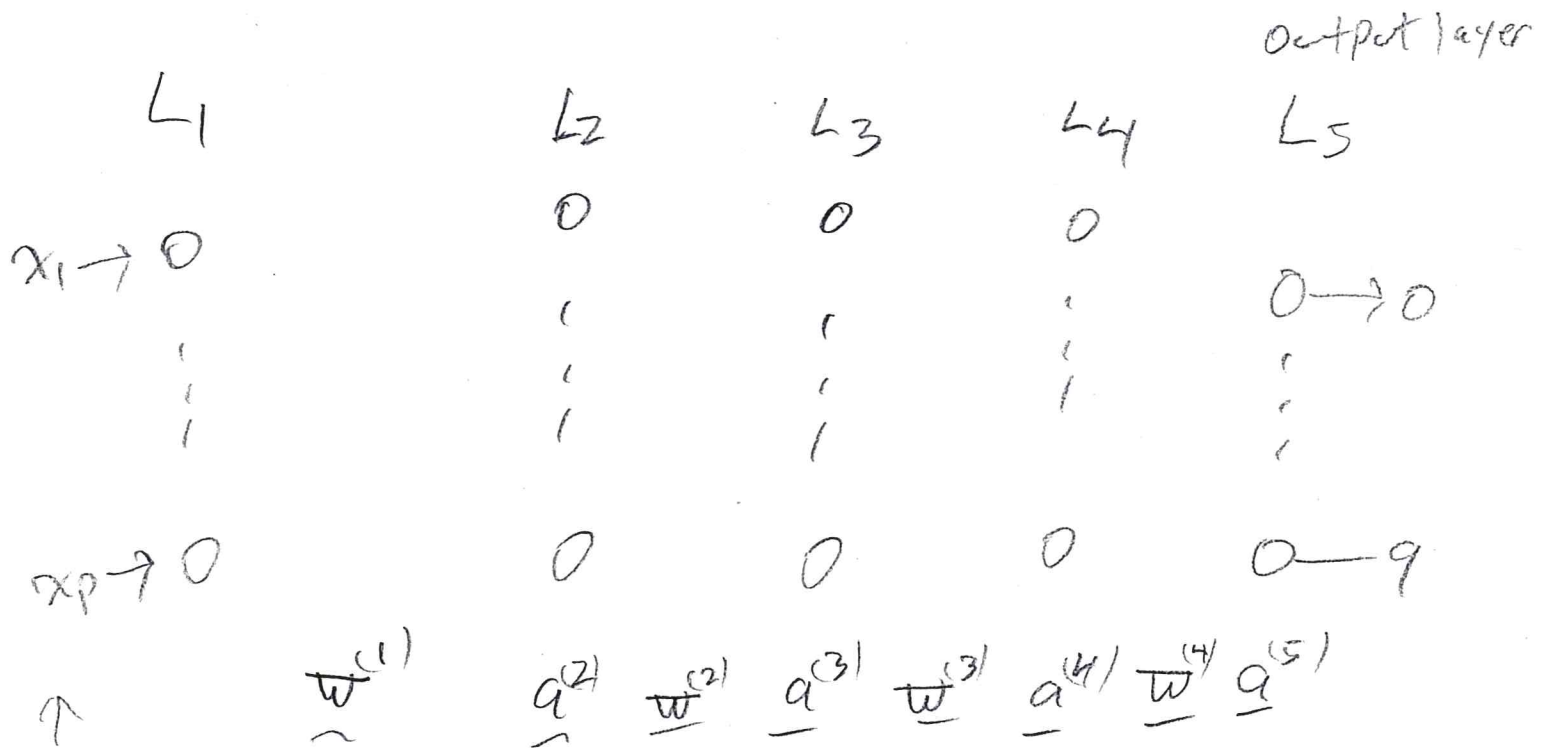
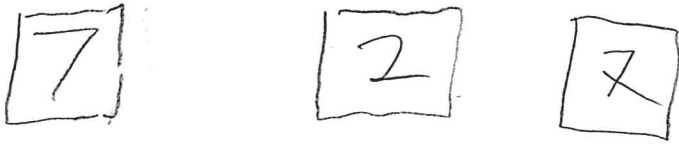
5) NNs are good at optical character recognition, for example, automatic

reading of handwritten digits  $\in \{0, 1, \dots, 9\}$

Build a classifier  $\hat{C}(x) \in \{0, 1, \dots, 9\}$

based on  $x \in \mathbb{R}^{28 \times 28}$  where each digit

is represented by a  $28 \times 28$  grayscale image



$n = 784$   
digits  
training data

This NN has close to 4 million weights and needs to be heavily regularized.

5) Let  $\underline{w}^{(k)}$  be the matrix of weights that go from layer  $L_{k-1}$  to  $L_k$ . Let  $\underline{a}^{(k)}$  be the  $\leq$  including bias = intercept



vector of activations at layer  $L_k$ , SL 64

$$\text{Let } \underline{z}^{(k)} = \underline{w}^{(k-1)} \underline{a}^{(k-1)} \quad \text{and} \quad \underline{a}^{(k)} = g^{(k)}(\underline{z}^{(k)})$$

where  $g^{(k)}$  operates elementwise on  $\underline{z}^{(k)}$  to

$$\text{produce vector } \underline{a}^{(k)} \quad (g^{(k)}: \mathbb{R}^{n_k} \rightarrow \mathbb{R}^{n_k})$$

7) For  $G$ -class classification, the softmax function  $g^{(G)}(\underline{z}_m; \underline{z}) = \frac{e^{z_m^{(G)}}}{\sum_{i=1}^G e^{z_i^{(G)}}}$

for  $m=1, \dots, G$ . These are  $G$  probabilities that sum to one. Then

$$\hat{C}(x) = d \quad \text{where } d = \max_{m=1, \dots, G} g^{(G)}(z_m^{(G)}, \underline{z}^{(G)})$$

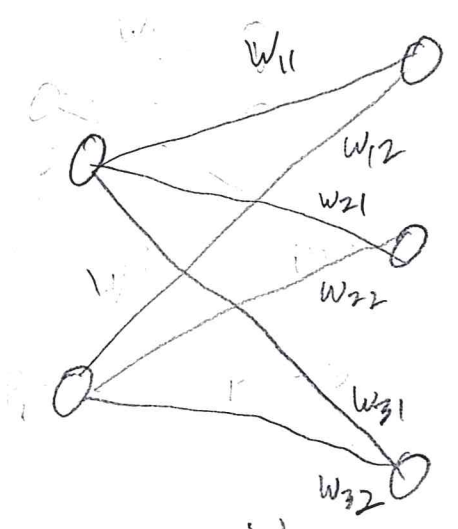
8) Let  $w$  be the collection of weights.

$$\min_w \frac{1}{n} \sum_{i=1}^n L[\bar{y}_i, \hat{M}(x_i, w)] + \lambda J(w)$$

not convex so not fast

$$J(w) = \frac{1}{2} \sum_{i=1}^{k-1} \sum_{j=1}^{p_i} \sum_{l=1}^{p_{i+1}} (w_{lj}^{(i)})^2 \quad \text{is often used}$$

for regression.   
 $\uparrow$  no intercepts.   
 $p_i = \# \text{ units in } i\text{th layer}$    
 $k = \# \text{ layers}$



$$w_{ij} \quad a^{(i+1)}$$

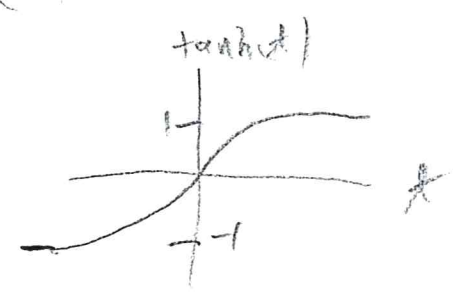
$$P_i = 2 \quad P_{i+1} = 3$$

9) Deeper neural networks have more hidden layers. With one hidden layer, the number of hidden units determines the number of parameters. Two hidden layers works well for digit recognition.

10) In addition to the sigmoid

$$g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \tanh(x) \text{ is used.}$$

The logistic curve (sigmoid) and  $\tanh(x)$  look similar.



The rectified linear  $g(x) = x_+ = \begin{cases} x & x \geq 0 \\ 0 & x \leq 0 \end{cases}$  is also used.

11) Deep learning  $\approx$  neural networks with better computers. NN are useful (classifying natural images (data base of flowers, snakes, mammals, insects, people) using what is known as convolutional architecture)

12) An autoencoder is a NN used to solve a nonlinear principal component decomposition. A single layer autoencoder solves

$$\min_{\underline{w}} \sum_{i=1}^n \| \underline{x}_i - \underline{w}' g(\underline{w} \underline{x}_i) \|^2$$

$P \times P$

13) PCA:  $\underline{x}_i \in \mathbb{R}^P$ ,  $\underline{z}_i = V^T \underline{x}_i \in \mathbb{R}^g$ ,  $g \leq P$  and columns of  $V$  are orthonormal. Taking  $g(\underline{x}) = \underline{x}$  (identity function) gives PCA.



deep learning network

15) Gradient descent for each layer is used since the NN criterion is not convex

