

# Linear Regression

Recap:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i=1, \dots, n$$

Using the LS principle, we tried to find estimates  $(\hat{\beta}_0, \hat{\beta}_1)$  of the true parameters  $(\beta_0, \beta_1)$ .

Specifically the LS estimators are:

$$\begin{aligned} (\hat{\beta}_0, \hat{\beta}_1) &= \underset{\beta_0, \beta_1}{\operatorname{argmin}} Q(\beta_0, \beta_1) \\ &= \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

We found a closed form:

$$\text{(random)} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \stackrel{\text{HW}}{=} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{(random)} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \left( \begin{array}{l} \text{b/c } \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \\ \hookrightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{array} \right)$$

## Properties of LS Estimators:

Are these tests any good?

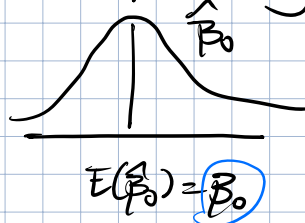
## Gauss - Markov Theorem:

Under the SLR model assumptions, the LS estimates  $\hat{\beta}_0$  &  $\hat{\beta}_1$  are:

i. unbiased for  $\beta_0$  &  $\beta_1$  respectively,

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

$$E(\hat{\beta}_1) = \beta_1$$



ii. the LS estimators are BLUE.

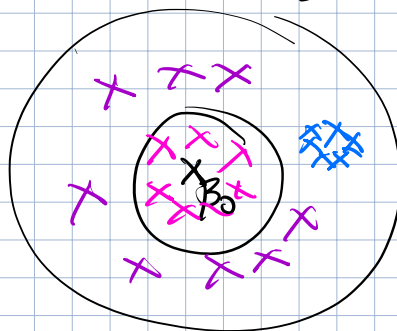
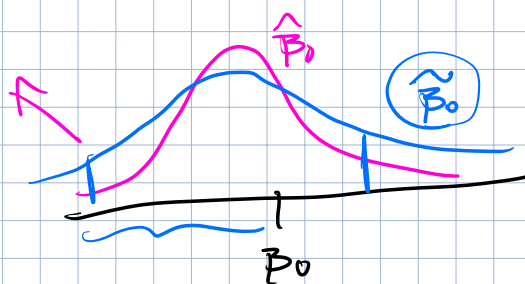
"best linear unbiased estimators"

i.e.

$$\text{Var}(\hat{\beta}_0) \leq \text{Var}(\tilde{\beta}_0)$$

$$\text{Var}(\hat{\beta}_1) \leq \text{Var}(\tilde{\beta}_1)$$

some other linear unbiased est



How do I know this?

i. how can I show  $E(\hat{\beta}_1) = \beta_1$ ?

WTS:

$$E(\hat{\beta}_1) = E\left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) = \beta_1$$

Def:  $SSX = \sum_{i=1}^n (x_i - \bar{x})^2$  (denom)  
 $SSXY = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

Lemma:

claim:  $SSXY = \sum_{i=1}^n (x_i - \bar{x}) y_i$

Why?

$$SSXY = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n [(x_i - \bar{x}) y_i - (x_i - \bar{x}) \bar{y}]$$

$$= \sum_{i=1}^n (x_i - \bar{x}) y_i - \sum_{i=1}^n (x_i - \bar{x}) \bar{y}$$

WTS = 0

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x}) \bar{y} &= \bar{y} \sum_{i=1}^n (x_i - \bar{x}) \\ &= \bar{y} \left( \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} \right) \\ &= \bar{y} (n\bar{x} - n\bar{x}) \\ &= \bar{y} (0) = 0 \end{aligned}$$

$$\left\{ \begin{array}{l} \sum_{i=1}^n a c_i = a \sum_{i=1}^n c_i \\ \bar{x} = \frac{\sum x_i}{n} \\ n\bar{x} = \sum x_i \end{array} \right.$$

$$\hat{\beta}_1 = \frac{SS_{XY}}{SS_X} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{SS_X}$$

$$= \frac{1}{SS_X} \sum_{i=1}^n (x_i - \bar{x}) y_i$$

$$= \sum_{i=1}^n \underbrace{\frac{(x_i - \bar{x})}{SS_X}}_{= k_i} y_i = \sum_{i=1}^n k_i y_i \quad (\text{linearity})$$

$$E(\hat{\beta}_1) = E\left(\sum_{i=1}^n k_i y_i\right) = \sum_{i=1}^n E(k_i y_i)$$

$$= \sum_{i=1}^n k_i E(y_i)$$

$$= \sum_{i=1}^n k_i E(\beta_0 + \beta_1 x_i + \varepsilon_i)$$

$$= \sum_{i=1}^n k_i [E(\beta_0) + E(\beta_1 x_i) + E(\varepsilon_i)]$$

$$= \sum_{i=1}^n k_i [\beta_0 + \beta_1 x_i + 0]$$

$$= \sum_{i=1}^n k_i \beta_0 + \sum_{i=1}^n \beta_1 k_i x_i$$

$$= \beta_0 \underbrace{\left(\sum_{i=1}^n k_i\right)}_{\substack{= 0 \\ \textcircled{A}}} + \beta_1 \underbrace{\left(\sum_{i=1}^n k_i x_i\right)}_{\substack{= 1 \\ \textcircled{B}}} = \beta_1$$

$$\textcircled{A} \sum_{i=1}^n k_i = \sum_{i=1}^n \frac{x_i - \bar{x}}{SSX} = \frac{1}{SSX} \sum_{i=1}^n (x_i - \bar{x}) = \frac{1}{SSX} (0) = 0 \quad \checkmark$$

$$\textcircled{B} \sum_{i=1}^n k_i x_i = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{SSX} \right) x_i$$

$$= \frac{1}{SSX} \sum_{i=1}^n (x_i - \bar{x}) x_i = \textcircled{SSX}$$

$$\frac{SSX}{SSX} = 1.$$

$$SSX = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})$$

$$= \sum_{i=1}^n [(x_i - \bar{x}) x_i - (x_i - \bar{x}) \bar{x}]$$

$$= \sum_{i=1}^n (x_i - \bar{x}) x_i - \sum_{i=1}^n (x_i - \bar{x}) \bar{x}$$

$$= \bar{x} \left( \sum_{i=1}^n (x_i - \bar{x}) \right)$$

$$= \bar{x} (\sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x})$$

$$= \bar{x} (n\bar{x} - n\bar{x})$$

$$= \bar{x} (n\bar{x} - n\bar{x}) = 0$$

That's the slope unbiasedness.

Y'all try the intercept.

## Part 2 of G-M Thm:

How can I show that

$$\underline{\text{Var}(\hat{\beta}_1)} \leq \text{Var}(\tilde{\beta}_1) \text{ for some other linear unbiased est } \tilde{\beta}_1?$$

Sketch:

$$\underline{\text{Var}(\hat{\beta}_1)} = \text{Var}\left(\sum_{i=1}^n k_i y_i\right) \quad \text{by covariance assumption.}$$

$$= \sum_{i=1}^n \text{Var}(k_i y_i)$$

$$= \sum_{i=1}^n k_i^2 \text{Var}(y_i)$$

$$= \sum_{i=1}^n k_i^2 \sigma^2$$

$$= \sigma^2 \sum_{i=1}^n k_i^2 = \frac{\sigma^2}{SSX}$$

Aside

$$\text{Var}(W+Z) = \text{Var}(W) + \text{Var}(Z) + 2\text{Cov}(W, Z)$$

$$\text{If } \text{Cov}(W, Z) = 0 \Rightarrow$$

$$\text{Var}(W+Z) = \text{Var}(W) + \text{Var}(Z)$$

What is  $\sum k_i^2$ ?

$$\sum_{i=1}^n \left( \frac{x_i - \bar{x}}{SSX} \right)^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{(SSX)^2}$$

$$= \frac{1}{(SSX)^2} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]$$

$$= \frac{SSX}{(SSX)^2} = \frac{1}{SSX}$$

I want to show  $\text{Var}(\tilde{\beta}_1) > \sigma^2 / \text{SSX} = \sigma^2 \sum_{i=1}^n k_i^2$

Write

$$\tilde{\beta}_1 = \sum_{i=1}^n \tilde{k}_i y_i = \sum_{i=1}^n (k_i + d_i) y_i$$

where

$d_i = \tilde{k}_i - k_i$   
for at least one  
 $i = 1, \dots, n$

$$\text{Var}(\tilde{\beta}_1) = \text{Var}\left(\sum_{i=1}^n (k_i + d_i) y_i\right)$$

$$= \sum_{i=1}^n \text{Var}((k_i + d_i) y_i)$$

$$= \sum_{i=1}^n (k_i + d_i)^2 \text{Var}(y_i)$$

$$= \sigma^2 \sum_{i=1}^n (k_i + d_i)^2$$

$$= \sigma^2 \left[ \sum_{i=1}^n (k_i^2 + 2k_i d_i + d_i^2) \right]$$

$$= \sigma^2 \left[ \sum_i k_i^2 + \underbrace{2 \sum_i k_i d_i}_{> 0} + \sum_i d_i^2 \right]$$

Ⓐ  $\sum_i k_i d_i = 0$  (try to show)

Ⓑ  $\sum_i d_i^2 > 0$  (easy to see)

not all  $d_i = 0 \Rightarrow$  one  $d_i^2 > 0 \Rightarrow \sum_i d_i^2 > 0$

## Sampling Distribution of $\hat{\beta}_0$ & $\hat{\beta}_1$ .

$$\left. \begin{array}{l} \text{Slope: } E(\hat{\beta}_1) = \beta_1 \\ \text{Var}(\hat{\beta}_1) = \sigma^2 / SSX \end{array} \right\} \text{under Li-iii)}$$

if I assume (iv)  $\Leftrightarrow$  normality of errors:

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / SSX) \quad \text{with (iv)}$$

$$? (\beta_1, \sigma^2 / SSX) \quad \text{without (iv).}$$

Intercept:

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

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SSX} \right) \quad (\text{see Cady notes / textbook})$$

Under Normality assumption:

$$\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SSX} \right)\right)$$

$$\textcircled{\text{or}} \quad ? \left( \beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SSX} \right) \right) \quad (\text{without iv.})$$



Consider the fitted values:

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

What is the distribution of  $\hat{y}_i$ ?

(under Normality is ok)

$$\hat{y}_i \sim N\left(\frac{\quad}{?}, \frac{\quad}{?}\right)$$