

Class 4a: Simple Linear Regression

Business Forecasting

Roadmap

This set of classes

- What is a simple linear regression?
- How to estimate it?
- How to test hypothesis in the regression?

Motivation

1. Suppose you are a consultant working for Ecobici
2. Your boss is worried about the impact of global warming on bike use
3. She wants to know: how the bike use will change when the temperature increases by 1 degreee
4. This is exactly what the linear regression will tell us!

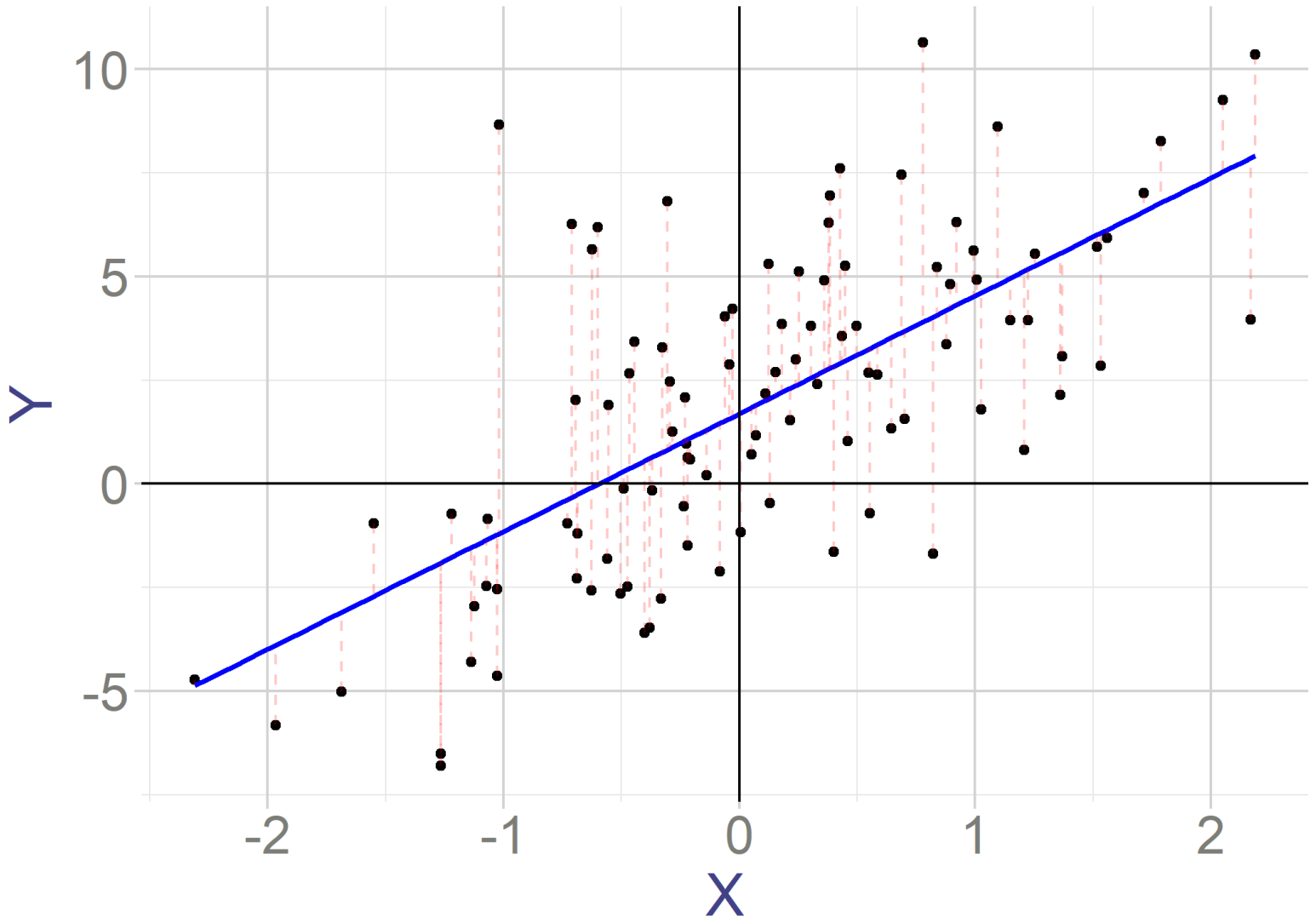
Simple linear regression

1. Suppose you have paired data: $\{(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)\}$
2. In the population, there exists a linear relationship between x_i and y_i of the form:

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

Where:

- y_i is a dependent variable
- x_i is a independent variable, or regressor, or predictor
 - (suppose non-random)
- β_0 and β_1 are parameters
- β_1 tells you how y_i changes (on average) when we change x_i by one unit
- β_0 is intercept, where the line cuts y axis
- u_i is a random error term (unknown)



Assumptions

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

Assumptions:

1. Model is linear in the parameter and with additive error term
2. $E(u_i) = 0 \rightarrow E(y_i|x = x_0) = \beta_0 + \beta_1 x_0$
3. $Var(u_i) = \sigma^2 \rightarrow var(y_i|x = x_0) = \sigma^2$
4. $cov(u_i, u_j) = 0$

Model is linear in the parameter and with additive error term

- Linear models

- $y_i = \beta_0 + \beta_1 x_i + e_i$
- $y_i = \beta_0 + \beta_1 x_i^2 + e_i$
- $y_i = \beta_0 + \beta_1 \log(x)_i + e_i$
- $y_i = \beta_0 + \beta_1 c^{x_i} + e_i$

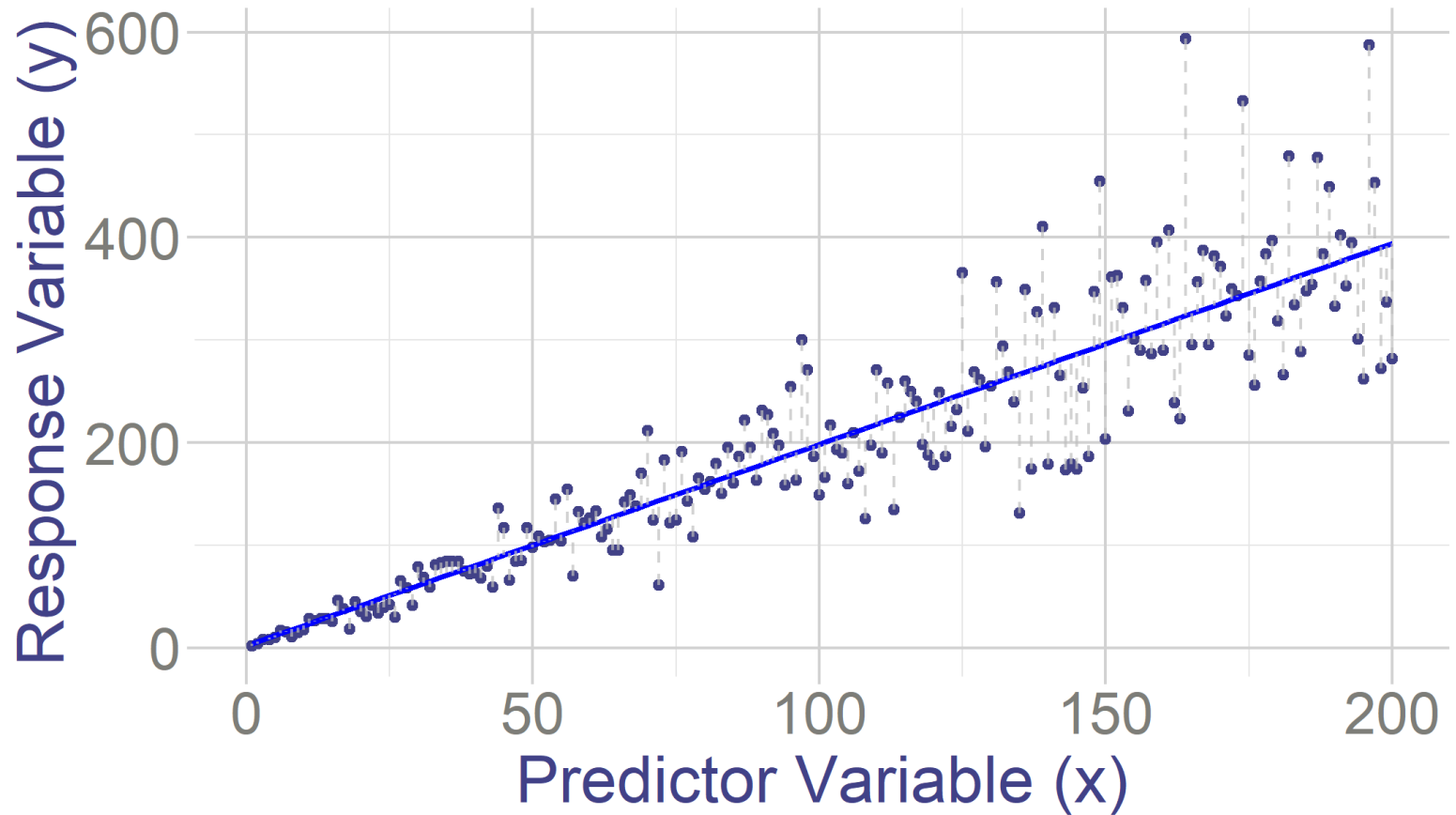
- Not linear models

- $y_i = (\beta_0 + \beta_1 x_i) * e_i$
- $y_i = \beta_0 + x_i^{\beta_1} + e_i$
- $y_i = \log(\beta_0 + \beta_1 x_i + e_i)$
- $y_i = \beta_0 + (\beta_1 x_i + e_i)^2$

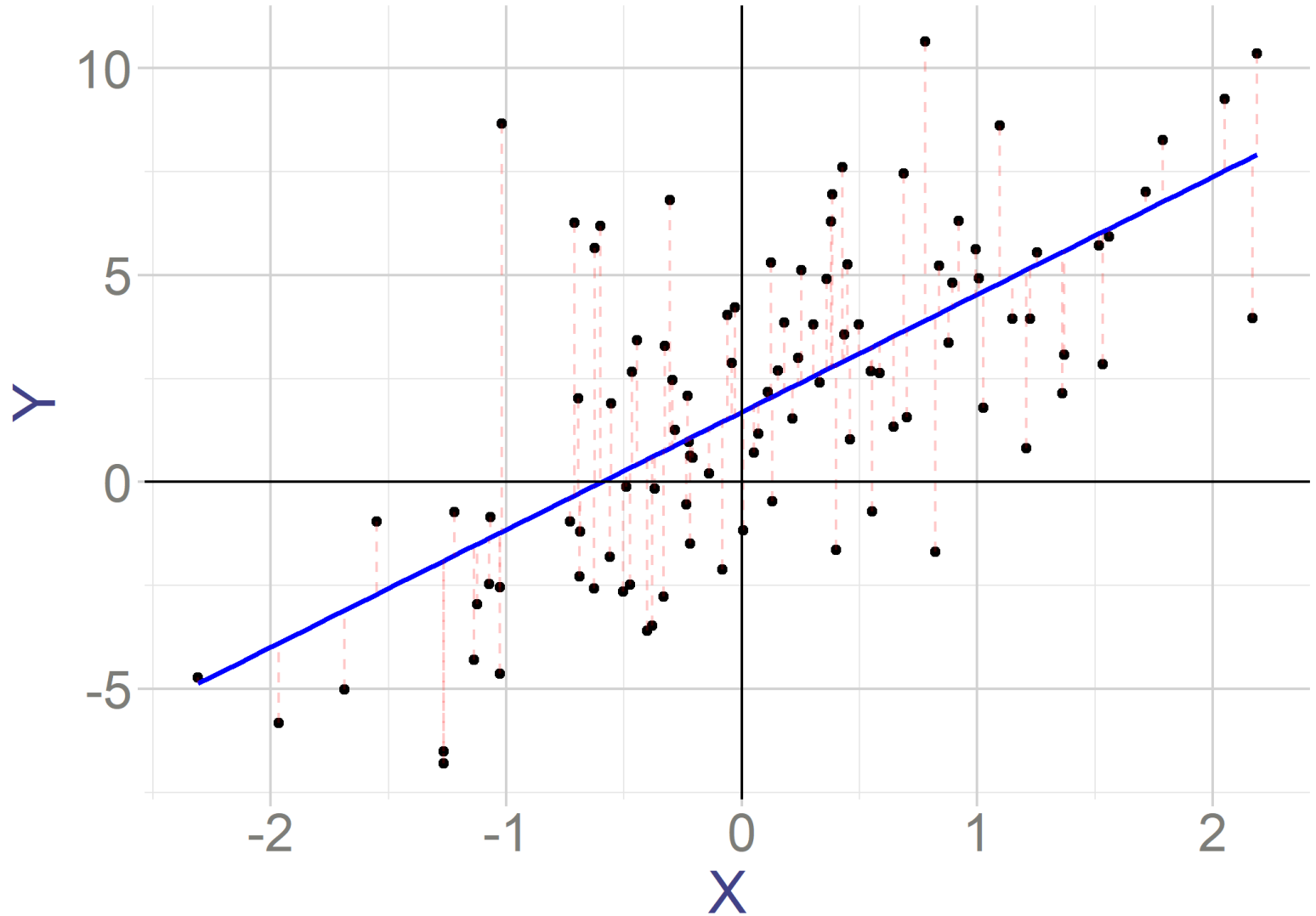
2 is in the app

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

What happens if this is not true?



Let's go back to our regression line



We want to estimate the parameters in this linear relationship based on our sample.

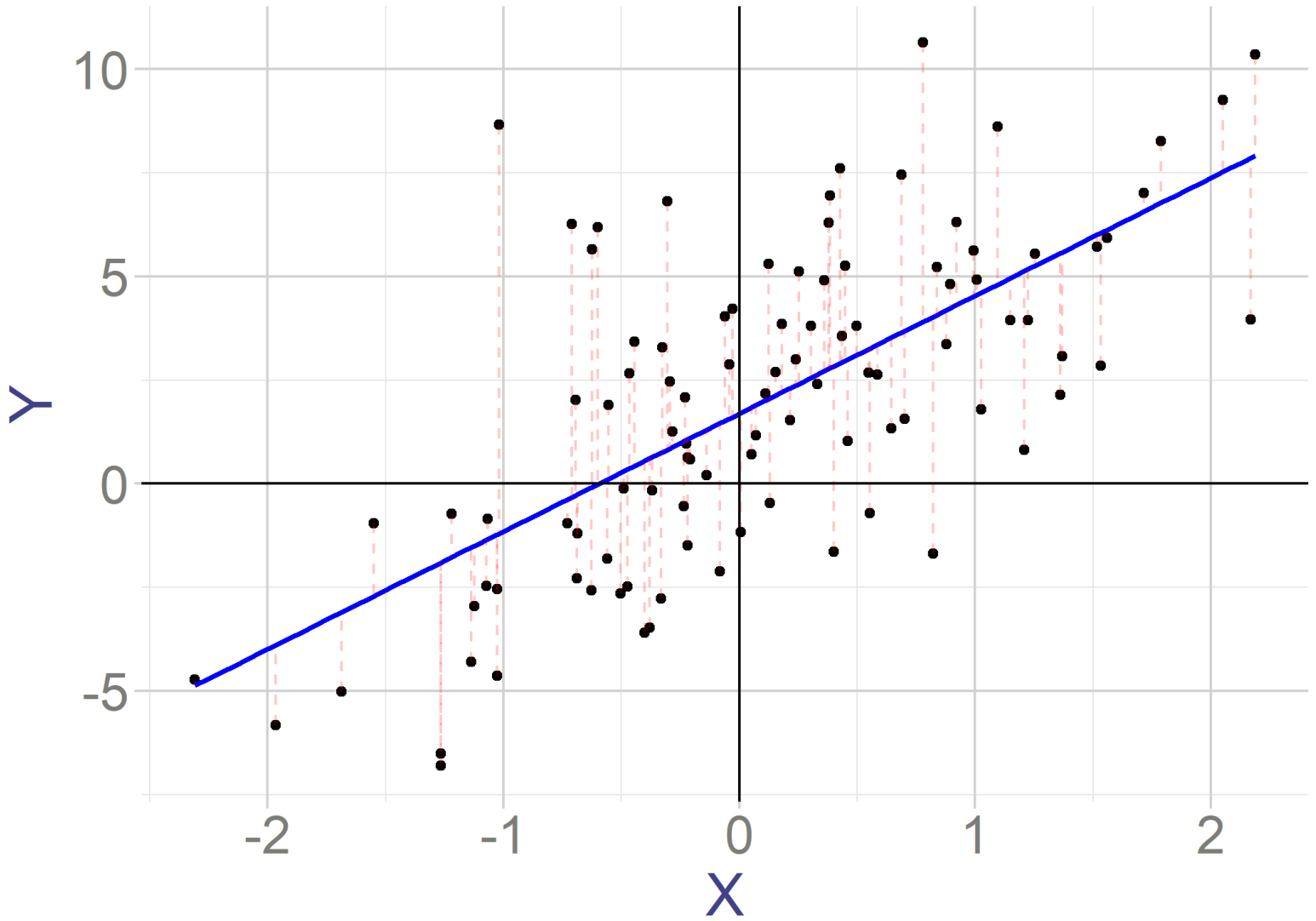
Once estimated, we can write y_i as

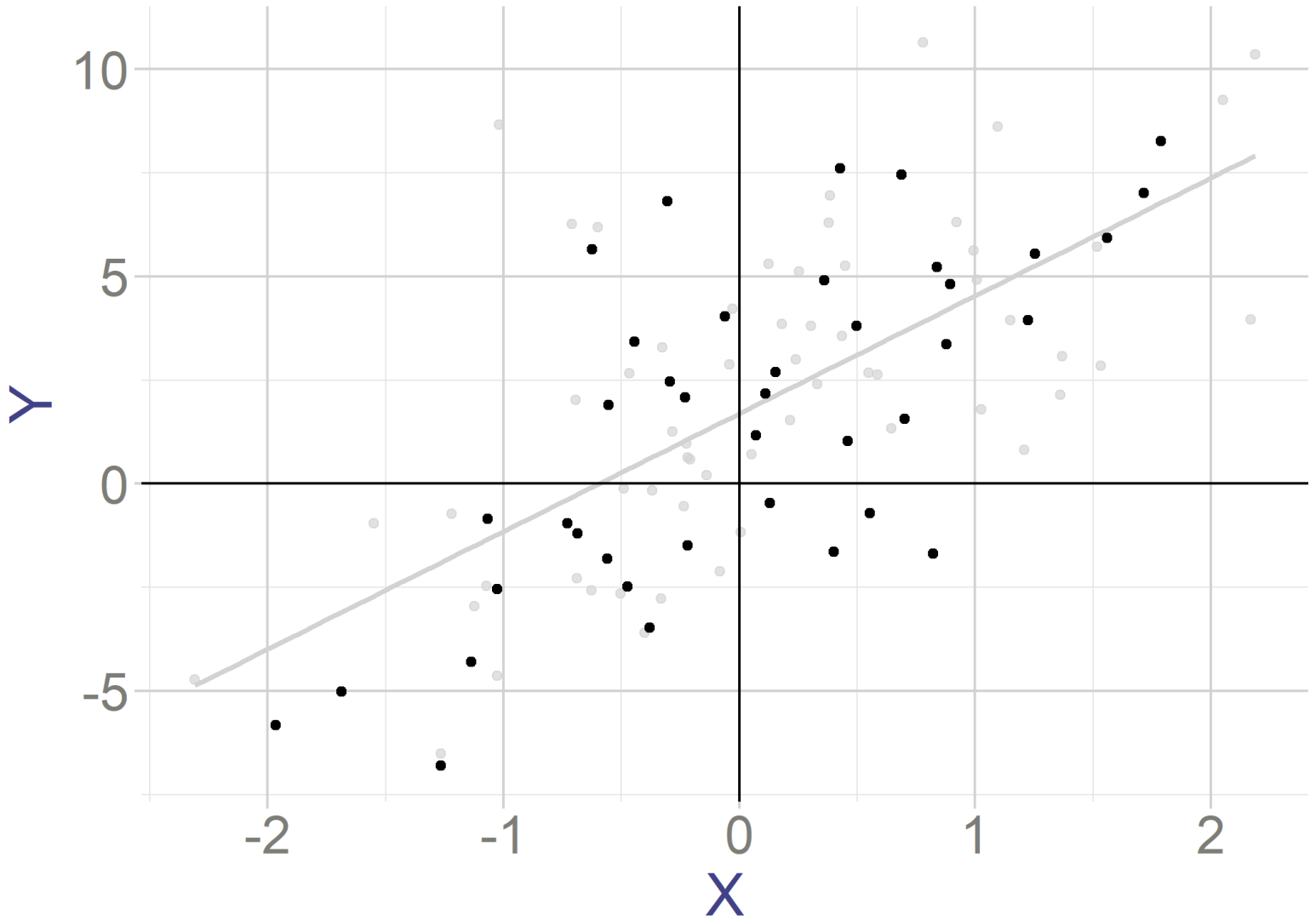
$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i$$

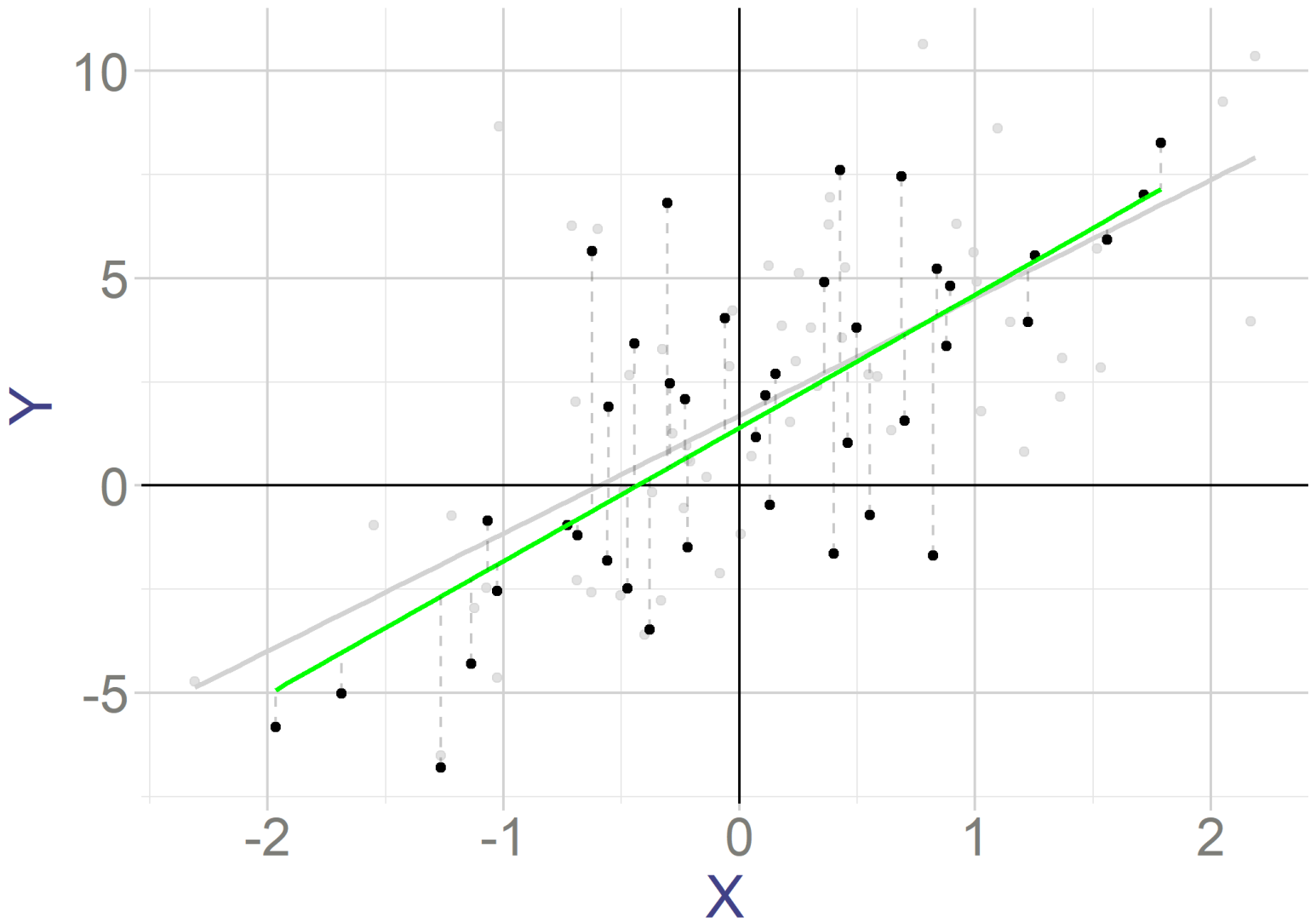
Error term here reflects both uncertainty about parameters and the random part present in population model

We can predict y_i for any x_i using our estimates

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







But how do we find $\hat{\beta}_0$ and $\hat{\beta}_1$?

Best fit line

The best fitting line will minimize the sum of squared residuals $SSE = \sum_i^n e_i^2$

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin}_{b_0, b_1} SSE = \operatorname{argmin}_{b_0, b_1} \sum_i^n e_i^2$$

$$\begin{aligned} SSE &= \sum_{i=1}^n e_i^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2 \end{aligned}$$

So effectively we are minimizing:

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin}_{b_0, b_1} SSE = \operatorname{argmin}_{b_0, b_1} \sum_i^n (y_i - (b_0 + b_1 x_i))^2$$

OLS

We called this estimator **OLS** - ordinary least squares

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin}_{b_0, b_1} SSE = \operatorname{argmin}_{b_0, b_1} \sum_i^n (y_i - (b_0 + b_1 x_i))^2$$

Best fit line 1

To find the minimum of SSE, we take partial derivatives with respect to β_0 and β_1 and set them equal to zero:

Partial derivative with respect to β_0 :

$$\frac{\partial SSE}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n \left(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)$$

Setting this derivative to zero:

$$-2 \sum_{i=1}^n \left(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right) = 0$$

$$\hat{\beta}_0 n + \hat{\beta}_1 \sum x_i = \sum y_i$$

Best fit line 2

Partial derivative with respect to $\hat{\beta}_1$:

$$\frac{\partial SSE}{\partial \hat{\beta}_1} = 2 \sum_{i=1}^n x_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))$$

Setting this derivative to zero:

$$2 \sum_{i=1}^n x_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)) = 0$$

$$\hat{\beta}_0 \sum x_i + \hat{\beta}_1 \sum x_i^2 = \sum x_i y_i$$

Best fit line

Putting it all together:

$$\hat{\beta}_0 n + \hat{\beta}_1 \sum x_i = \sum y_i$$

$$\hat{\beta}_0 = \frac{\sum y_i - \hat{\beta}_1 \sum x}{n} = \bar{y} - \hat{\beta}_1 \bar{x}$$

And plugging this here:

$$\hat{\beta}_0 \sum x_i + \hat{\beta}_1 \sum x_i^2 = \sum x_i y_i$$

We get:

$$\hat{\beta}_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\widehat{cov}(x_i, y_i)}{\widehat{var}(x_i)}$$

Or

$$\hat{\beta}_1 = \frac{\widehat{cov}(x_i, y_i)}{\widehat{var}(x_i)} = \frac{\widehat{cov}(x_i, y_i)}{\sqrt{\widehat{var}(x_i)} \sqrt{\widehat{var}(x_i)}} \frac{\sqrt{\widehat{var}(y_i)}}{\sqrt{\widehat{var}(y_i)}} = \widehat{\rho}(x, y) \frac{\sqrt{\widehat{var}(y_i)}}{\sqrt{\widehat{var}(x_i)}}$$



Source: [<https://observablehq.com/@yizhe-ang/interactive-visualization-of-linear-regression>)]

Back to Motivating example

Show entries

fecha_retiro	Trips	TMP	PM2.5
2017-01-02	20797	14.49	23.03
2017-01-03	26040	15.22	31.5
2017-01-04	27551	16.89	26.61
2017-01-05	28444	15.99	35.02
2017-01-06	26191	17.85	47.21
2017-01-09	31350	10.91	42.24
2017-01-10	33228	12.85	29.42

Showing 1 to 7 of 781 entries

Previous

1

2

3

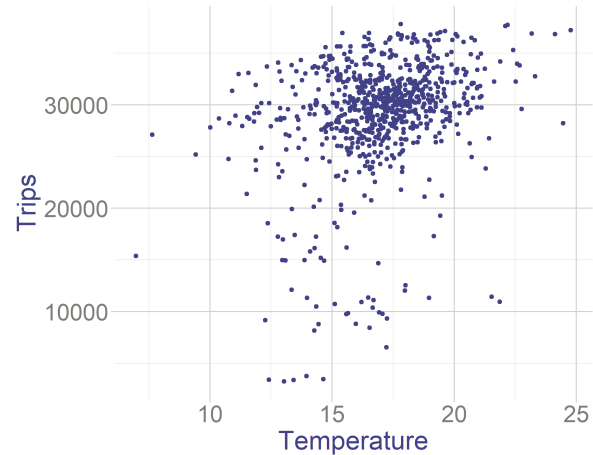
4

5

...

112

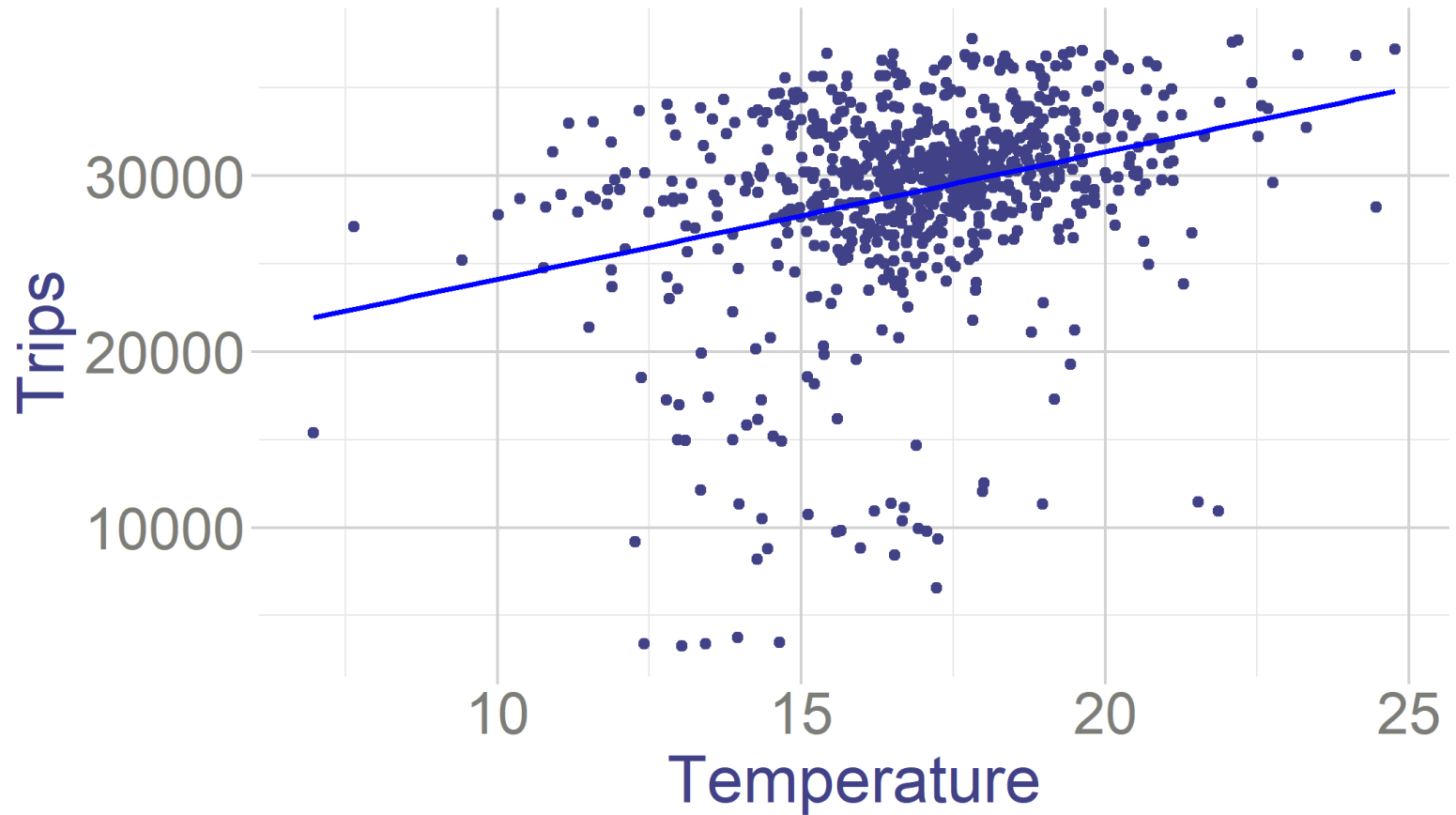
Next



We want to estimate the following relationship:

$$Trips_i = \beta_0 + \beta_1 Temperature_i + u_i$$

Best Fit Line



Regression output in R

```
# Fit a linear regression model
lm_model <- lm(Trips ~ TMP, data = Data_BP)
# Display the summary of the linear regression model
summary(lm_model)
```



```
##
## Call:
## lm(formula = Trips ~ TMP, data = Data_BP)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-24010.5	-1508.4	774.5	2920.5	8900.2


```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	16892.66	1427.32	11.835	<2e-16 ***
TMP	723.55	83.37	8.679	<2e-16 ***


```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5302 on 779 degrees of freedom
## Multiple R-squared:  0.08817,    Adjusted R-squared:  0.087
## F-statistic: 75.32 on 1 and 779 DF,  p-value: < 2.2e-16
```


2. [34 puntos] You have been hired to analyse the relationship between campaign spending and vote share for the forthcoming presidential elections using a simple linear regression model in the form:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad ; \quad i = 1, \dots, 150.$$

where

- y_i represents the share of votes received by the incumbent in the i^{th} election, that is, the candidate who has run for another charge in past elections (could be a mayor or another position that is elected by popular vote). Note that this is operationalised as a proportion of total votes obtained that it may take values between 0 and 1.
- x_i represents the share of total campaign spending by the incumbent in the i^{th} election who has been elected for a political position before. Note that this is operationalised as a proportion of total spending by all candidates and that it may take values between 0 and 1.
- ϵ_i is the i^{th} random error which satisfies Gauss–Markov’s assumptions.

You have been provided with some statistics for data from 150 past elections such as

$$\bar{x} = 0.40 \quad ; \quad \bar{y} = 0.50 \quad ; \quad s_X = 0.20 \quad ; \quad s_Y = 0.15 \quad ; \quad r_{XY} = 0.60$$

Answer the following questions with the information provided:

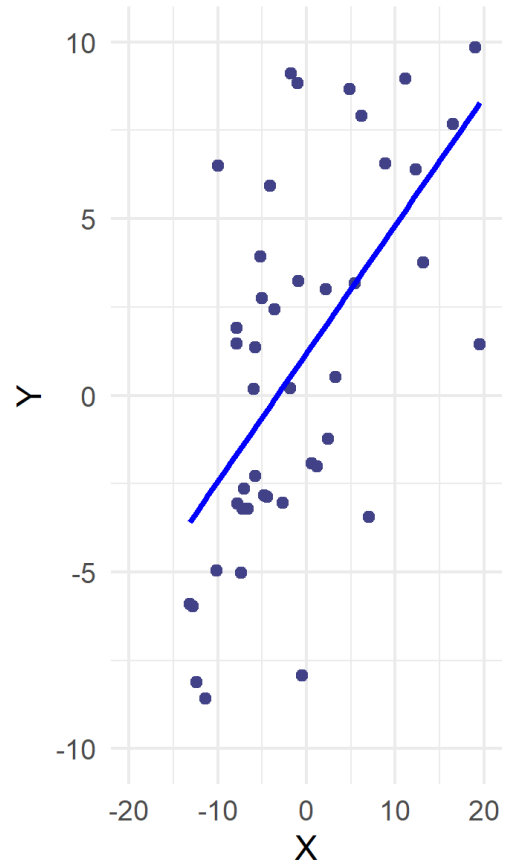
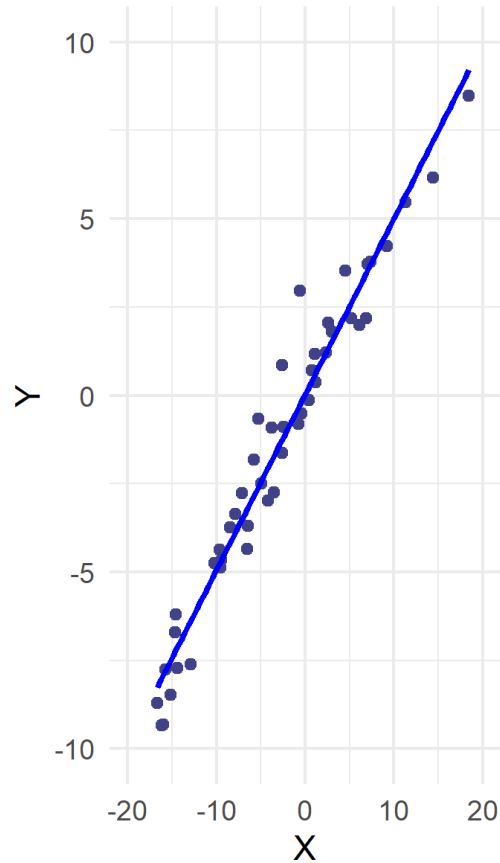
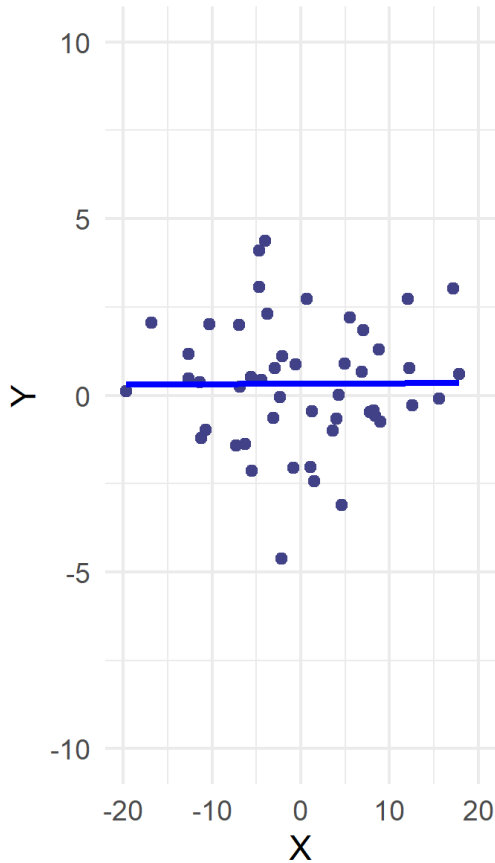
- a) [6 puntos] Calculate the estimates for the model’s parameters.
- b) [6 puntos] Without making any formal inferential process, interpret the coefficients estimated.
- c) [5 puntos] Determine how much campaign spending is needed to obtain at least 40 % of the total vote share.

Properties of this estimator

Here is a couple of cool useful properties of OLS. Let's derive them:

- $\sum e_i = \sum (y_i - \hat{y}_i) = 0$
- $\sum y_i = \sum \hat{y}_i$
- $\hat{y}_i | (x_i = 0) = 0 * \hat{\beta}_1 + \hat{\beta}_0 = \hat{\beta}_0$
- $\sum x_i e_i = 0$
- $\sum \hat{y}_i e_i = 0$
- $var(e_i) = \frac{\sum_i (y_i - \hat{y}_i)^2}{n-2} = \frac{SSE}{n-2}$

Fit of linear regression



Measure of fit - R squared

How much we managed to explain with our regression?

- SST= total sum of squares = $S_{yy} = \sum (y_i - \bar{y})^2 = \sum y_i^2 - n\bar{y}^2$
- SSR= regression sum of squares = $\sum (\hat{y}_i - \bar{y})^2 = \sum \hat{y}_i^2 - n\bar{y}^2$

Measure of fit is:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

Intuition:

- How much variation in y can we explain with our model
- It is always between 0 and 1
 - In fact $SST = SSR + SSE = \sum (\hat{y}_i - \bar{y})^2 + \sum (\hat{y}_i - y_i)^2$
- SSE/SST is proportion that cannot be explained with the model
- so 1-SSE/SST is the variation that we can explain with the model

Illustration in the app

Measure of fit: R squared

If we have just one regressor, the R^2 is related to correlation between x and y .

$$R^2 = (\rho(x, y))^2$$

Moreover, we can show that:

$$R^2 = (\rho(x, y))^2 = \hat{\beta}_1^2 \frac{S_{xx}}{S_{yy}} = \hat{\beta}_1^2 \frac{\sum (x_i - \bar{x})^2}{\sum (y_i - \bar{y})^2}$$

And

$$SSR = \hat{\beta}_1 * S_{xy} = \hat{\beta}_1 * \sum (x_i - \bar{x})(y_i - \bar{y})$$

How much of bike usage does the temperature explains?

- $\beta_1 = 723.55$
- $S_{xx} = \text{var}(x) * (n - 1) = 4043.965$
- $S_{yy} = \text{var}(y) * (n - 1) = 24012556582$

```
##
## Call:
## lm(formula = Trips ~ TMP, data = Data_BP)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -24010.5  -1508.4    774.5   2920.5   8900.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 16892.66    1427.32   11.835  <2e-16 ***
## TMP          723.55     83.37    8.679   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5302 on 779 degrees of freedom
## Multiple R-squared:  0.08817,    Adjusted R-squared:  0.087
## F-statistic: 75.32 on 1 and 779 DF,  p-value: < 2.2e-16
```

Question

Midterm 2, fall 2022, Long question 2, a) and b)

Scaling of variables:

- You built a linear regression explaining how one more peso spent on training improves the performance of the employee.
- You will present this regression to a client from US, who has no idea what a peso is.
- You need to translate it to dollars

Suppose that we used x and y in our sample to estimate $\hat{\beta}_1$ and $\hat{\beta}_0$.

- Let's say that the scale of x changed. New $z = ax + c$.
 - How do $\hat{\beta}_1$ and $\hat{\beta}_0$ change?
- Let's say that the scale of y changed. New $y' = by + d$.
 - How do $\hat{\beta}_1$ and $\hat{\beta}_0$ change?
- Suppose that $\bar{y} = 0$ and $\bar{x} = 0$. What is $\hat{\beta}_0$?

Scaling of variables:

Effect on slope is easiest derived using the definition with correlation:

$$\begin{aligned}\hat{\beta}'_1 &= \text{cor}(z, y') \cdot \frac{\text{sd}(y')}{\text{sd}(z)} \\ &= \text{cor}(ax + c, by + d) \cdot \frac{\text{sd}(by + d)}{\text{sd}(ax + c)} \\ &= \text{cor}(x, y) \cdot \frac{b \cdot \text{sd}(y)}{a \cdot \text{sd}(x)} \\ &= \frac{b}{a} \hat{\beta}_1\end{aligned}$$

- correlation does not change when we scale variables
- adding constants does not matter for the slope
- multiplication of y or x changes the slope

Scaling of variables:

Effect on the intercept is easiest seen through its formula:

$$\begin{aligned}\hat{\beta}'_0 &= \bar{y}' - \hat{\beta}'_1 \bar{z} \\ &= (b\bar{y} + d) - \left(\frac{b}{a} \hat{\beta}_1 \right) (a\bar{x} + c) \\ &= b\bar{y} + d - b\hat{\beta}_1 \bar{x} - \frac{b}{a} \hat{\beta}_1 c \\ &= b(\bar{y} - \hat{\beta}_1 \bar{x}) + d - \frac{b}{a} \hat{\beta}_1 c \\ &= b\hat{\beta}_0 + d - \frac{b}{a} \hat{\beta}_1 c\end{aligned}$$

- multiplying y changes the intercept
- adding a constant to y changes the intercept
- adding a constant to x changes the intercept
- multiplying x only changes the intercept if we also add a constant to x

6. [5 puntos] A group of experts used data relating weekly spending on food delivery through an *app* (Y) and reported monthly income (X), both measured in dollars, obtaining estimates in a regression:

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

with $\hat{\beta}_j$ being least squares estimators for $j = 0, 1$. The analysis revealed that even when the reported income is zero, there was on average positive spending in the app. Additionally, it was found that income had a positive impact on spending in the app. Now, suppose you want to perform the same analysis but with both variables measured in pesos at an exchange rate of \$17.93 pesos per dollar, and you obtain new least squares estimations $\hat{\beta}_0^*$ and $\hat{\beta}_1^*$. Then, it is true that:

$$a) \hat{\beta}_1^* > \hat{\beta}_1 ; \quad b) \hat{\beta}_1^* < \hat{\beta}_1 ; \quad c) \hat{\beta}_0^* \geq \hat{\beta}_0 ; \quad d) \hat{\beta}_0^* < \hat{\beta}_0$$

Regression through the origin

Suppose the following model:

$$y_i = \beta_1 x_i + u_i$$

- What is the least square estimator for β_1 ?
- What happens if we use this estimator when it's not going through the origin?

4. [5 points] Suppose a linear regression model of the form:

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

where ϵ_i is the random error term with the usual assumptions. If we define:

$$v_i = y_i - \bar{y}, \quad u_i = x_i - \bar{x},$$

and intend to adjust a new regression model given by:

$$v_i = \delta_0 + \delta_1 u_i + \epsilon_i; \quad i = 1, \dots, n,$$

then:

- (a) the new modeled line must pass through the origin.
- (b) the new modeled line will have a strictly positive y-intercept.
- (c) the new modeled line will have a strictly negative y-intercept.
- (d) the new modeled line will have a y-intercept different from zero, i.e., either positive or negative indiscriminately.

Regression with a categorical variable

- What if x_i is a categorical variable?
- **Example:** $x_i = 1$ if female, $x_i = 0$ if male
- We called it a binary variable, or a dummy variable

$$\hat{\beta}_0 = \bar{y}_{x_i=0}$$

and

$$\hat{\beta}_1 = \bar{y}_{x_i=1} - \bar{y}_{x_i=0}$$

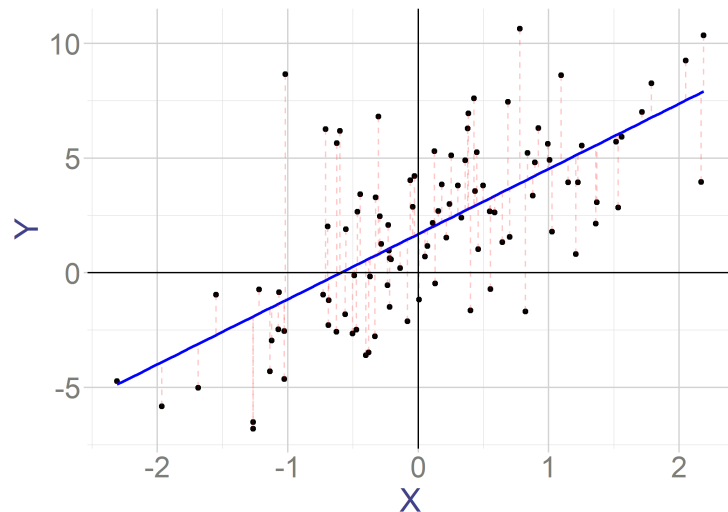
$$wage_i = \beta_0 + \beta_1 female_i + u_i$$

- $female_i = 1$ if i female and 0 if male
- What is β_0 ?
- What is β_1 ?

Statistical Properties of OLS

Uncertainty in the Estimate

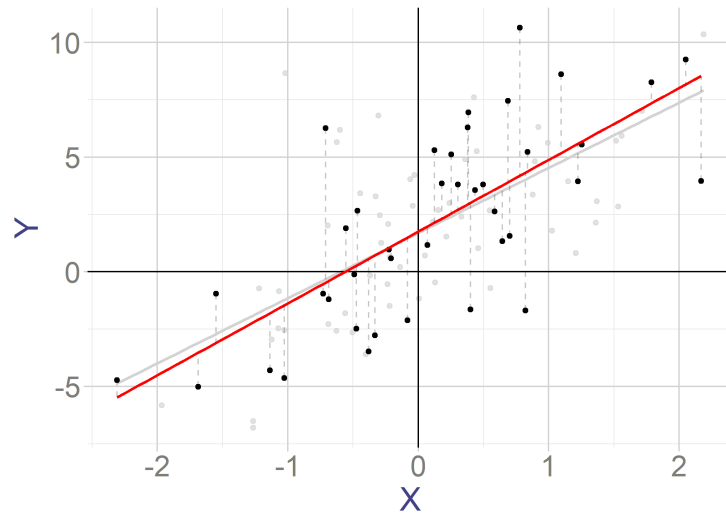
We only have samples, and yet we want to learn something about the population parameters



Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Uncertainty in the Estimate



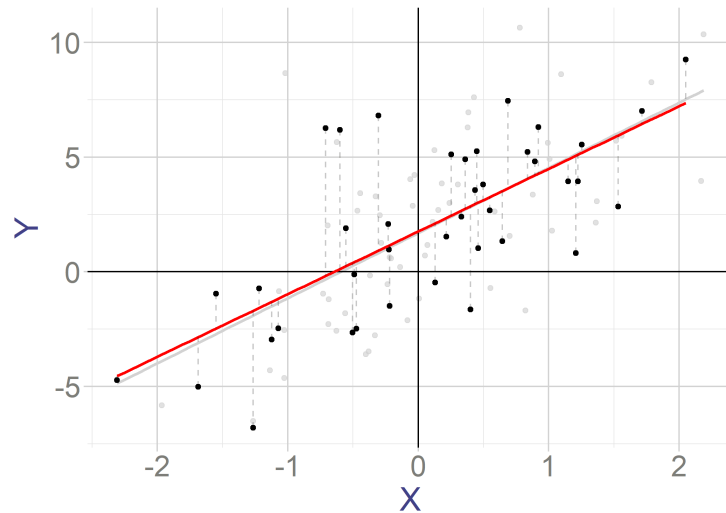
Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Sample Estimate

$$\hat{y}_i = 1.75 + 3.13x_i$$

Uncertainty in the Estimate



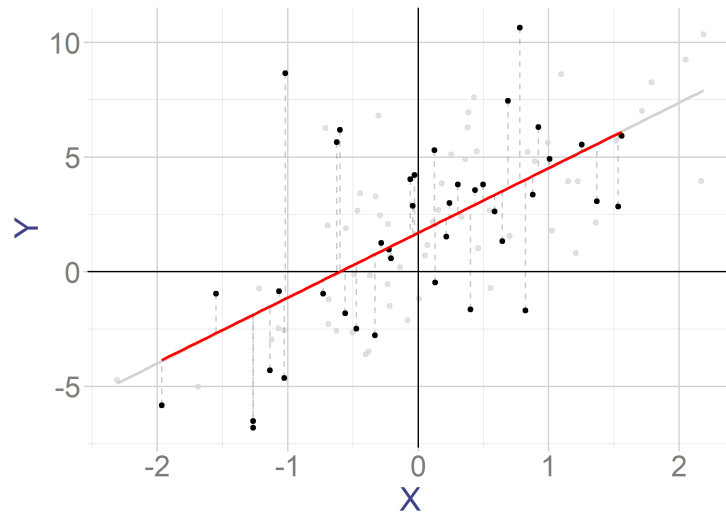
Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Sample Estimate

$$\hat{y}_i = 1.76 + 2.73x_i$$

Uncertainty in the Estimate



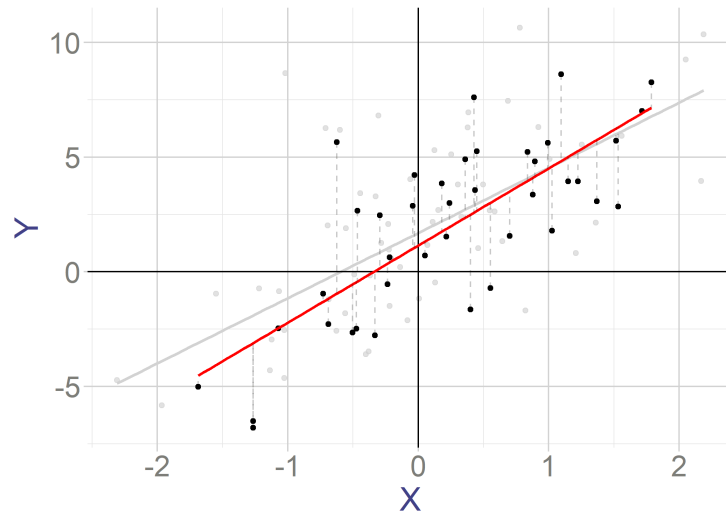
Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Sample Estimate

$$\hat{y}_i = 1.7 + 2.82x_i$$

Uncertainty in the Estimate



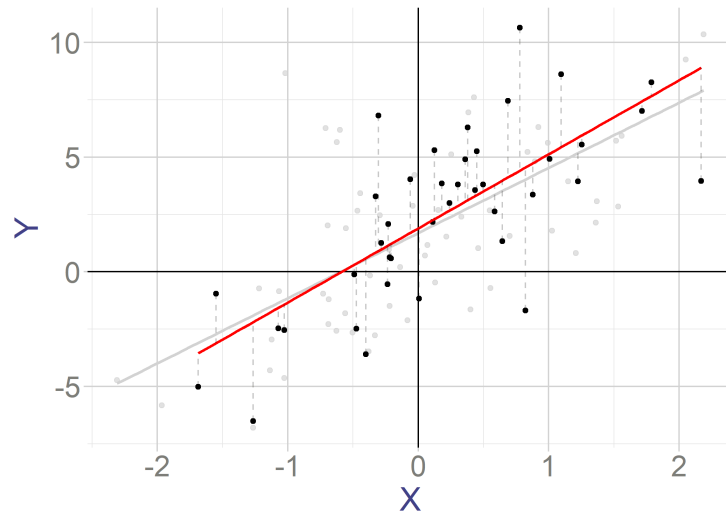
Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Sample Estimate

$$\hat{y}_i = 1.15 + 3.36x_i$$

Uncertainty in the Estimate



Population Regression

$$y_i = 1.69 + 2.84x_i + u_i$$

Sample Estimate

$$\hat{y}_i = 1.89 + 3.23x_i$$

Uncertainty in the Estimate

- $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimators
- And they are random variables
 - Because their values depend on the random samples
- Are they good estimators?
 - Are they unbiased?
 - Do they have small variance?

Uncertainty in the Estimate

Under these assumptions:

1. Relationship is linear in parameters with linear disturbance
2. $E(u_i) = 0$
3. $Var(u_i) = \sigma^2$
4. $cov(u_i, u_j) = 0$

- OLS is unbiased

$$E(\hat{\beta}_1) = E\left(\frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}\right) = \beta_1 \quad \text{and} \quad E(\hat{\beta}_0) = \beta_0$$

- Assumption 1 is enough for being unbiased $E(u_i) = 0$

Uncertainty in the Estimate

- What is the variance of $\hat{\beta}_1$ and $\hat{\beta}_0$?

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var} \left(\frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \right) \\ &= \text{Var} \left(\sum_i \frac{(x_i - \bar{x})y_i}{\sum_i (x_i - \bar{x})^2} \right) = \sum_i \left(\frac{(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right)^2 \text{Var}(y_i) \\ &= \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}\end{aligned}$$

Because x_i don't change: $\text{var}(y_i) = \text{var}(\beta_0 + \beta_1 x_i + u_i) = \text{var}(u_i) = \sigma^2$

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

Standard error is standard deviation of the estimator: $SE(\hat{\beta}) = \sqrt{\text{Var}(\hat{\beta})}$

Uncertainty in the Estimate

- How to estimate the σ^2 ?

$$\hat{\sigma}^2 = \frac{\sum_i e_i^2}{n - 2}$$

- Is unbiased for σ^2 :

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

Regression Output

```
##
## Call:
## lm(formula = Trips ~ TMP, data = Data_BP)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -24010.5  -1508.4    774.5   2920.5   8900.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 16892.66    1427.32  11.835  <2e-16 ***
## TMP          723.55     83.37   8.679  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5302 on 779 degrees of freedom
## Multiple R-squared:  0.08817,    Adjusted R-squared:  0.087
## F-statistic: 75.32 on 1 and 779 DF,  p-value: < 2.2e-16
```

Problem:

Suppose that instead of measuring TMP in celcius, we measure it in *Farenheits*
Practically: $F = 1.8C + 32$

- How would β_1 and $SE(\hat{\beta}_1)$ change?

Gauss Markov Theorem

Under assumptions 1-4, among all linear and unbiased estimators, OLS has the smallest variance.

$$\text{var}(\hat{\beta}_1) \leq \text{var}(\hat{\beta}'_1) \quad \text{and} \quad \text{var}(\hat{\beta}_0) \leq \text{var}(\hat{\beta}'_0)$$

Where $\hat{\beta}'_1$ $\hat{\beta}'_0$ are any linear and unbiased estimators of β_1 and β_0 respectively.

It's **BLUE** - Best, Linear, Unbiased Estimator

Linear estimator basically means it's a weighted sum of y_i s:

$$\hat{\beta}'_1 = \sum_i c_i y_i$$

where c_i are some weights, usually function of x_i

In OLS:

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{\sum_i (x_i - \bar{x})y_i}{\sum_i (x_i - \bar{x})^2} \quad \text{so} \quad c_i^{OLS} = \frac{(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}$$

UPDATE on Gauss Markov

- Science is in progress
- A new paper in 2022 by Hansen shows linearity is not needed
- OLS, under our assumptions, is BUE (Best Unbiased Estimator)

Question 6 [5 points]:

Consider the linear model of the form:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

with $E[\epsilon_i] = 0$; $var(\epsilon_i) = \sigma_i^2 \neq \sigma^2$; $cov(\epsilon_i, \epsilon_j) = 0$ for all $i \neq j$, and the estimation of the model by Least Squares. Now consider the following statements:

A: The Least Squares estimators will no longer be unbiased.

B: The Least Squares estimators will no longer have minimum variance.

Then:

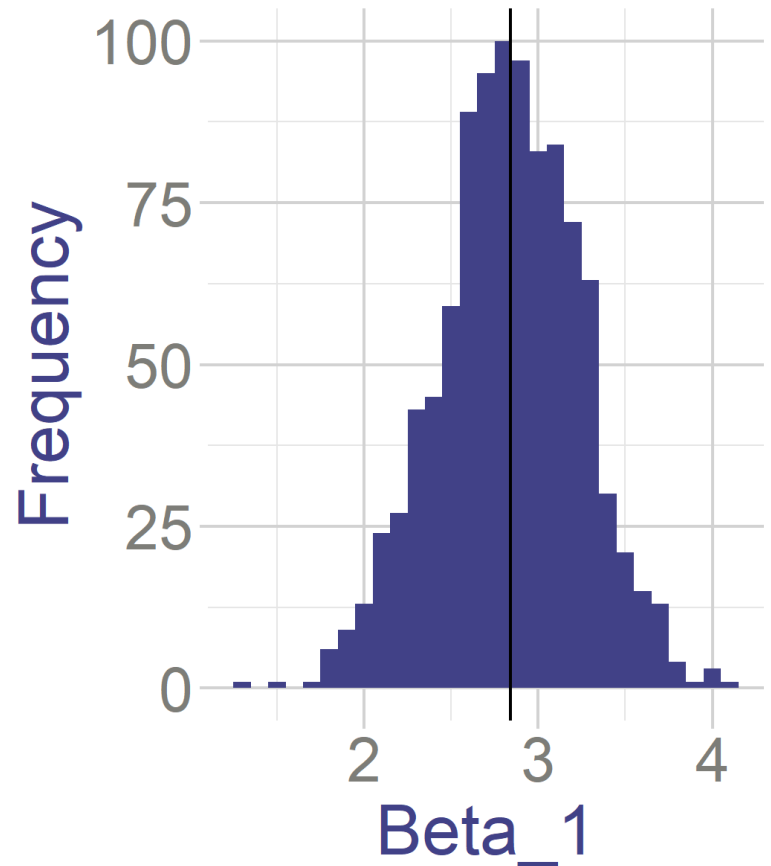
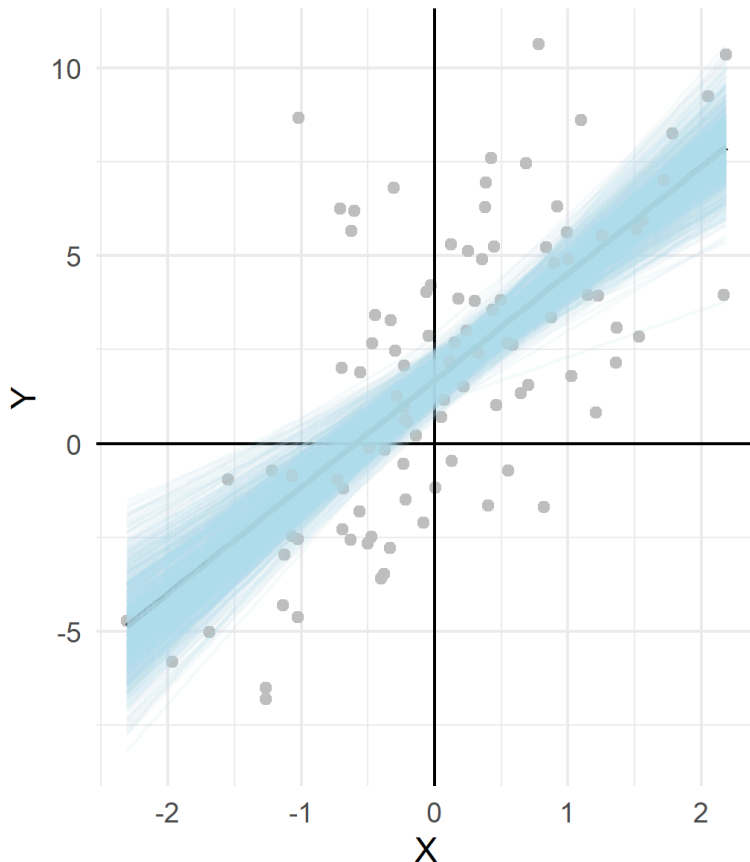
Inference

- Until now, we haven't made any assumptions about the **distributions** of the underlying data or β
 - We don't need it for calculating coefficients $\hat{\beta}_0$ or $\hat{\beta}_1$
 - We don't need it for making predictions $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
 - We don't need it to calculate variance or expectation of coefficients
 - We don't need it for Gauss-Markov Theorem
- However, to make **inference** (confidence intervals, hypothesis testing), we need to know something about distribution of $\hat{\beta}$
 - In particular, we will assume that population errors are normally distributed: $u_i \sim N(0, \sigma)$
 - This will help us to determine the distribution of β
 - y_i or x_i does not need to be normally distributed
 - But if $u_i \sim N(0, \sigma)$, then conditional on x_i : $y_i | x_i \sim N(\beta_0 + \beta_1 x_i, \sigma)$

Suppose I take 1000 samples of size 40 from the population where $u_i \sim N(0, 2)$:

$$y_i = 1.69 + 2.84x_i + u_i$$

And I estimate the β_1 and β_0 for each sample.



Distributions

Given that

- $u_i \sim N(0, \sigma)$
- linear combination of normal variables is normal

We can derive the following distributions:

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma}{\sqrt{S_{xx}}}\right) \quad \text{and} \quad \hat{\beta}_0 \sim N\left(\beta_0, \sigma \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}\right)$$

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2$$

Hypothesis Testing

Our **test statistic** for β_1 and its distribution under the null hypothesis:

$$H_0 : \beta_1 = b_1$$

$$T = \frac{\hat{\beta}_1 - b_1}{SE(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - b_1}{\frac{\hat{\sigma}}{\sqrt{S_{xx}}}} \sim t_{n-2}$$

Similarly, for β_0 the null hypothesis: $H_0 : \beta_0 = b_0$

$$T = \frac{\hat{\beta}_0 - b_0}{SE(\hat{\beta}_0)} = \frac{\hat{\beta}_0 - b_0}{\hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}} \sim t_{n-2}$$

With that, we can use usual hypothesis testing procedures

Example:

Does temperature predicts bike rides? Let's test it at $\alpha = 0.05$

$$H_0 : \beta_1 = 0 \quad H_A : \beta_1 \neq 0$$

$$T_{test} = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{723.55}{83.37} = 8.679$$

We can compare it to critical value (n=781):

$$t_{779, \frac{\alpha}{2}} \approx z_{\frac{\alpha}{2}} = 1.96 < 8.679 = T_{test}$$

We confidently reject the the null that the temperature does not predict bike rides.

P-Value

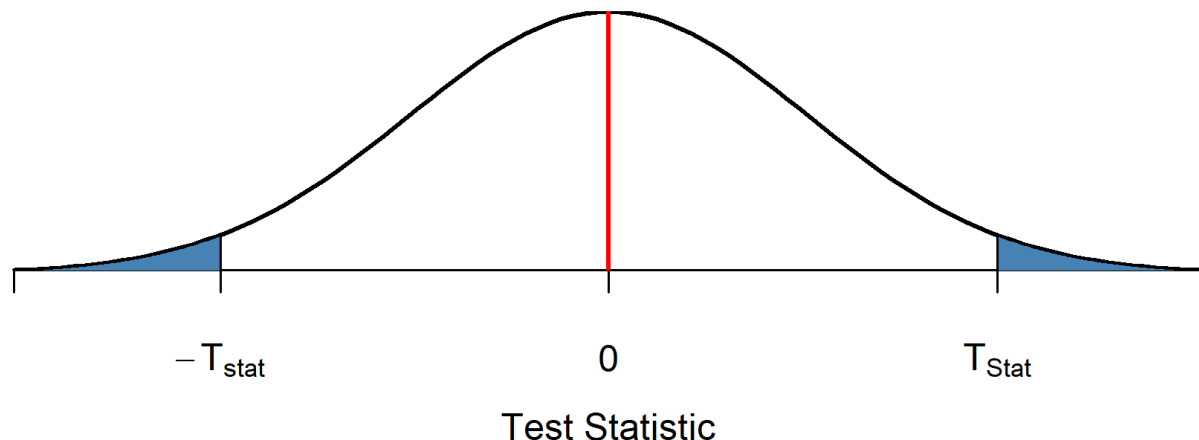
Alternatively, calculate **p-value**: the probability of seeing our test statistic or a more extreme test statistic if the null hypothesis were true.

In regressions we usually use two-sided tests. Hence the p-value is:

$$p - value = 2 * P(t_{n-2, \frac{\alpha}{2}} > T_{test})$$

Small p-values mean that it would be unlikely to see our results if the null hypothesis were really true.

Distribution of the statistic under the null



Regression Output

```
##
## Call:
## lm(formula = Trips ~ TMP, data = Data_BP)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -24010.5  -1508.4    774.5   2920.5   8900.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 16892.66    1427.32  11.835  <2e-16 ***
## TMP          723.55     83.37   8.679  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5302 on 779 degrees of freedom
## Multiple R-squared:  0.08817,    Adjusted R-squared:  0.087
## F-statistic: 75.32 on 1 and 779 DF,  p-value: < 2.2e-16
```

Confidence Intervals

Using the distributions, we can figure out confidence intervals for our estimates:

$$P(-t_{n-2, \frac{\alpha}{2}} < \frac{\hat{\beta}_1 - \beta}{SE(\hat{\beta}_1)} < t_{n-2, \frac{\alpha}{2}}) = 1 - \alpha$$

$$CI_{\beta_1} = \left(\hat{\beta}_1 - t_{n-2, \frac{\alpha}{2}} \underbrace{\frac{\hat{\sigma}}{\sqrt{S_{xx}}}}_{SE(\hat{\beta}_1)}, \hat{\beta}_1 + t_{n-2, \frac{\alpha}{2}} \underbrace{\frac{\hat{\sigma}}{\sqrt{S_{xx}}}}_{SE(\hat{\beta}_1)} \right)$$

And Similarly for β_0

$$CI_{\beta_0} = \left(\hat{\beta}_0 - t_{n-2, \frac{\alpha}{2}} \underbrace{\hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}}_{SE(\hat{\beta}_0)}, \hat{\beta}_0 + t_{n-2, \frac{\alpha}{2}} \underbrace{\hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}}_{SE(\hat{\beta}_0)} \right)$$

Confidence Intervals

What's the confidence 95% interval for the effect on temperature?

$$CI_{\beta_1} = (723.55 - 1.96 * 83.37, 723.55 + 1.96 * 83.37)$$

$$CI_{\beta_1} = (560.87, 886.23)$$

Question

You have data on monthly performance-based bonuses (X) measured in thousands of pesos and job satisfaction ratings (Y) on a scale of 0 to 100 for a sample of 14 employees. To evaluate if the bonus allocation policy is effective in promoting employee job satisfaction, a linear model was estimated, and the following information was obtained:

$$y_i = \underset{(11.2303)}{33.7083} + \underset{(1.6623)}{\hat{\beta}_1} x_i$$

where the numbers contained inside parentheses denote the standard error of the estimates. Additionally, you have:

$$s^2 = 202.3062 \quad ; \quad r_{XY} = 0.6574 \\ \bar{x} = 6.3571 \quad ; \quad \bar{y} = 65.6429 \quad ; \quad s_X^2 = 5.6319 \quad ; \quad s_Y^2 = 328.8626$$

With the above information, please answer the following two questions:

1. **[5 puntos]** Based on the information above and considering a 90 % confidence level, the average increase in job satisfaction for every thousand pesos of bonus earned falls within the interval:
a) (2.061, 7.9861) ; b) (1.402, 8.645) ; c) (9.240, 58.177) ; d) (13.693, 53.724)

Confidence Intervals

Suppose we instead want to estimate the impact of pollution (PM10) on bike trips.

```
##
## Call:
## lm(formula = Trips ~ PM10, data = Data_BP)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -27079.4  -1298.2    947.1   3155.8   8938.6
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  28382.98     576.49   49.235  <2e-16 ***
## PM10          16.99       11.68    1.455   0.146
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5544 on 779 degrees of freedom
## Multiple R-squared:  0.002709,    Adjusted R-squared:  0.001429
## F-statistic: 2.116 on 1 and 779 DF,  p-value: 0.1462
```

- Can we reject null of no impact at 10%?
- What's the 90% confidence interval?

Confidence Intervals

Average response: What would be average number of rides on days with temperature of 30C?

$$(\bar{y}|x = x_0) = \hat{\beta}_0 + \hat{\beta}_1 x$$

What's the expectation?

$$E(\bar{y}|x = x_0) = E(\hat{\beta}_0 + \hat{\beta}_1 x_0) = \beta_0 + \beta_1 x_0$$

What's the variance?

$$var(\bar{y}|x = x_0) = Var(\hat{\beta}_0 + \hat{\beta}_1 x_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

What's the distribution:

$$(\bar{y}|x = x_0) \sim N \left(\beta_0 + \beta_1 x_0, \sigma \sqrt{\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \right)$$

Confidence Intervals

We can build the confidence intervals as before:

$$CI_{(\bar{y}|x=x_0)} = \hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t_{n-2, \frac{\alpha}{2}} \underbrace{\hat{\sigma} \sqrt{\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}}_{SE}$$

Confidence Intervals

What would be 95% CI for average number of rides if temperature is 30C?

- $\hat{\beta}_0 = 16892.66$ and $\hat{\beta}_1 = 723.55$
- $n=781$
- $\bar{x} = 16.96$
- $S_{xx} = 4044$
- $\sum_i e^2 = 21895427100$
- $\hat{\sigma} = \sqrt{\frac{\sum_i e^2}{n-2}} = 5301.613$

$$CI_{(\bar{y}|x=x_0)} = 16892.66 + 723.55 * 30 \pm 1.96 * \underbrace{5301.613 \sqrt{\left(\frac{1}{781} + \frac{(30 - 16.96)^2}{4044}\right)}}_{SE}$$

$$CI_{(\bar{y}|x=x_0)} = 38599.16 \pm 2161.588$$

- Interpretation?
 - If we take a lot of samples, and calculate confidence interval using data from each, 95% of them would contain the true value
 - We are 95% confident, true value is in the interval

Confidence Intervals

R code

```
lm_model <- lm(Trips ~ TMP, data = Data_BP)
new_data<- data.frame(TMP= c(30))
predict(lm_model, newdata = new_data, interval = "confidence", level = 0.95)
```

```
## $fit
##           fit          lwr          upr
## 1 38599.23 36434.32 40764.14
##
## $se.fit
## [1] 1102.851
##
## $df
## [1] 779
##
## $residual.scale
## [1] 5301.613
```

Mean response vs New response

- Suppose you are checking how people react to a new drug for balding. You estimated the following regressions:

$$\text{Number of hairs/cm}^2 = \hat{\beta}_0 + \hat{\beta}_1 \text{Amount of drug in mg}$$

- For now, you were only giving doses between 1-25mg. You want to increase dosage to 30mg.
- You can have two types of confidence intervals

- For **Mean Response**

- Suppose you give 30mg to many, many people, and you are interested in average Number of hairs/ cm^2 among those who got 30mg
- Since you average among many people, the u_i individual error terms does not play a role ($E(u_i) = 0$)
- The uncertainty comes from whether you did a good job estimating β s

- For **New Response**

- Suppose you give 30mg to one person, and you are interested in their outcome.
- Since there is only one person, u_i will play a role
- Maybe you picked someone who naturally has a lot of hair, or who will be on other medication which makes him lose hair
- Those factors average out in mean response, so don't play a role
- There will be more uncertainty about this new response, hence wider CI
- In particular, $var(\text{new response}) = var(\text{mean response}) + var(u_i)$

Confidence Intervals

New response: What would be the number of rides on some day with temperature 30C?

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

What's the expectation?

$$E(\hat{y}|x = x_0) = E(\hat{\beta}_0 + \hat{\beta}_1 x_0) = \beta_0 + \beta_1 x_0$$

How much true value varies around this prediction?

$$\text{var}(y_0 - \hat{y}|x = x_0) = \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0) + \text{Var}(u_i) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

What's the distribution:

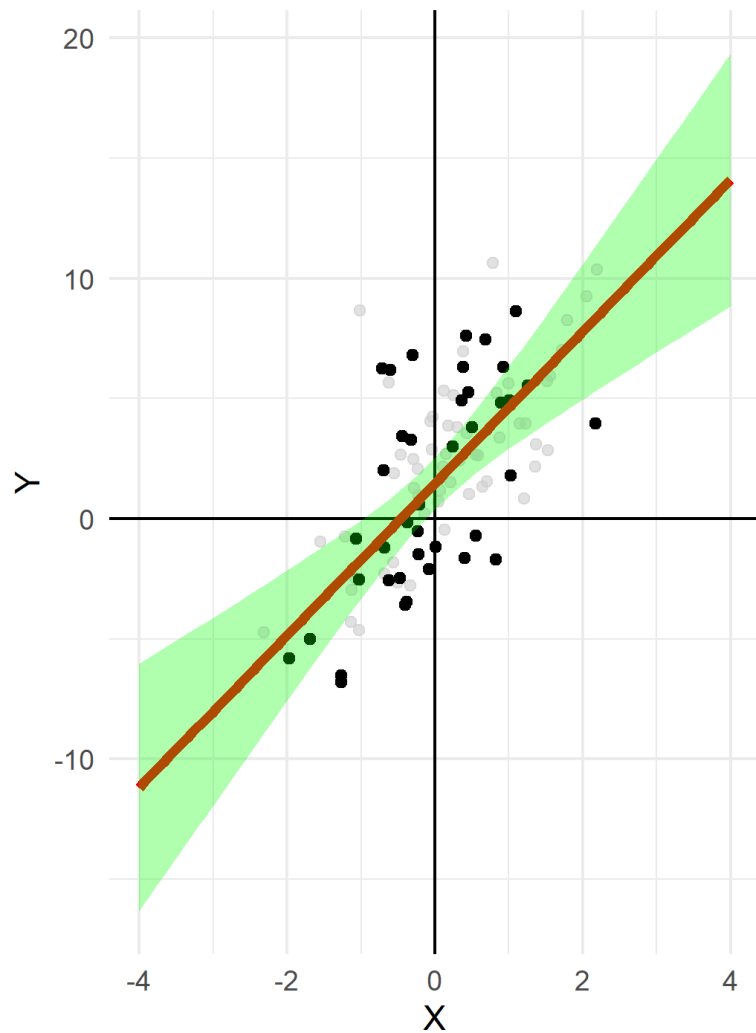
$$(\bar{y}|x = x_0) \sim N \left(\beta_0 + \beta_1 x_0, \sigma \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \right)$$

Confidence Intervals

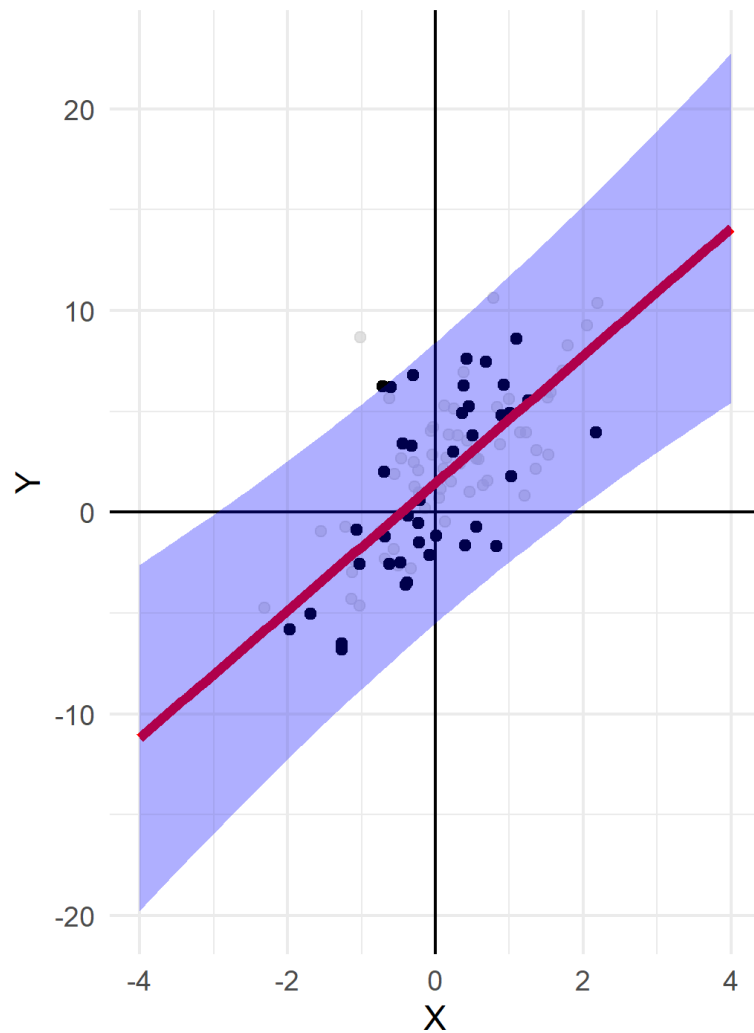
We can build the confidence intervals as before:

$$CI_{(\bar{y}|x=x_0)} = \hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t_{n-2, \frac{\alpha}{2}} \underbrace{\hat{\sigma} \sqrt{\left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}\right)}}_{SE}$$

Mean Response Interval



New Response Interval



Confidence Intervals

What would be 95% CI for number of rides on some day with 30C?

R code

```
lm_model <- lm(Trips ~ TMP, data = Data_BP)
new_data<- data.frame(TMP= c(30))
predict(lm_model, newdata = new_data, interval = "predict", level = 0.95)
```

```
## $fit
##           fit      lwr      upr
## 1 38599.23 27969.3 49229.16
##
## $se.fit
## [1] 1102.851
##
## $df
## [1] 779
##
## $residual.scale
## [1] 5301.613
```

Question

Midterm 2 2022, fall, Long question 1

Question

Suppose a model where we have employee's salary and their years of education. Predictor variable is education, response variable is salary. We try to establish the relationship between education and salary.

- What type of factors may affect the stochastic error u_i ?
- Are they correlated with education?
- Would the estimator be unbiased?

Practice

- Lista 02