

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

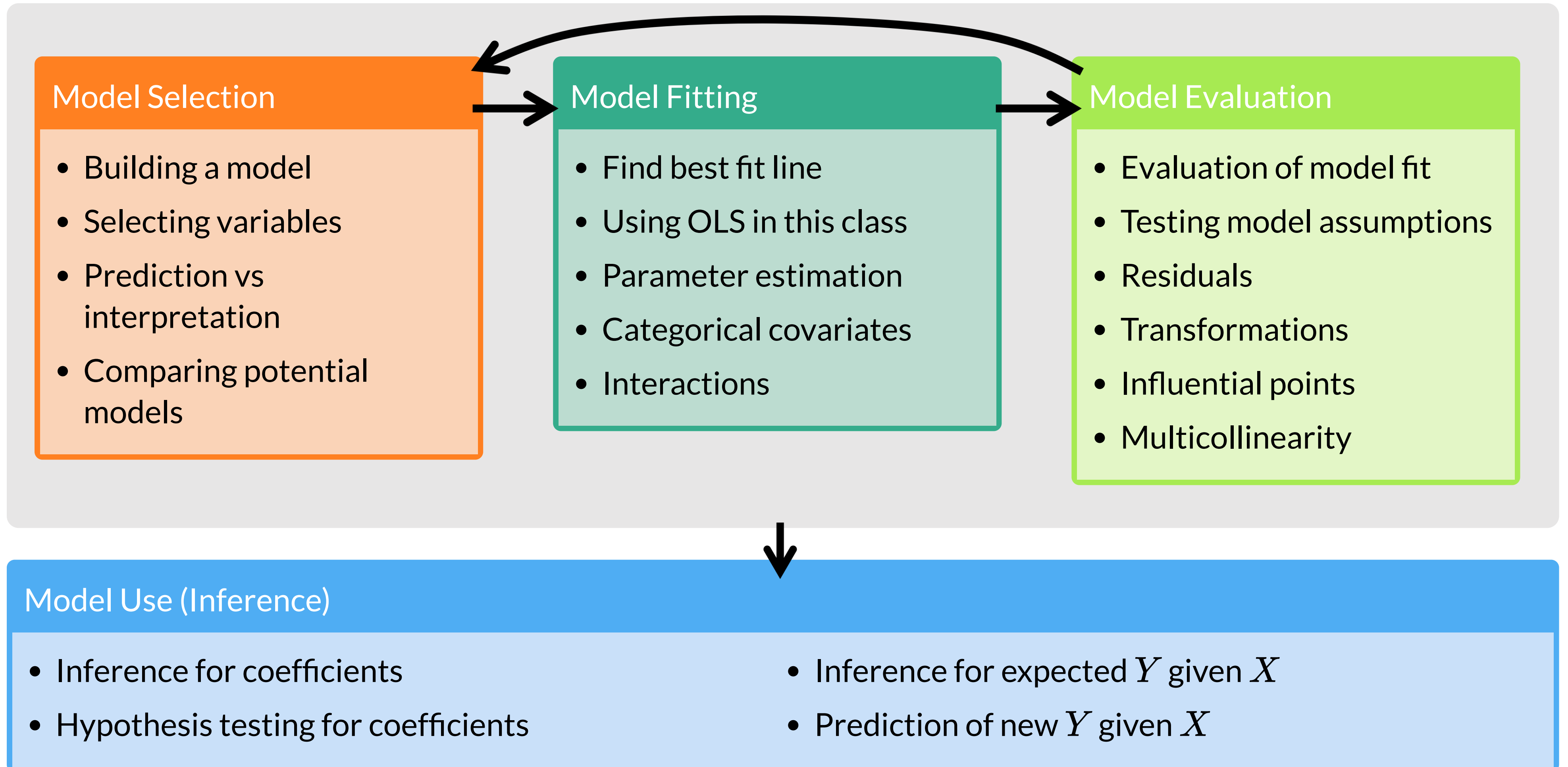
Nicky Wakim

2026-01-12

# Learning Objectives

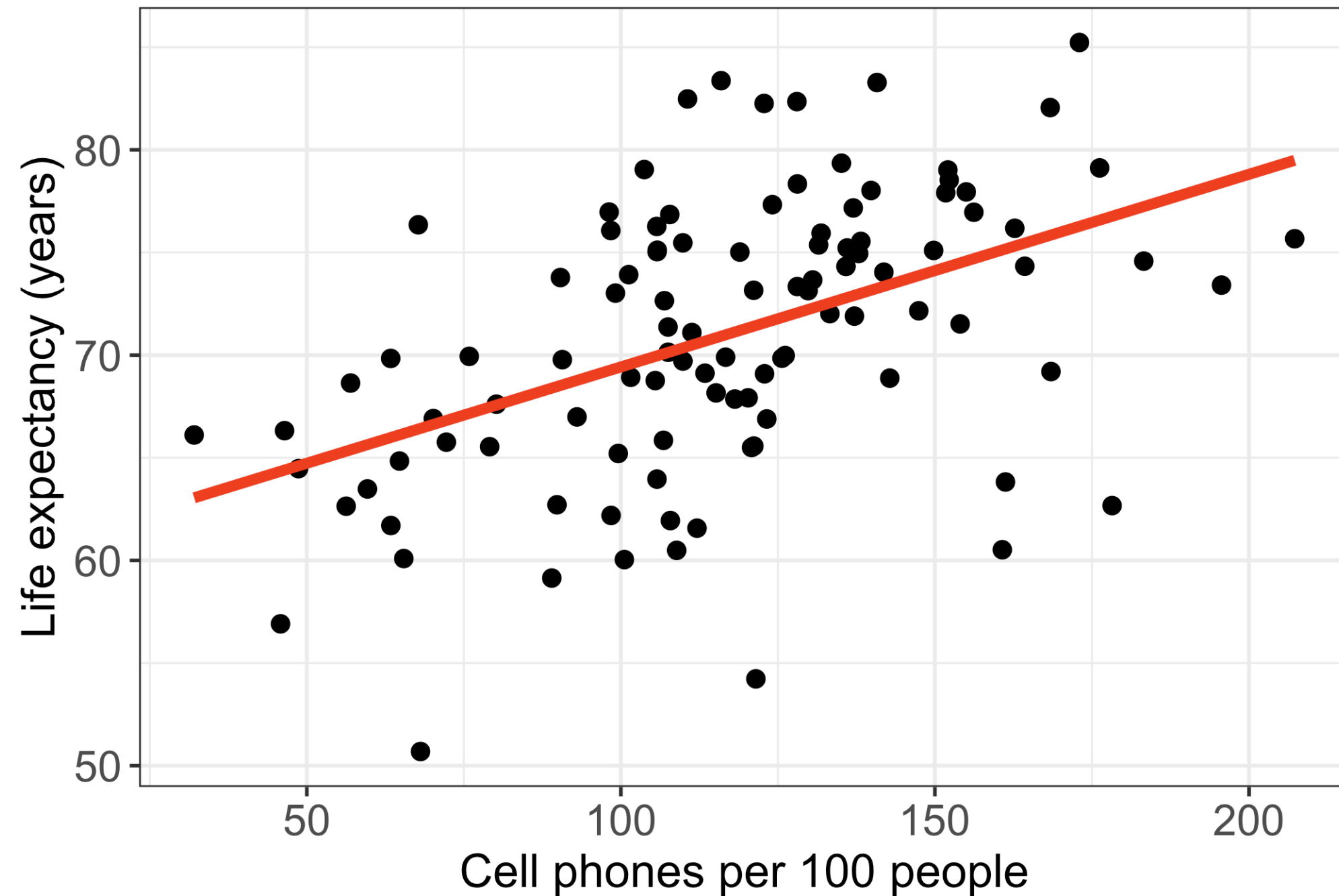
1. Identify the aims of your research and see how they align with the intended purpose of simple linear regression
2. Identify the simple linear regression model and define statistics language for key notation
3. Illustrate how ordinary least squares (OLS) finds the best model parameter estimates
4.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



# Let's start with an example

Relationship between life expectancy and cell phones



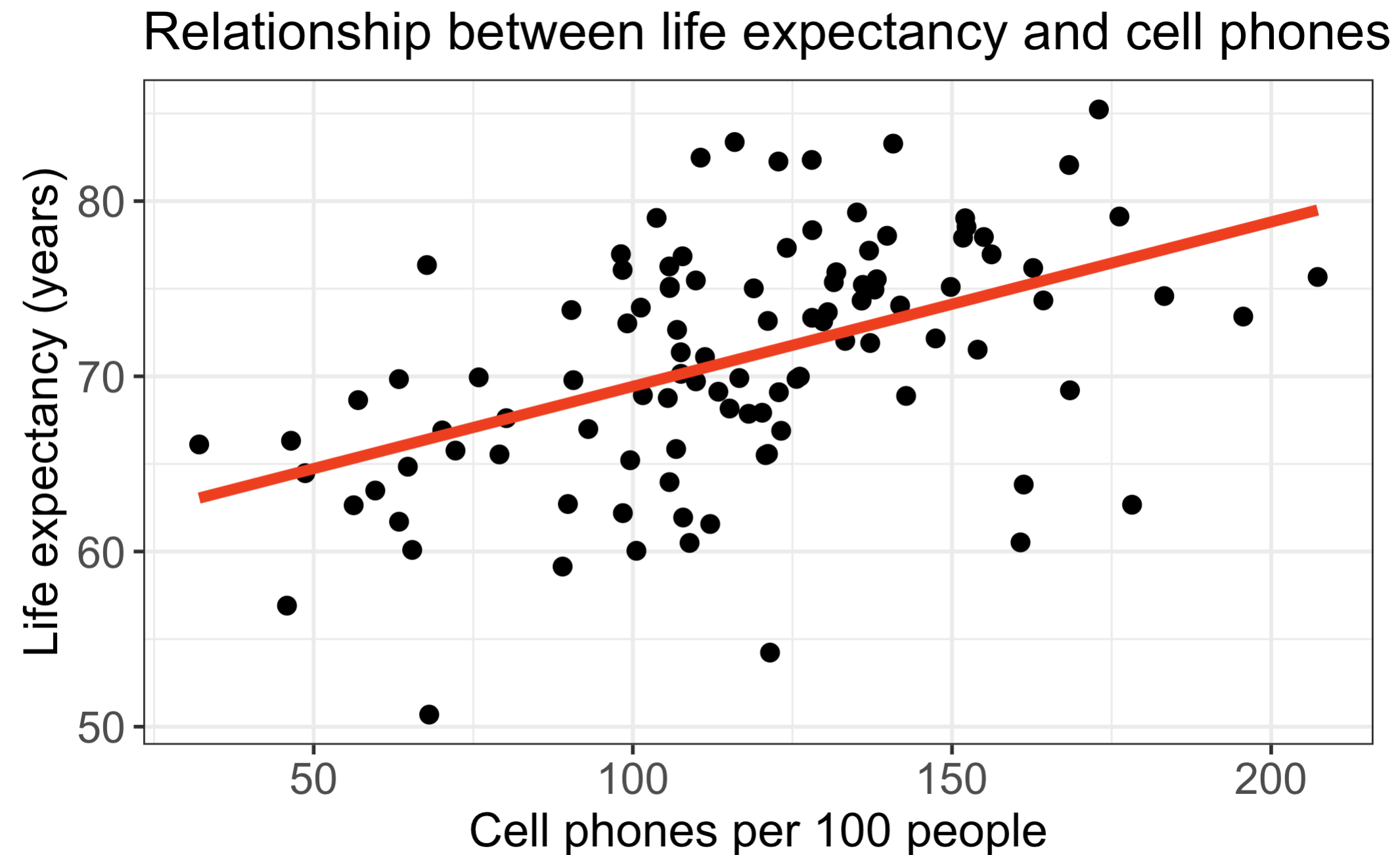
Life expectancy vs. cell phones

- Each point on the plot is for a different country/territory
- $X$  = country's number of cell phones per 100 people
- $Y$  = country's life expectancy (years)

$$\widehat{\text{life expectancy}} = 60.04 + 0.094 \cdot \text{cell phones}$$

# Reference: How did I code that?

```
1 gapm %>%
2   ggplot(aes(x = cell_phones_100,
3             y = life_exp)) +
4   geom_point(size = 4) +
5   geom_smooth(method = "lm", se = FALSE, size = 3, colour="#F14124") +
6   labs(x = "Cell phones per 100 people",
7        y = "Life expectancy (years)",
8        title = "Relationship between life expectancy and cell phones") +
9   theme(axis.title = element_text(size = 27),
10        axis.text = element_text(size = 25),
11        title = element_text(size = 25))
```



# Research and dataset description

Research question: Is there an association between life expectancy and number of cell phones?

- Data file: `gapminder.Rdata`
- [You can access my codebook here](#)
- Data were downloaded from [Gapminder](#)
  - 2022 is the most recent year with the most complete data
  - Observational study measuring different characteristics of countries/territories, 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.
- **Cell phones per 100 people** is the number of cell phone subscriptions per 100 people in a given population, indicating the level of mobile phone use and accessibility.

# Poll Everywhere Question 1

# Get to know the data (1/3)

- Load data

```
1 load(here("data", "gapminder2022.RData"))
```

# Get to know the data (2/3)

- Glimpse of the data

```
1 glimpse(gapm)
```

```
Rows: 105
```

```
Columns: 11
```

```
$ geo      <chr> "afg", "alb", "are", "arg", "arm", "aze", "ben", "...
$ territory <chr> "Afghanistan", "Albania", "UAE", "Argentina", "Arm...
$ life_exp <dbl> 62.64, 76.07, 73.41, 75.37, 73.66, 71.37, 63.96, 7...
$ freedom_status <chr> "NF", "PF", "NF", "F", "PF", "NF", "PF", "PF", "F"...
$ vax_rate  <dbl> 69, 99, 98, 94, 98, 96, 89, 99, 96, 96, 76, 88, 98...
$ co2_emissions <dbl> 211455404, 294574910, 5324389134, 8574249437, 4512...
$ basic_sani <dbl> 70.39219, 99.30948, 98.97272, 98.46960, 100.00000,...
$ happiness_score <dbl> 12.81, 52.12, 67.38, 62.61, 53.82, 45.76, 42.17, 3...
$ income_level_4 <chr> "Low income", "Upper middle income", "High income"...
$ cell_phones_100 <dbl> 56.2655, 98.3950, 195.6250, 131.4840, 130.5400, 10...
$ basic_sani_80_above <chr> "Low access", "High access", "High access", "High ...
```

# Get to know the data (3/3)

- Get a sense of the summary statistics

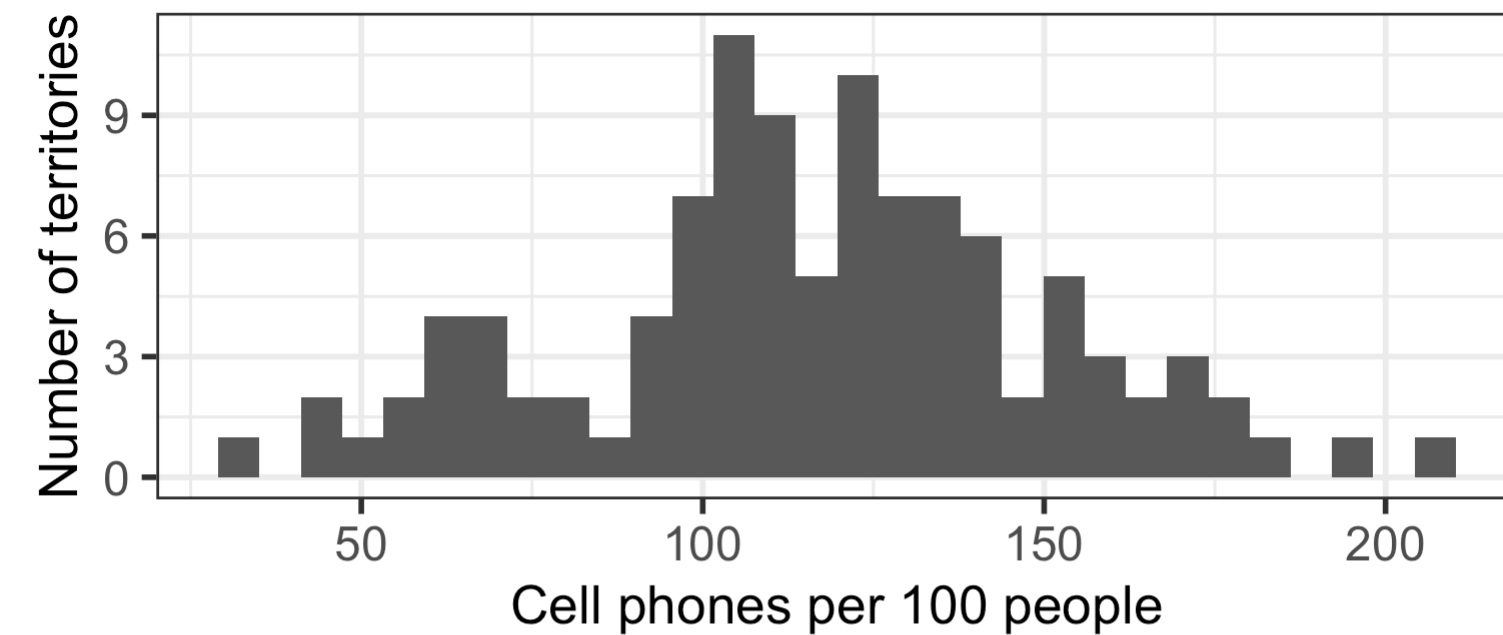
```
1 gapm %>%  
2   select(cell_phones_100,  
3         life_exp) %>%  
4   summary()
```

cell_phones_100	life_exp
Min. : 32.06	Min. : 50.69
1st Qu.: 99.13	1st Qu.: 66.11
Median : 116.68	Median : 71.52
Mean : 116.52	Mean : 70.97
3rd Qu.: 137.18	3rd Qu.: 75.67
Max. : 207.28	Max. : 85.23

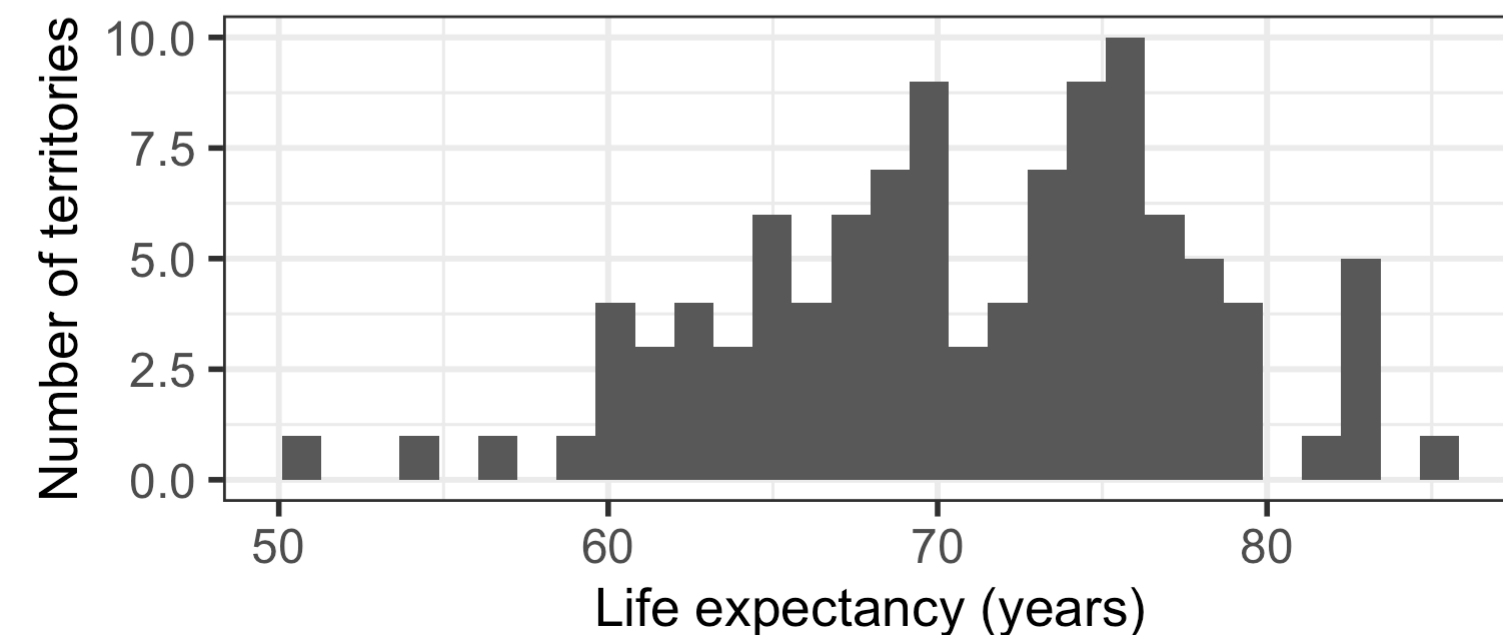
- Plot the individual variables

## ► Code

Distribution of cell phones per 100 people



Distribution of life expectancy



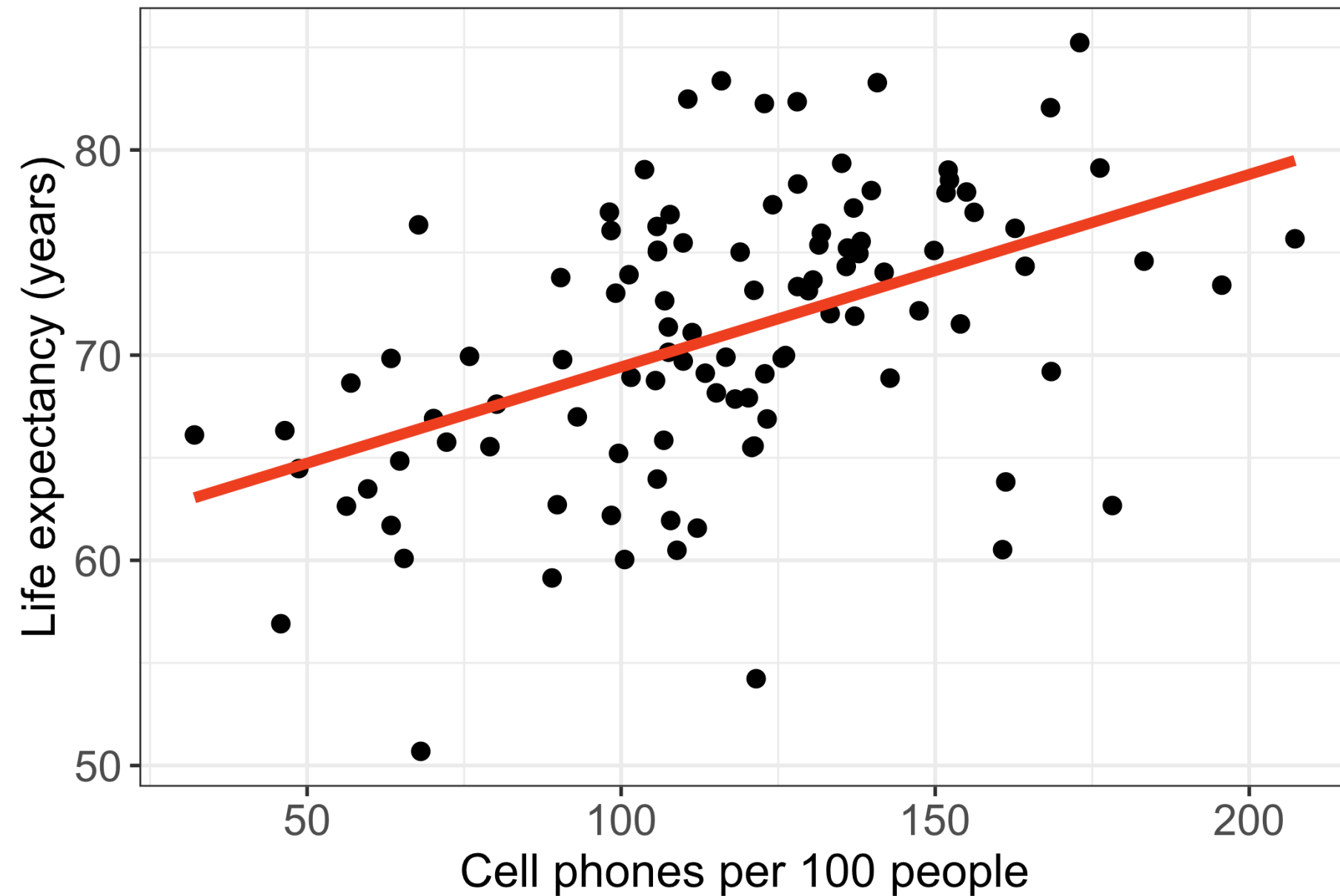
# Poll Everywhere Question 2

# 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 cell phones



- 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

$$\widehat{\text{life expectancy}} = 60.04 + 0.094 \cdot \text{cell phones}$$

# Association vs. prediction

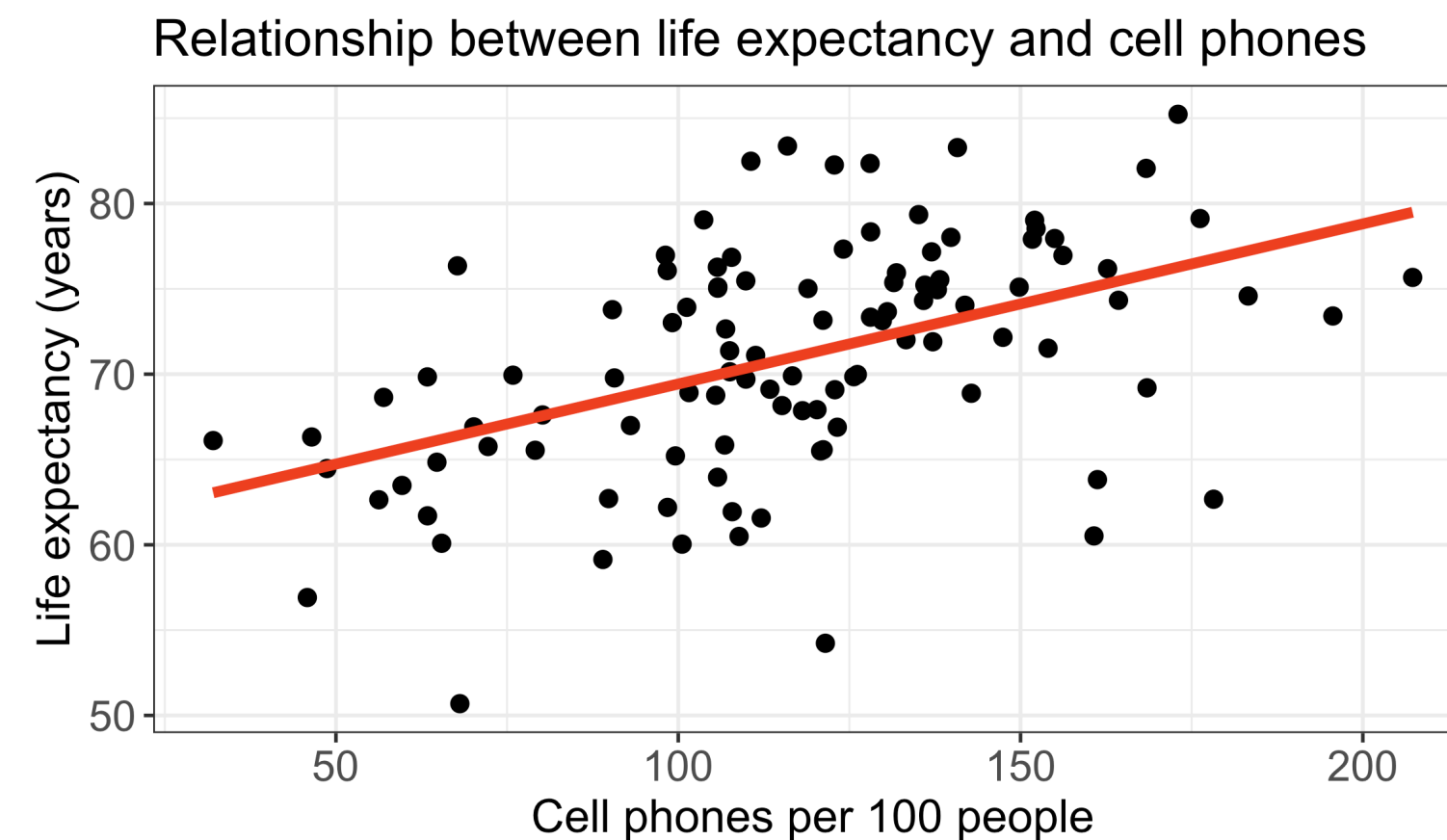
## Association

- What is the association between countries' life expectancy and cell phones?
- Use the slope of the line or correlation coefficient

## Prediction

- What is the expected life expectancy for a country with a specified number of cell phones per 100 people?

$$\widehat{\text{life expectancy}} = 60.04 + 0.094 \cdot \text{cell phones}$$



# 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

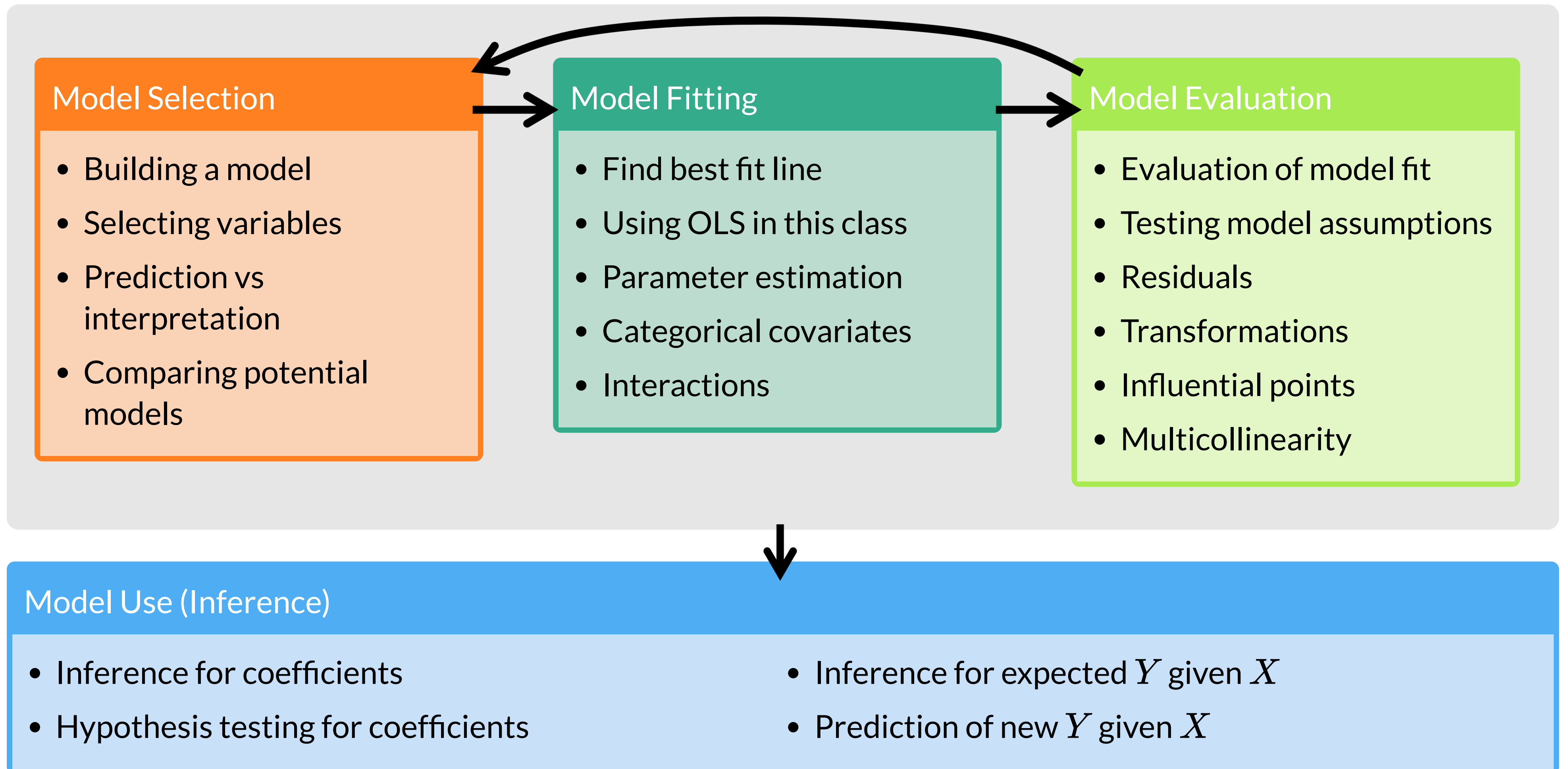
## Quasi-experiment

- 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 casual inference methods...

# Let's revisit the regression analysis process



# Poll Everywhere Question 3

# 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

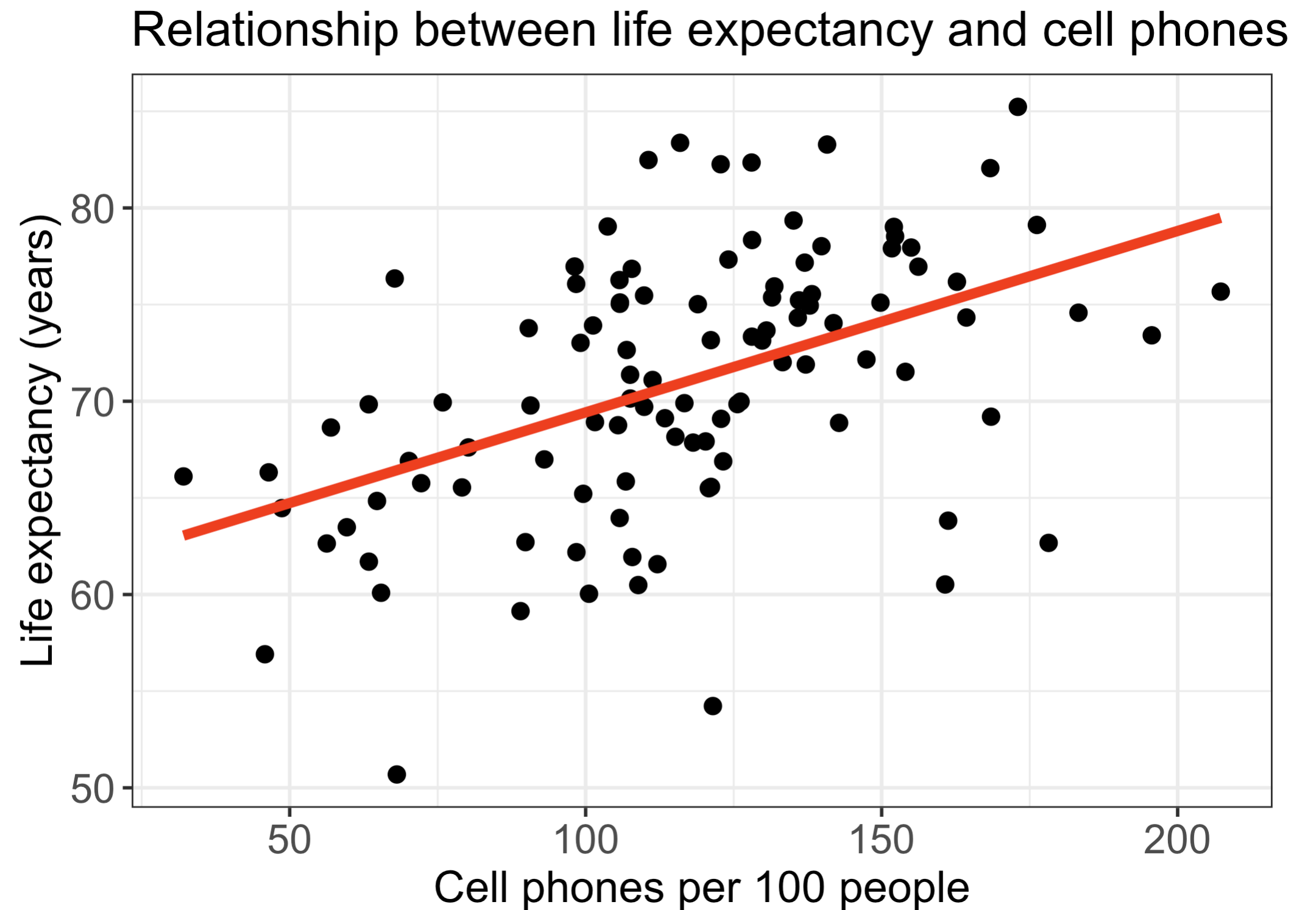
---

$\epsilon$  residuals, error term

# 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  $\beta_1$



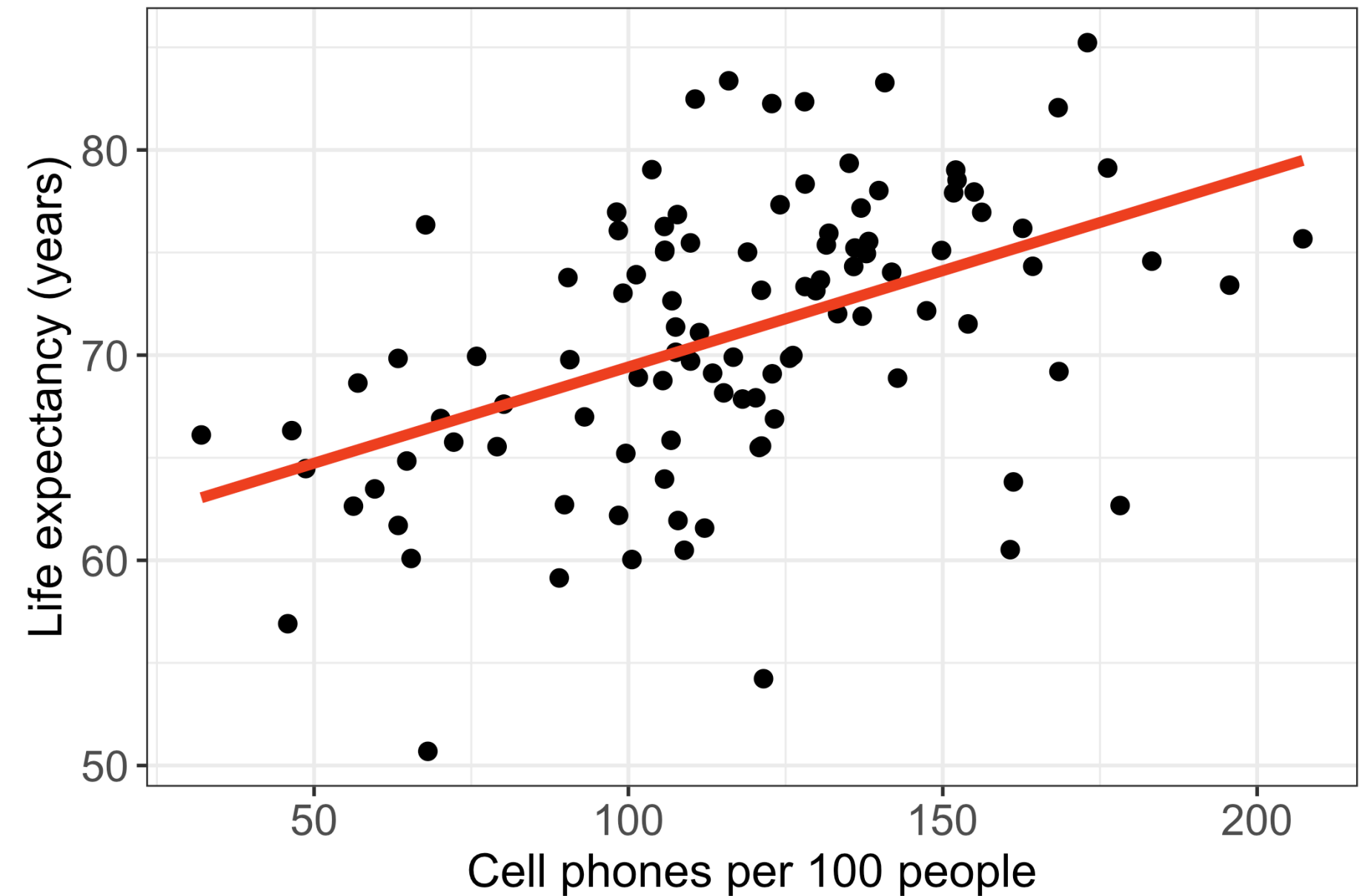
# Poll Everywhere Question 4

# 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  $X$
- $\hat{\beta}_0$  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 cell phones



# Simple Linear Regression Model

## Population regression *model*

Think of this as proposed model **before** we fit the data

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

## Components

$Y$	response, outcome, dependent variable
-----	---------------------------------------

$\beta_0$	intercept
-----------	-----------

$\beta_1$	slope
-----------	-------

$X$	predictor, covariate, independent variable
-----	--

$\epsilon$	residuals, error term
------------	-----------------------

## Estimated regression *line*

Think of this as the actualized model **after** we fit data

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

## Components

$\hat{Y}$	<i>estimated expected</i> response given predictor $X$
-----------	--

$\hat{\beta}_0$	<i>estimated</i> intercept
-----------------	----------------------------

$\hat{\beta}_1$	<i>estimated</i> slope
-----------------	------------------------

$X$	predictor, covariate, independent variable
-----	--

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

First let's take a break!!

# 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:

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

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

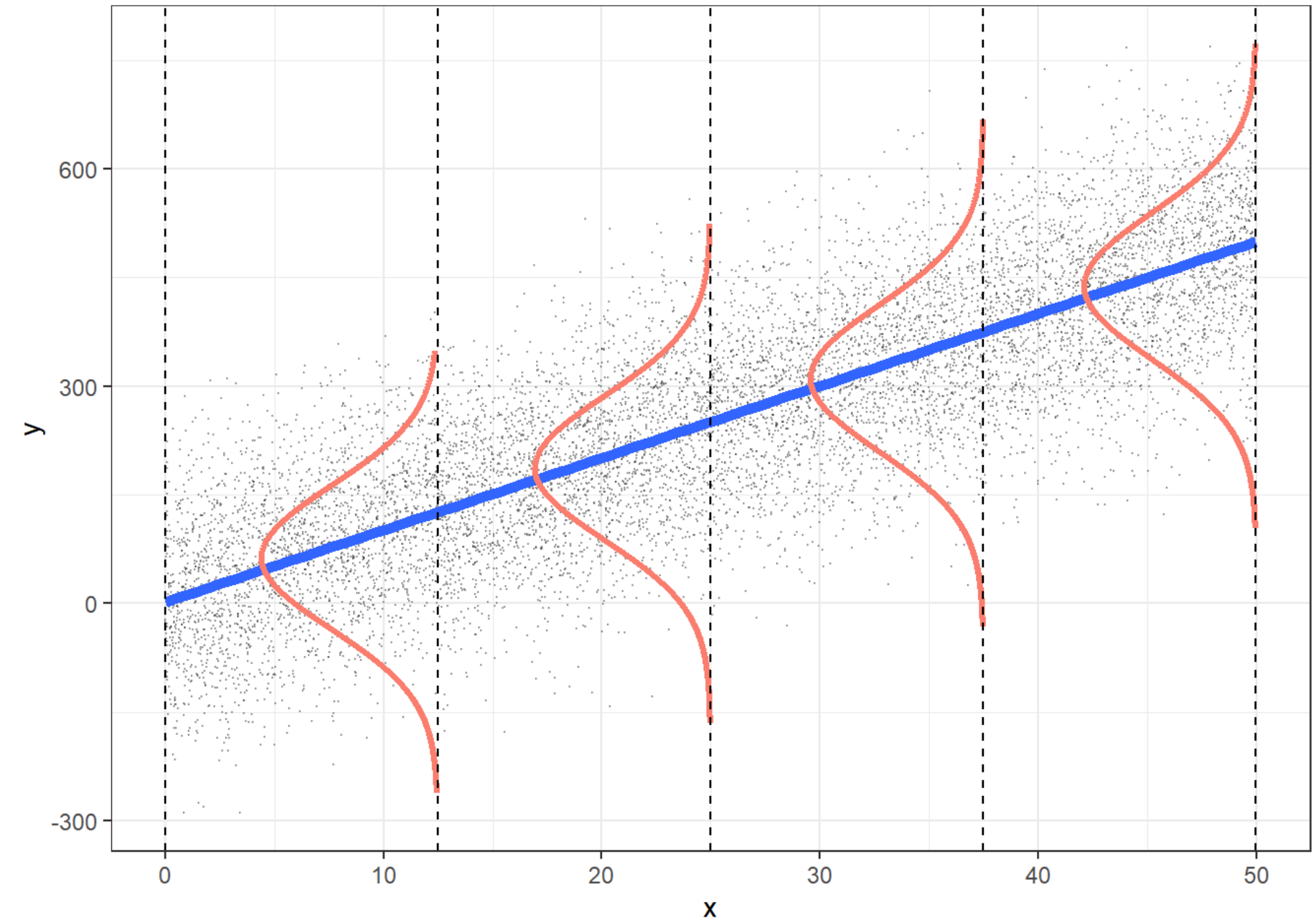
$$E[Y] = E[\beta_0 + \beta_1 X + \epsilon]$$

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

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

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

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



# So now we have two representations of our population model

With observed  $Y$  values and residuals:

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

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:

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

$$Y = E[Y|X] + \epsilon$$

$$Y - E[Y|X] = \epsilon$$

$$\epsilon = Y - E[Y|X]$$

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:

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

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

$$\hat{E}[Y|X] = \hat{\beta}_0 + \hat{\beta}_1 X$$

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

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

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

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

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

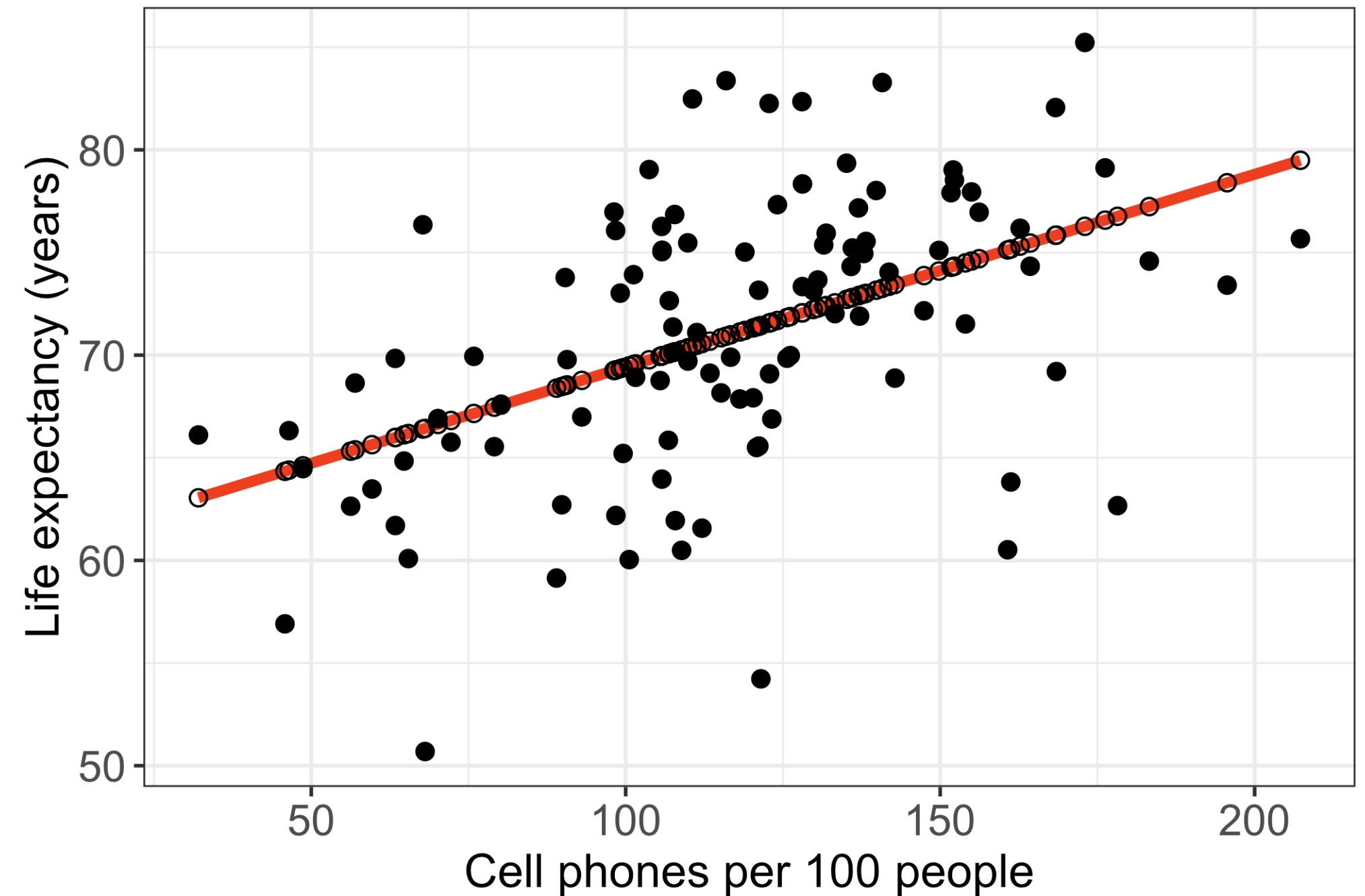
$$\hat{\epsilon} = Y - \hat{Y}$$

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

# Residuals for observation $i$ in the estimated/fitted model

- **Observed values for each country  $i$ :  $Y_i$** 
  - Value in the dataset for country  $i$
- **Fitted value for each country  $i$ :  $\hat{Y}_i$** 
  - 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  $\hat{Y}_i$

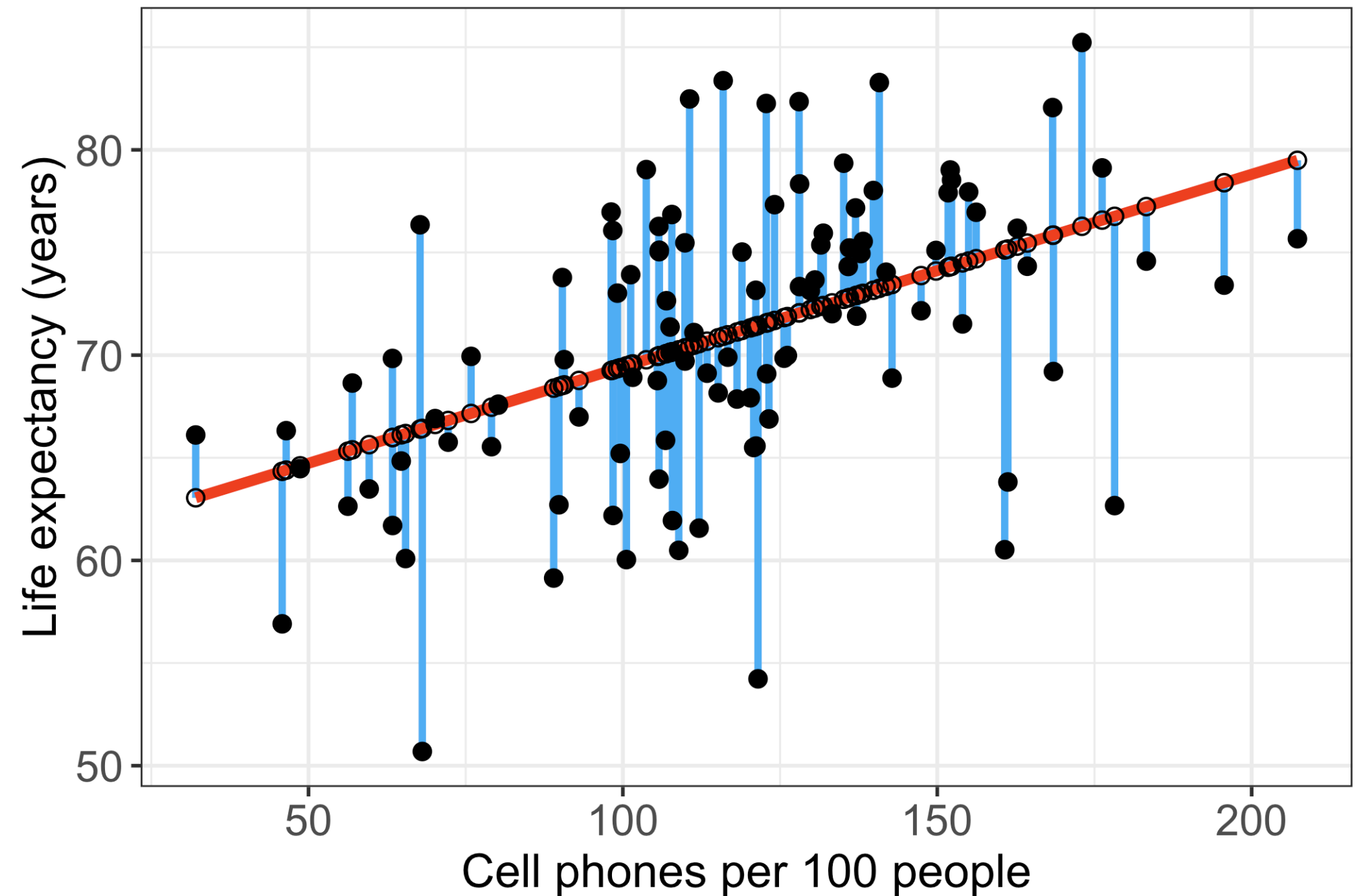
Relationship between life expectancy and cell phones



# Residuals for observation $i$ 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 cell phones



# Poll Everywhere Question 5

# 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**
  - 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)

$$SSE = \sum_{i=1}^n \hat{\epsilon}_i^2$$

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

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

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

## 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
  - 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 $\hat{\beta}_0$

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

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

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

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

#### Things to use

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

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

$$\begin{aligned}\frac{\partial SSE}{\partial \hat{\beta}_0} &= 0 \\ -2 \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) &= 0 \\ \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) &= 0 \\ \sum_{i=1}^n Y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n X_i &= 0 \\ \frac{1}{n} \sum_{i=1}^n Y_i - \hat{\beta}_0 - \hat{\beta}_1 \frac{1}{n} \sum_{i=1}^n X_i &= 0 \\ \bar{Y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{X} &= 0 \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X}\end{aligned}$$

### 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 $\hat{\beta}_1 = 0$

$$\begin{aligned}\frac{\partial SSE}{\partial \hat{\beta}_1} &= 0 \\ \sum_{i=1}^n (X_i Y_i - \hat{\beta}_0 X_i - \hat{\beta}_1 X_i^2) &= 0 \\ \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \hat{\beta}_0 - \sum_{i=1}^n X_i^2 \hat{\beta}_1 &= 0 \\ \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i (\bar{Y} - \hat{\beta}_1 \bar{X}) - \sum_{i=1}^n X_i^2 \hat{\beta}_1 &= 0 \\ \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \bar{Y} + \sum_{i=1}^n \hat{\beta}_1 X_i \bar{X} - \sum_{i=1}^n X_i^2 \hat{\beta}_1 &= 0 \\ \sum_{i=1}^n X_i (Y_i - \bar{Y}) + \sum_{i=1}^n (\hat{\beta}_1 X_i \bar{X} - X_i^2 \hat{\beta}_1) &= 0 \\ \sum_{i=1}^n X_i (Y_i - \bar{Y}) + \hat{\beta}_1 \sum_{i=1}^n X_i (\bar{X} - X_i) &= 0\end{aligned}$$

### 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

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



# Regression in R: `lm()`

- Let's discuss the syntax of this function

```
1 model1 <- gapm %>% lm(formula = life_exp ~ cell_phones_100)
```

In the general form:

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

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

```
1 model1 <- gapm %>% lm(formula = life_exp ~ cell_phones_100)
2 summary(model1)
```

Call:

```
lm(formula = life_exp ~ cell_phones_100, data = .)
```

Residuals:

Min	1Q	Median	3Q	Max
-17.211	-3.268	0.615	3.818	12.449

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	60.04051	2.05567	29.207	< 2e-16 ***
cell_phones_100	0.09384	0.01692	5.546	2.27e-07 ***

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

Residual standard error: 5.964 on 103 degrees of freedom

Multiple R-squared: 0.23, Adjusted R-squared: 0.2225

F-statistic: 30.76 on 1 and 103 DF, p-value: 2.271e-07

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

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

term	estimate	std.error	statistic	p.value
(Intercept)	60.04051297	2.05566959	29.207278	1.215444e-51
cell_phones_100	0.09383818	0.01691978	5.546063	2.271176e-07

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

$$\widehat{\text{life expectancy}} = 60.04 + 0.094 \cdot \text{cell phones}$$

# How do we interpret the coefficients?

$$\widehat{\text{life expectancy}} = 60.04 + 0.094 \cdot \text{cell phones}$$

- **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 60.4 years for a country with 0 cell phones per 100 people.

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

- For every increase of 1 unit in the  $X$ -variable (if continuous), there is an expected increase of  $\hat{\beta}_1$  units in the  $Y$ -variable.
- We only say that there is an expected increase and not necessarily a causal increase.
- **Example:** For every 1 additional cell phone per 100 people, life expectancy increases, on average, 0.09 years.
  - Can also say “...average life expectancy increases 0.09...” or “... expected life expectancy increases 0.09...”

# Next time

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