

Handout 2

RECAP

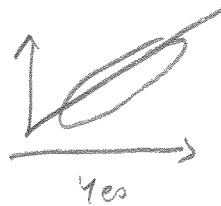
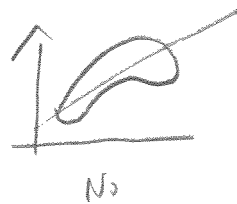
Summary of $y \sim x$ relationship : $y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{linear}} + \varepsilon_i$, ε_i uncorrelated
 $E(\varepsilon_i) = 0$
 $V(\varepsilon_i) = \sigma^2$ (constant)

\Rightarrow Fit model (β_0, β_1) to data using Least Squares

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

Basic assumptions

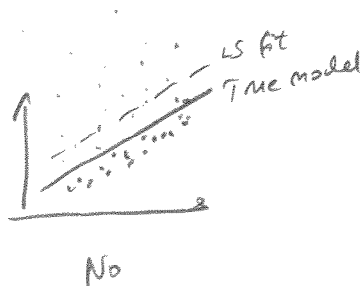
① Model sufficiency



Tools: graph

Fix: data transformations

② Symmetry of errors



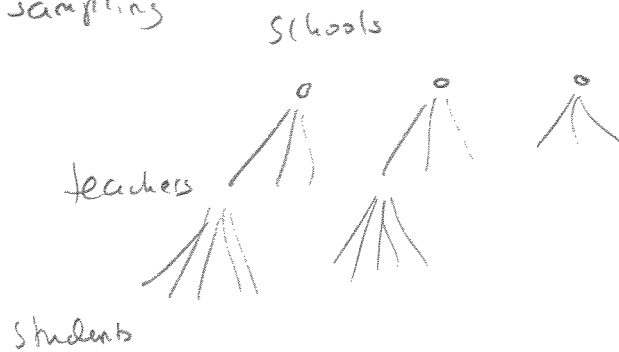
Tools: scatter, residual plots
QQ plot of residuals

Fix: data transformations
or
different fitting criteria
 \rightarrow median line
ML

③ Uncorrelated errors

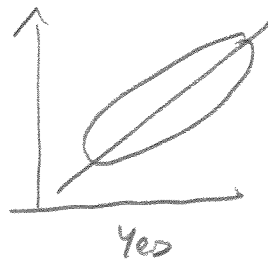
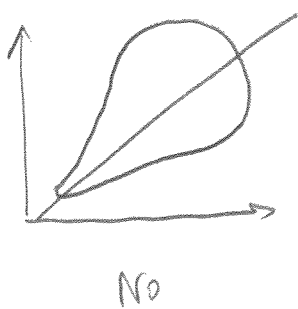
Think about the sampling process

- is there a time component
- are observations followed over time?
- Same person, unit of observation measured repeatedly
- Clustered sampling



If any of these → need to use other methods like time series, mixed effects modeling, ...

④ Constant error variance



Tools: scatter, residual plots

Fix: data transformations like $\log(y)$, \sqrt{y} , $y^{1/3}$

or

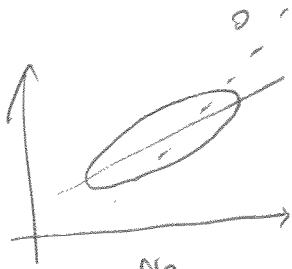
Weighted Least Squares (WLS)

$$\sum_{i=1}^n w_i (y_i - (\beta_0 + \beta_1 x_i))^2$$

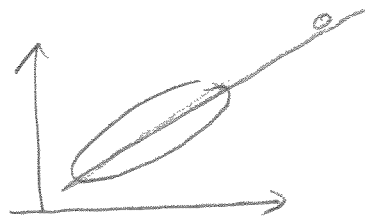
(more later)

⑤ No outliers

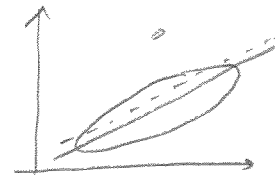
③



No
↓
bad fit



No
↓
fit looks "too good"



No
↓
you will overestimate
the error scatter
→ impact on inference
(more later)

But, careful about dropping too many observations.

Try to figure out why unusual.

Tools: scatter and residual plots.

- Other issues:
- groups in data
 - missing values
 - omitted variables
 - extrapolation
 - causal interpretation
 - ;
 - ,

More about least squares

(4)

$$Q = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \text{minimize wrt } (\beta_0, \beta_1)$$

$$\Rightarrow \begin{cases} \frac{\partial Q}{\partial \beta_0} = \sum_{i=1}^n -2 \cdot (y_i - (\beta_0 + \beta_1 x_i)) = 0 \\ \frac{\partial Q}{\partial \beta_1} = \sum_{i=1}^n -2x_i \cdot (y_i - (\beta_0 + \beta_1 x_i)) = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \sum_{i=1}^n y_i = n\beta_0 + \left(\sum_{i=1}^n x_i\right)\beta_1 \\ \sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i\right)\beta_0 + \left(\sum_{i=1}^n x_i^2\right)\beta_1 \end{cases}$$

The normal equations

Solve for β_0 : $\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$ $\left(\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i\right)$
eq (1)

Solve for β_1 : $\sum x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2$
eq (2)

w. β_0 solution
plugged in

$$\Rightarrow \hat{\beta}_1 = \frac{\sum x_i y_i - \frac{1}{n} (\sum x_i) (\sum y_i)}{\sum x_i^2 - \frac{1}{n} (\sum x_i)^2}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\equiv \frac{\text{cov}(x, y)}{\text{cov}(x, x)} = \underbrace{\text{correlation}(x, y)}_{\substack{\text{linear dependency} \\ \text{measure} \\ \in [-1, 1]}} \underbrace{\sqrt{\frac{V(y)}{V(x)}}}_{\substack{\text{scale} \\ \text{factor}}}$$

So, the slope parameter $\hat{\beta}_1 \sim \text{Correlation}(x, y)$

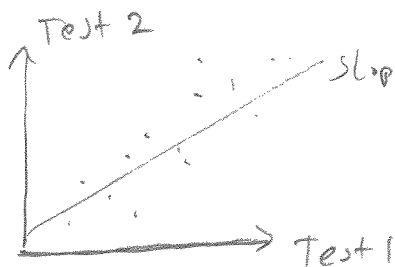
(5)

As a consequence;

- regression - linear relationship summary

- If y and x are standardized, $|\hat{\beta}_1| \leq 1$

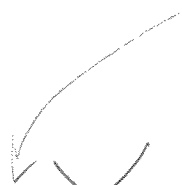
(classical misinterpretation)



Slope < 1 : meaning Test 2 score $<$ Test 1 score
on average

\Rightarrow "Best student - scores worse
on Test 2 - slacking off"

"Worst student - scores better
on Test 2 - got scared"



Wrong! "Regression towards the mean"

Best - nowhere to go but down

Worst - nowhere to go but up

More about the slope estimate

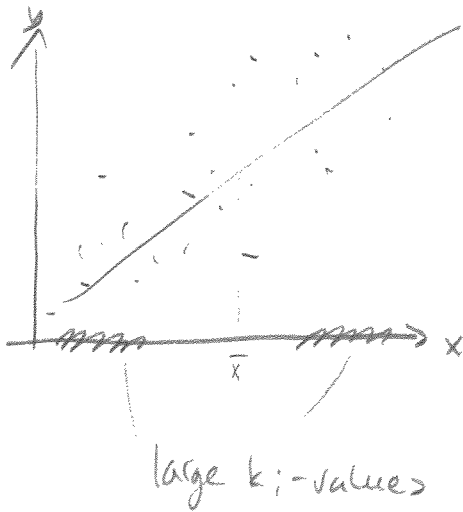
6

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n \left[\frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2} \right] y_i = \sum_{i=1}^n k_i \cdot y_i$$

• That is, $\hat{\beta}_1$ is a linear combination of y -values where weights $k_i = \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2}$ are large for y -values with

extreme x -values.

Note, the weights depend only on x so can be pre-computed (not data-dependent (y)) = a.k.a. linear model



Properties $\hat{\beta}_1 = \sum k_i \cdot y_i$

$$E(\hat{\beta}_1) = E(\sum k_i \cdot y_i) = \sum k_i E(y_i) = \sum k_i (\beta_0 + \beta_1 x_i)$$

$$= \sum k_i \cdot \beta_0 + \beta_1 \sum k_i \cdot x_i = \beta_1 \quad \checkmark$$

$$\sum k_i = \sum_i \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2} = 0$$

$$\sum k_i x_i = \frac{\sum (x_i - \bar{x}) x_i}{\sum_j (x_j - \bar{x})^2} = 1$$

So $\hat{\beta}_1$ is an unbiased estimate.

(7)

What about the variance?

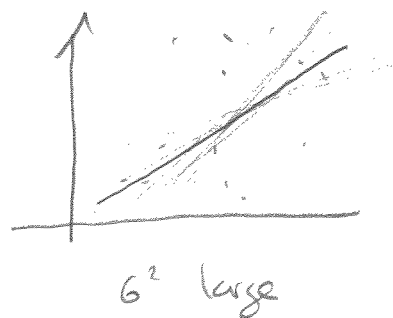
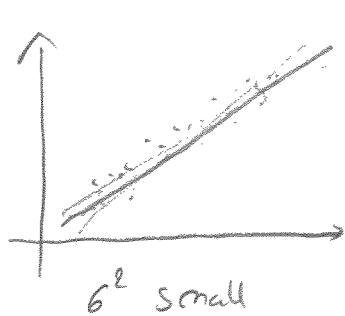
$$V(\hat{\beta}_1) = V(\sum k_i y_i) = \sum k_i^2 V(y_i) = \sum k_i^2 \sigma^2$$

↑
since y_i 's
are uncorrelated
(because ϵ_i are)

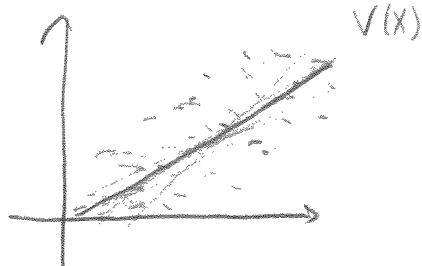
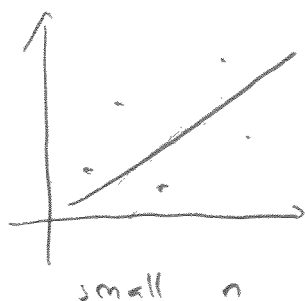
$$= \sum_i \frac{(x_i - \bar{x})^2}{(\sum_j (x_j - \bar{x})^2)^2} \sigma^2 = \frac{\sigma^2}{\sum_j (x_j - \bar{x})^2}$$

Meaning?

① $V(\hat{\beta}_1)$ increases with noise level in data (σ^2)

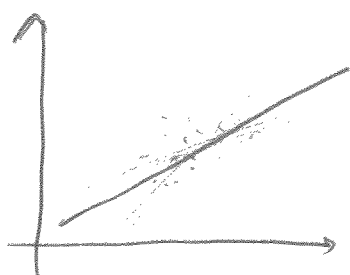


② $V(\hat{\beta}_1)$ decreases $\sim \frac{1}{n}$ since $V(\hat{\beta}_1) = \frac{\sigma^2}{n \left(\frac{1}{n} \sum_j (x_j - \bar{x})^2 \right)} = \frac{\sigma^2}{(n \cdot c)}$

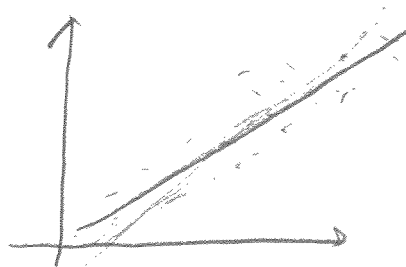


③ $V(\hat{\beta}_1)$ decreases w. variance of x

⑧



$\sum(x_i - \bar{x})^2$ small



$\sum(x_i - \bar{x})^2$ large

So, easier to estimate a $y \sim x$ relationship if range of x is wide - opportunity to see how y is affected.

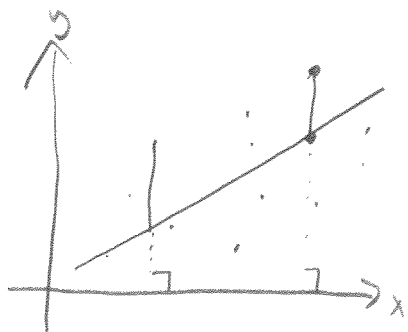
More properties & Diagnostics

• Residuals sum to 0 if intercept included in the model

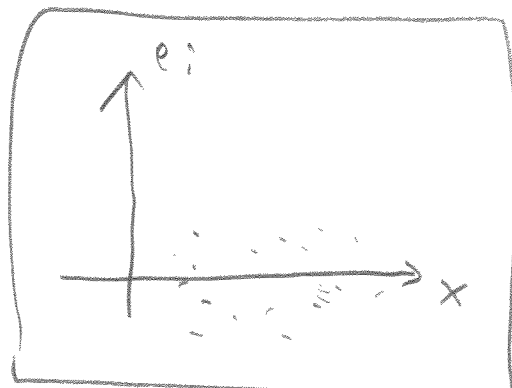
$$\sum e_i = \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)) = 0$$

• $\sum x_i e_i = 0$: "orthogonal projection"

Meaning - we 'used up' the linear information in x about y w. our LS fit.



→ diagnostic plot



Should see no pattern ↑ here if linear model sufficient

'proof'

$$e_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i = y_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x})$$

$$= y_i - \bar{y} - \frac{\sum_j (x_j - \bar{x}) y_j}{\sum_j (x_j - \bar{x})^2} (x_i - \bar{x})$$

$$\sum x_i e_i = \sum x_i y_i - n \bar{x} \bar{y} - \frac{\sum_j (x_j - \bar{x}) y_j}{\sum_j (x_j - \bar{x})^2} \sum x_i (x_i - \bar{x}) = 0$$

• $\sum \hat{y}_i e_i = 0$, fitted values \perp the residuals

\hat{y}_i = a linear function of x_i - part of y_i explainable from x

e_i = orthogonal to x_i - part of y_i not explainable from x

Side note - what about $\hat{\beta}_0$?

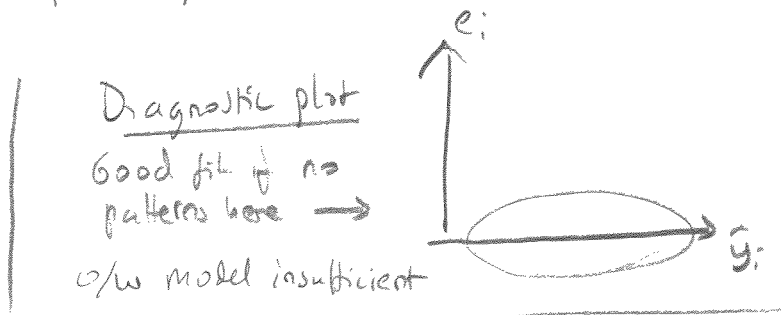
$$E(\hat{\beta}_0) = \beta_0$$

$$V(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_j - \bar{x})^2} \right)$$

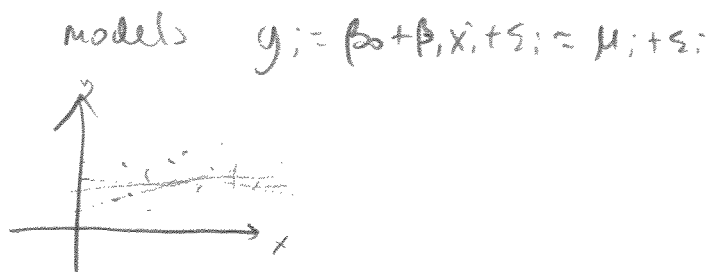
$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \bar{x}}{\sum (x_j - \bar{x})^2}$$

↑ if \bar{x}^2 large compared w spread $\sum (x_j - \bar{x})^2$, i.e. if x near constant

; large negative if x near constant



Why? Because if x near constant so we don't need 2 parameters



Fitted values

$$\begin{aligned}\hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x}) = \bar{y} + \sum_j k_j y_j (x_i - \bar{x}) \\ &= \sum_j \left(\frac{1}{n} y_j + k_j (x_i - \bar{x}) y_j \right) = \sum_j \left(\frac{1}{n} + k_j (x_i - \bar{x}) \right) y_j = \sum_j h_{ij} y_j \\ &= \text{a weighted average of } y\text{-values.}\end{aligned}$$

Note, if x is constant, all $h_{ij} = \frac{1}{n}$ and $\hat{y}_i = \bar{y}$!

Let's look at how much y_i contributes to its own fitted value, \hat{y}_i :

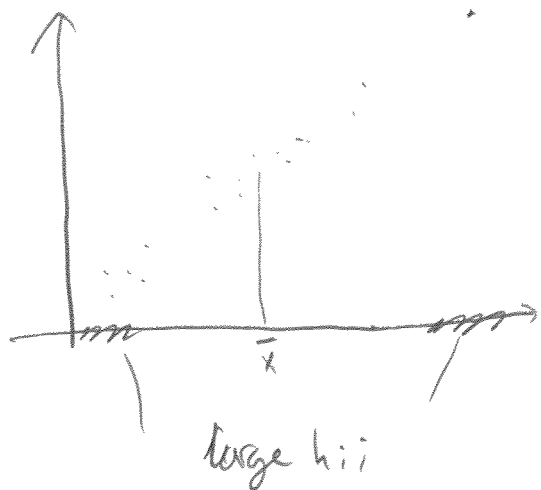
$$h_{ii} = \frac{1}{n} + k_i (x_i - \bar{x}) = \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_j (x_j - \bar{x})^2} \right] \quad \underline{\underline{\text{Leverage}}}$$

$$h_{ii} \sim \theta \left(\frac{1}{n} \right)$$

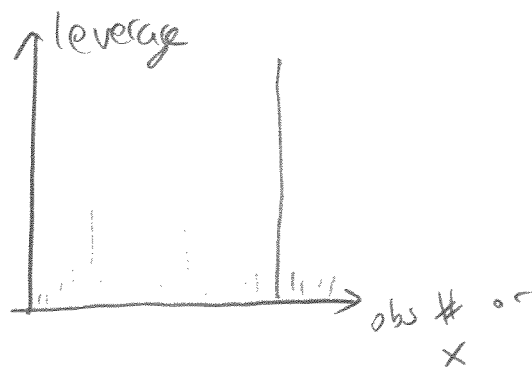
\sim is large if x_i 'extreme' - far from \bar{x}

\sim if $h_{ii} \gg h_{ij}$, \hat{y}_i dominated by y_{ii} .

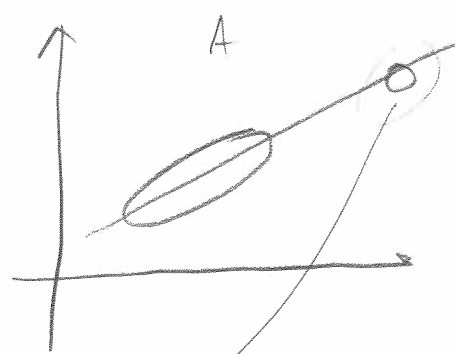
Discuss



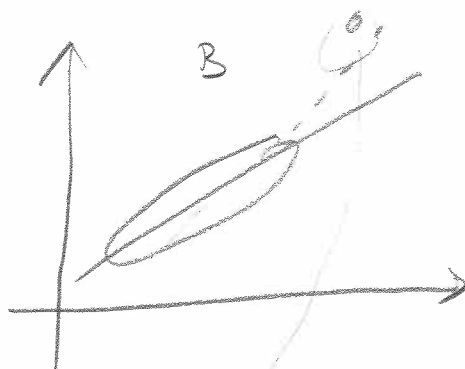
Diagnostic plot



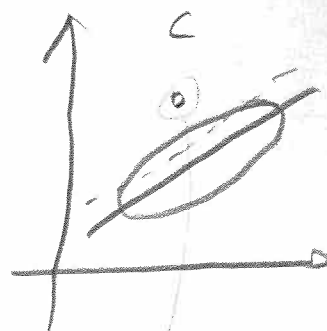
Problem If observation i has a high leverage value, a contamination (outlier) at this location can have a huge impact on the analysis (11)



high leverage,
but not
influential



high leverage
and influential



low leverage
but creates a
bias in fit
"pure outlier"

What do do?

(A) Do analysis w & w/o outlier - discuss impact

(B)

(C) - - - - - check if inference affected

Spotting the problem

(A) - leverage plot, but does not show up in residual!

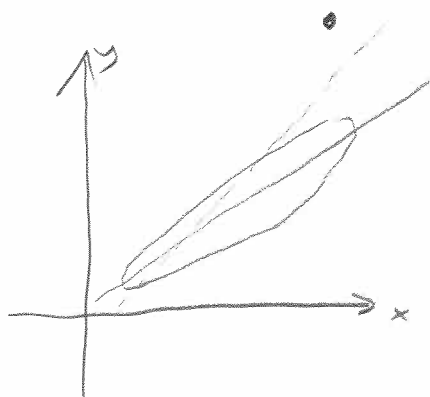
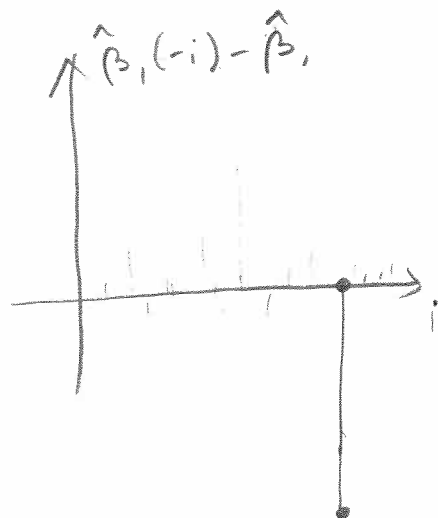
(B) - leverage & residual plot

(C) - residual plot

Other diagnostic tools

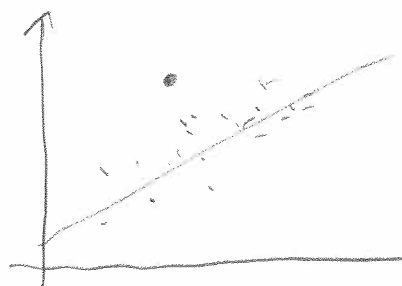
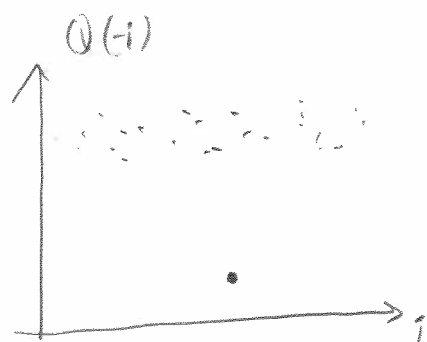
Dropping an observation - examine refit.

① Impact on slope



② Impact on LS criterion $\sum (y_i - (\hat{\beta}_0 - \hat{\beta}_1(x_i)))^2 = LS = Q$

vs $\sum (y_i - (\hat{\beta}_0(-i) - \hat{\beta}_1(-i)x_i))^2 = LS(-i) = Q(-i)$



③ Combining residuals, leverage & re-fitting ...

example Cook's D = $D_i = \left(\frac{e_i}{\hat{\sigma} \sqrt{1-h_{ii}}} \right)^2 \left(\frac{h_{ii}}{1-h_{ii}} \right) \cdot \frac{1}{2}$

(more later)