

Handout 2

①

RECAP

Summary $y \sim x$ relationship $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, ε_i uncorrelated, $E(\varepsilon_i) = 0$

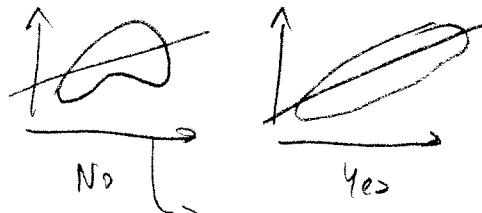
LS criterion $\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$

↑
Equal
Contribution of
all terms

squared error

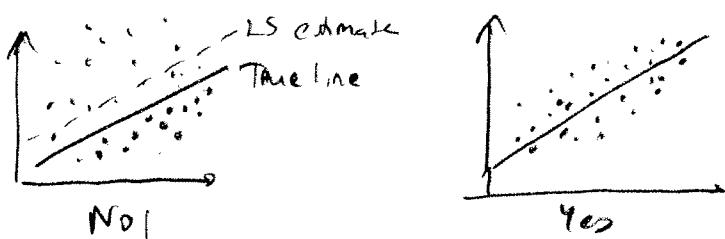
BASIC ASSUMPTIONS

① Model sufficiency



(check? Plot, plot, plot!)

② Symmetry of error scatter



(check? • QQplot of residuals
from LS fit.)

Solution? Transformation may help
→ also use a different fitting
criterion

Scatter plot y on x
may already show a
clear asymmetry of course.

③ Unrelated errors

(check? Think about the sampling process.)

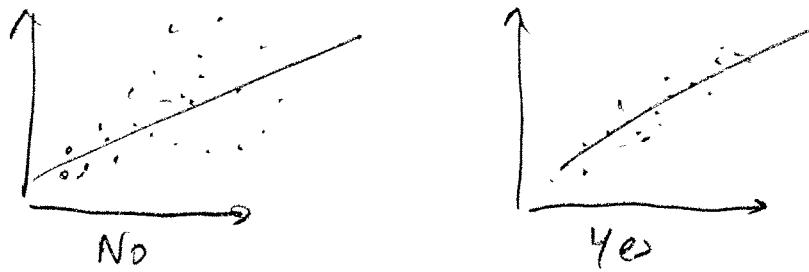
- Is the data sampled sequentially over time? \Rightarrow Time Series
 - Or samples nested? (e.g. students from the same class)
- Solution \Rightarrow Mixed effects or hierarchical linear model

Solution:

patients from the same hospital

④

(4) Constant Variance of error scatter

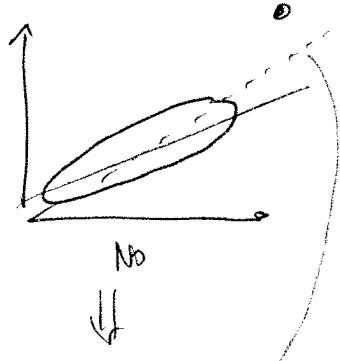


Check? Plot, plot, plot; both y vs x and residuals

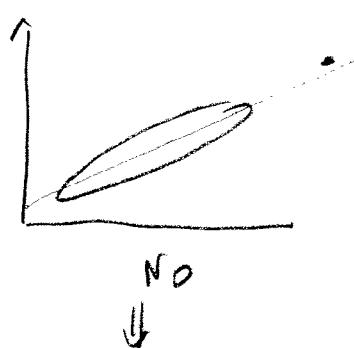
Solution \rightarrow Try some data transformations $\log(y), \sqrt{y}, \sqrt[3]{y}, \frac{1}{y} \dots$
 \Rightarrow If that doesn't work \Rightarrow Weighted Least Squares

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

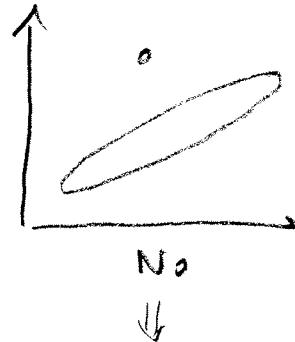
(5) No outliers



Effect: bad fit



Effect: fit looks
"too good"



Effect: you will
overestimate the error
scatter (more on this later)

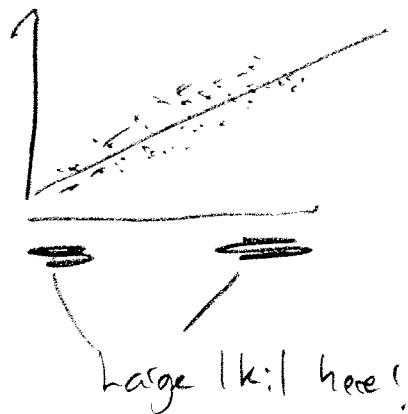
LS estimates \rightarrow Properties

③

Have $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$, $\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_j (x_j - \bar{x})^2} = \sum_i k_i y_i$

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

Which observations contribute? The ones w. x_i far from \bar{x} !



Note, LS estimates are unbiased

$$E(\hat{\beta}_1) = E\left(\sum_i k_i y_i\right) = \sum_i k_i E(y_i) = \sum_i k_i (\beta_0 + \beta_1 x_i) \quad \left[\text{since } E(y_i) = \beta_0 + \beta_1 x_i + E(\epsilon_i) \right]$$

$$\stackrel{?}{=} \left(\sum_i k_i \right) \cdot \beta_0 + \beta_1 \left(\sum_i k_i x_i \right) = \beta_1 \quad \checkmark$$

$$\sum_i k_i = \sum_i \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2} = 0 \qquad \sum_i k_i x_i = \sum_i \frac{(x_i - \bar{x})^2}{\sum_j (x_j - \bar{x})^2} = 1$$

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

So, if we obtain many new data sets (x, y) , on average we'll obtain the true regression line.

What about the variance?

since $V(y_i) \cdot V(\varepsilon_i) = \sigma^2$

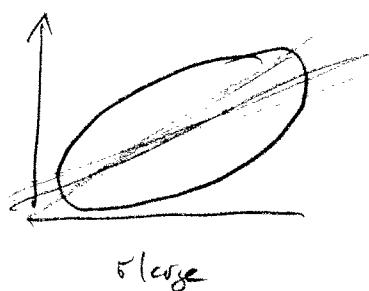
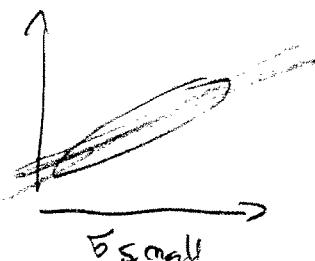
$$V(\hat{\beta}_1) = V\left(\sum b_i y_i\right) = \sum b_i^2 V(y_i) = \sum b_i^2 \sigma^2$$

↑
Since y_i 's
are independent
(because ε_i 's 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 the noise level

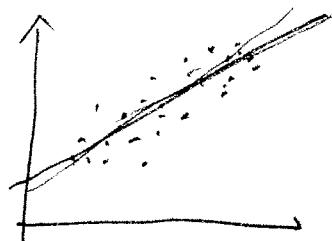
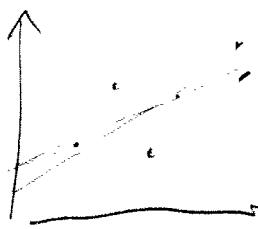
→ more uncertainty, more difficult to identify the slope



② $V(\hat{\beta}_1)$ decreases $\sim \frac{1}{n}$ w. the sample size n

$\sum_{j=1}^n (x_j - \bar{x})^2 = \text{sum of } n \text{ terms} \approx n \cdot c$, c some constant

⇒ More data, easier to identify the slope.



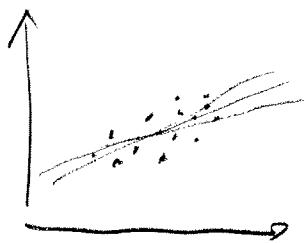
n small

n large

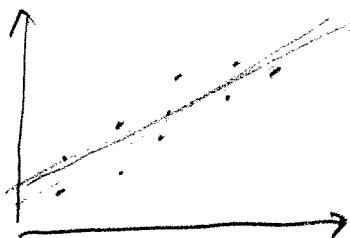
③ $V(\hat{\beta}_i)$ decreases w. the variance of x

(5)

\Rightarrow The more spread in x we have, the easier it will be to identify the $y|x$ relationship



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



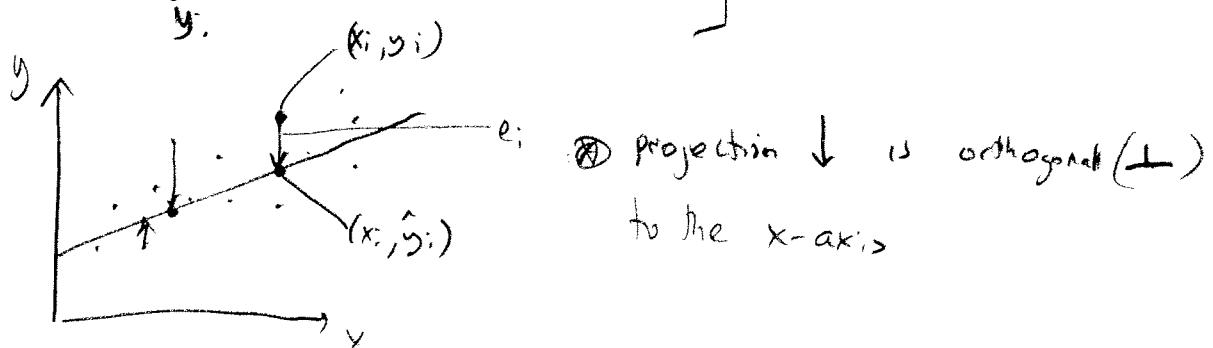
$$\sum(x_i - \bar{x})^2 \text{ large}$$

More properties

- Residuals sum to 0, $\sum e_i = 0$ (if β_0 included in the model)
- $\sum x_i e_i = 0$, "orthogonal projection"

Meaning \Rightarrow We used up all the linear information in x about y w. our LS fit

$$\left[e_i = y_i - \hat{y}_i - \underbrace{\left(\frac{\sum(x_j - \bar{x}) y_j}{\sum(x_j - \bar{x})^2} \right) (x_i - \bar{x})}_{\hat{y}_i}, \sum x_i e_i = \sum x_i y_i - n \bar{x} \bar{y} - \sum x_i \bar{y} - \underbrace{\sum x_i (x_i - \bar{x})^2}_{\sum(x_i - \bar{x})^2} \right] = 0$$



More properties

(6)

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

\hat{y}_i = linear function of x_i , the part that's explainable from x

e_i = orthogonal to x_i , the part of y_i that's not explainable from x

Side note $V(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_j - \bar{x})^2} \right)$ increased if \bar{x}^2 large compared w. spread $\sum(x_j - \bar{x})^2$
 \rightarrow meaning if x near constant

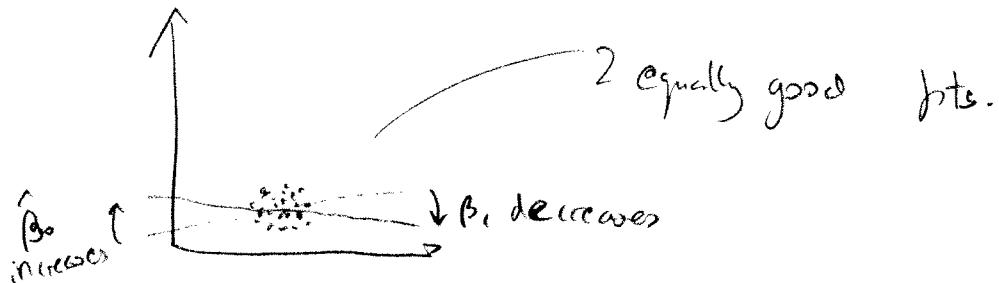
Also $\text{cor}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\bar{x} \sigma^2}{\sum(x_j - \bar{x})^2}$, large negative if x near constant

So - i) x near constant, the model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \approx \mu + \varepsilon_i$$

and we don't need 2 parameters β_0, β_1 .

Equally good fits are obtained w. lots of regression line, meaning there is a lot of uncertainty in the estimates of β_0, β_1 .



What about the fitted values?

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1(x_i - \bar{x}) = \bar{y} + \left(\sum_j h_{ij} y_j \right) (x_i - \bar{x})$$
$$= \sum_j \left(\frac{1}{n} y_j + h_{ij}(x_i - \bar{x}) y_j \right) = \sum_j \left(\frac{1}{n} + h_{ij}(x_i - \bar{x}) \right) y_j = \sum_j h_{ij} y_j$$

= a weighted average of y -values.

Note, i) x is constant $h_{ii} = \frac{1}{n}$ since $x_i = \bar{x} \forall i$

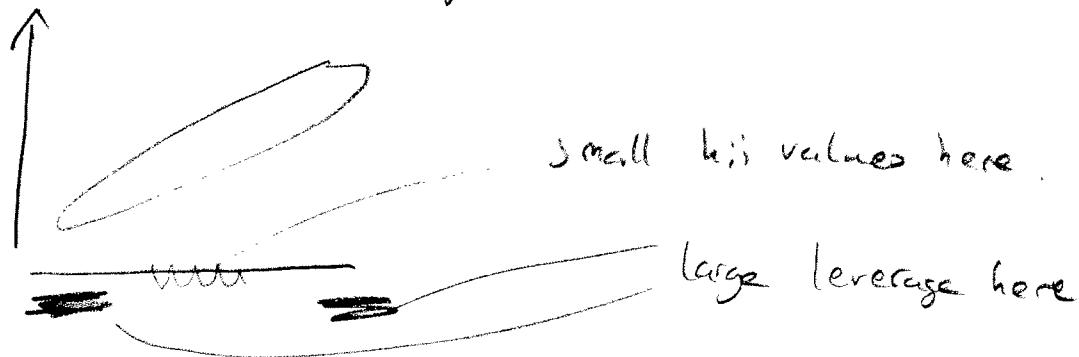
$\Rightarrow \hat{y}_i = \bar{y}$ if x has no information about y .

Observation i contributes h_{ii} to its own fitted value

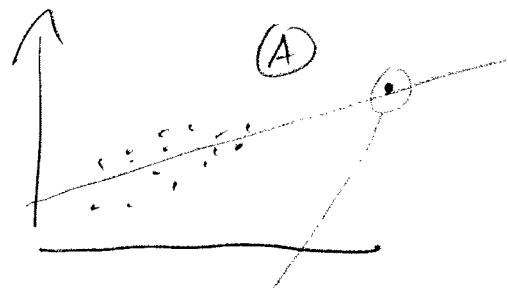
\rightarrow Extreme $h_{ii} = 1$ $h_{ij} = 0 \forall j \neq i \Rightarrow \hat{y}_i = y_i$

Now we use no part of the other observations
to help learn the $y \sim x$ relationship

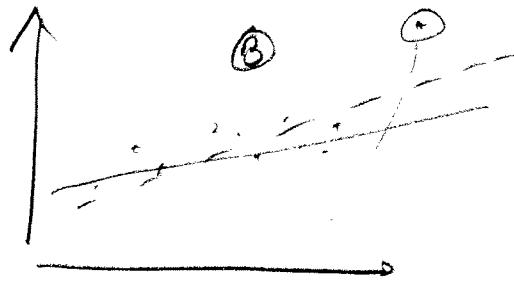
$$h_{ii} \text{ "the leverage"} = \left(\frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} + \frac{1}{n} \right)$$



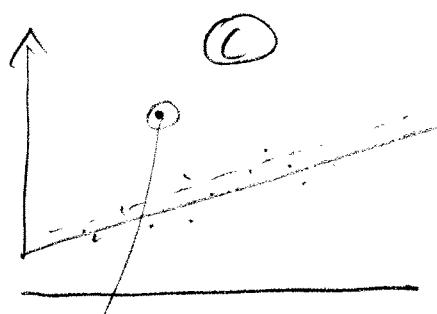
Problem: ②) observation i has high leverage, a contamination (or outlier) at this location can have a huge impact on the analysis.



High leverage, but
not influential



High leverage, and
influential



Low leverage, but
creates a bias in the
fit. "Pure outlier"

↑

Case A: Do the analysis with and without the observation ①
& B included \rightarrow discuss impact

Case C: Do the analysis with & without ... ③

\rightarrow check if inferences are affected

Go back to the data source and see if you can figure out
why the outlier is present.

More about fitted values ;

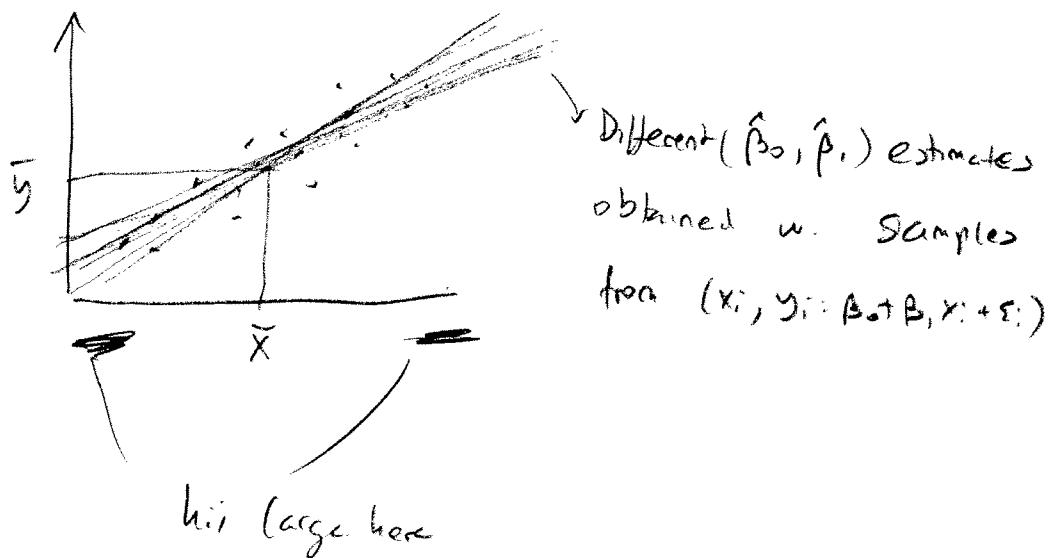
$$\boxed{V(\hat{y}_i) = \sigma^2 h_{ii}}$$

$$\left[V(\hat{y}_i) = V(\hat{\beta}_0 + \hat{\beta}_1 x_i) = V(\bar{y} + \hat{\beta}_1(x_i - \bar{x})) = V(\bar{y}) + (x_i - \bar{x})^2 V(\hat{\beta}_1) \right]$$

$$\left(\begin{array}{l} \text{Since } Cov(\bar{y}, \hat{\beta}_1) = Cov\left(\sum_i y_i, \sum_j k_j \beta_j\right) \\ = \sum_j \frac{1}{n} k_j V(y_j) = \frac{\sigma^2}{n} \sum k_j = 0 \end{array} \right)$$

$$= \frac{\sigma^2}{n} + \frac{\sigma^2 (x_i - \bar{x})^2}{\sum_j (x_j - \bar{x})^2} = \sigma^2 h_{ii}$$

Meaning ; At locations of high leverage , the regression line can fluctuate a lot from sample to sample



What about the residuals e_i ?

$$V(e_i) = \sigma^2(1-h_{ii})$$

$$\left[\begin{aligned} e_i &= y_i - \hat{y}_i = y_i - \sum h_{ij} y_j \\ \text{cov}(e_i, e_j) &= \text{cov}(y_i - \hat{y}_i, y_j - \hat{y}_j) = \text{cov}(y_i, y_j) + \text{cov}(\hat{y}_i, \hat{y}_j) - 2\text{cov}(y_i, \hat{y}_j) \\ &= \sigma^2 + \sigma^2 h_{ii} - 2\text{cov}(\hat{y}_i + e_i, \hat{y}_j) = \sigma^2 + \sigma^2 h_{ii} - 2\sigma^2 h_{ij} = \sigma^2(1-h_{ij}) \end{aligned} \right]$$

$\swarrow \searrow$
 $\left(\begin{array}{l} \perp \text{cov}(e_i, \hat{y}_j) = 0 \\ \text{LS properties} \end{array} \right)$

$$\begin{aligned} \text{Corr}(e_i, e_j) &= \text{cor}(y_i - \hat{y}_i, y_j - \hat{y}_j) = \text{cor}(y_i, y_j) + \text{cor}(\hat{y}_i, \hat{y}_j) - \text{cor}(y_i, \hat{y}_j) \\ &\stackrel{=} {=} -\text{cor}(y_i, \hat{y}_j) \\ &= \text{cor}\left(\sum h_{ik} y_k, \sum h_{jk} y_k\right) = \text{cor}(y_i, \sum h_{jk} y_k) - \text{cor}(y_j, \sum h_{ik} y_k) \\ &\stackrel{=} {=} \sum h_{ik} h_{jk} \sigma^2 - h_{ij} \sigma^2 - h_{ij} \sigma^2 = \sigma^2 h_{ij} - 2\sigma^2 h_{ij} = -\sigma^2 h_{ij} \end{aligned}$$

- MEANING :
- ① Residuals are correlated (but errors ε_i are not)
 - ② residuals have nonconstant variance $V(e_i) = \sigma^2(1-h_{ii})$ (since $V(\varepsilon_i) = \sigma^2$)
 - ③ residual variance is small in high leverage regions, since observations here drive the fit

Note residuals are not directly comparable since their variances differ

⇒ better to look @ and use

Standardized residuals

$$\tilde{e}_i = \frac{e_i}{\sqrt{1-h_{ii}}}$$

Note $V(\tilde{e}_i) = \sigma^2$

So far

$$\left\{ \begin{array}{l} E(\hat{\beta}_i) = \beta_i, \quad V(\hat{\beta}_i) = \frac{\sigma^2}{\sum(x_j - \bar{x})^2} \\ E(\hat{y}_i) = \beta_0 + \beta_1 x_i, \quad V(\hat{y}_i) = \sigma^2 h_{ii} \\ E(e_i) = 0, \quad V(e_i) = \sigma^2 (1-h_{ii}) \end{array} \right.$$

Note, all the variances involve one extra parameter

$\boxed{\sigma^2}$

If we don't know it, we can't make inferences

⇒ We fit the data to estimate this "nuisance" parameter as well.

(12)

- RSS = residual sum of squares (or SSE error sum of squares)

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- MSE = mean squared error $= \frac{\text{RSS}}{n-2} = \hat{\sigma}^2$

is an unbiased estimate of the error variance σ^2

Note denominator $n-2$; we've already used the data for two parameter estimates ($\hat{\beta}_0$ & $\hat{\beta}_1$) so only $n-2$ degrees of freedom remains

Next lecture $\Rightarrow R^2$, and the F & t-tests.