

Statistical Properties of LS Estimators

Lecture 3

Alexandra Chronopoulou



COLLEGE OF LIBERAL ARTS & SCIENCES

Department of Statistics
101 Illini Hall, MC-374
725 S. Wright St.
Champaign, IL 61820-5710

Learning objectives

In this lecture we will:

- study properties of $(\hat{\beta}_0, \hat{\beta}_1)$ as an estimate of the true coefficient vector (β_0, β_1) .
- construct confidence/prediction intervals for (β_0, β_1) .

- Uppercase letters are normally used for **Random Variables**, and lowercase letters for **observed values** of the random variables.
- Uppercase letter will be also reserved for matrices.
- In some occasions lowercase letter will also be used for random variables.

LS Estimators Properties

$$\text{Model: } y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

Assumptions

The **errors** $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ are assumed to

- have **mean zero**: $\mathbb{E}(\varepsilon_i) = 0$
 - be **uncorrelated**: $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j$,
 - be **homoscedastic**: $\text{Var}(\varepsilon_i) = \sigma^2$ does not depend on i .
- ⇒ We can combine the last two and write it as

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij}$$

where $\delta_{ij} = 0$ if $i \neq j$.

Moments of $(Y|X)$

Based on the SLR model moment assumptions on the error terms, we have the following assumptions for the moments of Y conditioning on X :

$$\mathbb{E}(y_i|x_i) = \beta_0 + \beta_1 x_i$$

$$\text{Var}(y_i|x_i) = \sigma^2$$

$$\text{Cov}(y_i, y_j|x_i, x_j) = 0, \quad i \neq j$$

Remark:

When we evaluate expectation, only y_i 's are random and x_i 's are treated as known, non-random constants.

Proposition

Both LS estimators $\hat{\beta}_1$, $\hat{\beta}_0$ are *unbiased*, i.e. $\mathbb{E}(\hat{\beta}_1) = \beta_1$, $\mathbb{E}(\hat{\beta}_0) = \beta_0$.

Proof for Slope

Recall that $\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{\sum_i (x_i - \bar{x}) \cdot y_i}{\sum_i (x_i - \bar{x})^2}$. So, we have

$$\begin{aligned}\mathbb{E}(\hat{\beta}_1) &= \mathbb{E}\left[\frac{\sum_i (x_i - \bar{x}) y_i}{\sum_i (x_i - \bar{x})^2}\right] = \frac{\sum_i (x_i - \bar{x}) \cdot \mathbb{E}(y_i)}{\sum_i (x_i - \bar{x})^2}, && \text{since the } x_i\text{'s are known} \\ &= \frac{\sum_i (x_i - \bar{x}) \cdot \mathbb{E}(\beta_0 + \beta_1 x_i)}{\sum_i (x_i - \bar{x})^2} = \sum_i c_i (\beta_0 + \beta_1 x_i), && \text{where } c_i = \frac{(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \\ &= \beta_0 \sum_i c_i + \beta_1 \sum_i c_i x_i = \beta_1\end{aligned}$$

where the last result is true since $\sum_i c_i = 0$, and $\sum_i c_i x_i = 1$.

Proof for Intercept

Recall that $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$. So, we have

$$\begin{aligned}\mathbb{E}(\hat{\beta}_0) &= \mathbb{E}(\bar{y} - \hat{\beta}_1 \bar{x}) \\ &= \mathbb{E}(\bar{y}) - \bar{x} \cdot \mathbb{E}(\hat{\beta}_1) = \frac{1}{n} \sum_i \mathbb{E}(y_i) - \bar{x} \cdot \beta_1 \\ &= \frac{1}{n} \sum_i \mathbb{E}(\beta_0 + \beta_1 x_i) - \bar{x} \cdot \beta_1 \\ &= \beta_0 + \bar{x} \cdot \beta_1 - \bar{x} \cdot \beta_1 = \beta_0\end{aligned}$$

(*) Note that since both estimators are unbiased $\Rightarrow MSE = \text{Variance}$.

MSE for Slope

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var}\left[\frac{\sum_i (x_i - \bar{x})y_i}{\sum_i (x_i - \bar{x})^2}\right] = \text{Var}\left(\sum_i c_i y_i\right) \quad (c_i \text{ as before}) \\ &= \sum_i c_i^2 \cdot \text{Var}(y_i) = \sum_i c_i^2 \sigma^2 \quad (\text{from model assumption}) \\ &= \sigma^2 \cdot \left(\frac{\sum_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}\right)^2 = \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2} = \sigma^2 \frac{1}{S_{xx}}\end{aligned}$$

MSE for Intercept

$$\text{Var}(\hat{\beta}_0) = \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)$$

SLR Model

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

Normality Assumption

Additionally, we assume that

$$\varepsilon_i \sim^{iid} \mathcal{N}(0, \sigma^2)$$

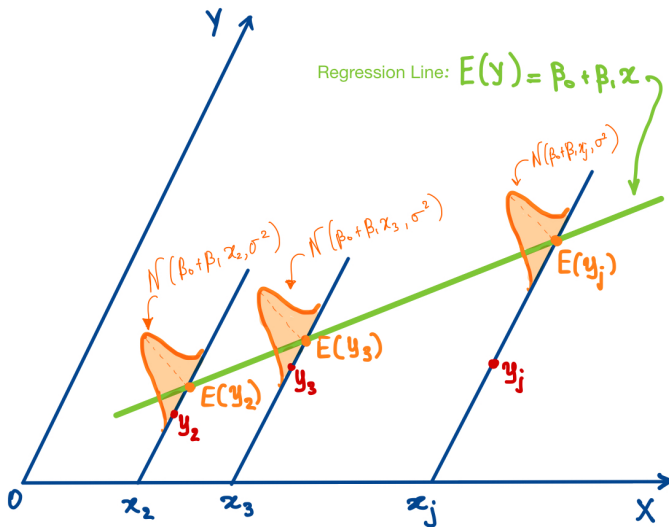
Equivalently, $y_i \sim^{iid} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$. (Why?)

Recall that the error terms ε_i are *independent, normally distributed* with *mean 0* and *variance σ^2* . Based on that, we can prove the following properties for the y_i 's.

Properties of y_i

- $\mathbb{E}(y_i) = \beta_0 + \beta_1 x_i$, since the ε_i 's have mean zero.
- y_i 's are independent, since ε_i 's are independent.
- $\text{Var}(y_i) = \text{Var}(\varepsilon_i) = \sigma^2$.
- y_i 's are a linear shift of the ε_i 's, so they are also *normally distributed*.
- The y_i 's are *jointly normal*, and so are linear combinations of the y_i 's, since the errors are normally distributed and uncorrelated/independent.

Normal Regression Model Illustration



- $\hat{\beta}_1$ and $\hat{\beta}_0$ are **jointly normally distributed** with

$$\mathbb{E}(\hat{\beta}_1) = \beta_1 \quad \text{Var}(\hat{\beta}_1) = \sigma^2 \frac{1}{S_{xx}}$$

$$\mathbb{E}(\hat{\beta}_0) = \beta_0 \quad \text{Var}(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_0) = -\sigma^2 \frac{\bar{x}}{S_{xx}}.$$

- $RSS = \sum_i (y_i - \hat{y}_i)^2 \sim \sigma^2 \chi_{n-2}^2$ which implies that

$$\mathbb{E}(\hat{\sigma}^2) = \mathbb{E} \left(\frac{RSS}{n-2} \right) = \frac{\sigma^2(n-2)}{n-2} = \sigma^2$$

- $(\hat{\beta}_0, \hat{\beta}_1)$ and RSS are **independent**.

Testing for the Slope

$$\begin{cases} H_0 : \beta_1 = c \text{ (null)} \\ H_\alpha : \beta_1 \neq c \text{ (alternative)} \end{cases}$$

where c is an known constant.

- The **test statistics** is

$$t = \frac{\hat{\beta}_1 - c}{\sqrt{\text{Var}(\hat{\beta}_1)}} = \frac{\hat{\beta}_1 - c}{\hat{\sigma} / \sqrt{S_{xx}}}$$

- The **distribution** of t *under the null* is T_{n-2} .
- The **p-value** is twice the area under the T_{n-2} distribution more extreme than the observed statistic t .

R outputs the p -value for testing β_1 against 0, i.e. $c = 0$.

Testing for the Intercept

$$\begin{cases} H_0 : \beta_0 = c \text{ (null)} \\ H_a : \beta_0 \neq c \text{ (alternative)} \end{cases}$$

- The **test statistics** is

$$t = \frac{\hat{\beta}_0 - c}{\sqrt{\text{Var}(\hat{\beta}_0)}}$$

- The **distribution** of t *under the null* is T_{n-2} .
- The **p-value** is twice the area under the T_{n-2} distribution more extreme than the observed statistic t .

R outputs the p -value for testing β_0 against 0, i.e. $c = 0$.

ANOVA Table & F -Test

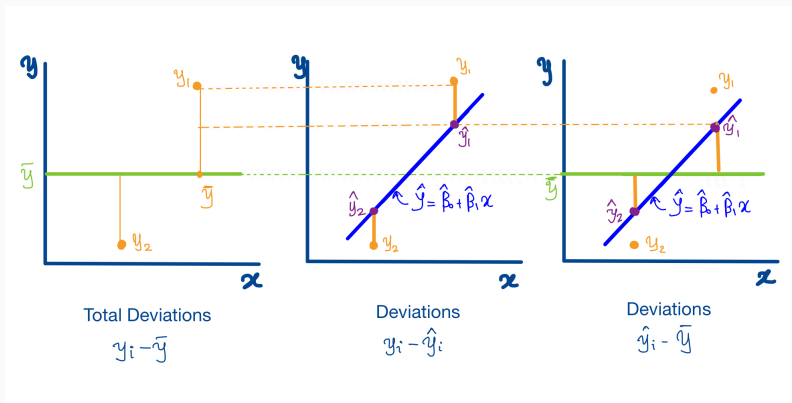
Partitioning the Total Variation (Revisited)

Recall the decomposition of the **Total Sum of Squares** (TSS)

$$\begin{array}{ccccc} TSS & = & FSS & + & RSS \\ \uparrow & & \uparrow & & \uparrow \\ \text{Total variation} & & \text{Variation accounted} & & \text{Variation left} \\ \text{in } y & & \text{through the model} & & \end{array}$$

- **TSS** = $\sum_i (y_i - \bar{y})^2$ is the measure of the total variation in y s: the greater **TSS** is, the more variation there is in the y values.
- **RSS** = $\sum_i (y_i - \hat{y}_i)^2$ measures the variation in the data using the stated regression model: the larger **RSS** is, the more y_i s vary around the *estimated regression line*.
- **FSS** = $\sum_i (\hat{y}_i - \bar{y})^2$ measures how far the predicted center of each probability distribution is from the overall center of all y 's together.

Partitioning the Total Variation (Revisited)



Breakdown of Degrees of Freedom

- $df_{TSS} = n - 1$: one df is lost, because the sample mean is used to estimate the population mean.
- $df_{RSS} = n - 2$: two df are lost, because the two parameters are estimated in obtaining the fitted values \hat{y} .
- $df_{FSS} = 1$: there are n deviations $\hat{y}_i - \bar{y}$, but all the fitted values are associated with the same regression line.

$$df_{TSS} = df_{RSS} + df_{FSS}$$

Sum of Squares	Expression	df
TSS	$\sum_i (y_i - \bar{y})^2$	$n - 1$
FSS	$\sum_i (\hat{y}_i - \bar{y})^2$	1
RSS	$\sum_i (y_i - \hat{y}_i)^2$	$n - 2$

ANOVA Table

Source	SS	df	MS	F
Regression (model)	FSS	1	$MS_{Reg} = \frac{FSS}{1}$	$F = \frac{MS_{Reg}}{MSE}$
Error	RSS	$n - 2$	$MSE = \frac{RSS}{n-2}$	
Total	TSS	$n - 1$		

(*) The Mean Squares are *not* additive.

An alternative way to test for the model parameters is using the F test:

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_a : \beta_1 \neq 0 \end{cases}$$

- Under H_0 , the F -test statistic is

$$F = \frac{MS_{\text{Reg}}}{MSE} = \frac{FSS}{RSS/(n-2)} \sim F_{1,n-2}$$

- It can be shown that the F -test statistic is equal to the square of the t -test statistic and their p -values are the same. So, this test is *equivalent* to the t -test before.

Estimation and Prediction

The LS line can be used to obtain values of the response (Y^*) for given values of the predictor ($X = x^*$). There are two variants of this problem¹:

1. **Estimation**: We want to estimate the mean response at x^* . This is equivalent to estimate: $\beta_0 + \beta_1 x^*$
2. **Prediction** of an outcome of random variable Y^* at a given value x^* , where

$$Y^* \sim N(\beta_0 + \beta_1 x^*, \sigma^2)$$

The fitted value (or point estimate) for estimation and prediction are the same: $\hat{\beta}_0 + \hat{\beta}_1 x^*$. However the accuracy for estimation and the one for prediction are different.

¹Estimation looks to get information from the data about a fixed but parameter, while prediction looks to get information about a random variable

- **Accuracy of the estimation** is measured by the expected value of the squared difference between the point estimate and the target.
- For **estimation** the target is $\beta_0 + \beta_1 x^*$:

$$\begin{aligned}& \mathbb{E} \left(\hat{\beta}_0 + \hat{\beta}_1 x^* - \beta_0 - \beta_1 x^* \right)^2 \\&= \text{Var} \left(\hat{\beta}_0 + \hat{\beta}_1 x^* \right) \\&= \text{Var} \left(\hat{\beta}_0 \right) + (x^*)^2 \text{Var} \left(\hat{\beta}_1 \right) + 2x^* \text{Cov} \left(\hat{\beta}_0, \hat{\beta}_1 \right) \\&= \sigma^2 \left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_i (x_i - \bar{x})^2} \right) \\&= \sigma^2 \left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right)\end{aligned}$$

2

²Recall that all our calculations are done *conditionally* on x^*

- A confidence interval is always reported for a parameter. An $(1 - \alpha)100\%$ Confidence Interval for the *Mean Response* when $x = x_*$ is given by

$$\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

- For **prediction** the target is $Y^* = \beta_0 + \beta_1 x^* + e^*$, where $e^* \sim N(0, \sigma^2)$. This new error e^* is independent of the previous n data points, i.e. is independent of $(\hat{\beta}_0, \hat{\beta}_1)$

$$\begin{aligned} & \mathbb{E}[(\hat{\beta}_0 + \hat{\beta}_1 x^* - Y^*)^2] \\ &= \mathbb{E}[(\hat{\beta}_0 + \hat{\beta}_1 x^* - \beta_0 - \beta_1 x^* - e^*)^2] \\ &= \mathbb{E}[(\hat{\beta}_0 + \hat{\beta}_1 x^* - \beta_0 - \beta_1 x^*)^2] + \mathbb{E}[(e^*)^2] \\ &= \sigma^2 \left(1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

- A prediction interval is reported for the value of a random variable, for example, Y^* . An $(1 - \alpha)100\%$ Prediction Interval for \hat{Y}^* when $x = x^*$ is given by

$$\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

- Based upon the $\text{Var}(\hat{Y}^*)$, the prediction interval is wider than the interval used to estimate the mean response at fixed $x = x^*$.
- So far, we have assumed that the x -levels are known constants. So, all the previous results hold if:
 1. $f(y|x)$ are independent and normally distributed with mean $\beta_0 + \beta_1 x$ and variance σ^2 conditionally on x .
 2. x are independent with distribution $g(x_i)$ that does not depend on β_0 , β_1 , or σ^2 .

- The statement “ X causes Y ” means that changing the value of X will change the distribution of Y . When X causes Y , X and Y will be associated, but the reverse is not, in general, true.

Association does not necessarily imply causation.

- If the data are from a *randomized study*, then the causal interpretation is correct.
- If the data are from a *observational study*, then the association interpretation is correct.