

# Lesson 3: Introduction to Simple Linear Regression (SLR)

Nicky Wakim

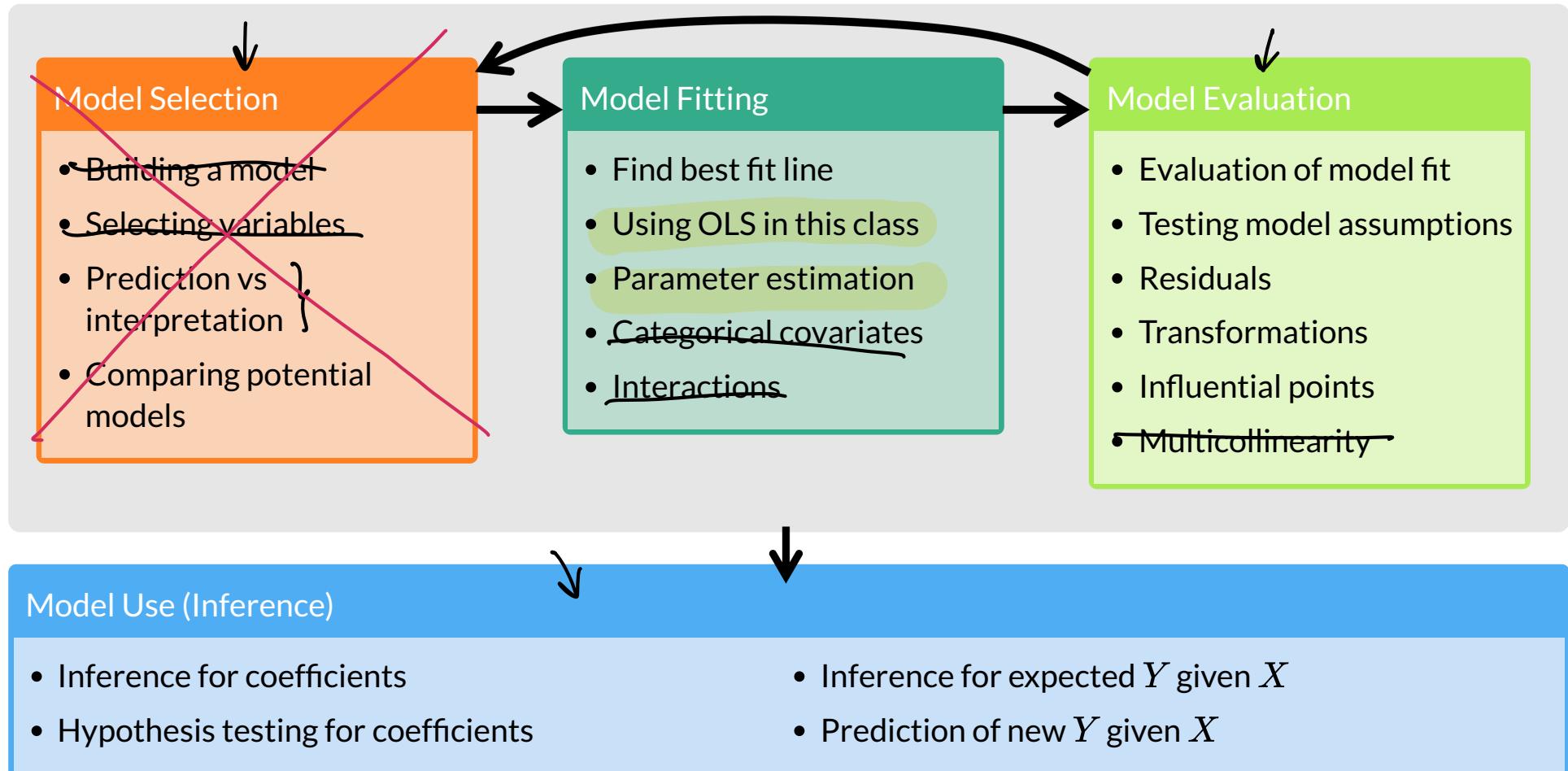
2025-01-13

# Learning Objectives

1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4. Solve the optimal coefficient estimates for simple linear regression using OLS
5. Apply OLS in R for simple linear regression of real data

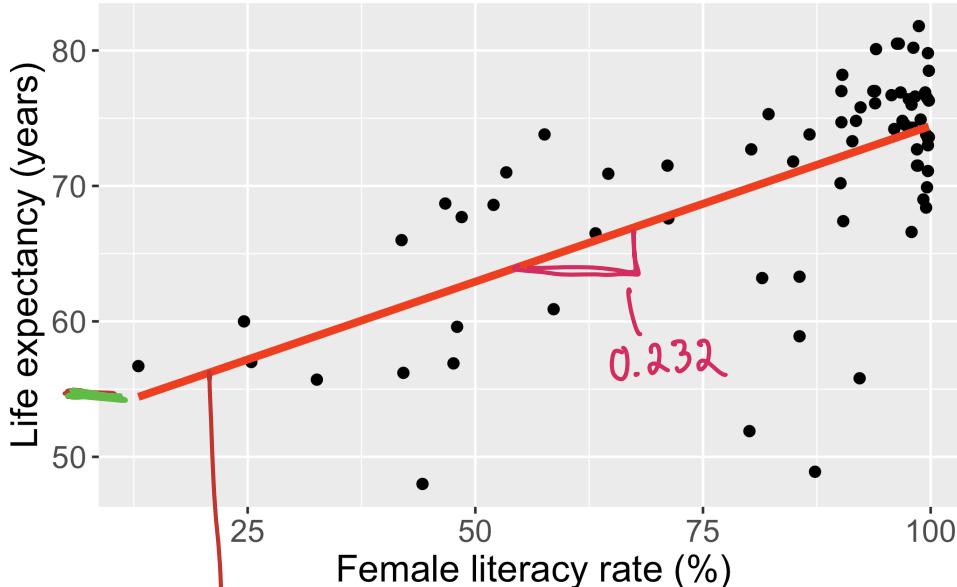
# Process of regression data analysis

SLR



# Let's start with an example

Relationship between life expectancy and the female literacy rate in 2011



Y vs X

~~What~~ life expectancy vs. female literacy rate

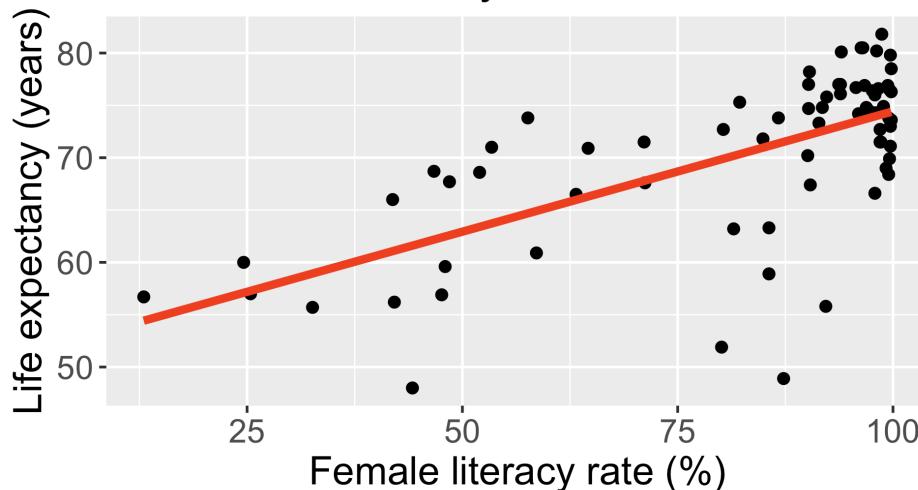
- Each point on the plot is for a different country
- $X$  = country's adult female literacy rate
- $Y$  = country's ~~average~~ life expectancy (years)

$$\text{life expectancy} = 50.9 + 0.232 \cdot \text{female literacy rate}$$

# Reference: How did I code that?

```
1 gapm %>%  
2   ggplot(aes(x = FemaleLiteracyRate,  
3               y = LifeExpectancyYrs)) +  
4     geom_point(size = 4) +  
5     geom_smooth(method = "lm", se = FALSE, size = 3, colour="#F14124") +  
6     labs(x = "Female literacy rate (%)",  
7           y = "Life expectancy (years)",  
8           title = "Relationship between life expectancy and \n the female literacy rate in 2011") +  
9     theme(axis.title = element_text(size = 30),  
10        axis.text = element_text(size = 25),  
11        title = element_text(size = 30))
```

Relationship between life expectancy and the female literacy rate in 2011



# Research and dataset description

**Research question:** Is there an association between life expectancy and female literacy rates?

- Data file: [Gapminder\\_vars\\_2011.xlsx](#)
- Data were downloaded from [Gapminder](#)
  - 2011 is the most recent year with the most complete data
  - Observational study measuring different characteristics of countries, including population, health, environment, work, etc.
- **Life expectancy** = the average number of years a newborn child would live if current mortality patterns were to stay the same.
- **Adult literacy rate** is the percentage of people ages 15 and above who can, with understanding, read and write a short, simple statement on their everyday life.
  
- **National Literacy Trust** in England has studied the link between these two variables
  - Please note that they clearly state that literacy is linked to life expectancy **through many socioeconomic and health factors**

# Poll Everywhere Question 1

13:22 Mon Jan 13 ... 98%

X Join by Web [PollEv.com/nickywakim275](https://PollEv.com/nickywakim275)

From how I described the data and research question, what is our observational unit?

Individual 0%

Hospital 0%

County 4%

Country  96%

Powered by  Poll Everywhere



# Get to know the data (1/3)

- Load data

```
1 library(readxl)  
2 gapm1 <- read_excel(here("data/Gapminder_vars_2011.xlsx"), na = "NA")
```

"missing"  
na = "missing"

na = "NA"

"NA"

that anything  
w/ "NA" is a  
missing value

# Get to know the data (2/3)

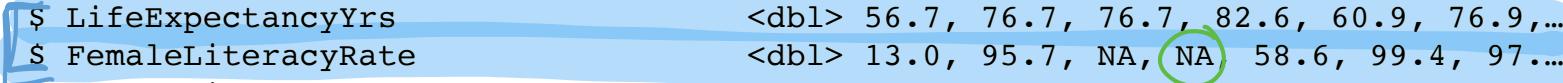
- Glimpse of the data

```
1 glimpse(gapm1)
```

Rows: 195

Columns: 18

\$ country  
\$ CO2emissions  
\$ ElectricityUsePP  
\$ FoodSupplykcPPD  
\$ IncomePP  
\$ LifeExpectancyYrs  
\$ FemaleLiteracyRate  
\$ population  
\$ WaterSourcePrct  
\$ geo  
\$ four\_regions  
\$ eight\_regions  
\$ six\_regions  
\$ members\_oecd\_g77  
\$ Latitude  
\$ Longitude  
\$ `World bank region`  
\$ `World bank, 4 income groups 2017`


``Ilama''

"Afghanistan", "Albania", "Algeria"...,  
0.412, 1.790, 3.290, 5.870, 1.250, ...  
NA 2210.0, 1120.0, NA 207.0, NA, ...  
2110, 3130, 3220, NA 2410, 2370, 3...  
1660, 10200, 13000, 42000, 5910, 18...  
56.7, 76.7, 76.7, 82.6, 60.9, 76.9,...  
13.0, 95.7, NA, NA 58.6, 99.4, 97....  
2.97e+07, 2.93e+06, 3.68e+07, 8.38e...  
52.6, 88.1, 92.6, 100.0, 40.3, 97.0...  
"afg", "alb", "dza", "and", "ago", ...  
"asia", "europe", "africa", "europe...  
"asia\_west", "europe\_east", "africa...  
"south\_asia", "europe\_central\_asia"...  
"g77", "others", "g77", "others", "...  
33.00000, 41.00000, 28.00000, 42.50...  
66.00000, 20.00000, 3.00000, 1.5210...  
"South Asia", "Europe & Central Asi...  
"Low income", "Upper middle income"...

- Note the missing values for our variables of interest

## Get to know the data (3/3)

- Get a sense of the summary statistics

```
1 gapm1 %>%
2   select(LifeExpectancyYrs,
3           FemaleLiteracyRate) %>%
4   summary()
```

	LifeExpectancyYrs	FemaleLiteracyRate
Min.	:47.50	Min. :13.00
1st Qu.	:64.30	1st Qu.:70.97
Median	:72.70	Median :91.60
Mean	:70.66	Mean :81.65
3rd Qu.	:76.90	3rd Qu.:98.03
Max.	:82.90	Max. :99.80
NA's	:8	NA's :115

# Remove missing values (1/2)

- Remove rows with missing data for life expectancy and female literacy rate

```
1 gapm <- gapm1 %>% drop_na(LifeExpectancyYrs, FemaleLiteracyRate)  
2 glimpse(gapm)
```

Rows: 80

Columns: 18

```
$ country                                <chr> "Afghanistan", "Albania", "Angola",...  
$ CO2emissions                            <dbl> 0.4120, 1.7900, 1.2500, 5.3600, 4.6...  
$ ElectricityUsePP                         <dbl> NA, 2210.0, 207.0, NA, 2900.0, 1810...  
$ FoodSupplykcPPD                          <dbl> 2110, 3130, 2410, 2370, 3160, 2790,...  
$ IncomePP                                 <dbl> 1660, 10200, 5910, 18600, 19600, 70...  
$ LifeExpectancyYrs                        <dbl> 56.7, 76.7, 60.9, 76.9, 76.0, 73.8,...  
$ FemaleLiteracyRate                        <dbl> 13.0, 95.7, 58.6, 99.4, 97.9, 99.5,...  
$ population                               <dbl> 2.97e+07, 2.93e+06, 2.42e+07, 9.57e...  
$ WaterSourcePrct                           <dbl> 52.6, 88.1, 40.3, 97.0, 99.5, 97.8,...  
$ geo                                      <chr> "afg", "alb", "ago", "atg", "arg", ...  
$ four_regions                             <chr> "asia", "europe", "africa", "americ...  
$ eight_regions                            <chr> "asia_west", "europe_east", "africa...  
$ six_regions                             <chr> "south_asia", "europe_central_asia"...  
$ members_oecd_g77                          <chr> "g77", "others", "g77", "g77", "g77...  
$ Latitude                                  <dbl> 33.00000, 41.00000, -12.50000, 17.0...  
$ Longitude                                 <dbl> 66.00000, 20.00000, 18.50000, -61.8...  
$ `World bank region`                     <chr> "South Asia", "Europe & Central Asi...  
$ `World bank, 4 income groups 2017`       <chr> "Low income", "Upper middle income"...
```

only countries with both life exp & FLR observed

## Remove missing values (2/2)

- And no more missing values when we look only at our two variables of interest

```
1 gapm %>%
2   select(LifeExpectancyYrs,
3         FemaleLiteracyRate) %>%
4   get_summary_stats()

# A tibble: 2 × 13
  variable      n    min    max median     q1     q3    iqr    mad mean     sd     se
  <fct>    <dbl> <dbl>
1 LifeExpect...    80    48   81.8   72.4   65.9   75.8   9.95   6.30   69.9   7.95  0.889
2 FemaleLite...    80    13   99.8   91.6   71.0   98.0  27.0   11.4   81.7  22.0   2.45
# i 1 more variable: ci <dbl>
```

### Note

- Removing the rows with missing data was not needed to run the regression model.
- I did this step since later we will be calculating the standard deviations of the explanatory and response variables for *just the values included in the regression model*. It'll be easier to do this if we remove the missing values now.

# Poll Everywhere Question 2

13:34 Mon Jan 13 95%

X

Join by Web [PollEv.com/nickywakim275](https://PollEv.com/nickywakim275)



What are other ways you would get to know your data? (Hint: What else have we learned to visualize or summarize the data?)

Summary stats, a scatter plot, histogram.

graph it

0 0

0 0

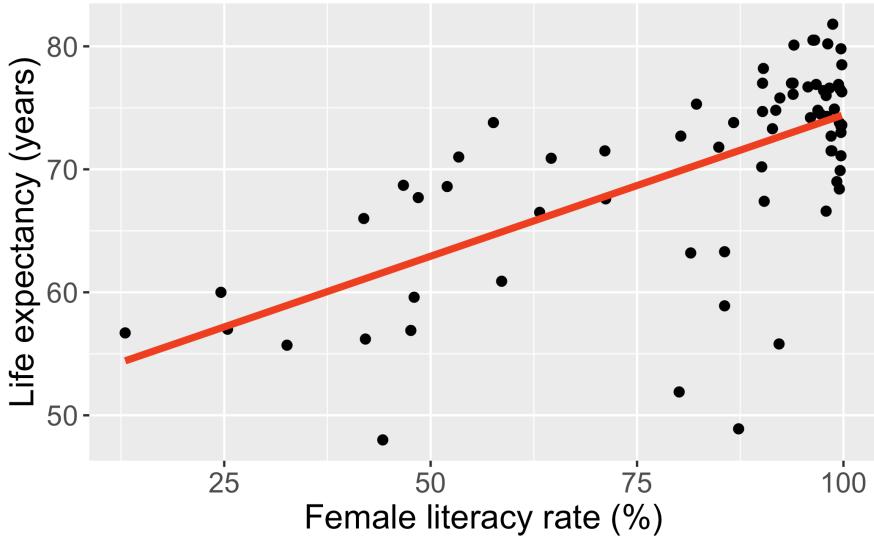
Powered by  Poll Everywhere

# Learning Objectives

1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4. Solve the optimal coefficient estimates for simple linear regression using OLS
5. Apply OLS in R for simple linear regression of real data

# Questions we can ask with a simple linear regression model

Relationship between life expectancy and the female literacy rate in 2011



- How do we...
  - calculate slope & intercept? ✓
  - interpret slope & intercept? ✓
  - do inference for slope & intercept?
    - CI, p-value
  - do prediction with regression line?
    - CI for prediction?
- Does the model fit the data well?
  - Should we be using a line to model the data?
- Should we add additional variables to the model?
  - multiple/multivariable regression

$$\widehat{\text{life expectancy}} = 50.9 + 0.232 \cdot \text{female literacy rate}$$

# Association vs. prediction

## Association

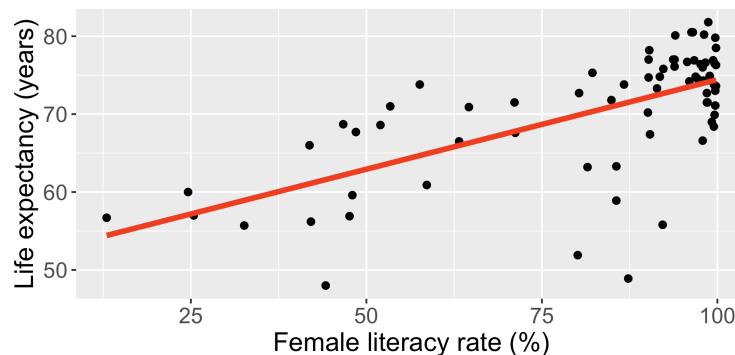
- What is the association between countries' life expectancy and female literacy rate?
- Use the slope of the line or correlation coefficient

## Prediction

- What is the expected ~~average~~ life expectancy for a country with a specified female literacy rate?

$$\widehat{\text{life expectancy}} = 50.9 + 0.232 \cdot \text{female literacy rate}$$

Relationship between life expectancy and the female literacy rate in 2011



# Three types of study design (there are more)



## Experiment

- Observational units are randomly assigned to important predictor levels
  - Random assignment controls for confounding variables (age, gender, race, etc.)
  - “gold standard” for determining causality
  - Observational unit is often at the participant-level

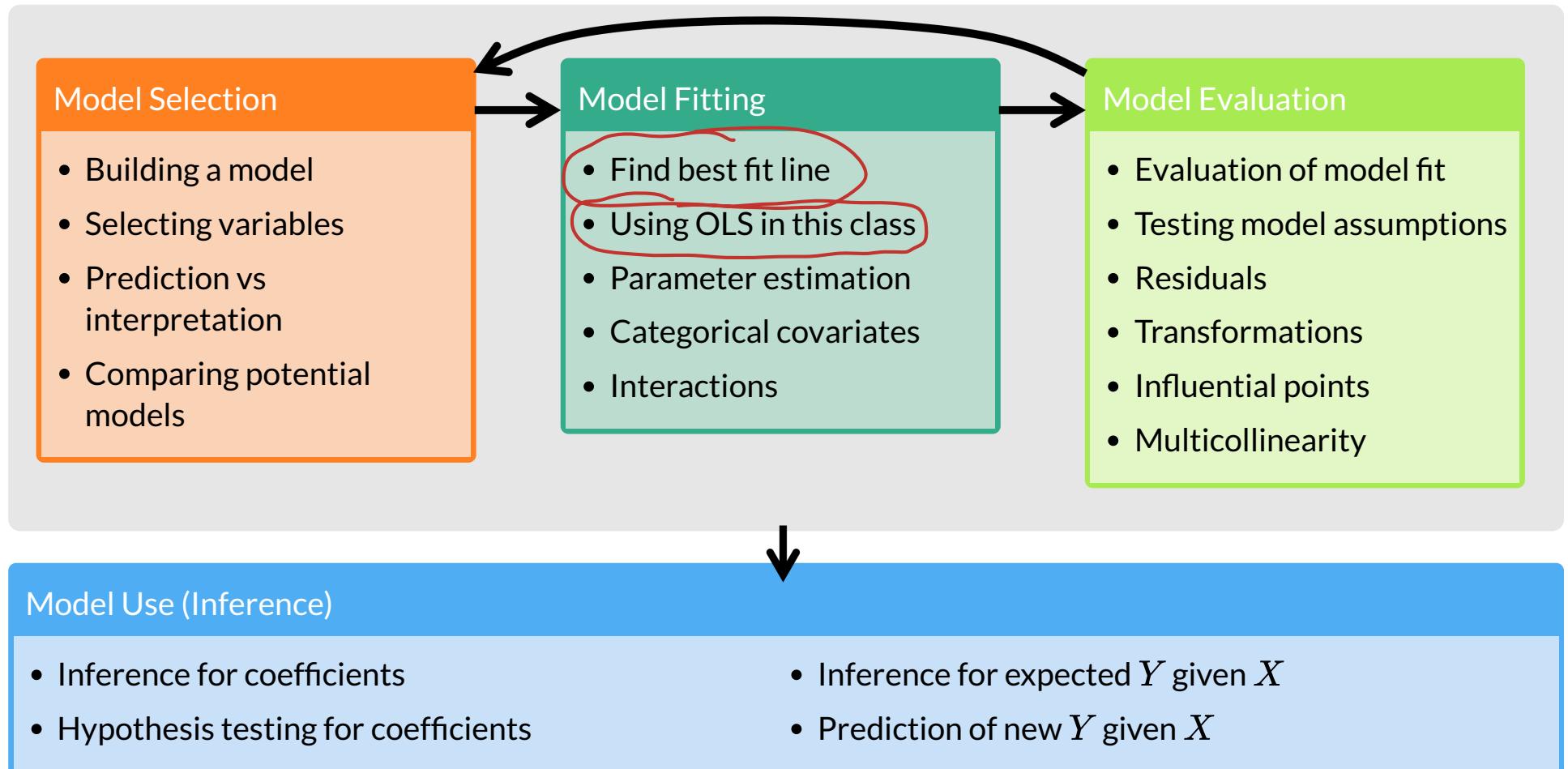
## Quasi-experiment

- Participants are assigned to intervention levels without randomization
- Not common study design

## Observational

- No randomization or assignment of intervention conditions
- In general cannot infer causality
  - However, there are causal inference methods...

# Let's revisit the regression analysis process



# Poll Everywhere Question 3

13:42 Mon Jan 13

Join by Web [PollEv.com/nickywakim275](https://PollEv.com/nickywakim275)

Why is there no model selection in simple linear regression?

We're just trying to keep things simple! 4%

SLR only has one predictor by definition ✓ 71%

There is! We can pick which predictor to use! ✓ 17%

We only want to learn model fitting right now! 8%

Powered by  Poll Everywhere



# Learning Objectives

1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4. Solve the optimal coefficient estimates for simple linear regression using OLS
5. Apply OLS in R for simple linear regression of real data

# Simple Linear Regression Model

The (population) regression model is denoted by:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

## Observable sample data

- $Y$  is our dependent variable
  - Aka outcome or response variable
- $X$  is our independent variable
  - Aka predictor, regressor, exposure variable

## Unobservable population parameters

- $\beta_0$  and  $\beta_1$  are **unknown** population parameters
- $\epsilon$  (epsilon) is the error about the line
  - It is assumed to be a random variable with a...
  - Normal distribution with mean 0 and constant variance  $\sigma^2$
  - i.e.  $\epsilon \sim N(0, \sigma^2)$

# Simple Linear Regression Model (another way to view components)

The (population) regression model is denoted by:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

## Components

$Y$  response, outcome, dependent variable

$\beta_0$  intercept

$\beta_1$  slope

$X$  predictor, covariate, independent variable

*exposure variable*

$\epsilon$  residuals, error term

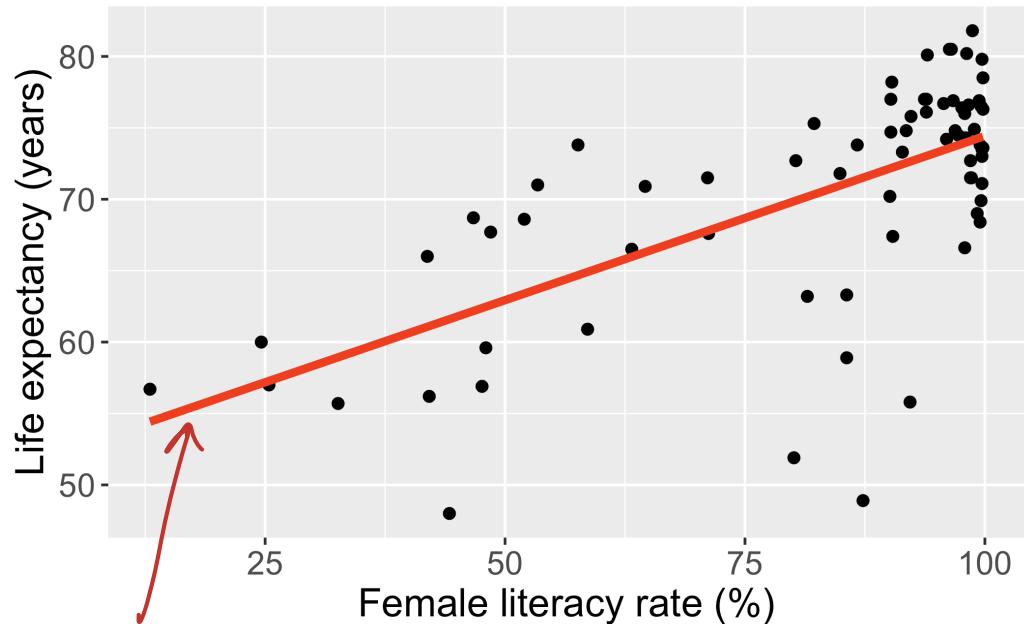
*epsilon*

# If the population parameters are unobservable, how did we get the line for life expectancy?

Note: the **population model** is the true, underlying model that we are trying to estimate using our sample data

- Our goal in simple linear regression is to estimate  $\beta_0$  and  $\beta_1$

Relationship between life expectancy and the female literacy rate in 2011



estimate  
(not population model)

# Poll Everywhere Question 4

13:48 Mon Jan 13

Join by Web [PollEv.com/nickywakim275](https://PollEv.com/nickywakim275)

Sample  
we use data  
to fit the red  
line

What do we label as the slope of the red line?

Relationship between life expectancy and the female literacy rate in 2011

Life expectancy (years)

Female literacy rate (%)

Powered by Poll Everywhere

$\beta_0$  7% → population intercept

$\hat{\beta}_0$  0% → estimated intercept (from sample)

$\beta_1$  66% → population slope

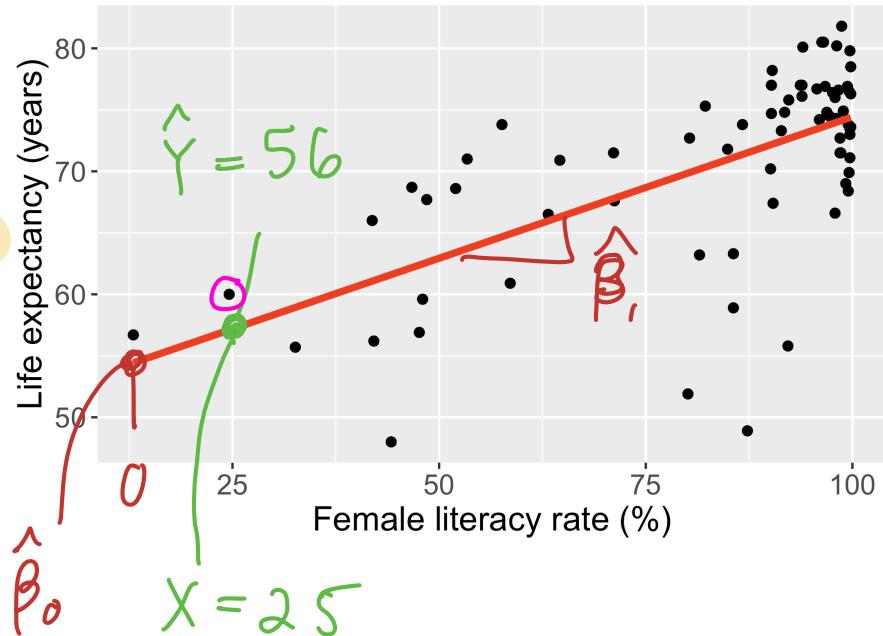
$\hat{\beta}_1$  28% → estimated slope (from sample)

# Regression line = best-fit line

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

- $\hat{Y}$  is the predicted outcome for a specific value of  $X$   
*or avg outcome*
- $\hat{\beta}_0$  is the intercept of the best-fit line
- $\hat{\beta}_1$  is the slope of the best-fit line, i.e., the increase in  $\hat{Y}$  for every increase of one (unit increase) in  $X$ 
  - slope = rise over run

Relationship between life expectancy and the female literacy rate in 2011



# Simple Linear Regression Model

Population regression model

$$Y = \beta_0 + \beta_1 X + \epsilon$$

## Components

$Y$  response, outcome, dependent variable

$\beta_0$  intercept

$\beta_1$  slope

$X$  predictor, covariate, independent variable

$\epsilon$  residuals, error term

Once we fit data:

Estimated regression line

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

~~$\hat{\beta}_0$~~

## Components

$\hat{Y}$  estimated expected response given predictor  $X$

$\hat{\beta}_0$  estimated intercept

$\hat{\beta}_1$  estimated slope

$X$  predictor, covariate, independent variable

$$\begin{aligned}\hat{Y} &= 59.2 + 0.232 X \\ &\cdot 10 \\ Y &= mx + b\end{aligned}$$

# We get it, Nicky! How do we estimate the regression line?

First let's take a break!!



Lat: looking for the  
truth, but only  
estimating it

detective beta

# Learning Objectives

1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4. Solve the optimal coefficient estimates for simple linear regression using OLS
5. Apply OLS in R for simple linear regression of real data

# It all starts with a residual...

- Recall, one characteristic of our population model was that the residuals,  $\epsilon$ , were Normally distributed:  $\epsilon \sim N(0, \sigma^2)$

- In our population regression model, we had:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

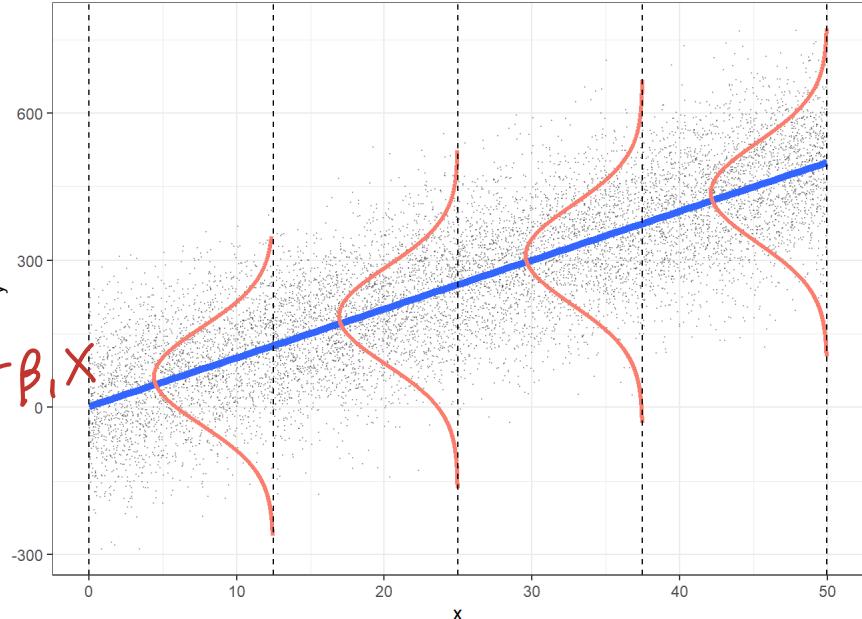
pop line:  
 $E(Y|X) = \beta_0 + \beta_1 X$

- We can also take the average (expected) value of the population model
- We take the expected value of both sides and get:

$$\begin{aligned} E[Y] &= E[\beta_0 + \beta_1 X + \epsilon] \\ E[Y] &= E[\beta_0] + E[\beta_1 X] + E[\epsilon] \\ E[Y] &= \beta_0 + \beta_1 X + E[\epsilon] \\ E[Y|X] &= \beta_0 + \beta_1 X \end{aligned}$$

given

- We call  $E[Y|X]$  the expected value (or average) of  $Y$  given  $X$



b/c Normally dist w/ mean of 0  
mean =  $E(\epsilon) = 0$

# So now we have two representations of our population model

With observed  $Y$  values and residuals:

pop  
model

$$Y = \beta_0 + \beta_1 X + \epsilon$$

With the population expected value of  $Y$  given  $X$ :

$$E[Y|X] = \beta_0 + \beta_1 X$$

Using the two forms of the model, we can figure out a formula for our residuals:

$$\begin{aligned} &\rightarrow Y = (\beta_0 + \beta_1 X) + \epsilon \\ &Y = E[Y|X] + \epsilon \\ &Y - E[Y|X] = \epsilon \\ &\epsilon = Y - E[Y|X] \end{aligned}$$

And so we have our **true, population model**, residuals!

This is an important fact! For the **population model**, the residuals:  $\epsilon = Y - E[Y|X]$

# Back to our estimated model

We have the same two representations of our estimated/fitted model:

With observed values:

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\epsilon}$$

estimated residual

With the estimated expected value of  $Y$  given  $X$ :

$$\begin{aligned}\hat{E}[Y|X] &= \hat{\beta}_0 + \hat{\beta}_1 X \\ \widehat{E[Y|X]} &= \hat{\beta}_0 + \hat{\beta}_1 X \\ \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X\end{aligned}$$

Using the two forms of the model, we can figure out a formula for our estimated residuals:

$$Y = (\hat{\beta}_0 + \hat{\beta}_1 X) + \hat{\epsilon}$$

$$Y = \hat{Y} + \hat{\epsilon}$$

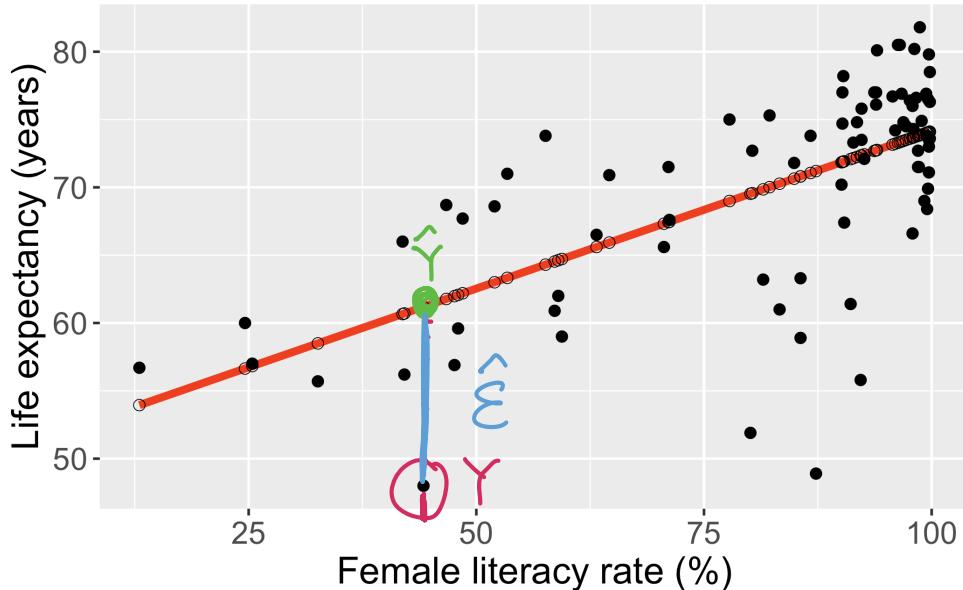
$$\hat{\epsilon} = Y - \hat{Y}$$

This is an important fact! For the **estimated/fitted model**, the residuals:  $\hat{\epsilon} = Y - \hat{Y}$

# *Individual $i$ residuals in the estimated/fitted model*

- Observed values for each individual  $i$ :  $Y_i$ 
  - Value in the dataset for individual  $i$
- Fitted value for each individual  $i$ :  $\hat{Y}_i$ 
  - Value that falls on the best-fit line for a specific  $X_i$
  - If two individuals have the same  $X_i$ , then they have the same  $\hat{Y}_i$

Relationship between life expectancy and the female literacy rate in 2011

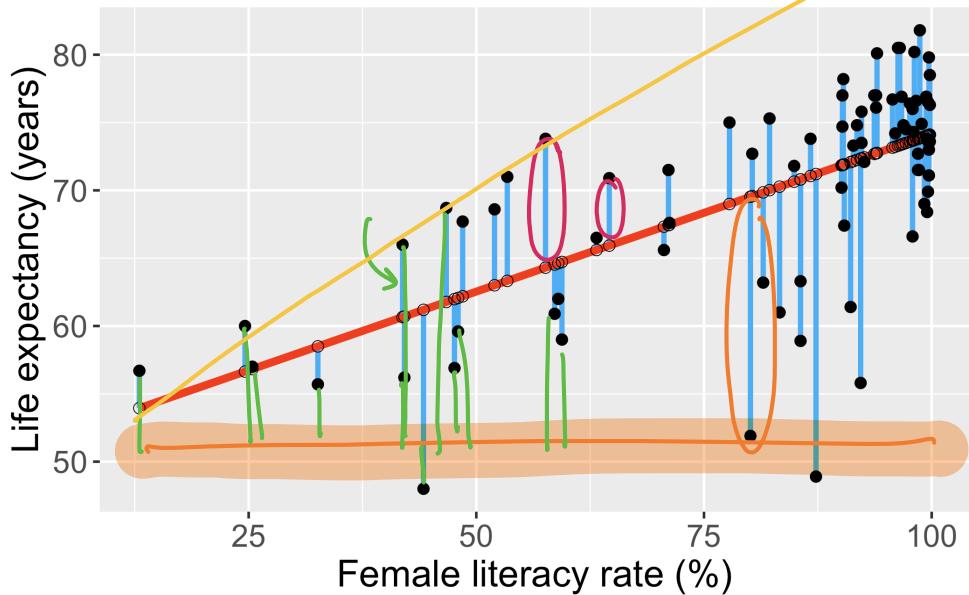


residual is diff b/w observed  $Y$  & fitted  $Y$  (est. expected val of  $Y$  given  $X$ )

# *Individual $i$ residuals in the estimated/fitted model*

- Observed values for each individual  $i$ :  $Y_i$ 
  - Value in the dataset for individual  $i$
- Fitted value for each individual  $i$ :  $\hat{Y}_i$ 
  - Value that falls on the best-fit line for a specific  $X_i$
  - If two individuals have the same  $X_i$ , then they have the same  $\hat{Y}_i$
- Residual for each individual:  $\hat{\epsilon}_i = Y_i - \hat{Y}_i$ 
  - Difference between the observed and fitted value

Relationship between life expectancy and the female literacy rate in 2011



# Poll Everywhere Question 5

14:23 Mon Jan 13

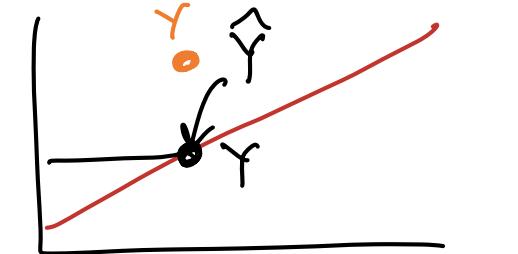
Join by Web [PollEv.com/nickywakim275](https://PollEv.com/nickywakim275)

QR code:

If our observed  $\hat{Y}$  value fell exactly on the best-fit line, what would the residual be?

0  
0

Powered by  Poll Everywhere


$$\hat{Y}_i - Y_i = 0$$

b/c  $\hat{Y}_i = Y_i$

# So what do we do with the residuals?

- We want to minimize the residuals
  - Aka minimize the difference between the observed  $Y$  value and the estimated expected response given the predictor ( $\hat{E}[Y|X]$ )
- We can use ordinary least squares (OLS) to do this in linear regression!
- Idea behind this: reduce the total error between the fitted line and the observed point (error between is called residuals)
  - Vague use of total error: more precisely, we want to reduce the sum of squared errors
  - ~~Think back to my R Shiny app!~~
  - We need to mathematically define this!
- Note: there are other ways to estimate the best-fit line!!
  - Example: Maximum likelihood estimation

# Learning Objectives

1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4. Solve the optimal coefficient estimates for simple linear regression using OLS
5. Apply OLS in R for simple linear regression of real data

# Setting up for ordinary least squares

- Sum of Squared Errors (SSE)

$$\Rightarrow SSE = \sum_{i=1}^n \hat{\epsilon}_i^2 \quad \text{sum of errors}$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$SSE = \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

Things to use

- $\hat{\epsilon}_i = Y_i - \hat{Y}_i$
- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

Then we want to find the estimated coefficient values that minimize the SSE!

# Steps to estimate coefficients using OLS

1. Set up SSE (previous slide)
2. Minimize SSE with respect to coefficient estimates →  $\hat{\beta}_0, \hat{\beta}_1$  are coefficient estimates
  - Need to solve a system of equations
3. Compute derivative of SSE wrt  $\hat{\beta}_0$
4. Set derivative of SSE wrt  $\hat{\beta}_0 = 0$
5. Compute derivative of SSE wrt  $\hat{\beta}_1$
6. Set derivative of SSE wrt  $\hat{\beta}_1 = 0$
7. Substitute  $\hat{\beta}_1$  back into  $\hat{\beta}_0$

## 2. Minimize SSE with respect to coefficients

- Want to minimize with respect to (wrt) the potential coefficient estimates ( $\hat{\beta}_0$  and  $\hat{\beta}_1$ )
- Take derivative of SSE wrt  $\hat{\beta}_0$  and  $\hat{\beta}_1$  and set equal to zero to find minimum SSE

$$\frac{\partial SSE}{\partial \hat{\beta}_0} = 0 \text{ and } \frac{\partial SSE}{\partial \hat{\beta}_1} = 0$$

- Solve the above system of equations in steps 3-6

### 3. Compute derivative of SSE wrt $\hat{\beta}_0$

$$\rightarrow SSE = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\begin{aligned}\frac{\partial SSE}{\partial \hat{\beta}_0} &= \frac{\partial \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\partial \hat{\beta}_0} = \sum_{i=1}^n \frac{\partial (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\partial \hat{\beta}_0} \\ &= \sum_{i=1}^n 2(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)(-1) = \sum_{i=1}^n -2(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)\end{aligned}$$

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

Things to use

- Derivative rule: derivative of sum is sum of derivative
- Derivative rule: chain rule

## 4. Set derivative of SSE wrt $\hat{\beta}_0 = 0$

$$\frac{\partial SSE}{\partial \hat{\beta}_0} = 0$$

$$-2 \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0$$

$$\sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0$$

$$\sum_{i=1}^n Y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n X_i = 0$$

$$\frac{1}{n} \sum_{i=1}^n Y_i - \hat{\beta}_0 - \hat{\beta}_1 \frac{1}{n} \sum_{i=1}^n X_i = 0$$

$$\bar{Y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{X} = 0$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Things to use

- $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$
- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

## 5. Compute derivative of SSE wrt $\hat{\beta}_1$

$$SSE = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\begin{aligned}\frac{\partial SSE}{\partial \hat{\beta}_1} &= \frac{\partial \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\partial \hat{\beta}_1} = \sum_{i=1}^n \frac{\partial (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{\partial \hat{\beta}_1} \\ &= \sum_{i=1}^n 2(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)(-X_i) = \sum_{i=1}^n -2X_i(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) \\ &= -2 \sum_{i=1}^n X_i(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)\end{aligned}$$

Things to use

- Derivative rule: derivative of sum is sum of derivative
- Derivative rule: chain rule

## 6. Set derivative of SSE wrt $\hat{\beta}_1 = 0$

$$\frac{\partial SSE}{\partial \hat{\beta}_1} = 0$$

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

$$\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \hat{\beta}_0 - \sum_{i=1}^n X_i^2 \hat{\beta}_1 = 0$$

$$\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \left( \bar{Y} - \hat{\beta}_1 \bar{X} \right) - \sum_{i=1}^n X_i^2 \hat{\beta}_1 = 0$$

$$\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \bar{Y} + \sum_{i=1}^n \hat{\beta}_1 X_i \bar{X} - \sum_{i=1}^n X_i^2 \hat{\beta}_1 = 0$$

$$\sum_{i=1}^n X_i (Y_i - \bar{Y}) + \sum_{i=1}^n (\hat{\beta}_1 X_i \bar{X} - X_i^2 \hat{\beta}_1) = 0$$

$$\sum_{i=1}^n X_i (Y_i - \bar{Y}) + \hat{\beta}_1 \sum_{i=1}^n X_i (\bar{X} - X_i) = 0$$

Things to use

- $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$
- $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$
- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$



$$\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i (Y_i - \bar{Y})}{\sum_{i=1}^n X_i (X_i - \bar{X})}$$

## 7. Substitute $\hat{\beta}_1$ back into $\hat{\beta}_0$

### Final coefficient estimates for SLR

Coefficient estimate for  $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i(Y_i - \bar{Y})}{\sum_{i=1}^n X_i(X_i - \bar{X})}$$

Coefficient estimate for  $\hat{\beta}_0$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\hat{\beta}_0 = \bar{Y} - \frac{\sum_{i=1}^n X_i(Y_i - \bar{Y})}{\sum_{i=1}^n X_i(X_i - \bar{X})} \bar{X}$$

# Poll Everywhere Question 6

14:36 Mon Jan 13 ... 76% 

X Join by Web PollEv.com/nickywakim275

What do  $\hat{\beta}_0$  and  $\hat{\beta}_1$  mean for our model?

They are the coefficient estimates that minimize every residual value 0%

They are the coefficient estimates that are closest to the population parameters 14%

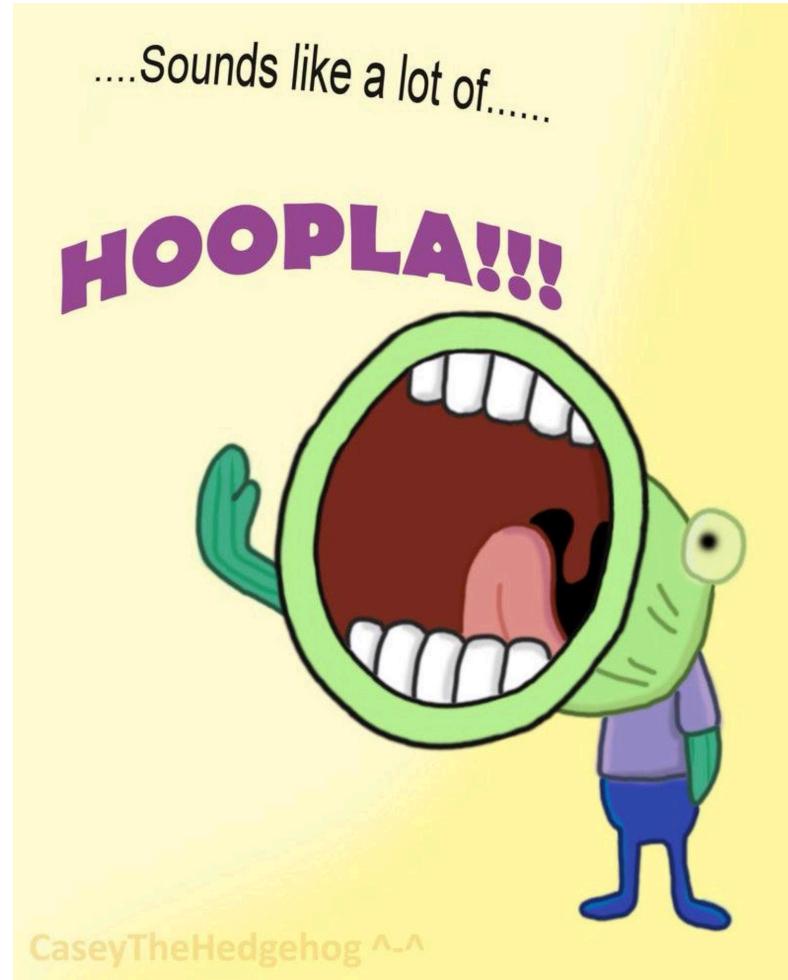
They are the coefficient estimates that perfectly fit our data 3%

They are the coefficient estimates that minimize the sum of the squared residuals 83%

Powered by  Poll Everywhere



# Do I need to do all that work every time??



# Regression in R: `lm()`

- Let's discuss the syntax of this function

```
1 modell <- gapm %>% lm(formula = LifeExpectancyYrs ~ FemaleLiteracyRate)
```

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

In the general form:

```
1 lm( Y ~ X, data = dataset_name)
2 dataset_name %>% lm( formula = Y ~ X )
```

female\_lit\_rate

# Regression in R: `lm()` + `summary()`

```
1 modell <- gapm %>% lm(formula = LifeExpectancyYrs ~ FemaleLiteracyRate)  
2 summary(modell)
```

Call:  
`lm(formula = LifeExpectancyYrs ~ FemaleLiteracyRate, data = .)`

Residuals:

Min	1Q	Median	3Q	Max
-22.299	-2.670	1.145	4.114	9.498

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept) $\hat{\beta}_0$	50.92790	2.66041	19.143	< 2e-16	***
FemaleLiteracyRate $\hat{\beta}_1$	0.23220	0.03148	7.377	1.5e-10	***
---					

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.142 on 78 degrees of freedom

Multiple R-squared: 0.4109, Adjusted R-squared: 0.4034

F-statistic: 54.41 on 1 and 78 DF, p-value: 1.501e-10

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

life expectancy =

50.9 +

0.232 ·

femal e  
literacy  
rate

# Regression in R: `lm()` + `tidy()`

```
1 tidy(model1) %>%
2   gt() %>%
3     tab_options(table.font.size = 45)
```

more presentable in slides

term	estimate	std.error	<u>t</u> statistic	p.value
(Intercept)	50.9278981	2.66040695	19.142898	3.325312e-31
FemaleLiteracyRate	0.2321951	0.03147744	7.376557	1.501286e-10

*estimated (not pop)*

- Regression equation for our model (which we saw a looong time ago):

$$\widehat{\text{life expectancy}} = 50.9 + 0.232 \cdot \text{female literacy rate}$$

# How do we interpret the coefficients?

$$\text{life expectancy} = 50.9 + 0.232 \cdot \text{female literacy rate}$$

↑ 1%

- Intercept ( $\hat{\beta}_0$ )

- The expected outcome for the  $Y$ -variable when the  $X$ -variable (if continuous) is 0

→ Example: The expected/average life expectancy is 50.9 years for a country with 0% female literacy.

- Slope ( $\hat{\beta}_1$ )

- For every increase of 1 unit in the  $X$ -variable (if continuous), there is an expected increase of, ~~0.232~~,  $\hat{\beta}_1$  units in the  $Y$ -variable.
- We only say that there is an expected increase and not necessarily a causal increase.

→ Example: For every 1 percent increase in the female literacy rate, life expectancy increases, on average, 0.232 years.

- Can also say "...average life expectancy increases 0.232..."

avg OR expected

$$\left\{ \begin{array}{l} E(Y|X) \\ \hat{Y} \end{array} \right.$$
$$\hat{E}(Y|X)$$

## Next time

- More on interpreting the estimate coefficients
- Inference of our estimated coefficients
- Inference of estimated expected  $Y$  given  $X$
- Prediction
- Hypothesis testing!

