

Last lecture  
 Sec. 11.4 Properties on  $r_{X,Y}$   $\begin{cases} r_{X,Y} = r_{Y,X} \\ r_{X,Y} \in [-1, 1] \end{cases}$   
 Sec. 11.5 the Error in Regression  $MSE = (1-r^2) \sigma_Y^2$   
 Sec. 12.1 Simple linear regression model  $\leftarrow$  "mathematical model"  
 (Least square linear regression)  $\leftarrow$   $\begin{cases} \text{define "good"} \\ \text{give best ans.} \end{cases}$

Today.  
 Sec. 12.2 Estimated slope in simple linear regression model

~~Sec. 12.3~~

Next week:

- Monday - Review
- Wednesday - Q & A
- Friday - Final, during lecture time.  
(almost same rule as MT)
- During MT revision "what to do before & during exam?"

Sec. 12.2 Distr. of estimated slope

Input  $\rightarrow$   $\boxed{\phantom{000}}$   $\rightarrow$  signal + error noise  $\rightarrow$  output

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, 2, 3, \dots, n$$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

From least squares  $(Y = aX + b)$   
 $\hat{\beta}_1$  (as  $a$  in least square)  
 $= \frac{\sum D_x D_y}{\sum D_x^2}$  (given by least square)

Use empirical distr. given by sample pts  $\{(x_i, Y_i) : i=1, 2, \dots, n\}$

In particular, the marginal of  $X$  is  $Unif(\mathbb{R}^n : i=1, 2, \dots, n)$

Under the empirical distr.

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$S_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

$$S_y^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\mathbb{E}(D_x D_y) = \mathbb{E}(X - \bar{X})(Y - \bar{Y})$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})$$

Note:  $x_i$ 's are constants

Estimated slope  $\hat{\beta}_1 = \frac{\mathbb{E}(D_x D_y)}{S_x^2}$

$$= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2}$$

Note:  $x_i$ 's are i.i.d. Normals

Weighted average of the slopes, with weights equal to squared horizontal distance from  $\bar{X}$

1 professional  $\times 5$   
 2 amateurs  $\times 1$

$$\frac{5 \times 4 + 2 \times 3}{5 + 1 + 1} = \frac{5}{3+1+1} \times 4 + \frac{1}{3+1+1} \times 2 + \frac{1}{3+1+1} \times 3$$

| n=5 |   |
|-----|---|
| X   | Y |
| 1   | 4 |
| 2   | 6 |
| 3   | 5 |

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$$

define  $S_x = \sum_{i=1}^n (x_i - \bar{X})^2$  known constants  
 $X_i^* = \frac{x_i - \bar{X}}{S_x}$

$\hat{\beta}_1$  is also a Normal.

Q Expectation of  $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$$

$$\mathbb{E} Y_i = \beta_0 + \beta_1 x_i$$

$$\mathbb{E} \bar{Y} = \mathbb{E} \frac{1}{n} \sum_{i=1}^n Y_i = \frac{1}{n} \sum_{i=1}^n \mathbb{E} Y_i$$

$$= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) = \beta_0 + \beta_1 \bar{X}$$

$$\mathbb{E}(Y_i - \bar{Y}) = (\beta_0 + \beta_1 x_i) - (\beta_0 + \beta_1 \bar{X}) = \beta_1 (x_i - \bar{X})$$

$$\mathbb{E} \hat{\beta}_1 = \frac{\mathbb{E} \sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{X}) \beta_1 (x_i - \bar{X})}{\sum_{i=1}^n (x_i - \bar{X})^2} = \beta_1$$
 - unbiased estimator

Q Var of  $\hat{\beta}_1$

Recall  $S_x = \sum_{i=1}^n (x_i - \bar{X})^2$  is const.

$$Var(\hat{\beta}_1) = Var\left(\frac{\sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})}{S_x}\right)$$

$$= \frac{1}{S_x^2} \cdot Var\left(\sum_{i=1}^n (x_i - \bar{X})(Y_i - \bar{Y})\right)$$

$$= \frac{1}{S_x^2} Var\left(\sum_{i=1}^n (x_i - \bar{X}) Y_i\right)$$

(Additivity of indep. R.V.s)

$$= \frac{1}{S_x^2} \sum_{i=1}^n Var((x_i - \bar{X}) Y_i)$$

$$= \frac{1}{S_x^2} \left(\sum_{i=1}^n (x_i - \bar{X})^2\right) \sigma^2$$

$$= \frac{\sigma^2}{S_x}$$

Note on  $Var(\hat{\beta}_1)$ :  $\sigma^2$  unknown parameter  
 $S_x = \sum_{i=1}^n (x_i - \bar{X})^2$  known fixed number  
 will be larger when # data pts gets larger  $\Rightarrow$  SE will be smaller

Estimate  $Var(\hat{\beta}_1)$ : need an estimation to  $\sigma^2$  - Variance of error/residual  
 $\hat{\sigma}^2 = MSE$  is empirical variance of error/residual

$$SD(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{S_x}}$$

$$SE(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{S_x}} \quad (S_x = \sum_{i=1}^n (x_i - \bar{X})^2)$$

When the SD of  $\hat{\beta}_1$  is estimated from data, we call it the SE (standard error)

Standardized slope - statistic that helps us give CI / perform test

$$T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}$$

Claim (w/o proof)  $T \sim N(0, 1)$  when sample size  $n$  is large

In this case, 95% CI?

$$P(-2 \leq T \leq 2) \approx 95\%$$

$$P(-) \leq \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \leq 2 \approx 95\%$$

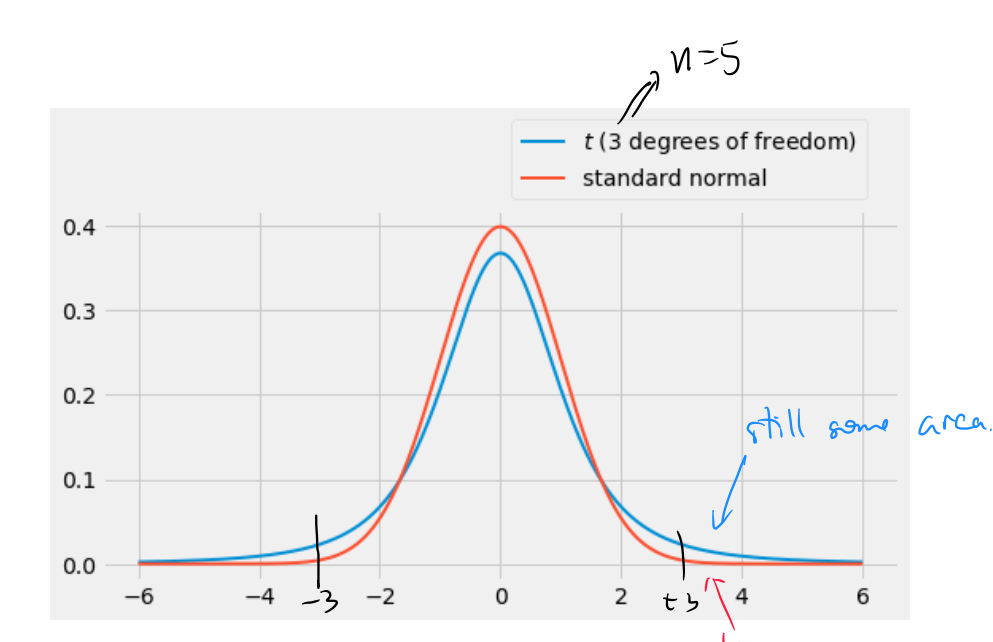
$$P(\hat{\beta}_1 - 2SE \leq \beta_1 \leq \hat{\beta}_1 + 2SE) \approx 95\%$$

$\Rightarrow$  95% CI for  $\beta_1$  is given by  $\hat{\beta}_1 \pm 2SE(\hat{\beta}_1)$ .

Test: (if  $X$  &  $Y$  are linearly associated?)  
 $H_0: \beta_1 = 0$   
 $H_a: \beta_1 \neq 0$  ( $H_a: \beta_1 > 0$ )  $\rightarrow$  Cannot use 95% CI  
 Quick way: use 95% CI we obtained  
 $\leftarrow$  if  $0 \in$  95% CI, accept  $H_0$   
 o/w. reject  $H_0$   
 Test statistic:  $T = \frac{\hat{\beta}_1 - \beta_1}{SE}$   
 Under  $H_0$ , and  $n$  large:  $T \sim N(0, 1)$

What would happen when  $n$  is not large enough?

For sample size  $n$   $T \sim t_{n-2}$   $\leftarrow$  t distr. with degree of freedom  $n-2$



Example

We see the following from computer program

|           |                       |   |
|-----------|-----------------------|---|
| $(x, y)$  | $n = 272$             |   |
| slope     | $= 1$                 | $\hat{\beta}_1$                         |
| intercept | $= 13$                | $\hat{\beta}_0$                         |
| $r$       | $= 0.6$               | correlation coefficient                 |
| $p$       | $= 2 \times 10^{-24}$ | p-value for the test $H_0: \beta_1 = 0$ |
| se slope  | $= 0.1$               | $SE(\hat{\beta}_1)$                     |

1) 95% CI of  $\hat{\beta}_1$  slope  $\pm 2 \cdot se\text{-slope} = [0.8, 1.2]$

2) SE of residual?

$$MSE = (1-r^2) \sigma_Y^2$$

$$\Rightarrow \hat{\sigma} = \sqrt{1-r^2} \sigma_Y = \sqrt{1-0.6^2} \cdot \sigma_Y$$

Alt. skls:

$$SE(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{S_x}}$$

SE-slope

$$\hat{\sigma} = se\text{-slope} \cdot \sqrt{S_x}$$

obtainable from program  
 i.e. SD of all  $y$ 's  
 Var of all  $x$ 's:  $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$