

STA 302 / 100)

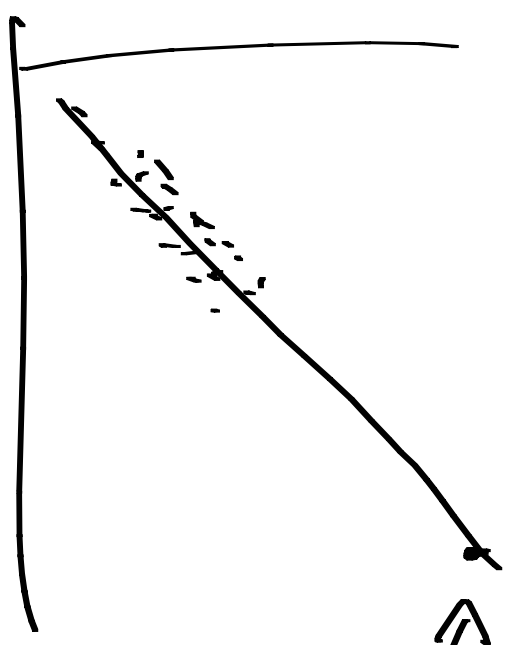
Tuesday Oct 18: Andriy will be here at 11:10 to

- return your assignments
- answer any questions you have about course material, practice problems, old tests, etc.

He'll also have some more office hours before the test.

Influential Points / Outliers / Leverage Points

Influential - take it out and makes a substantial difference in the est. of slope and intercept



This makes a large difference in R^2

\Rightarrow makes a substantial difference in fitted values

Point is not influential (a "good" leverage point)

Outliers - in x direction, the y -direction or
the y -residuals

If residuals \Rightarrow model failure

Outliers with respect to the explanatory variable
are called leverage points

They may be influential and possibly
cause the s.e. of est. of coefficients to be smaller than
they would be if point was removed
Sneeaker calls leverage points that follow

The pattern of the data "good" leverage points and leverage points that are influential are "bad" leverage points.



Quantifying leverage
" h_{ij} "

$$\hat{y}_i = b_0 + b_1 x_i$$

$$= \bar{y} + b_1 (x_i - \bar{x}) \quad b_1 = \frac{\sum_j (x_j - \bar{x}) y_j}{S_{xx}}$$

$$= \bar{y} + \frac{\sum_j [(x_j - \bar{x}) y_j]}{S_{xx}} (x_i - \bar{x})$$

$$= \frac{1}{n} \sum_j y_j + \frac{(x_i - \bar{x})}{S_{xx}} \sum_j (x_j - \bar{x}) y_j$$

$$= \sum_j \left\{ \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}} \right\} y_j$$

$$= \sum_{j=1}^n h_{ij} y_j \quad , \quad h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}}$$

why "h_{ij}" , h is for "hat"

Exercise: Show $\sum_{j=1}^n h_{ij} = 1$

$$\hat{y}_i = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$$

h_{ii} = leverage of i^{th} data point

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \quad \text{If } h_{ii} \text{ is large, } x_i \text{ is far from } \bar{x}$$

$$\text{Suppose } h_{ii} \approx 1, \sum_{j \neq i} h_{ij} \approx 0 \Rightarrow \hat{y}_i \approx y_i$$

$$\begin{aligned} \text{Average of } h_{ii} &= \frac{1}{n} \sum_i h_{ii} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right\} \\ &= \frac{1}{n} + \frac{1}{n} = \frac{2}{n} \end{aligned}$$

Rule of thumb for what is large for leverage part is $h_{ii} > 2 * \text{average value} = \frac{4}{n}$

OK large gap between h_{ii} for i^{th} observation and i^s value for all other observations

Dealing with leverage points

- Remove only if invalid; there's a good reason to say that i^s unlike other points
- Is model correct? Maybe a curvilinear model would be more appropriate.

Bad leverage points are an example of model failure
Good leverage points - decrease $s.e.$ of est of β_1 's because increase s_{xy}

- increase R^2
→ make you think you have a better fit than you actually do

Measuring Influence of i^{th} Observation

Notation: subscript (i) indicates i^{th} observation
_____ has been deleted from calculation

DFBETAS - difference in est. of β 's with and without i^{th} observation

$$DFBETAS_{ik} = \frac{b_k - b_{k(i)}}{s.e. \text{ of } b_{k(i)}}$$

$i = 1, \dots, n$
 $k = 0, 1 \rightarrow \text{slope}$
Intercept \leftarrow

$b_1(a_i)$ is est of slope for $n-1$ observations
(1st observation removed)

Rule of thumb | $DF_{\text{BERTS}} = k$ | $> 2/\sqrt{n}$ for
large data sets

> 1 for small
data sets.

OR separated by large gap
from the other values
of DF_{BERTS}

Game example: $2/\sqrt{23} = .417$

For NYC, $DEFFITS_{1,23} = .98$ \leftarrow Influential

DEFFITS_i - how ith predicted value changes with and without the inclusion of ith observation

$$DEFFITS_i = \frac{\hat{y}_i - \hat{y}_{i(i)}}{S.E. \text{ of } \hat{y}_{i(i)}}$$

$\hat{y}_{i(i)}$ is predicted at x_i calculated from $n-1$ observations (ith observation excluded)

Rule of Thumb $|DEFFITS_i| > 2\sqrt{\frac{2}{n}}$ for large set
> 1 for small data sets

DL group between DIFFITS_i and
DIFFITS for all other
observations

Crime Example $2\sqrt{\frac{1}{23}} = .59$, $\frac{2}{\sqrt{23}}$
(or $\frac{2}{\sqrt{23}}$)

Cobb's Distance

: Measure of influence of i^{th}
observation

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_{j(i)} - \hat{y}_j)^2}{2s^2}$$

$s^2 = \text{MSE}$

Can show $D_i = \frac{r_i^2 h_{ii}}{2(1-h_{ii})}$

r_i is the i^{th} standardized residual

Large D_i can be due to
 long residual or
 leverage or both

Rule of thumb $D_i > 4/n-2$

or gaps between D_i & other D_i 's

Some example $4/23-2 = .19$ ($D_i = .53$ for 24^{th})

Gauss-Markov assumptions:

$$E(e_i) = 0$$

$$\text{Var}(e_i) = \sigma^2 \text{ all } i$$

e_i, e_j uncorrelated, $i \neq j$

e_i 's are normally distributed.

Residuals $\hat{e}_i = y_i - \hat{y}_i = y_i - b_0 - b_1 x_i = y_i - \hat{y}_i$

What do we know about \hat{e}_i ?

$E(\hat{e}_i) = 0$ since $E(\hat{y}_i) = b_0$ and $E(y_i) = \beta_1$

Now $\sum_{i=1}^n \hat{e}_i = 0$, $\sum_{i=1}^n \hat{e}_i x_i = 0$, $\sum_{i=1}^n \hat{e}_i y_i = 0$

$\hat{\epsilon}_i$ as estimators of ϵ_i

$$E(\hat{\epsilon}_i) = 0$$

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2 (1 - h_{ii})$$

Calcⁿ on p. 61 of

- high leverage points have lower variance
I'll show this later using matrices

$\hat{\epsilon}_i$'s not uncorrelated (know $\sum \hat{\epsilon}_i = 0$)

could show $\text{Cov}(\hat{\epsilon}_i, \hat{\epsilon}_j) = -h_{ij}\sigma^2, i \neq j$

$\hat{\epsilon}_i$'s are linear combinations of y_i 's so they are normally distributed

They can be more closely to normally distributed than the e_i 's because of the cut

If no leverage points, h_{ii} is small, variances of \hat{e}_i 's close, lack of equal variance sometimes ignored.
To adjust for unequal variance, standardizing

$$r_i = \frac{e_i}{\sqrt{1-h_{ii}}}$$

SAs call this "Student's"

~~points~~

This is internal studentization.
Sometimes people use external studentization

$$\frac{e_i}{s_{(i)} \sqrt{1-h_{ii}}}$$

s.e. of \hat{e}_i for observation i with i removed

externally studentized residuals
have a t distribution
SPSS calls externally studentized
residuals "student"

— FND of TEST
CONVERSE —