

Underfitting is a serious problem:

$$Y = \underline{X}_I^T \underline{\beta}_I + e_I \quad \text{var}(e_I) \rightarrow \text{var}(e) = \sigma^2$$

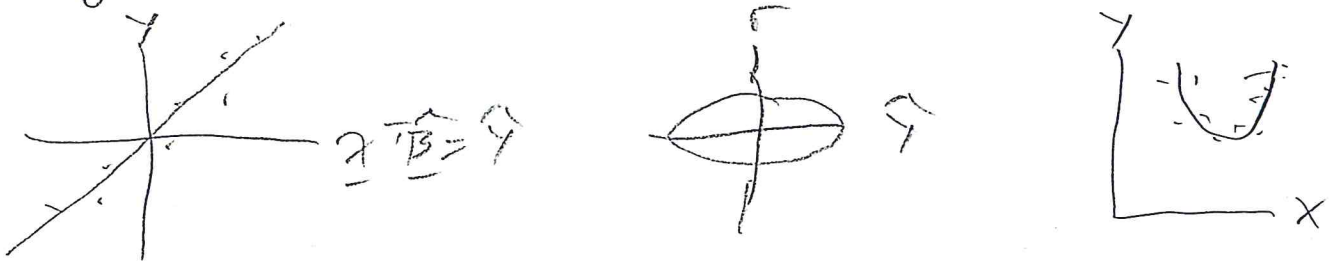
47

may not be constant, could depend on case  $i$

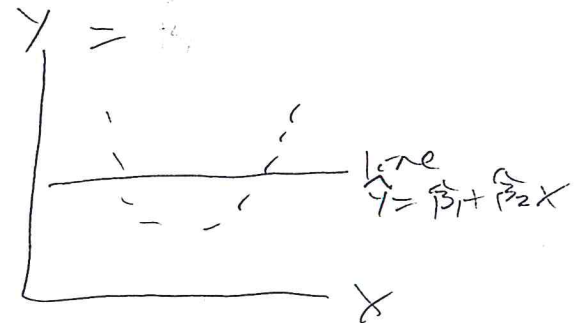
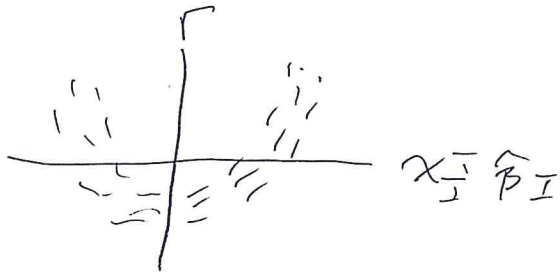
and the model may no longer be linear. Check with response and residual plots.

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

$SP = \underline{x}^T \underline{\beta}$  is a hyperplane that is a quadratic in  $x$ .



If  $SP(I) = \beta_1 + \beta_2 x = \text{line in } x$



expect  $\hat{\beta}_2 \approx 0$

8) Forward selection forms a sequence of submodels  $I_1, \dots, I_m$ , eg  $m = \min(p, \lceil \sqrt{\frac{n}{5}} \rceil)$ .

$I_1$  uses  $x_1^* = x_1$ , a constant but no nontrivial predictors. (47.9)

To form  $I_2$ , consider all models  $I$  with 2 predictors including  $x_1^*$ . Compute  $SSE(I)$   
 $= RSS(I) = \sum_{i=1}^n (y_i - \hat{y}_i(I))^2$ . Let  $I_2$  minimize

$SSE(I)$ , 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  $SSE(I)$  and use predictors  $x_1^*, \dots, x_j^*$ .

$I_j$		# models (at least) to find $I_j$
$I_1$	$x_1^* = x_1$	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

$$1 + 2 + \dots + p-1 = \frac{p(p-1)}{2}$$

10} Need to choose the final model from the sequence of  $m$  models  $I_1, \dots, I_m$ ,

Let  $\sigma^2$  and  $\sigma^2$  be ... for a given data set,  $p, n$ , and  $\underbrace{\hat{\sigma}^2 = \text{MSE}}_{\text{full model}}$  act as

(48)

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 Criterion

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

$C_p(I) = \text{AIC}_s(I)$  uses  $k_n = 2$  while  $\text{BIC}_s(I)$  uses  $k_n = \log(n)$ . Want  $n \geq 5p$ ,  $5 \geq s$ , preferably  $5 \geq 10$ . The following criterion still need  $\frac{n}{p}$  large.

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

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

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

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

$$\log \left[ \binom{p}{a} \right] = \underbrace{\log [p!]}_{\text{constant for a given data set}} - \log [(p-a)!] - \log [a!]$$

11) If  $\hat{\beta}_{US} = \hat{\beta}_{I_{\min}, 0}$  is a consistent estimator of  $\beta$ , (W.B.S)  
 then the probability that  $I_{\min}$  underfits  $\rightarrow 0$   
 as  $n \rightarrow \infty$ . Hence  $P(S \subseteq I_{\min}) \rightarrow 1$  as  $n \rightarrow \infty$ ,  
 a condition that holds for  $C_p$ , AIC, BIC  
 lasso variable selection and elastic net  
 variable selection.

12) If  $S \subseteq I$  then  $\hat{\beta}_{I, 0}$  is a  $\sqrt{n}$  consistent  
 estimator of  $\beta$  for  $Y = \underline{x}^T \beta + e = \underline{x}_S^T \beta_S + e$ .  
 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,  $\hat{\beta}_{US} = \hat{\beta}_{I_{\min}, 0}$  is a  $\sqrt{n}$  consistent  
 estimator of  $\beta$  if  $P(S \subseteq I_{\min}) \rightarrow 1$  as  $n \rightarrow \infty$ .

§11.7 13) A random vector  $\underline{U}$  has  
a mixture distribution of random vectors  $\underline{U}_j$   
with probabilities  $\pi_j$  if  $\underline{U}$  equals  $\underline{U}_j$  with  
 probabilities  $\pi_j$  for  $j = 1, \dots, J$  (and the  
 selection mechanism does not change the  
 distribution of the  $\underline{U}_j$ ). Let  $\underline{U}$  and  $\underline{U}_j$  be

random v ... .. then ... .. = ... ..

is  $F_U(x) = \sum_{j=1}^J \pi_j F_{U_j}(x)$  where  $0 \leq \pi_j \leq 1$

$\sum_{j=1}^J \pi_j = 1$ ,  $J \geq 2$ , and  $F_{U_j}(x)$  is the cdf of  $U_j$ .

Suppose  $E[\bar{h}(U)]$  and  $E[\bar{h}(U_j)]$  exist, Then

$E[\bar{h}(U)] = \sum_{j=1}^J \pi_j E[\bar{h}(U_j)]$  and  $E(U) = \sum_{j=1}^J \pi_j E(U_j)$ .

Hence  $\text{COV}(U) = E(UU^T) - [E(U)][E(U)]^T$

$= \sum_{j=1}^J \pi_j E(U_j U_j^T) - [E(U)][E(U)]^T =$

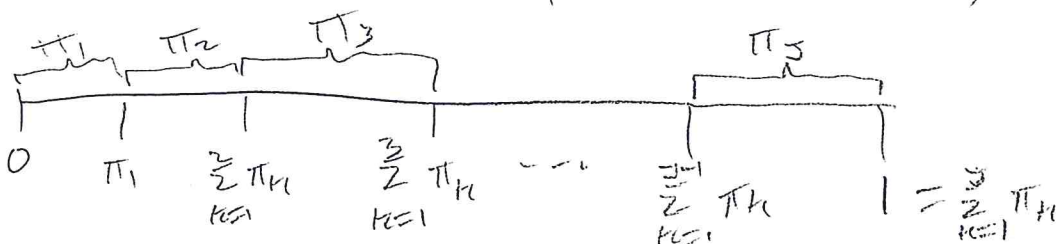
$\sum_{j=1}^J \pi_j \text{COV}(U_j) + \sum_{j=1}^J \pi_j E(U_j)E(U_j)^T - E(U)E(U)^T$ .

If  $E(U_j) = \underline{0}$  (for  $j=1, \dots, J$ ), then  $E(U) = \underline{0}$

and  $\text{COV}(U) = \sum_{j=1}^J \pi_j \text{COV}(U_j)$ .

14} Random selection works: generate a uniform (0,1) RV  $W \perp U_j$  and set

$U = U_j$  if  $W \in \left( \sum_{k=0}^{j-1} \pi_k, \sum_{k=0}^j \pi_k \right)$  with  $\pi_0 = 0$ .



15) Variable selection changes the  $\hat{\beta}_{OLS}$  of  $\underline{u}_{in}$  to  $\underline{w}_{in}$ , say.

16)  $\hat{\beta}_{VS} = \hat{\beta}_{I_K, 0}$  with prob  $\pi_{KN}$ . Let  $\hat{\beta}_{MIX} = \hat{\beta}_{I_{KN}, 0}$  with prob  $\pi_{KN}$  but  $\hat{\beta}_{MIX}$  uses random selection instead of variable selection, can't compute  $\hat{\beta}_{MIX}$  since the  $\pi_{KN}$  are unknown.

17) By OLS CLT, if  $S \subseteq I_j$ , then

$$\sqrt{n} (\hat{\beta}_{I_j} - \beta_{I_j}) \xrightarrow{D} N_{q_j}(\underline{0}, V_{j_1}).$$

Hence  $\sqrt{n} (\hat{\beta}_{I_{j,0}} - \beta) \xrightarrow{D} N_p(\underline{0}, V_{j0})$  where

$V_{j0}$  adds rows and columns of  $Q_S$  corresponding to  $X_i$  not in  $I_j$ . Thus  $V_{j0}$  is singular unless  $I_j$  is the full model.

ex)  $\beta = (\beta_1, \beta_2, \beta_3, \beta_4)^T = (\beta_1, \beta_2, 0, 0)^T$  with

$$\beta_S = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \quad \beta_E = \begin{pmatrix} \beta_3 \\ \beta_4 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}. \quad \beta_{I_1} = \beta_S \quad a_1=2$$

$$\beta_{I_2} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}, \quad \beta_{I_3} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_4 \end{pmatrix} \quad \text{and} \quad \beta_{I_4} = \beta \quad a_4=4$$

for  $I_j$  with  $S = \{1, 2\} \subseteq I_j$