

58) Let the sufficient predictor

$$E(\underline{Y}|\underline{X}_I) = SP = \underline{\underline{X}}_I^T \underline{\underline{\beta}}_I.$$

so  $E(\underline{Y}|\underline{X}) = SP = \underline{\underline{X}}^T \underline{\underline{\beta}}$  for the full model.

Let the estimated sufficient predictor

$$\hat{ESP} = \hat{\underline{Y}} = \underline{\underline{X}}^T \underline{\underline{\beta}}$$

$$\hat{ESP}(I) = \hat{\underline{Y}}_I = \underline{\underline{X}}_I^T \underline{\underline{\hat{\beta}}}_I.$$

59) For  $\underline{X}_S$  and  $\underline{X}_I$ , suppose  $S \subseteq I$

and that 42) holds:  $\underline{Y} = \underline{\underline{X}}^T \underline{\underline{\beta}} + e = \underline{\underline{X}}_S^T \underline{\underline{\beta}}_S + e$

$$\begin{aligned} \text{Then } SP &= \underline{\underline{X}}^T \underline{\underline{\beta}} = \underline{\underline{X}}_S^T \underline{\underline{\beta}}_S = \underline{\underline{X}}_S^T \underline{\underline{\beta}}_S + \underline{\underline{X}}_I^T \underline{\underline{\beta}}_{IS} + \underline{\underline{X}}_0^T \underline{\underline{\beta}}_0 \\ &= \underline{\underline{X}}_I^T \underline{\underline{\beta}}_I, \text{ and } \text{corr}(\underline{\underline{X}}^T \underline{\underline{\beta}}, \underline{\underline{X}}_I^T \underline{\underline{\beta}}_I) = 1. \end{aligned}$$

Here  $\underline{X}_{IS}$  is the vector of predictors in  $\underline{X}_I$  but not in  $\underline{X}_S$ .

So Submodel I is worth considering

if the full model  $\underline{Y} = \underline{\underline{X}}^T \underline{\underline{\beta}} + e$  is good  
and  $\text{corr}(\underline{\underline{X}}^T \underline{\underline{\hat{\beta}}}, \underline{\underline{X}}_I^T \underline{\underline{\hat{\beta}}}_I) \geq 0.95$ .

60) full model:  $\hat{\underline{Y}}$ ,  $r = \underline{Y} - \hat{\underline{Y}}$ ,  $\hat{\underline{\beta}}$

Submodel I:  $\hat{\underline{Y}}_I$ ,  $r_I = \underline{Y} - \hat{\underline{Y}}_I$ ,  $\hat{\underline{\beta}}_I$   
 $\hat{Y}_{i,I} = \hat{\underline{\beta}}_I^T \underline{\underline{X}}_I$ ,  $r_i = Y_i - \hat{\underline{\beta}}_I^T \underline{\underline{X}}_I$

61) \* P10D If you have one data set,  
build a full model for that data set, then

perform variable selection to obtain a final (sub)model, then an extreme amount of data snooping was used. (34.5)

The final model may be useful for exploratory data analysis and description, but the p-values for the t tests and Anova F test are likely too small while the p-value for the partial F test that uses the final model as the reduced model is likely too high. The final model tends to fit the data set from which it was built better than future obs's. PIs tend to be too short so coverage < nominal coverage e.g. 95%.

62) Inference such as tests and PI's are valid if the final MLR model is selected before gathering data, not after using variable selection to build a final model. Data splitting is useful.

63) p101 Leaving out important predictors is called underfitting. Such models tend to be very poor and may violate linearity ( $y = \underline{x}_I^T \beta_I + e$ ) and constant variance ( $V(y_i | \underline{x}_I) = \sigma^2$ ) assumptions.

64) After performing variable selection, the tests and CIs are useful for exploratory

purposes and description, but are not valid for inference.

65) P101 Having unnecessary predictors in the model I is called overfitting or fitting noise.

Let 42) hold with  $S \subseteq I$ .

$$\text{If } Y = X^T B + e = \underbrace{X^T B_I}_{p \times 1} + \underbrace{e}_{j \times 1} = \underbrace{X_S^T B_S}_{k \times 1} + e$$

$$\frac{1}{n} \sum V(\hat{Y}_{II}) = \frac{\sigma^2 j}{n} > \frac{\sigma^2 k}{n} = \frac{1}{n} \sum V(\hat{Y}_S).$$

↑ worst when  $j=p$  (consider  $n \approx p \gg k$ )

66) P109 \*Rules of thumb 3.7 Assume full model has good response and residual plots and that  $n > 5p$ . Let subset I have  $k$  predictors including a constant. Let

$I_{\min}$  be the min Cp model and let

$I_I$  be the model with the fewest predictors satisfying  $C_p(I_I) \leq C_p(I_{\min}) + 1$ .

\* Do not use more predictors than  $I_I$ .

Often can't have the following 10 rules hold simultaneously: submodel I is good if

i) the response and residual plots for the submodel look like the response and residual plots for the full model

ii)  $\text{corr}(\hat{Y}, \hat{Y}_I) \geq 0.95$

iii) The plotted points in the FF plot cluster tightly about the identity line

- 35.5
- iv)  $p\text{val} \geq 0.01$  (not 0.05) for the partial F test that uses  $I$  as the reduced model.
- v) want  $K \leq n/10$
- vi) the plotted points in the RR plot should cluster tightly about the identity line.
- vii) want  $R^2(I)$  close to  $R^2(\text{full})$  ( $R^2(I) \leq R^2(\text{full})$  since adding predictors does not decrease  $R^2$ )
- viii) want  $C_p(I_{\min}) \leq C_p(I) \leq \min(2k, p)$  with no big jumps in  $C_p$  (the increase should be less than 4) as variables (with 1df) are deleted.
- ix) want hardly any predictors with  $p\text{val} > 0.05$
- x) want few predictors with  $.01 \leq p\text{val} \leq .05$

### Common E2 Problem

ex problem 3.3

	L1	L2	L3	L4
# predictors	10	5	4	3
# with $.01 \leq p\text{val} \leq .05$	0	1	0	0
# with $p\text{val} > .05$	8	0	0	0
$R_I^2 = R^2(I)$	.655	.650	.648	.630
corr ( $\hat{Y}, \hat{Y}_I$ )	1.0	.996	.992	.981
$C_p(I)$	10.0	4.0	5.60	13.81
$p\text{val}$ for partial Ftest	1.0	.550	.272	.015

$L_1$  is the full model,  $L_2$  is the  $\min(C_p)$  model. a) which model is  $I$ ? b) which other models, if any, should be looked at?

Soln a)  $I_I = L_2$  since  $5.6 = C_p(L_3) > C_p(I_{\min}) + 1 = 5$   
 b)  $L_3$  since  $5.6 < C_p(I_{\min}) + 4$  (and  $C_p(L_4) > 2k$ ).

§ 2.5  
and § 3.4.1 Inference After Variable Selection 38  
certain  
82) Prediction intervals and Prediction regions

Skipped notes 36, 37 for now  
Math 428

applied to a bootstrap sample can  
be used as confidence intervals and  
confidence regions.

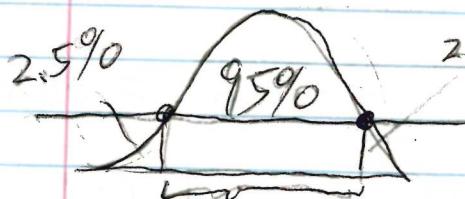
83) Let  $Y_{(1)}, \dots, Y_n$  be the data.

The order statistics are

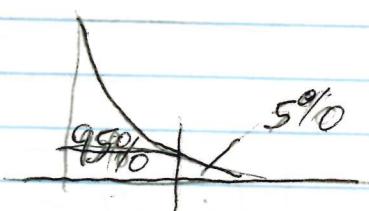
$$Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)}$$

↑ min ↑ max

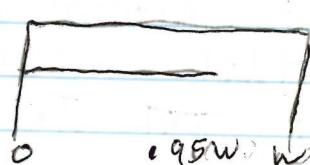
84) The highest  $100(1-\delta)\%$  density region  
of a pdf is found by moving  
a horizontal line down from the  
top of the pdf so that the  
line "intersects the pdf" at one  
or more intervals and the sum  
of the areas under the pdf  
corresponding to the intervals is  $1-\delta$



highest 95%  
region  
symmetric unimodal  
lop off  $100\frac{\delta}{2}\%$  from  
both sides

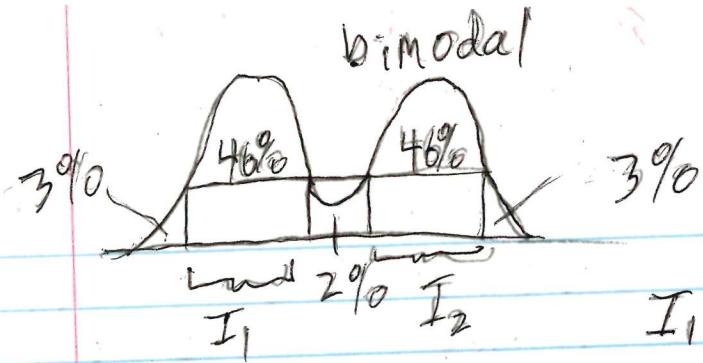


Unimodal right skewed  
lop off  
 $100\delta\%$  from the right  
tail



$U(0, w)$  any line  
within pdf of length  
 $0.95w$  works

38.5


 $I_1 \cup I_2 = \text{highest } 92\% \text{ region}$ 

- 85) Consider intervals that contain  $c$  cases:

$[Y_{(1)}, Y_{(c)}], [Y_{(2)}, Y_{(c+1)}], \dots, [Y_{(n-c+1)}, Y_{(n)}]$ . Compute  $Y_{(c)} - Y_{(1)}$ ,  $Y_{(c+1)} - Y_{(2)}$ , i.e.,  $Y_{(n)} - Y_{(n-c+1)}$ .

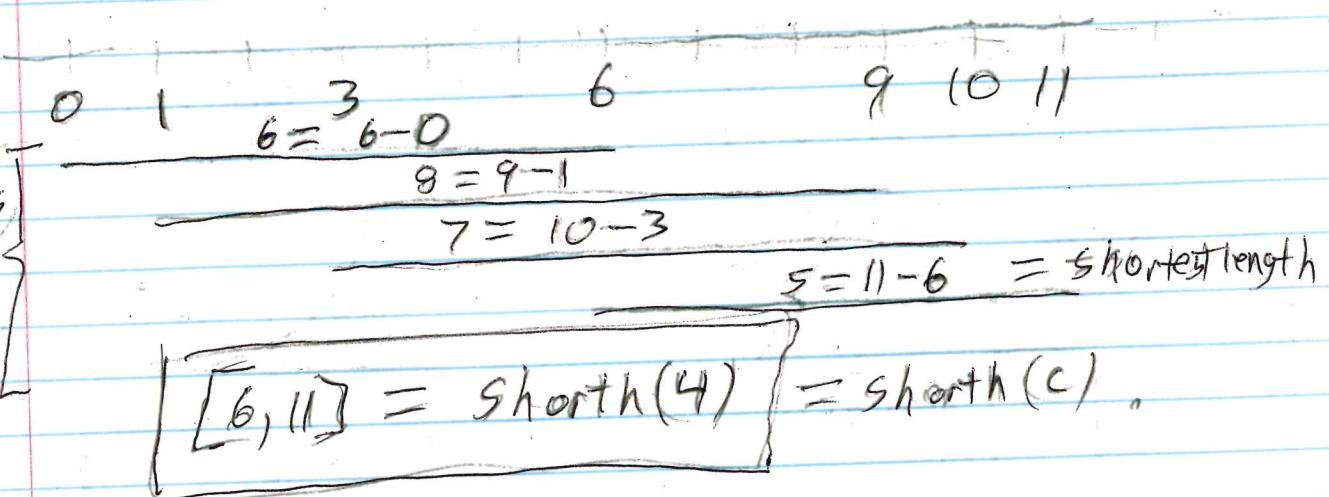
Then  $\text{Shorth}(c) = [Y_{(s)}, Y_{(sc-c)}]$  is the closed interval with the shortest length.

- ex) Know for E2. Let  $c = 4$ . Data below has  $n = 7$ . Find  $\text{Shorth}(4)$ .

see  
Hw7

6

intervals  
containing  
 $c=4$   
cases



- 86) If  $\frac{c}{n} \rightarrow 1-\delta$ , the  $\text{Shorth}(c)$  interval

estimates the highest density  $100(1-\delta)\%$  region if that region is an interval (unimodal pdf that decreases from the mode).