

HANDOUT 1

(1)

Simple linear regression

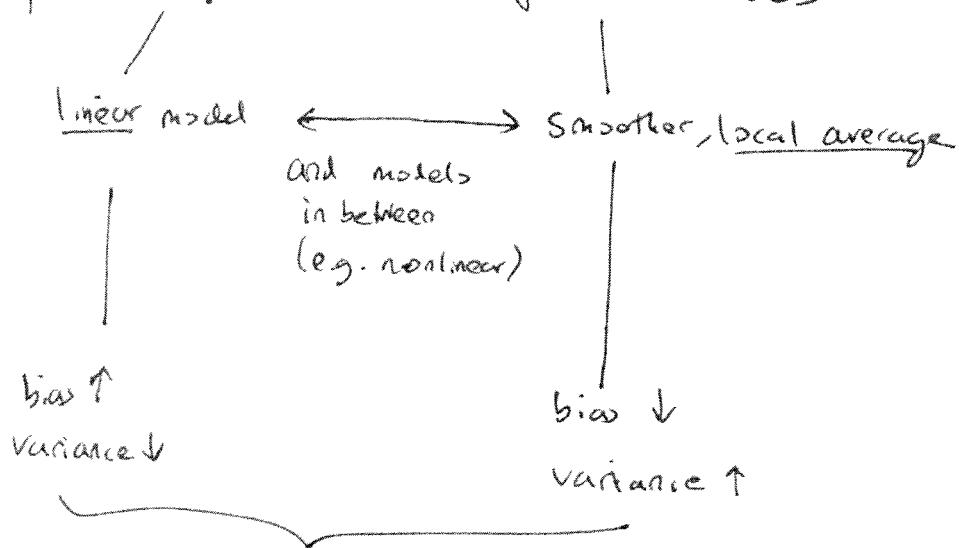
- Summarize relationship between predictor variable y and explanatory variable x :

$$y = f(x) + \varepsilon$$

/
'Explainable'
part of y

random error, scatter
'unexplainable' part

Trade-off: simple vs. flexible models



Why? Simple models use all of the data to generate the regression line/curve, whereas the local methods use only a subset of the data for each local decision: more data \Rightarrow lower variance.

(2)

Typically, we fit a model to the data using Least squares LS.

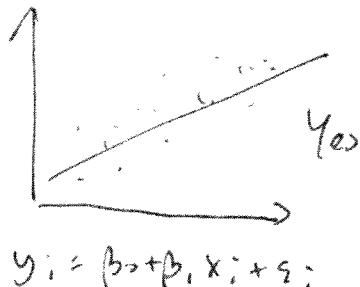
$$\min_{\beta_0, \beta_1} Q, Q = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

all observations contribute equally

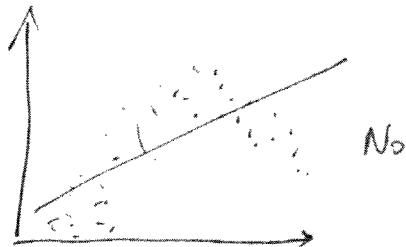
errors are squared, so positive & negative errors (above & below the regression line $\beta_0 + \beta_1 x$) are equally important.

BASIC ASSUMPTIONS BEHIND THE LS FIT

① Linear model is sufficient to describe the data.

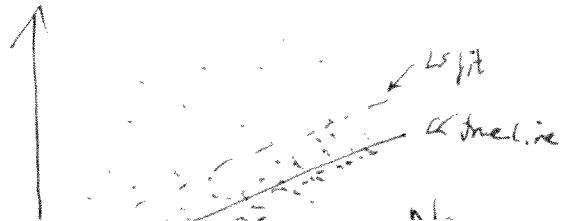
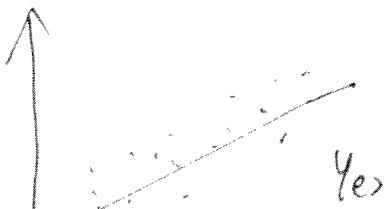


$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$



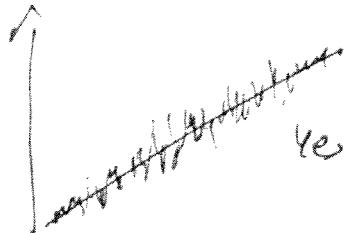
O/w — regression line is an inadequate summary of the $y \times x$ relationship.

② Symmetric errors (since we use $(\)^2$, letting positive & negative errors contribute equally)

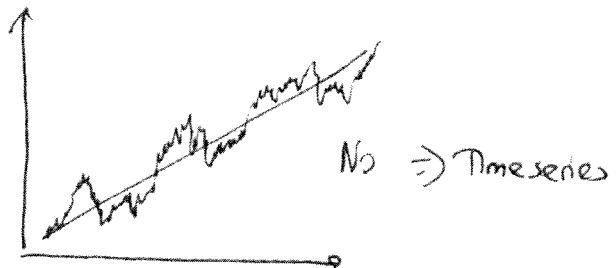


(3)

- ③ Uncorrelated errors (OLS shouldn't be using a simple sum $\sum_{i=1}^n$)

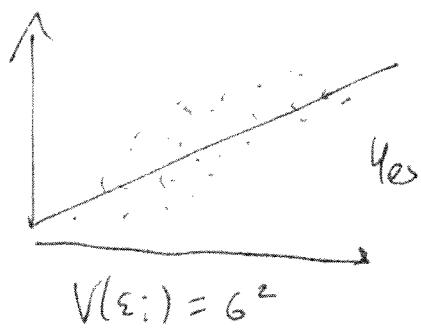


$$\text{Corr}(\varepsilon_i, \varepsilon_j) = 0$$

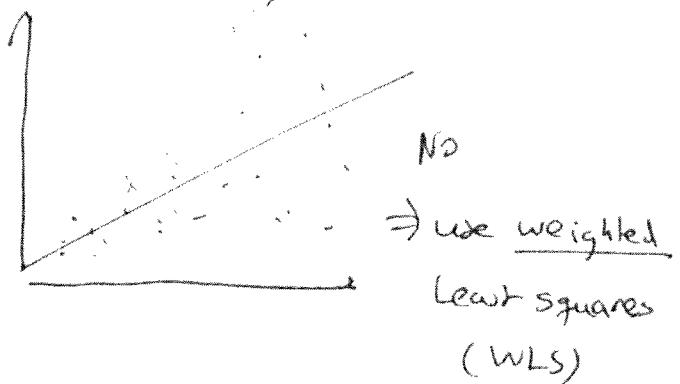


No \Rightarrow Time series

- ④ Constant variance (OLS shouldn't let all observations contribute equally; large variance = less information about line)

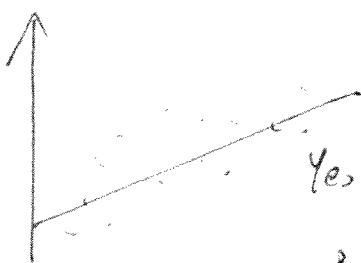


$$V(\varepsilon_i) = \sigma^2$$



No
 \Rightarrow use weighted least squares (WLS)

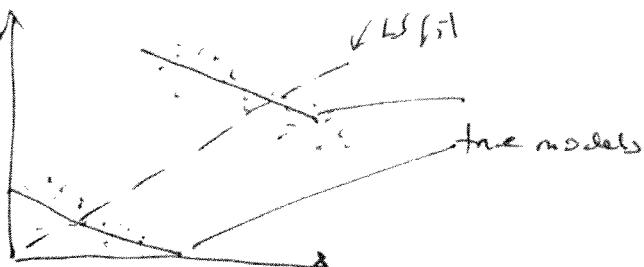
- ⑤ No outliers (OLS shouldn't let all observations contribute equally)



(4)

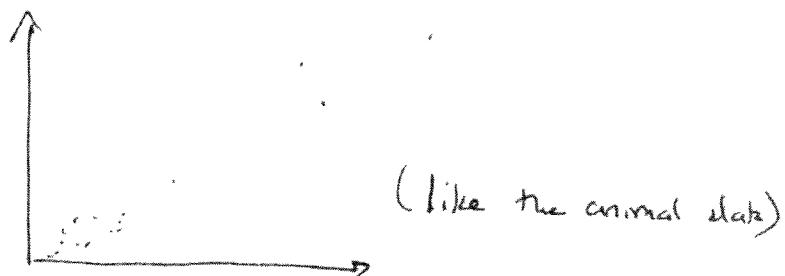
Other things to look out for

- groups in data



LM not sufficient.

- uneven spread in X



Sometimes [①-⑤] doesn't hold for the data, but a simple transformation (like $\log()$, $\sqrt{()}$, $\frac{1}{()}$) can sort this out.

Graph! Graph! Graph! Try different transformations & models to get as close as possible to ①-⑤ before proceeding with your analysis.

(5)

Note: in regression we model the conditional distribution $y|x$ (not (x,y)) jointly. x is assumed known & fixed, only ε is random $\Rightarrow y$ is random through ε .

Extra assumption on ε : $\varepsilon \sim N(0, \sigma^2)$

\rightarrow makes some inferences easier.

(more later)

WARNING

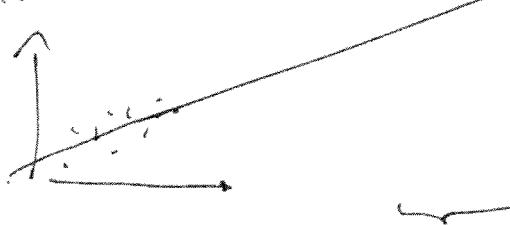
- Watch out for omitted variables.

• true: $y = f(z) + \varepsilon$ \Rightarrow model $y = h(x) + \varepsilon'$
 $x = g(z) + \eta$

~~obtains~~ \curvearrowleft artificial relationship

- Missing values - Why/how missing?

- extrapolation



\curvearrowleft (interpretation out here is invalid.)

- causal interpretation of regression line

- regression effect, regression toward the mean

(if an observation is the most extreme in x , it has nowhere to go but down in its rank in y).

(6)

LS estimates

$$Q = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad = \text{cost function to minimize}$$

$$\begin{cases} \frac{\partial Q}{\partial \beta_0} = \sum (y_i - (\beta_0 + \beta_1 x_i)) \cdot -2 = 0 \\ \frac{\partial Q}{\partial \beta_1} = \sum (y_i - (\beta_0 + \beta_1 x_i)) \cdot -2 x_i = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \sum y_i = n \beta_0 + (\sum x_i) \beta_1 & (1) \\ \sum x_i y_i = (\sum x_i) \beta_0 + (\sum x_i^2) \beta_1 & (2) \end{cases} \quad \text{"The Normal Equations"}$$

A) Solve for $\beta_0 \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ (where $\bar{y} = \frac{1}{n} \sum y_i$, $\bar{x} = \frac{1}{n} \sum x_i$)

B) Plug $\hat{\beta}_0$ into (2) $\Rightarrow \sum x_i y_i = (\sum x_i) (\bar{y} - \hat{\beta}_1 \bar{x}) + (\sum x_i^2) \hat{\beta}_1$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum x_i y_i - \sum x_i \sum y_i}{\sum x_i^2 - (\sum x_i)^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Note ; $\hat{\beta}_1 \equiv \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}} = \text{Correlation}(X, Y) \sqrt{\frac{\text{Var}(Y)}{\text{Var}(X)}}$
 covariance variance $\rho \in [-1, 1]$

(only $= 1 \rightarrow -1$ if $\Sigma = 0$ exactly)

Consider X, Y normalized such that $\text{Var}(Y) = \text{Var}(X) = 1$. Then $\boxed{\hat{\beta}_1 = \rho}$
 and $|\hat{\beta}_1| \leq 1$, slope always less than 1
 = regression effect.

Looking ahead - Matrix notation

(2)

$$\frac{\partial Q}{\partial \beta} = \begin{pmatrix} \frac{\partial Q}{\partial \beta_0} \\ \vdots \\ \frac{\partial Q}{\partial \beta_n} \end{pmatrix} \Rightarrow \begin{pmatrix} \sum y_i \\ \vdots \\ \sum x_i y_i \end{pmatrix} = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_n \end{pmatrix}$$

Normal Equations

B A

If A is invertible, solve for $\begin{pmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_n \end{pmatrix} = A^{-1} B$

$$= \frac{\begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix}}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

$$= \frac{\begin{pmatrix} \sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i \\ -\sum x_i \sum y_i + n \sum x_i y_i \end{pmatrix}}{n \sum x_i^2 - (\sum x_i)^2}$$

Now \Rightarrow Line 1: $\frac{\bar{y} \sum x_i^2 - \bar{x} \sum x_i y_i}{n \sum (x_i - \bar{x})^2} = \frac{\bar{y} \sum x_i^2}{n \sum (x_i - \bar{x})^2} - \hat{\beta}_0 \bar{x} - \bar{x} \bar{y} \frac{\sum x_i}{\sum (x_i - \bar{x})^2} = \bar{y} \left(\frac{\sum x_i^2 - \frac{1}{n} (\sum x_i)^2}{\sum (x_i - \bar{x})^2} \right) - \hat{\beta}_0 \bar{x}$

Line 2: $\frac{\sum x_i y_i - \sum x_i y_i}{n \sum (x_i - \bar{x})^2} = \frac{n(\bar{x} - \bar{y})(\bar{y} - \bar{y})}{n \sum (x_i - \bar{x})^2} = \hat{\beta}_1$ ✓ $= \bar{y} - \hat{\beta}_0 \bar{x}$ ✓

Matrix notation

(8)

Very important!

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i \quad i=1, \dots, n \Leftrightarrow \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\Leftrightarrow \underline{y} = \underline{\mathbf{X}} \underline{\beta} + \underline{\varepsilon}$$

$n \times 1 \quad n \times 2 \quad 2 \times 1 \quad n \times 1 \quad \text{vector / matrices}$

Multivariate case

$$y_{ij} = \beta_0 + \beta_1 x_{i1j} + \beta_2 x_{i2j} + \dots + \beta_{p-1} x_{ip-1j} + \varepsilon_{ij} \Leftrightarrow \begin{pmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1n} \\ y_{21} \\ y_{22} \\ \vdots \\ y_{2n} \\ \vdots \\ y_{m1} \\ y_{m2} \\ \vdots \\ y_{mn} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & x_{111} & x_{121} & \dots & x_{1p-11} \\ 1 & x_{211} & x_{221} & \dots & x_{2p-11} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m11} & x_{m21} & \dots & x_{mp-11} \end{pmatrix}}_{\mathbf{X}} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{1n} \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \vdots \\ \varepsilon_{2n} \\ \vdots \\ \varepsilon_{m1} \\ \varepsilon_{m2} \\ \vdots \\ \varepsilon_{mn} \end{pmatrix}$$

\mathbf{X} is called the
design matrix

$$\Leftrightarrow \underline{y} = \underline{\mathbf{X}} \underline{\beta} + \underline{\varepsilon}$$

$n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1$

$$\hat{Q} = \sum_{i=1}^n \left(\underbrace{y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{ip})}_{\text{error}} \right)^2 = \sum_{i=1}^n e_i^2 = (\underline{e}, \underline{e}_1 \dots \underline{e}_n) \begin{pmatrix} \underline{e}_1 \\ \underline{e}_2 \\ \vdots \\ \underline{e}_n \end{pmatrix}$$

@:

$$= \underline{e}' \underline{e}, \text{ where } \underline{e} = \begin{pmatrix} \underline{e}_1 \\ \underline{e}_2 \\ \vdots \\ \underline{e}_n \end{pmatrix}$$

and $\underline{e}' = \text{transpose of } \underline{e}$

(9)

Now, we also make errors $e_i = (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))$

In vector form $\underline{e} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \underline{y} - \underline{\mathbf{X}}\underline{\beta}$

$$\Rightarrow Q = \underline{e}' \underline{e} = (\underline{y} - \underline{\mathbf{X}}\underline{\beta})' (\underline{y} - \underline{\mathbf{X}}\underline{\beta})$$

We want to minimize Q with respect to $\underline{\beta}$.

\Rightarrow Matrix derivative

$$\frac{\partial Q}{\partial \underline{\beta}} = -2 \underline{\mathbf{X}}' (\underline{y} - \underline{\mathbf{X}}\underline{\beta}) = 0$$

//
inner derivative

$$\Rightarrow \boxed{(\underline{\mathbf{X}}'\underline{\mathbf{X}})\underline{\beta} = \underline{\mathbf{X}}'\underline{y}}$$

Normal equations on
matrix form.

"Interpretation" $(\underline{\mathbf{X}}'\underline{\mathbf{X}}) \sim \text{Cov}(\underline{\mathbf{X}})$, where diagonal of $\underline{\mathbf{X}}'\underline{\mathbf{X}}$ matrix is the $\text{Var}(x_1), \text{Var}(x_2), \dots$ and off-diagonal elements are the covariances $\text{Cov}(x_i, x_j), \dots$

So if some x 's are linearly dependent, or close to

$\rightarrow (\underline{\mathbf{X}}'\underline{\mathbf{X}})^{-1}$ may not exist, or be numerically unstable

$(\underline{\mathbf{X}}'\underline{y}) \sim \text{Cov}(\underline{\mathbf{X}}, \underline{y})$ — measure of linear dependence between the

$$\textcircled{1} \quad (\mathbf{X}'\mathbf{X})^{-1} \text{ exists } \Rightarrow \hat{\beta} = \underline{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'y}$$

\textcircled{10}

Problems: if $(\mathbf{X}'\mathbf{X})^{-1}$ numerically unstable, meaning x_i 's are related \Rightarrow difficult to interpret the meaning of the slopes $\hat{\beta}_1, \dots, \hat{\beta}_{p-1}$. They are not just measuring the contribution of one x_i , if x_i 's are related.

Back to simple linear regression \Rightarrow

Properties

- unbiased estimates

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

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

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

That is, $\hat{\beta}_1$ is a weighted average of y -values!

Note, extreme x_i 's, far from \bar{x} , contribute the most to the slope estimate

