

# Linear Regression

8/25/25

Recap:

$$y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{signal}} + \underbrace{\varepsilon_i}_{\text{noise}} \quad i=1, \dots, n$$

Using the LS principle, we tried to find estimates  $(\hat{\beta}_0, \hat{\beta}_1)$  of the 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 the closed form:

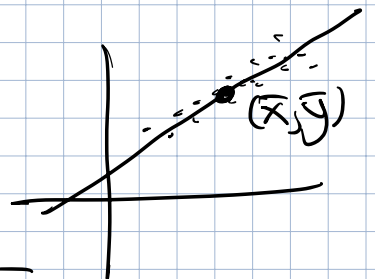
$$\overset{\text{random}}{\downarrow} \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \stackrel{\text{trick}}{=} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\overset{\text{random}}{\downarrow} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

↑ came from

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

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



## Properties of LS Ests:

Are these estimators any good?

## Gauss-Markov Theorem:

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

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

$$\Leftrightarrow E(\hat{\beta}_0) = \beta_0 \text{ \& } E(\hat{\beta}_1) = \beta_1$$

&

ii. the LS estimators are BLUE

= "best linear unbiased estimators" //

Focus: on  $\hat{\beta}_1$  for now

linear in  $y_i$

(linear)

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

(blue)

(best)

$$\text{Var}(\hat{\beta}_1) \leq \text{Var}(\tilde{\beta}_1) \quad \text{where } \tilde{\beta}_1 \text{ is some other linear unbiased estimator.}$$

How do I know this then is true?

Still focus on  $\hat{\beta}_1$  for now:

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

Want to show

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

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

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

Lemma:

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

Why?

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

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

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

want to show = 0.

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

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

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

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

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

$$= \sum_{i=1}^n \underbrace{\left[ \frac{(x_i - \bar{x})}{SS_X} \right]}_{=k_i} y_i = \sum_{i=1}^n k_i y_i$$

⚡ linearity in  $y_i$

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

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

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

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

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

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

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

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

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

(A) (B)

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

$$= \frac{1}{SS_X} (0) = 0$$

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

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

$$\sum x_i^2 - x_i \bar{x}$$

$$\sum x_i^2 - \bar{x} \sum x_i$$

$$\boxed{\sum x_i^2 - n\bar{x}^2}$$

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

$$\sum x_i^2 - 2\bar{x} \sum x_i + n\bar{x}^2$$

$$\sum x_i^2 - 2\bar{x} \sum x_i + n\bar{x}^2$$

$$\sum x_i^2 - 2n\bar{x}^2 + n\bar{x}^2$$

$$\boxed{\sum x_i^2 - n\bar{x}^2}$$

$$= \frac{SSX}{SSX} = \frac{\sum x_i^2 - n\bar{x}^2}{\sum x_i^2 - n\bar{x}^2} = 1$$

$$\Rightarrow E(\hat{\beta}_1) = \beta_0(0) + \beta_1(1) = \beta_1 \quad \square$$

Your turn:

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

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

but

$$\sum_{i=1}^n x_i^2 \neq n\bar{x}^2$$

## Part 2 of GM Thm:

Goal: If I consider some other linear unbiased estimator  $\tilde{\beta}_1$  then

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

How can I show this?

Sketch:

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

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

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

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

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

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

Consider other linear unbiased estimators  $\tilde{\beta}_1$ :

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

$\tilde{k}_i = k_i + d_i$   
 $\Leftrightarrow d_i = \tilde{k}_i - k_i$

$d_i \neq 0 \quad \forall i$   
 $\neq$  "for all"

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

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

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

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

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

if  $\geq 0$  then were done

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

$\text{(A)} = 0$        $\text{(B)} > 0$



Ⓐ  $\sum_i d_i k_i = 0$  try it. (Exercise)

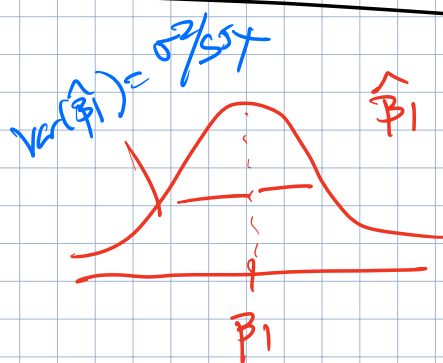
Ⓑ  $\sum_i d_i^2 > 0$  (sum of qty's that are 0 or 1)

↑ repeat for  $\beta_0$  ↓

Slope:

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

$$\text{Var}(\hat{\beta}_1) = \sigma^2 \sum_{i=1}^n k_i^2 = \sigma^2 / SS_X$$



If (iv):

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / SS_X)$$

(classical)

$$? (\beta_1, \sigma^2 / SS_X)$$

(i-iii. only)

Int:

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

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{SS_X} \right)$$

↖ see notes / textbook

If (iv), we know

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

(classical)

$$? (\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{SS_X} \right))$$

(i-iii. only)

Fitted values:

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

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

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

Think about it for next time!