

# Simple Linear Regression

Wednesday, September 1, 2021 6:23 PM

- Overview
- Model

$$Y \approx \beta_0 + \beta_1 X$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

- Estimation

- Training data:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$RSS = e_1^2 + e_2^2 + \dots + e_n^2 \quad e_i = y_i - \hat{y}_i$$

↑  
Residual sum of squares

$$\min_{\hat{\beta}_0, \hat{\beta}_1} RSS$$

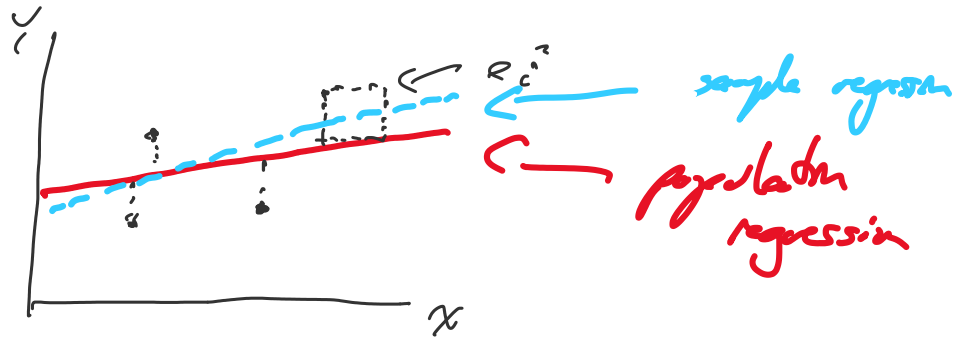
$$\therefore \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\therefore \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

- Analysis

↳ True relationship:

$$Y = \underbrace{\beta_0 + \beta_1 x}_{f(x)} + \varepsilon$$



\* Bias: Sample statistic vs Population value

e.g. mean  $\rightarrow \hat{\mu}$  vs.  $\mu$   
as  $n \rightarrow \infty$  in my sample, does  $\hat{\mu} \rightarrow \mu$ ?

• Variance of my estimate:

e.g.  $\text{Var}(\hat{\mu}) = SE(\hat{\mu})^2 = \frac{\sigma^2}{n}$   $\leftarrow$  population variance

$$SE(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

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

}  $\sigma^2 = \text{Var}(\epsilon)$

① Confidence Intervals:

-  $\alpha\%$  confidence interval  $\rightarrow$  "true" value is within a defined range with  $\alpha$  percent probability

$$\hat{\beta}_1 \pm \underline{2 SE(\hat{\beta}_1)} \quad \text{the } \underline{95\%} \text{ prob that}$$

actual  $\beta_1$  is within  $[\hat{\beta}_1 - 2SE(\hat{\beta}_1), \hat{\beta}_1 + 2SE(\hat{\beta}_1)]$

\* Assumption that  $\varepsilon$  is Gaussian / Normal  
 ↳ with low  $n$ , Student-t

↳ Hypothesis Tests

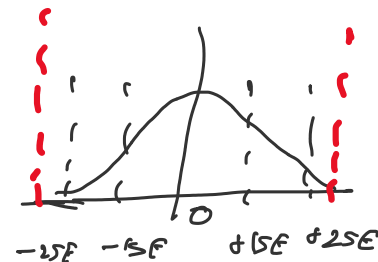
- $H_0$ : Null Hypothesis, No Relationship
- $H_a$ : Alternative, Is Relationship

$$H_0: \beta_1 = 0$$

$$H_a: \beta_1 \neq 0$$

\* T-stat:

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$



t-dist with  $n-2$  df  $\rightarrow$  probability of observing  $\hat{\beta}_1 \neq 0$

p-value  $\rightarrow$   $< 0.05$   $> 95\%$

↳ Accuracy

- RSE: Estimate of  $sd(\varepsilon)$

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

↑  
adjust df

\* measure of "lack of fit"

-  $R^2 \Rightarrow$  Proportion of variance explained

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

$$TSS = \sum (y - \bar{y})^2 \leftarrow \text{total sum of squares}$$

$$* R^2 = r^2 \rightarrow r = \text{cor}(X, Y)$$

↳ single var case!