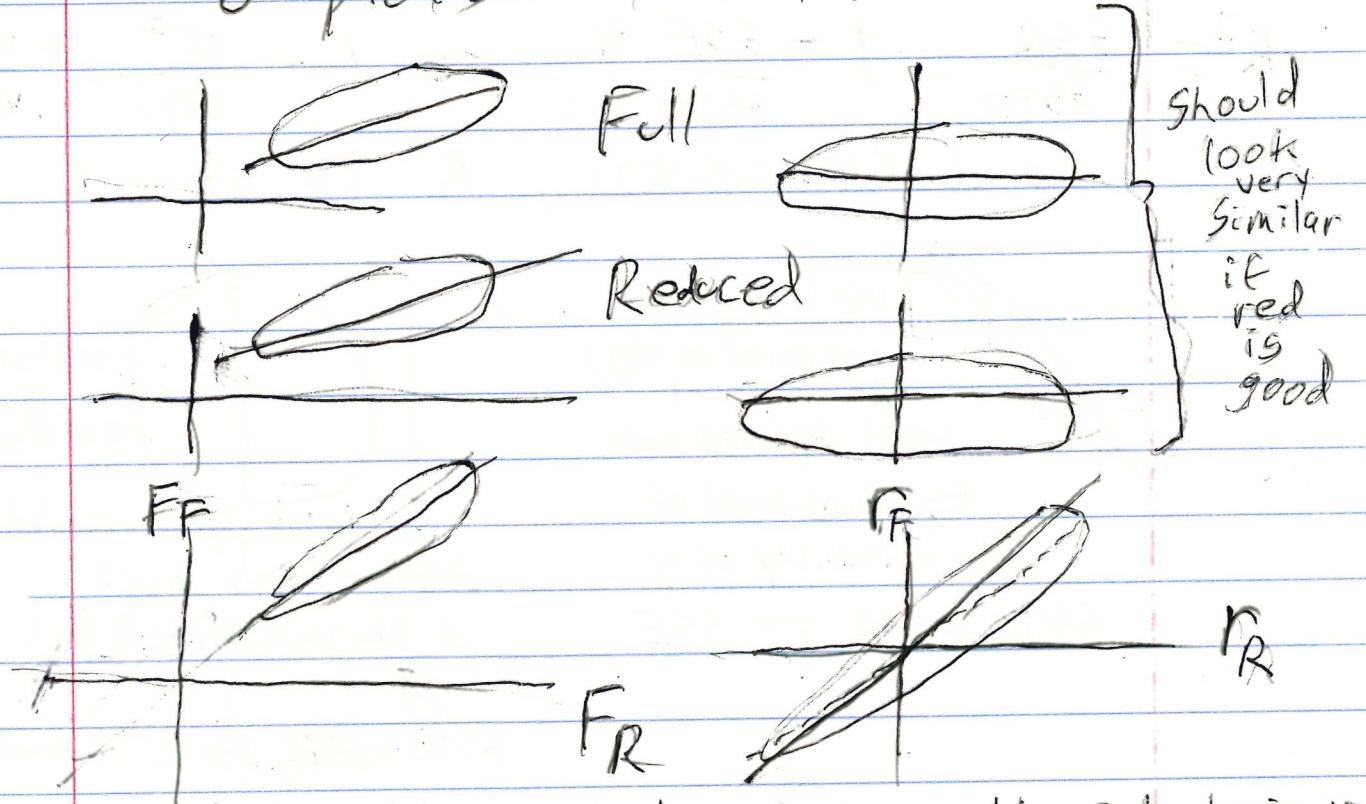


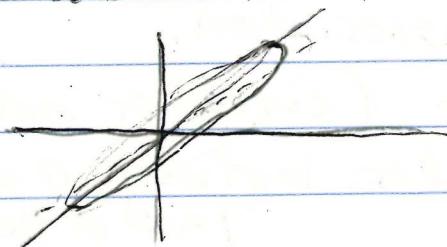
Math 484 15

67) P 47 Make the response and residual plots for the full and reduced model, the RR plot and the FF plot: 6 plots total:



Make these plots for variable selection in ch3.

68) P 34 For the Anova F test, if you fail to reject  $H_0$ , the reduced model residuals are  $\hat{r}_F = y_i - \bar{Y}$  and Full model residuals are  $r_F = y_i - \hat{Y}$



see RR plot on p 35

69) P 47 If the FF plot looks good, but the RR plot is poor, the reduced model

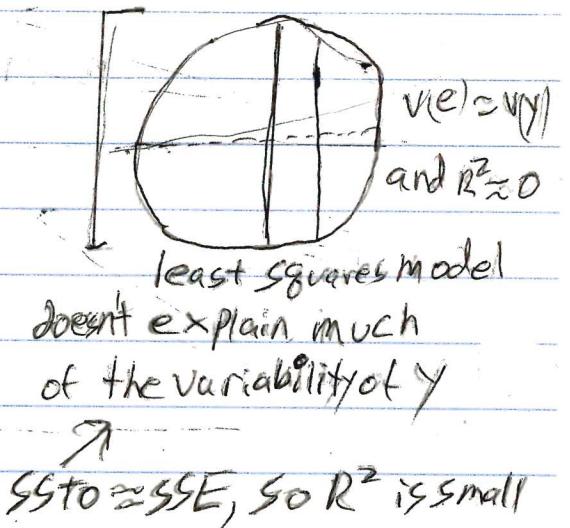
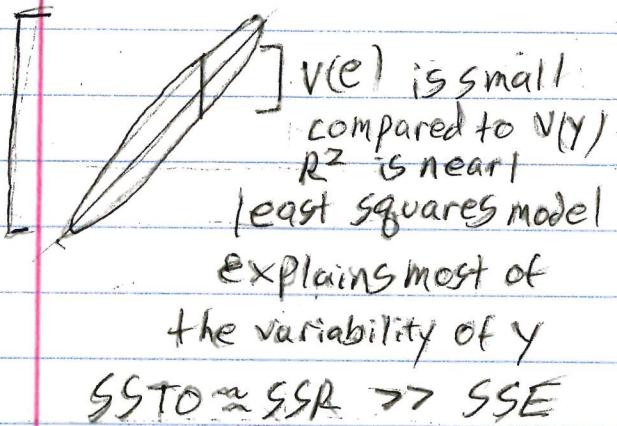
may be good if the main goal of the analysis is to predict  $Y$  (or description). 15.3

70) If a constant is in the model,

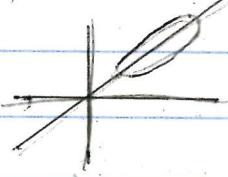
$$SSTO = SSE + SSR \text{ and}$$

$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO} \approx 1 - \frac{\sigma^2}{V(Y)} = \frac{1 - V(e)}{V(Y)}$$

So if  $V(e) \ll V(Y)$ ,  $R^2$  is near 1.



71) i)  $R^2$  is best for a good response plot

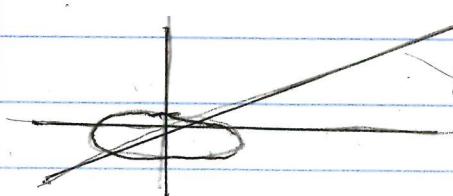


ii) If response plot is linear with gaps,  $R^2$  is usually too high and over emphasizes linear association.



iii) If plot is nonlinear (curved or outliers)

$R^2$  can be any number between 0 and 1



identity line

} most of the variability of  $Y$  is due to the cluster of outliers

If  $Y = X^2$  on an interval symmetric about  $X=0$ , could have  $R^2 = 0$ . But given  $X$ , you know  $Y$ .  $R^2$  measures linear association  $\neq$  association.  
end exam material

72) PGI Let  $X_1 \in \mathbb{R}$  and  $\tilde{X}_{(k)} = (x_1, x_2, \dots, \hat{x}_k, \dots, x_p)^T$

be the vector of predictors with  $x_k$  deleted.

Let  $\tilde{\epsilon}_{(k)}$  be the vector of residuals

from regressing  $Y$  on  $\tilde{X}_{(k)}$ . Let

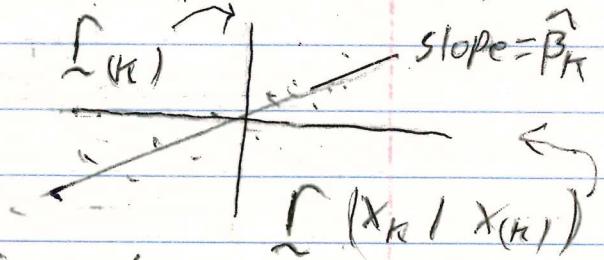
$\tilde{\epsilon}(x_k | \tilde{X}_{(k)})$  be the vector of

residuals from regressing  $x_k$  on  $\tilde{X}_{(k)}$ .

An added variable plot for  $x_k$

is a plot of  $\tilde{\epsilon}(x_k | \tilde{X}_{(k)})$  vs  $\tilde{\epsilon}_{(k)}$

for  $k=2, \dots, p$ .



73] This plot gives information about  $H_0: \beta_k = 0$ . The plotted points cluster about a line through the origin with slope  $\hat{\beta}_k$ . Small range on the horizontal axis means  $x_k$  is well explained by the other predictors.  $\tilde{\epsilon}(x_k | \tilde{X}_{(k)})$  represents

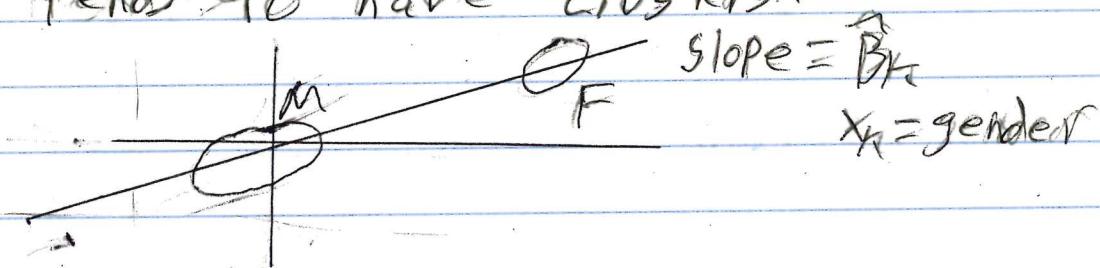
16.5

the part of  $x_k$  not explained by the remaining variables while  $\hat{e}_{(k)}$  represents the part of  $y$  that is not explained by the remaining variables.

74) p51 An added variable plot with clearly nonzero slope and tight clustering about some line through the origin implies  $x_k$  is needed in the MLR model for  $y$  given the other variables are in the model.

Slope near zero implies that  $x_k$  may not be needed, see Fig 2.4.

75) p53 For a categorical variable like gender, the plot tends to have clusters.



76) p25 If residual and response plots look good after several points are moved slightly, then the plots are OK; (If the points have to be moved a lot, they may be outliers.)

77) p27 The sample correlation between  $w_i$  and  $z_i$  is  $\text{corr}(w, z) = \hat{\rho}(w, z) = \frac{1}{n} \sum_{i=1}^n (w_i - \bar{w})(z_i - \bar{z})$

$$= \frac{\sum w_i z_i - n \bar{w}\bar{z}}{\sqrt{\sum (w_i - \bar{w})^2} \sqrt{\sum (z_i - \bar{z})^2}}$$

$$(n\bar{w} = \sum w_i z_i - \sum w_i \bar{z} - \sum z_i \bar{w} + n \bar{w}\bar{z})$$

$$= \sum w_i z_i - n \bar{w} \bar{z}$$

78) p26-7 Suppose  $\mathbf{X}$  is  $n \times p$  with full rank  $p$ ,  $n > p$ , and a constant is in the model. Fact if  $\mathbf{x} = \bar{\mathbf{x}}$ , then  $\hat{y} = \bar{y}$ .

i)  $\sum r_i = 0$  so  $\bar{r} = 0$

ii)  $\sum r_i \hat{y}_i = \underbrace{\mathbf{r}^\top \mathbf{y}}_{j\text{th column of } \mathbf{X}} = 0$  so  $\text{corr}(\mathbf{r}, \hat{y}) = 0$   
 $\downarrow$   $\downarrow$   
 $j\text{th predictor}$

iii)  $\mathbf{X}^\top \mathbf{L} = 0$  so  $\mathbf{X}^\top \mathbf{L} = 0$  so  $\text{corr}(\mathbf{r}, \mathbf{X}_j) = 0$

iv)  $\sum y_i = \sum \hat{y}_i = \sum \hat{y}_i + \sum r_i$

79) p25 i) The residual plot of  $\hat{y}$  vs  $r$  should always be made.

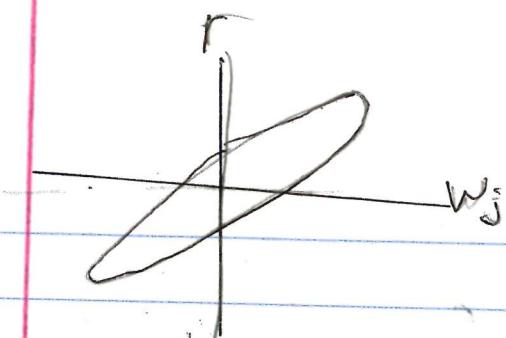
ii) Residual plots of the  $j$ th predictor  $x_j (= \bar{x}_j)$  vs  $r$  should look ellipsoidal about the  $r=0$  line for  $j=2, \dots, p$

if the errors are independent of the predictors and the predictors are quantitative.

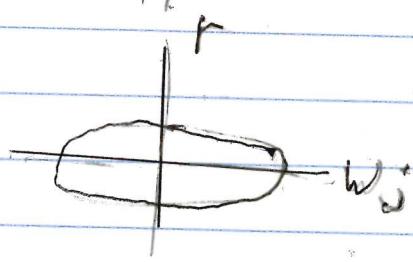
iii) Residual plots of potential predictors  $w_j$  vs  $r$  may be made. If  $w_j$  is not important, the plot should look ellipsoidal about  $r=0$  if  $w_j$  is quantitative.

iv) If the plot is parabolic, add  $x_j^2$  to the model or add  $w_j$  and  $w_j^2$  to the model.

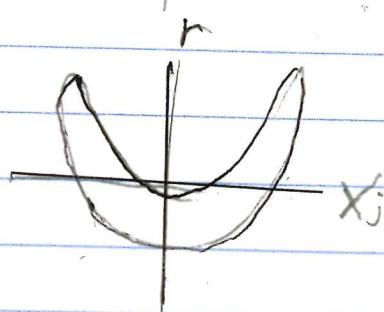
17.5



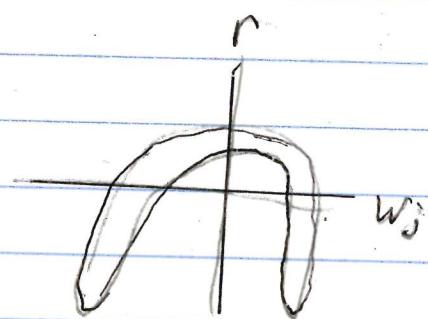
add  $w_j$  to model



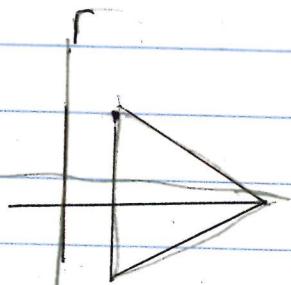
don't add  $w_j$  to model



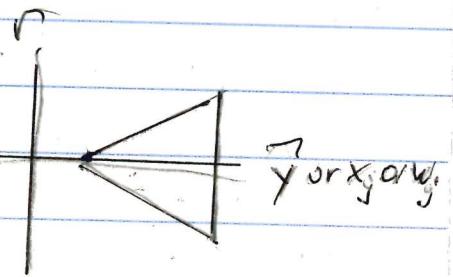
add  $x_j^2$  to model



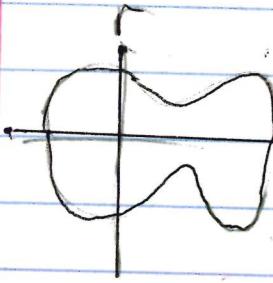
add  $w_j$  and  $w_j^2$  to model



$\gamma$  or  $x_j$  or  $w_j$



$\gamma$  or  $x_j$  or  $w_j$



$\gamma$ ,  $x_j$  or  $w_j$

check constant  
variance assumption  
if  $\gamma$  is used

\*now

- 80) Response plots are good for visualizing the MLR model ( $\gamma_1 x^T \beta$ ) in the background of the data, for checking goodness of fit, linearity, strength of relationship as measured by  $R^2$  or signal to noise ratio, for detecting