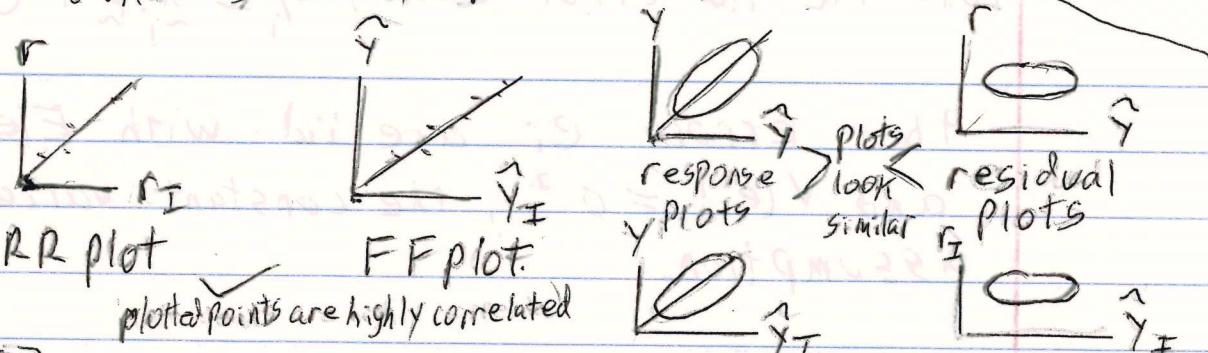


- 67] P105* Full model and candidate for final submodel \mathcal{I}



- 68] P104 For models with k predictors,

the model with the smallest $C_p(\mathcal{I})$

maximizes $\text{corr}(r, r_I)$,

$$C_p(\mathcal{I}) \leq 2k \Rightarrow \text{corr}(r, r_I) \geq \sqrt{1 - \frac{p}{n}},$$

and as $\text{corr}(r, r_I) \rightarrow 1$, $\text{corr}(\hat{y}, \hat{y}_I) \rightarrow 1$.

- 69] Prop 3.2 proves that for models with a constant, the plotted points scatter about the identity line for RR, FF and response plots.

(Go to 83)-86).

- 70] P144 Let R_k^2 be the R^2 from regressing x_k on the other predictors. If R_k^2 is high, multicollinearity is present and $SE(\hat{\beta}_k)$ or $\frac{1}{1-R_k^2}$ is large. Highly correlated linear plots in the scatterplot matrix also suggest multicollinearity. Multicollinearity makes CI for $\hat{\beta}_k$ very long, $(X^T X)$ nearly singular, but is "not a problem for predicting y ".

- §3.5 71] P129 Diagnostics are used to

Check that model assumptions are reasonable. (36.5)

For the iid error model, $y_i = \underline{x}_i^T \underline{\beta} + e_i$,

the errors e_i are iid with $E(e_i) = 0$ and $V(e_i) = \sigma^2$, the constant variance assumption.

72] P132 The best diagnostics are the response and residual plots, but numerical diagnostics are popular.

73] In matrix form $\underline{y} = \underline{X} \underline{\beta} + \underline{e}$

$$\underline{X} = [\underline{x}^1, \underline{x}^2, \dots, \underline{x}^p] = [\underline{I}, \underline{\Omega}]$$

where $\underline{\Omega} = [\underline{x}^2, \dots, \underline{x}^p]$ is the matrix of $n \times (p-1)$

nontrivial predictors. Let $\underline{x}_i = (1, x_{i1}, \dots, x_{ip})^T$
 $= (1, \underline{u}_i^T)^T$ be the i th case, the

i th row of \underline{X} .

74] P13D The sample mean and covariance matrix of the nontrivial predictors $(\underline{u}_1, \dots, \underline{u}_n)$ are $\underline{\bar{u}} = \frac{1}{n} \sum_{i=1}^n \underline{u}_i$

and $\hat{\Sigma} = \widehat{\text{cov}}(\underline{\Omega}) = \frac{1}{n-1} \sum_{i=1}^n (\underline{u}_i - \underline{\bar{u}})(\underline{u}_i - \underline{\bar{u}})^T$.
[$\widehat{\text{cov}}(\underline{\Omega})$] is closely related to $(\underline{X}^T \underline{X})^{-1}$.

75] $\hat{\underline{y}} = H \underline{y}$, $H = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T$

76) Leave one out or case diagnostics are computed by omitting the i th case.

Let $\tilde{Y}_{(i)} = \underline{\underline{X}} \hat{\beta}_{(i)}$ be the fitted values when the i th case is omitted.

$V(r_i) = \sigma^2 (I - h_i)$, $h_i = h_{ii} = H_{ii}$, $0 \leq h_i \leq 1$, is the i th leverage. The studentized residual $\hat{e}_i = \frac{r_i}{\hat{\sigma} \sqrt{1-h_i}}$ where $\hat{\sigma}^2 = \text{MSE}$.

77) * p130-1 $MD_i^2 = (\underline{\underline{v}}_i - \bar{\underline{\underline{v}}})^T \hat{\Omega}^{-1} (\underline{\underline{v}}_i - \bar{\underline{\underline{v}}})$ is the i th squared Mahalanobis distance.

The i th Cook's distance

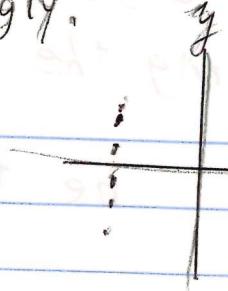
$$CD_i = \frac{(\tilde{Y}_{(i)} - \bar{Y})^T (\tilde{Y}_{(i)} - \bar{Y})}{p \hat{\sigma}^2}$$

$$= \frac{r_i^2}{p \hat{\sigma}^2 (1-h_i)} \frac{h_i}{1-h_i} = \frac{\hat{e}_i^2}{p} \frac{h_i}{1-h_i}.$$

CD_i measures "influence" of case i on $\hat{\beta}$.
 influence \approx leverage (distance of $\underline{\underline{v}}_i$ from $\bar{\underline{\underline{v}}}$)

78) Leverage h_i is large if $\underline{\underline{v}}_i$ is far from $\bar{\underline{\underline{v}}}$ where "far" depends on $[\text{Cov}(\underline{\underline{v}})]^{-1}$ or $(\underline{\underline{X}}^T \underline{\underline{X}})^{-1}$.
 $\sum_{i=1}^n h_i = p$, $0 < h_i \leq 1$, V_i , and $h_i = \frac{1}{n-1} MD_i^2 + \frac{1}{n}$.

Cases with high h_i and CD_i tend to effect $\hat{\beta}$ strongly.

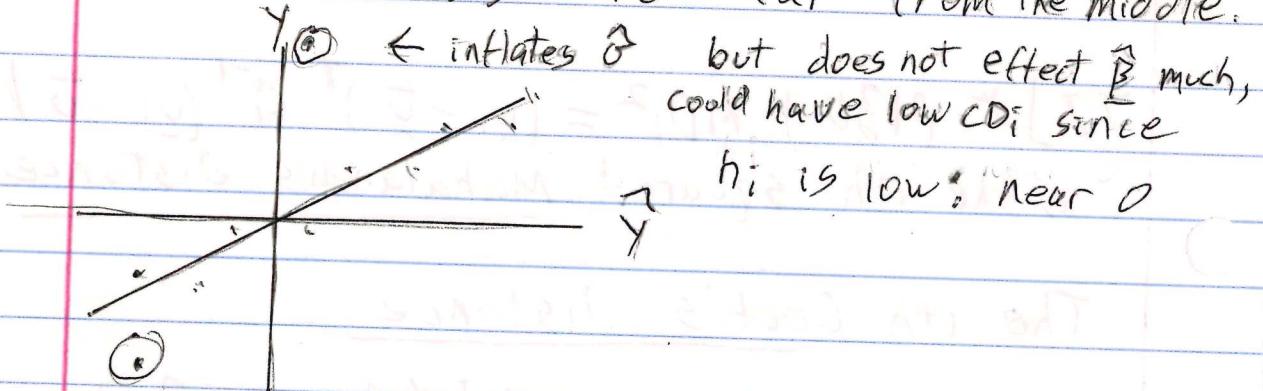


may have $r_i \approx 0$ and low CD_i

high leverage since

$X \quad \hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)$ changes
as this point is
moved vertically.

- 79] In the response plot, cases with large CD_i tend to have large $|r_i|$ if they are close to the middle of the plot. The $|r_i|$ need not be so large if the cases are far from the middle.



h_i is low: near 0

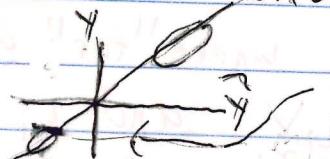
80] $h_i > \frac{2p}{n}$, $MD_i^2 > \chi_{p-1, 95}^2$

$CD_i > \min(1.5, \frac{2p}{n})$ are large

and should be checked.

- 81] These numerical diagnostics do not work well if 2 or more cases are close together but unusual.

~~influential group on $\hat{\beta}$~~
~~most cases in group have small h_i~~



many cases will not have
large CD_i , h_i or MD_i