

# Simple Linear Regression - Model

# Research Objective

**Research Question:** Is the adult height of a child determined by the height of the mother? In other words, what is the relationship between student's height and mother's height for all BYU students.

**Population:** All BYU students.

**Parameter of Interest:**

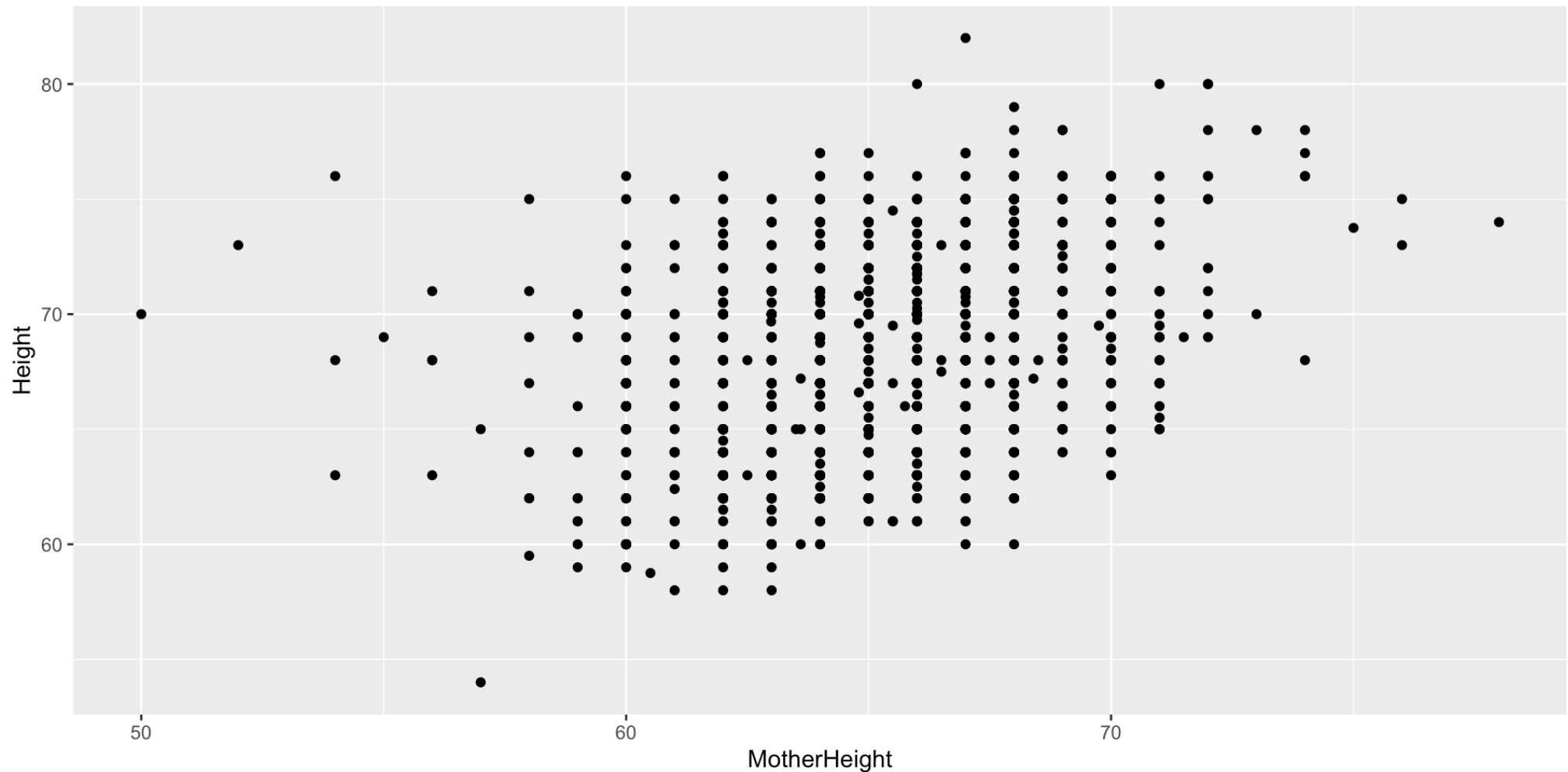
- Some number measurement of the “relationship” between student's height and mother's height.
- For this subunit we are going to focus on what a “relationship” means.

**Sample:** A convenience sample of 1575 BYU students who are in Stat 121.

Are there any issues with this study setup?

# Simple Linear Regression Model

- Main goal: Specify the population relationship between student's height and mother's height.



# Review: Equation of a Line

Equation you are probably used to:

$$y = mx + b$$

where:

- $m$  = slope
- $b$  = intercept

# Review: Equation of a Line

We are going to change notation to:

$$y = \beta_0 + \beta_1 x$$

where  $\beta$  is pronounced “beta”,

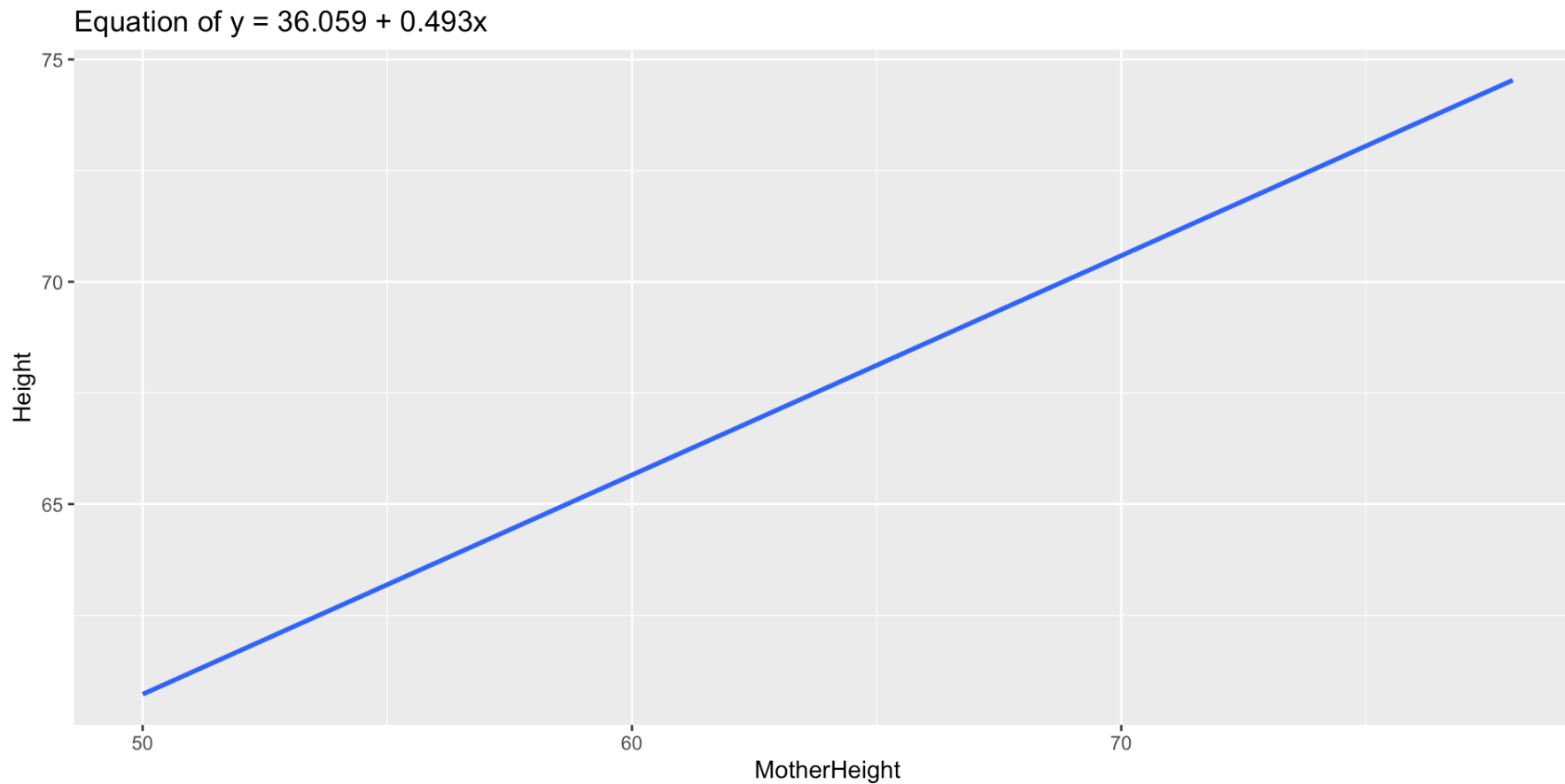
- $\beta_0$  = intercept
- $\beta_1$  = slope

Why? Remember that greek letters in here will always represent population parameters so this notation is more consistent with that standard (and this is the notation that everyone uses for regression).

# Review: Equation of a line

Equation of a line:

$$y = \beta_0 + \beta_1 x$$



# Review: Equation of a line

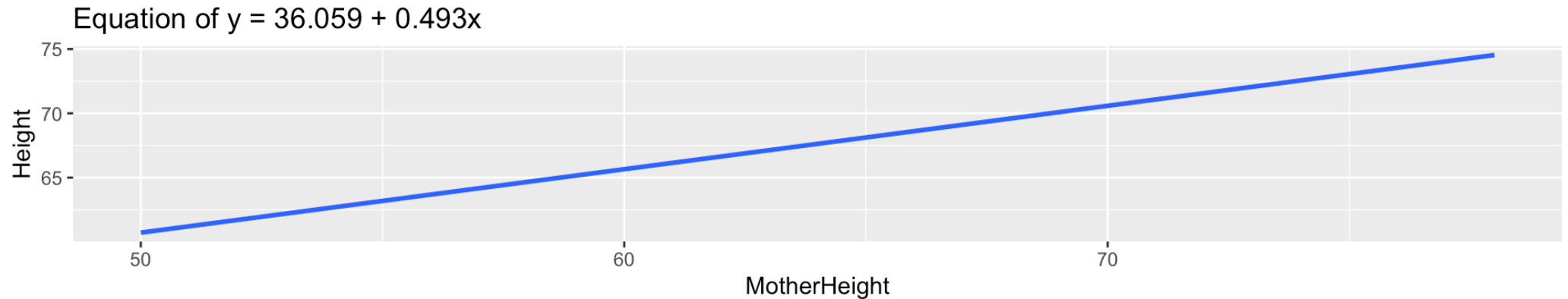
Equation of a line:

$$y = \beta_0 + \beta_1 x$$

Interpretations:

- Slope ( $\beta_1 =$  “rise over run”): As  $x$  increases/decreases by 1,  $y$  increases/decreases by  $\beta_1$ . If  $x$  “runs” by 1, then  $y$  “rises” by  $\beta_1$ .
- Intercept ( $\beta_0$ ): If  $x$  is 0, then  $y$  is  $\beta_0$ .

# Review: Equation of a line

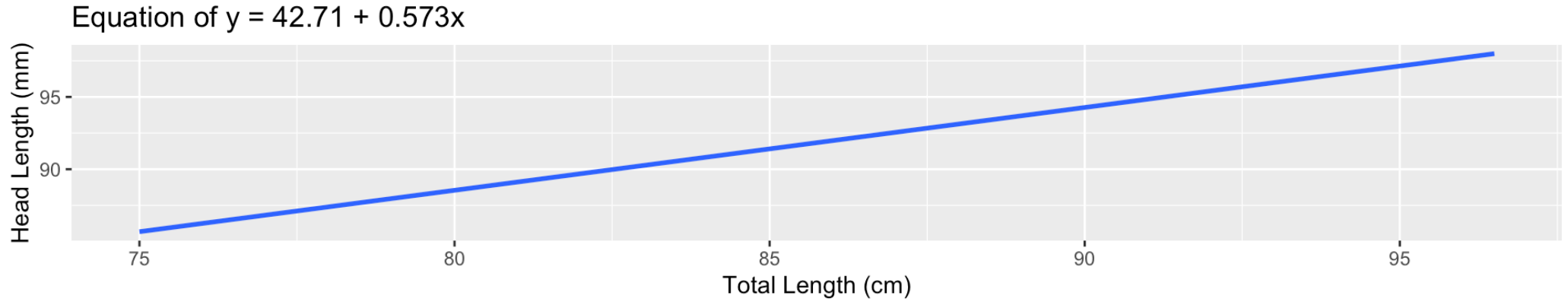


## Practice: Height vs. Mother's Height

- How would you interpret the intercept?
  - If the mother is zero inches tall ( $x = 0$ ) then the student height is 36.059.
- How would you interpret the slope?
  - If the mother's height increases by 1, then the student height goes up by 0.493.
- What is  $y$  when  $x = 64$ ?
  - Plug in  $y = 36.059 + 0.493 \times 64 = 67.611$



# Review: Equation of a line

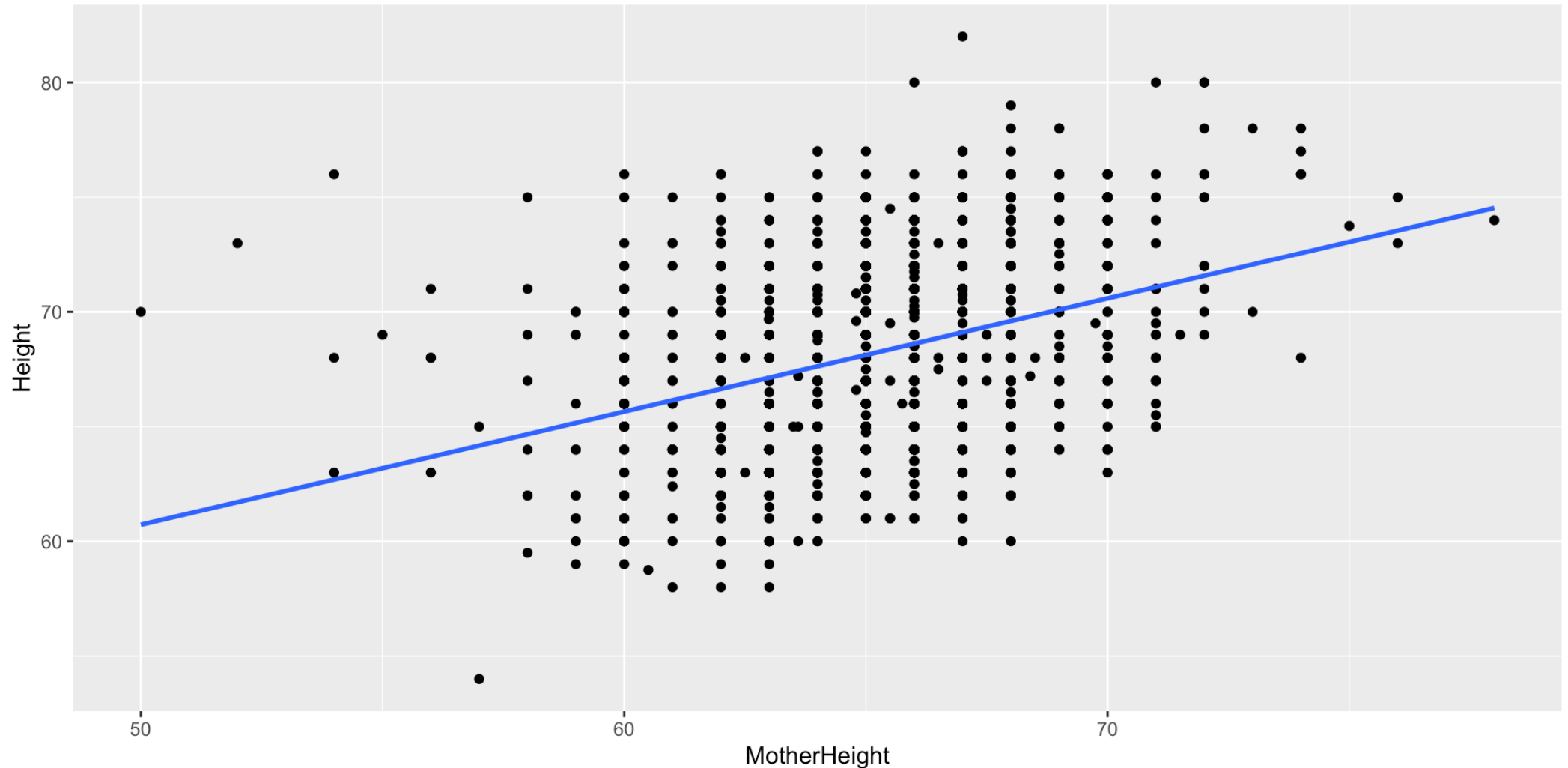


## Practice: Possum lengths

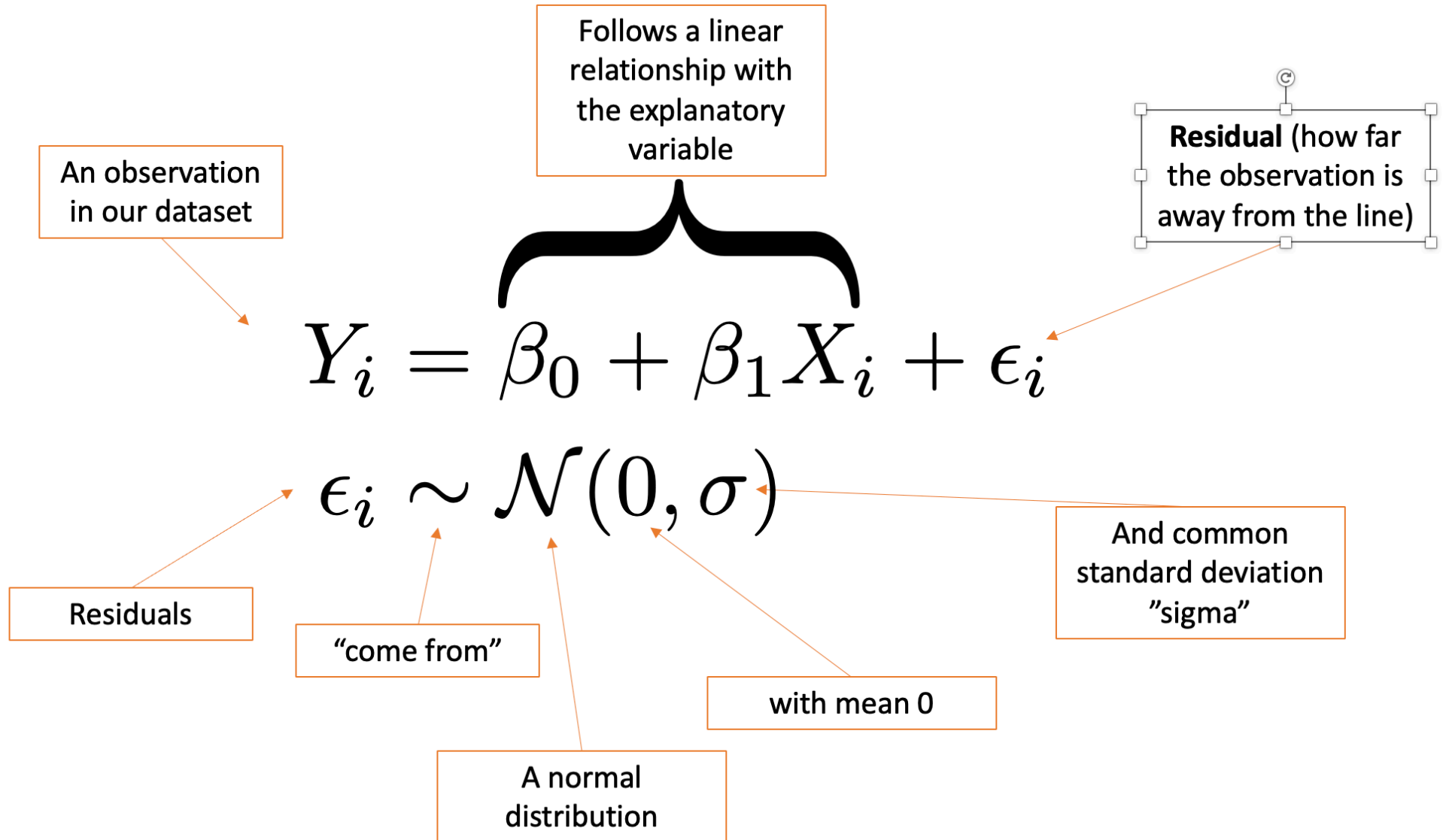
- How would you interpret the intercept?
  - If total length is zero ( $x = 0$ ) then the head length is 42.71.
- How would you interpret the slope?
  - If the total length goes up by 1, then the head length goes up by 0.573.
- What is head length when total length = 95?
  - Plug in  $y = 42.71 + 0.573 \times 95 = 97.145$

# Simple Linear Regression Model

Issue: When specifying a model for the relationship, the data do not perfectly follow a line:

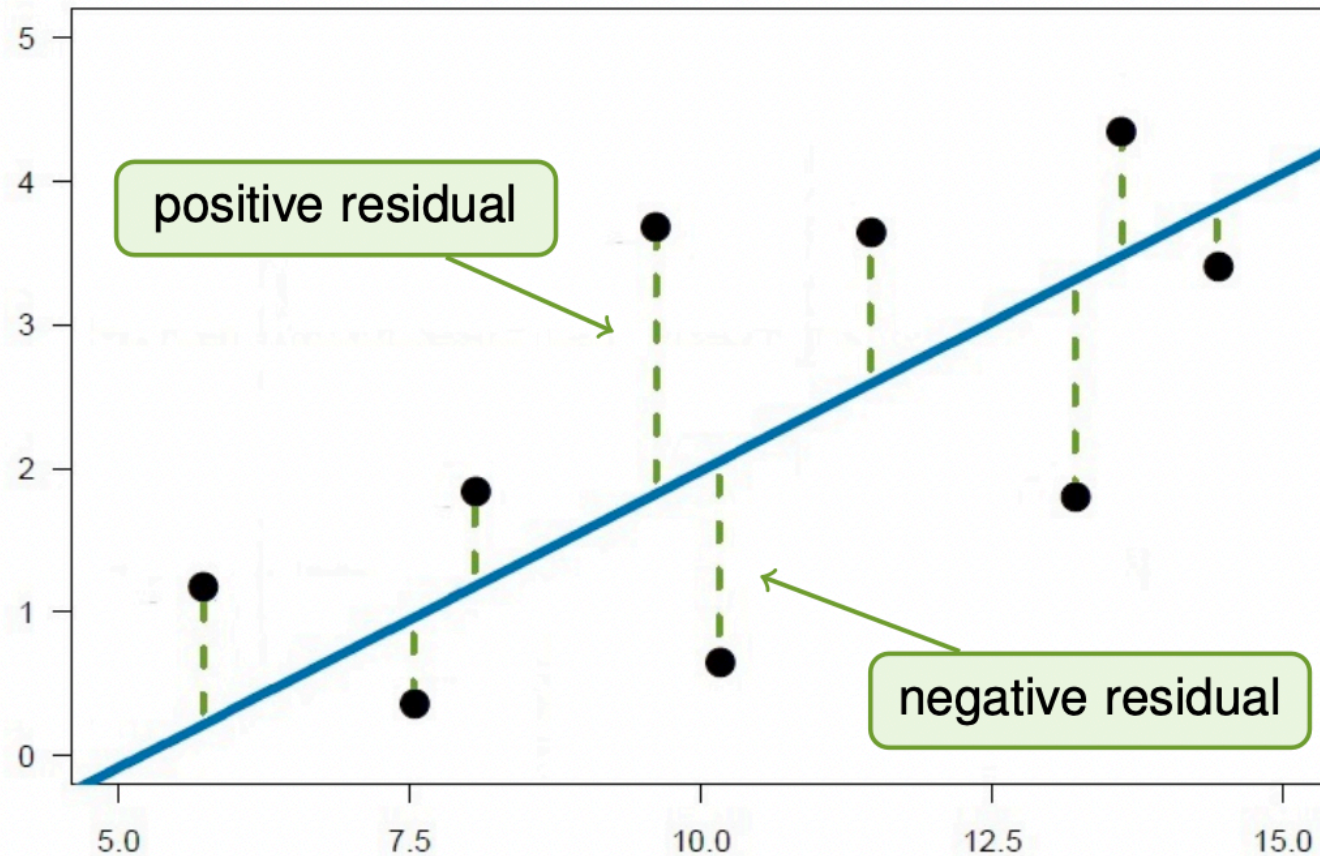


# Simple Linear Regression Model

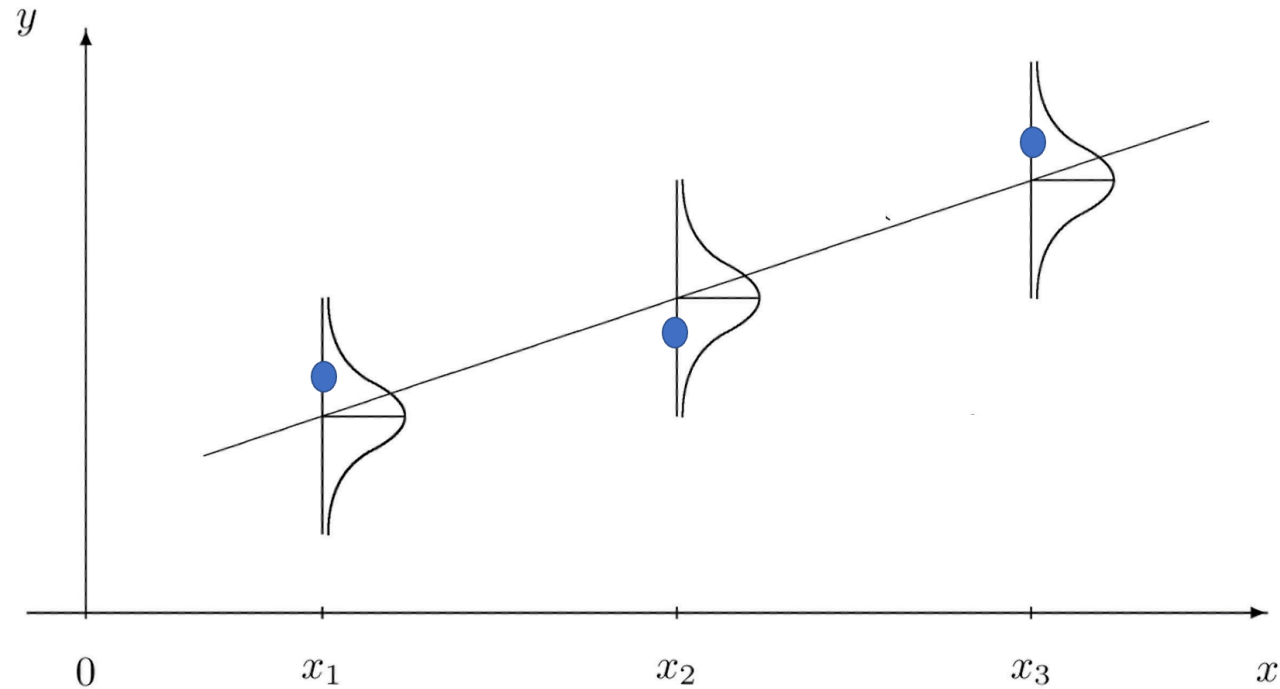


# Residuals

$$\begin{aligned}\text{Residual} = \epsilon_i &= \text{Observation} - \text{Predicted Value} \\ &= Y_i - (\beta_0 + \beta_1 X_i)\end{aligned}$$

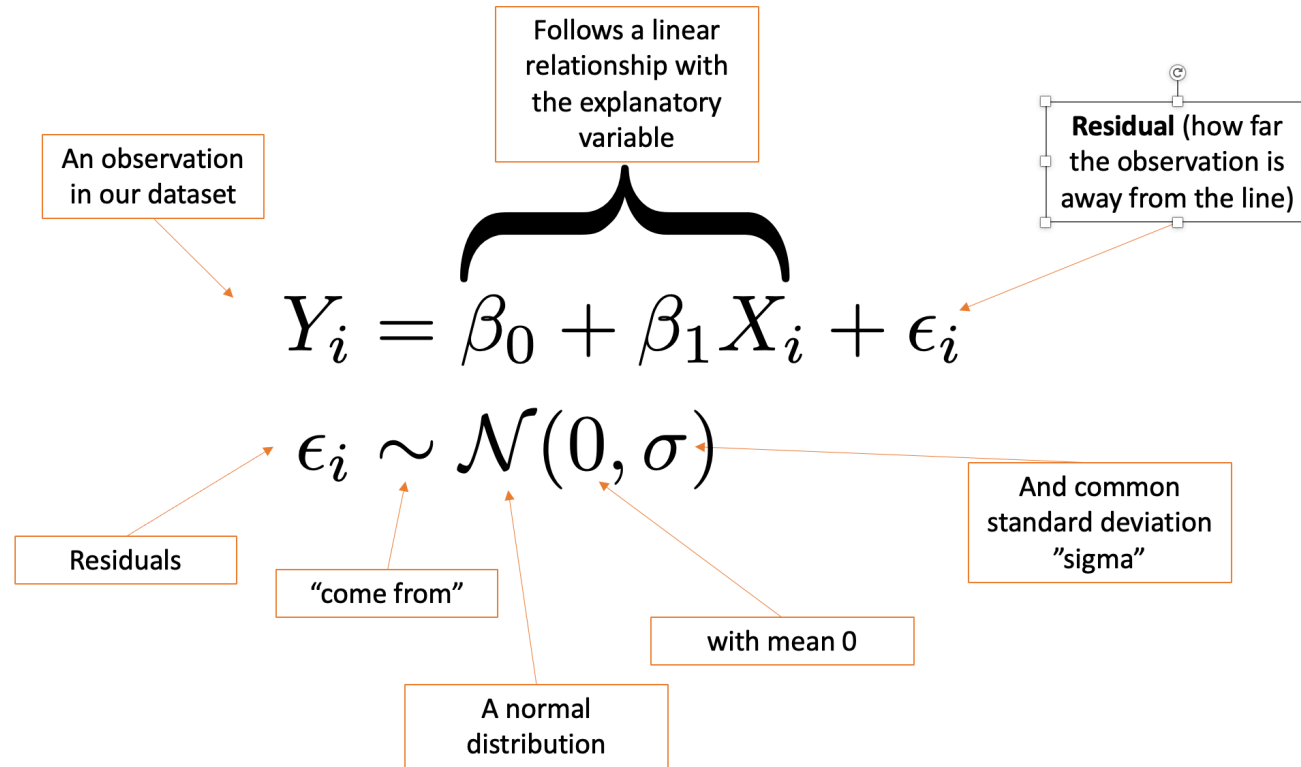


# Visualizing the SLR Model



- $\sigma$  is the standard deviation and controls the spread of the dots about the regression line. The bigger the  $\sigma$ , the farther the dots from the line.

# Interpreting the SLR Model



Slight change in interpretation:

- Intercept ( $\beta_0$ ): If  $X = 0$ , we expect  $Y$  to be  $\beta_0$ .
- Slope ( $\beta_1$ ): If  $X$  goes up by 1, we expect  $Y$  to go up by  $\beta_1$ .

# Assumptions of the SLR Model

Easy way to remember what we are assuming about the population in a simple linear regression model:

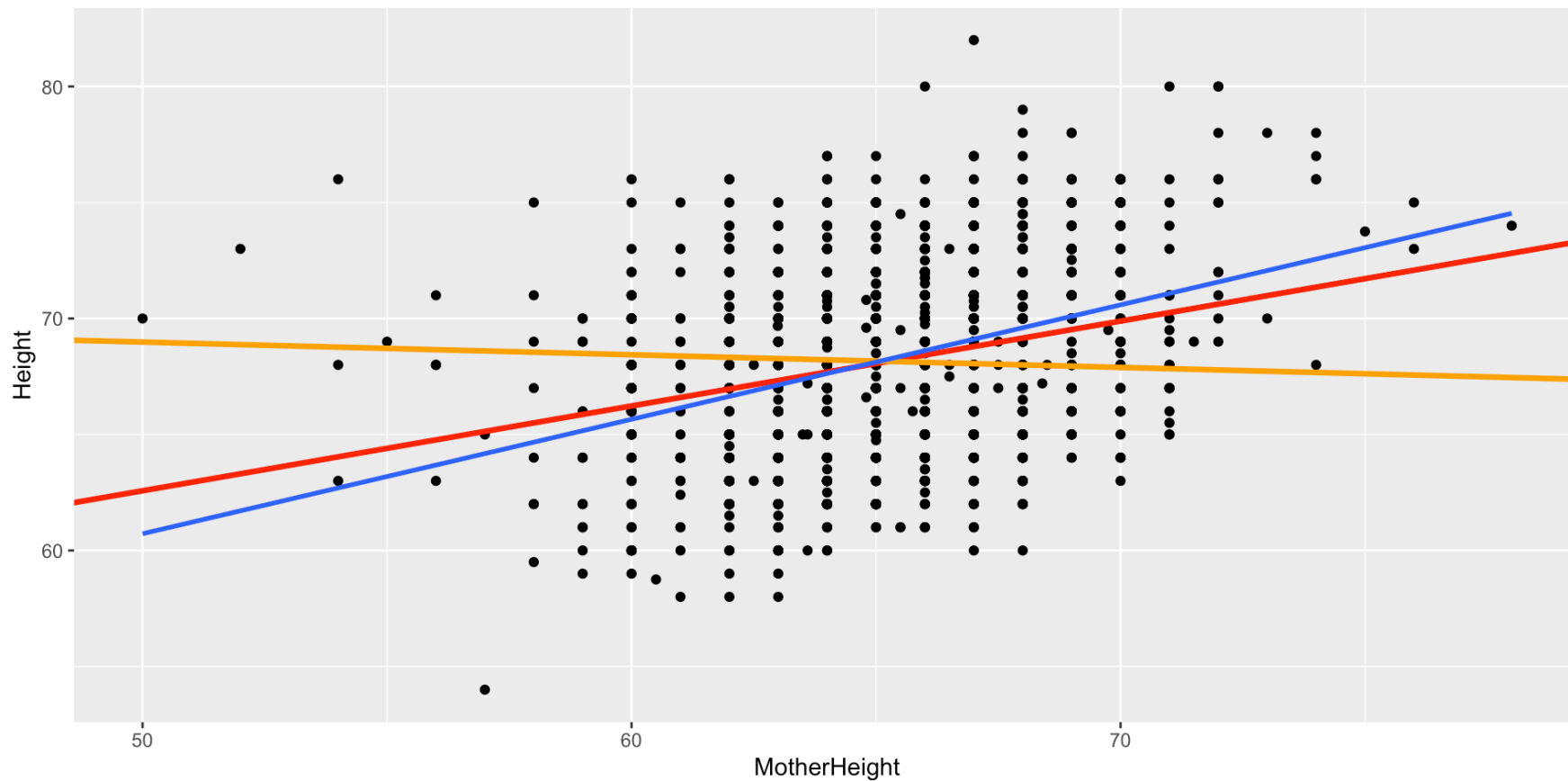
- L - Linear relationship between  $x$  and  $y$
- I - Independence (one obs. doesn't impact the other)
- N - Normal residuals (distance from line is normal)
- E - Equal spread of residuals around the line

More on why these assumptions are important and how to check these in the next subunit.

# Parameter Estimation

Parameters we want to estimate:  $\beta_0$  &  $\beta_1$  (which defines the line) and  $\sigma$  (so we know how spread out things are)

Goal: Find the line that goes “closest” to the data points.





# Parameter Estimation

What do we mean by “line closest to points”? We want to find  $\hat{\beta}_0$  and  $\hat{\beta}_1$  so that:

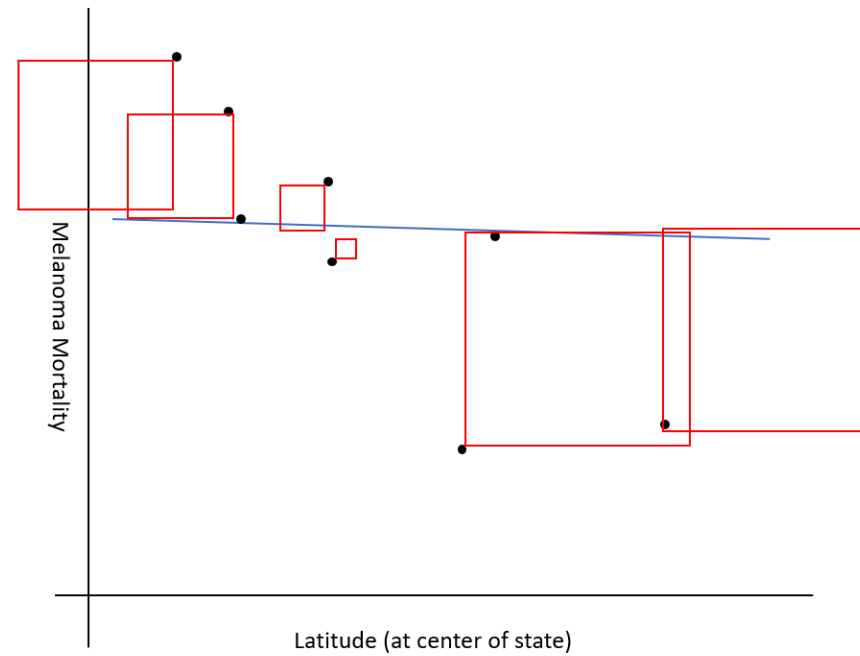
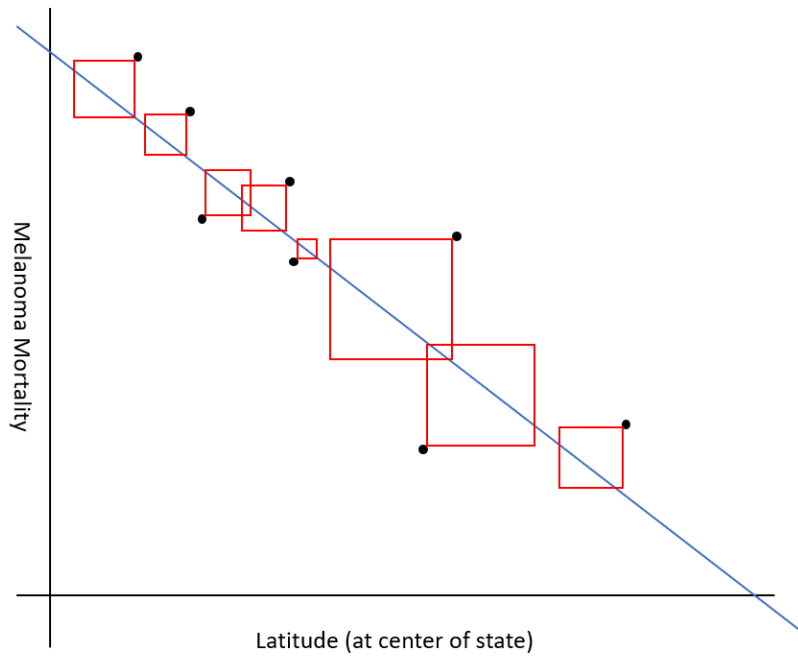
$$\sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2$$

is as small as possible. This is called the **least squares regression line**.

A few notes:

1. We “square” distances so that, for example, a 5 “above” and 5 “below” the line are the same “distance”.
2. We sum squared residuals because we look at all the data.
3. We use “hats” to denote estimates from sample (for example,  $\hat{\beta}_1$  is our estimate of  $\beta_1$ )

# Parameter Estimation



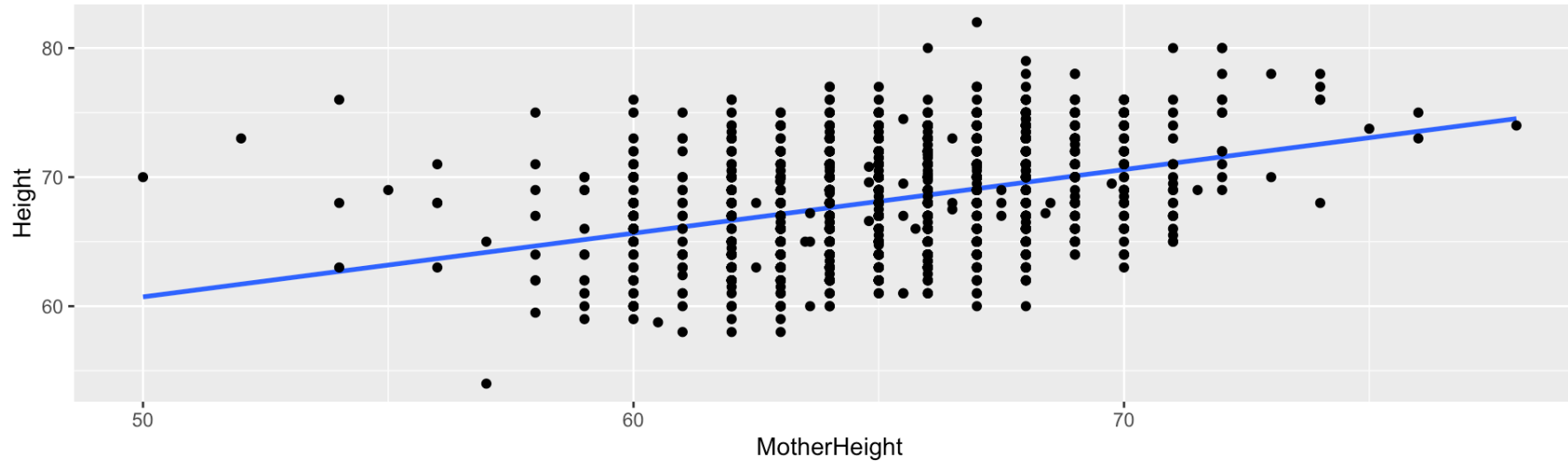
# Parameter Estimation

How do we find  $\hat{\beta}_0$  and  $\hat{\beta}_1$  that minimizes

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

1. Guess and check
  2. Use calculus
- In either case, we'll let the computer do the hard work for us

# The Fitted SLR Model



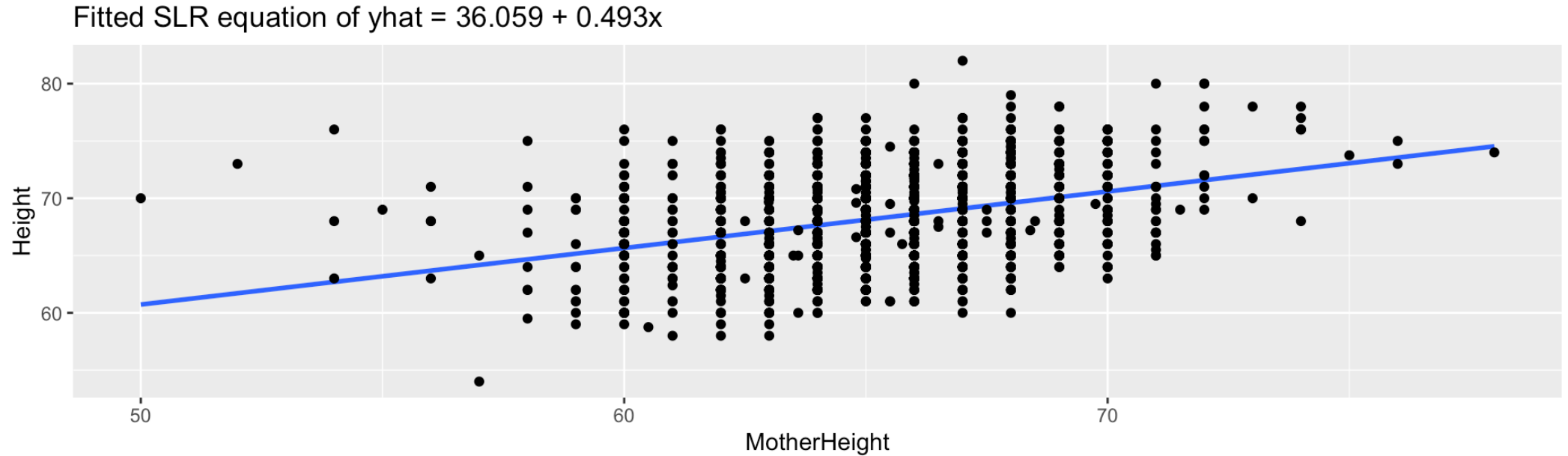
Fitted Regression Line Equation:

$$\hat{y} = 36.059 + 0.493 \times x$$

where:

- $\hat{y}$  is the fitted height value (the height value on the line)
- $\hat{y} \neq y_i$  because  $y_i$  is an observed height

# The Fitted SLR Model



A interesting point: The sign (positive/negative) of the correlation will always match the sign of the slope (positive/negative). Not the same number but the same sign.

# Parameter Estimation

An estimate of  $\sigma$  is more complicated to explain (take more stats courses), so for purposes of this class, the computer estimates it for us.

- $\hat{\sigma} = 3.896$

How do we interpret  $\hat{\sigma}$ ?

- On average, the actual student's heights are about 3.896 inches away from the estimated heights.

# Using the Analysis Tool

Stat 121 Analysis Tool

Exploratory Data Analysis

Normal Probability Calculator

Central Limit Theorem

Analysis for Means <

Analysis For Proportions <

Regression <

>> Simple Linear Regression

>> Multi Linear Regression

Use this section for Unit 6

## Simple Linear Regression

### 1) Dataset Selection

**Data Selection**

Use Preexisting Dataset

Upload Your Own Dataset

**Select Dataset**

Melanoma

Description: Melanoma mortality rates (per 10 million people) for each state in the continental US.

Sample size: 49

Display Dataset

Select This Dataset

Choose the dataset

# Using the Analysis Tool

## 2) Select Variables

Please select the explanatory variable. The explanatory variable should "explain" what happens to the response variable.

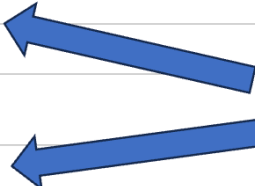
**Select Response Variable:**

Mort

**Select Explanatory Variable:**

Lat

Proceed to EDA



Make sure you get these right or everything below will be messed up



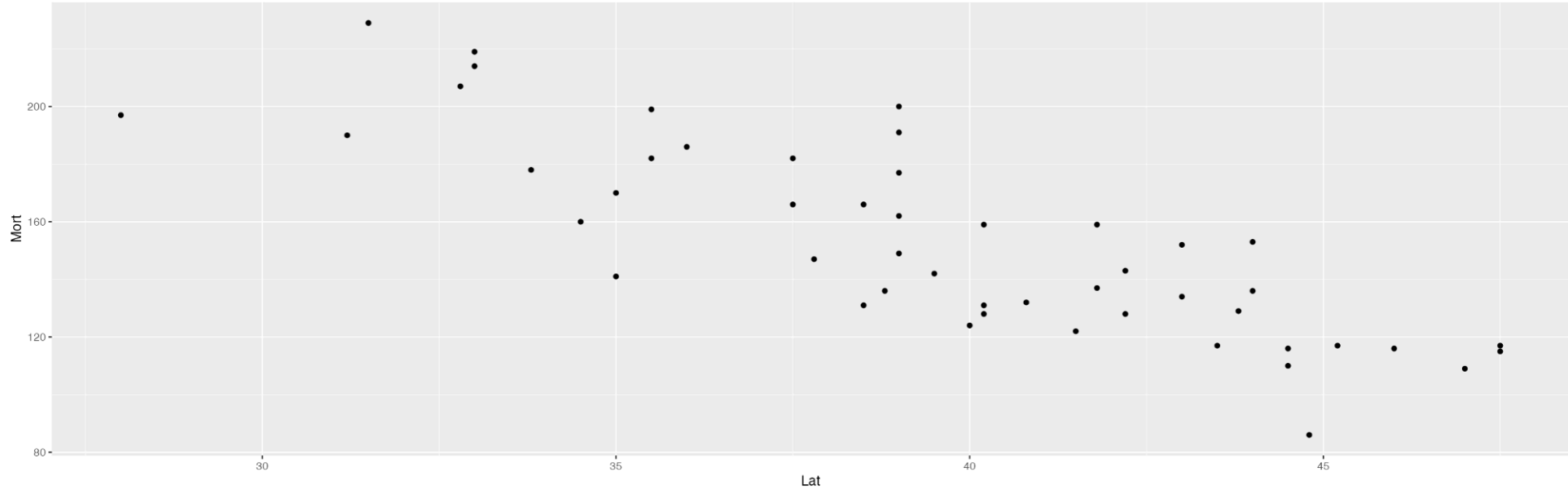
# Using the Analysis Tool

### 3) Exploratory Data Analysis

Choose the plot you want to draw  
(scatterplot is most useful)

Which plot would you like to draw?

Scatterplot



Which numerical summary do you want to calculate?

Correlation between Explanatory and Response Variable

Choose value you want to calculate  
(correlation and covariance are most useful)

Correlation ( $r$ ) =  $-0.8245$

Proceed to Checking Assumptions

# Using the Analysis Tool

4) Check Regression Assumptions

What regression assumption plot do you want to look at?

---

For now, you can ignore Part 4 and just proceed to regression analysis (we will come back to this section next unit when we talk about inference)

Proceed to Regression Analysis (Statistical Inference)

# Using the Analysis Tool

5) Regression Analysis

Confidence Level for Slope and Intercept:

0.5 0.95 0.99

Regression Analysis of Mort (Y) explained by Lat (X)  
Coefficient Table:

Show 5 entries

	Estimate	t value	p-value	CI Lower Bound	CI Upper Bound
(Intercept)	389.1894	16.344	0	341.2852	437.0936
Lat	-5.9776	-9.9898	0	-7.1814	-4.7739

Showing 1 to 2 of 2 entries Previous 1 Next

R-squared: 0.6798  
sigma: 19.115

Show Fitted Regression Line

Proceed to Predictions

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

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

Estimate of spread ( $\sigma$ )

Ignore this for now

# Assessing Model Fit

Coming back to the student height example, we had  $\hat{\sigma} = 3.896$  which we interpret to be the difference between the actual heights and the predicted heights. Does  $\hat{\sigma} = 3.896$  mean that the observations are “close” to the line or not?

- It’s hard to tell just from  $\hat{\sigma}$  if this is “good” or “bad” because it depends on the problem. A better measure would be a standardized measure that can be used for all regression problems.

# Assessing Model Fit

Mathematical formula:

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2}{\sum_{i=1}^n (Y_i - \bar{y})^2} = 0.12142$$

Intuition:

- Formal interpretation: The percent of variability in  $Y$  that is explained by  $X$ .
- $R^2$  is between 0 and 1 with 1 meaning the data perfectly follow a line and 0 meaning the data don't follow the line at all.
- Intuition: you can think of  $R^2$  as a “grade” for your regression line where  $R^2 = 1$  is a perfect line and  $R^2 = 0$  is a terrible line.

# Using the Analysis Tool

5) Regression Analysis

Confidence Level for Slope and Intercept:

0.5 0.95 0.99

0.5 0.55 0.6 0.65 0.7 0.75 0.8 0.85 0.9 0.95 0.99

Regression Analysis of Mort (Y) explained by Lat (X)  
Coefficient Table:

Show  entries

	Estimate	t value	p-value	CI Lower Bound	CI Upper Bound
(Intercept)	389.1894	16.344	0	341.2852	437.0936
Lat	-5.9776	-9.9898	0	-7.1814	-4.7739

Showing 1 to 2 of 2 entries Previous  Next

R-squared: 0.6798  
sigma: 19.115

Show Fitted Regression Line

Proceed to Predictions

# Additional SLR Practice

Does a higher GPA lead to better pay? Use the salary data and a simple linear regression model to answer the following questions:

1. What is the estimated pay for someone who completely fails college (0.0 GPA)?
2. For two people who differ by 1.0 GPA, how much higher (or lower) should the pay be for person with the higher GPA on average?
3. On average, how far away are pay amounts from estimated pay amounts?
4. How well does the GPA explain pay?



# Additional SLR Practice Answers

Does a higher GPA lead to better pay? Use a simple linear regression model (and the course app) to answer the following questions (Salary dataset):

1. What is the estimated pay for someone who completely fails college (0.0 GPA)?

- $\hat{\beta}_0 = 51135.68$

2. For two people who differ by 1.0 GPA, how much higher (or lower) should the pay be for person with the higher GPA on average?

- $\hat{\beta}_1 = 6510.04$

3. On average, how far away are pay amounts from estimated pay amounts?

- $\hat{\sigma} = 10353.03$

4. How well does the GPA explain pay?

- $R^2 = 0.1147$

# Key Terminology

- Least squares
- Simple linear regression model
- Slope
- Intercept
- $R^2$
- Relationship between correlation and slope
- Spread about regression line ( $\sigma$ )