## Chapter 10
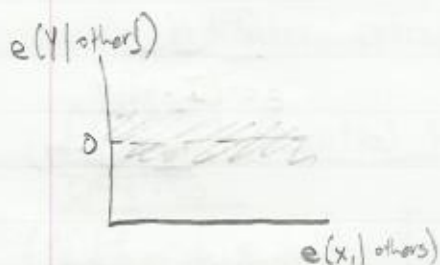
### Added Variable Plots
Show the maginal role of a predictor variable given all other predictors are in the model



$e(Y|others)$ ... $e(x_1|others)$

↑ $x_1$ contains no new information for predicting Y

↑ $x_1$ may be a helpful linear addition to the model

↑ $x_1$ may be a helpful factor to the model (maybe an interaction)

SAS

```
proc reg;
  model y = x_1 x_2 / partial;
run;
```

### Finding Outlying observations
sgscatter



① Outlying w/ respect to its y value, x value or both

②,③,④ Outlying w/ respect to x as they are larger than the other cases

①,② Look like they wont be too influential

③,④ Will be influential to the regression line

④

Residuals $\quad e_i = Y_i - \hat{Y}_i$

Semistudentized Residuals $\quad e_i^* = \frac{e_i}{\sqrt{MSE}}$

Hat Matrix $\quad H_{n \times n} = \underline{X}(\underline{X}'\underline{X})^{-1}\underline{X}'$ $\qquad$ Note $\quad \underline{\hat{Y}} = H\underline{Y}$

$\qquad 0 \leq h_{ii} \leq 1$ $\qquad\qquad\qquad\qquad\qquad\qquad e = (I - H)Y$

$\qquad \sum h_{ii} = p$ $\qquad\qquad\qquad\qquad\qquad\qquad \sigma^2(e) = \underline{\sigma}^2 (I - H)$

Studentized Residuals $\quad r_i = \frac{e_i}{s\{e_i\}}$ $\qquad\qquad$ (student residual)

Deleted Residual $\quad d_i = $ actual data $-$ predicted value from model w/o said data point

$\qquad\qquad\qquad d = Y_i - \hat{Y}_{i(i)}$

$\qquad\qquad\qquad\quad = \frac{e_i}{1 - h_{ii}}$

Studentized Deleted Residuals $\quad t_i = d_i / s\{d_i\}$ $\qquad\qquad$ (RStudent)

$\qquad\qquad\qquad\qquad = e_i / \sqrt{MSE_{(i)}(1 - h_{ii})}$

$\qquad\qquad\qquad\qquad = e_i \left[\frac{n - p - 1}{SSE(1 - h_{ii}) - e_i^2}\right]^{1/2}$

Bonferroni Critical Value $\quad t(1 - \alpha/2 ; n - p - 1) = qt(1 - \frac{\alpha}{2}, n - p - 1)$ in R

$\qquad\qquad$ Test: If $|$Student Deleted Residual$|$ > Bonferroni Critical Value $\Rightarrow$ outlier

$\qquad\qquad\qquad$ else we're okay

Hat Matrix for Outlying Observation $\qquad$ (Hat Diag H)

$\quad \cdot$ $h_{ii}$ is called the leverage, distance from the center

$\quad \cdot$ $h_{ii}$ is considered large if $h_{ii} > 2\bar{h} = \frac{2 \sum h_{ii}}{n} = \frac{2p}{n}$

$\quad \cdot$ Another guideline: $\quad h_{ii} \geq .5 \quad$ high leverage

$\qquad\qquad\qquad\qquad .2 \leq h_{ii} \leq .5 \quad$ moderate leverage $\Big\}$ large n

$\qquad\qquad\qquad\qquad .2 \leq h_{ii} \qquad$ low leverage

42

## Hat matrix for Extrapolation

$$h_{new, new} = X'_{new} (X'X)^{-1} X_{new}$$

*If $h_{new, new}$ is much larger than the leverage values it indicates extrapolation.

## Finding Influential Points

### DFFITS

$$(DFFITS)_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{MSE_{(i)} h_{ii}}} = e_i \left[\frac{n-p-1}{SSE(1-h_{ii}) - e_i^2}\right]^{1/2} \left(\frac{h_{ii}}{1-h_{ii}}\right)^{1/2} = t \left(\frac{h_{ii}}{1-h_{ii}}\right)^{1/2}$$

\# Considered influential if $\begin{cases} |DFFITS|_i > 1 & \text{for small data} \\ |DFFITS| > 2\sqrt{\frac{p}{n}} & \text{for large data sets} \end{cases}$

### Cook's D

$$D_i = \frac{\sum_j (\hat{Y}_j - \hat{Y}_{j(i)})^2}{p \, MSE} = \frac{(\hat{Y} - \hat{Y}_{(i)})'(\hat{Y} - \hat{Y}_{(i)})}{p \, MSE} = \frac{e_i^2}{p \, MSE}\left[\frac{h_{ii}}{(1-h_{ii})^2}\right]$$

\* Considered influential if $pF(D_i, p, n-p) \geq .5$

moderately if $.2 \leq pF(D_i, p, n-p) \leq .5$

### DFBetas

$$(DFBETAS)_{k(i)} = \frac{b_k - b_{k(i)}}{\sqrt{MSE_{(i)} \, c_{kk}}}$$

$c_{kk}$ is the $k^{th}$ diag of $(X'X)^{-1}$

\* Considered influential if $\begin{cases} > 1 & \text{for small data sets} \\ > 2/\sqrt{n} & \text{for large data sets} \end{cases}$

## VIF $\geq 10$ indicate multicollinearity

[SAS]

```
proc reg;
  model y= X1 X2 X3 / r influence vif;
run;
```

(43)