Linear Models

Nate Wells

Math 141, 2/16/22

Fitting a Line by Least-Squares Regression

Outline

In this lecture, we will...

Fitting a Line by Least-Squares Regression

Outline

In this lecture, we will...

- Investigate the linear model
- Discuss predictions and residuals
- Explore a formula for finding the line of best fit

Section 1

Introduction to Linear Regression

Fitting a Line by Least-Squares Regression 0000000000

Overview

"All models are wrong, but some are useful."

— George Box, statistician

Overview

"All models are wrong, but some are useful."

— George Box, statistician

• Linear regression is both an accessible and potent tool in statistical analysis.

Overview

"All models are wrong, but some are useful."

— George Box, statistician

Linear regression is both an accessible and potent tool in statistical analysis.

What is the Relationship between Income and Life Expectancy?



Overview

"All models are wrong, but some are useful."

— George Box, statistician

Linear regression is both an accessible and potent tool in statistical analysis.



What is the Relationship between Income and Life Expectancy?

• Quantitative variables, by nature, are amenable to algebraic manipulation.

- Quantitative variables, by nature, are amenable to algebraic manipulation.
- Given two quantitative variables X and Y, we construct a mathematical model that expresses the values of Y as a function of the values of X:

- Quantitative variables, by nature, are amenable to algebraic manipulation.
- Given two quantitative variables X and Y, we construct a mathematical model that expresses the values of Y as a function of the values of X:

Y=f(X)

- Quantitative variables, by nature, are amenable to algebraic manipulation.
- Given two quantitative variables X and Y, we construct a mathematical model that expresses the values of Y as a function of the values of X:

Y=f(X)

• Linear functions are the simplest of all mathematical functions, and so are the starting place for modeling

- Quantitative variables, by nature, are amenable to algebraic manipulation.
- Given two quantitative variables X and Y, we construct a mathematical model that expresses the values of Y as a function of the values of X:

Y=f(X)

• Linear functions are the simplest of all mathematical functions, and so are the starting place for modeling

 $Y = \beta_0 + \beta_1 X$ with β_0, β_1 fixed constants

- Quantitative variables, by nature, are amenable to algebraic manipulation.
- Given two quantitative variables X and Y, we construct a mathematical model that expresses the values of Y as a function of the values of X:

$$Y=f(X)$$

• Linear functions are the simplest of all mathematical functions, and so are the starting place for modeling

 $Y = \beta_0 + \beta_1 X$ with β_0, β_1 fixed constants

• Of course, in the wild, the observed values of Y will **not** be perfectly predicted by the values of X.

$$Y = \beta_0 + \beta_1 X + \epsilon$$
 where ϵ is the error

Fitting a Line by Least-Squares Regression 0000000000

Scatterplots and Linear Relationships I







Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships I

State-by-State Graduation and Poverty Rates



Poverty rate based 2020 US Census, obtained from US Census website Grad rate based 2018–19 school year, obtained from NCES website • Explanatory Variable:

Fitting a Line by Least-Squares Regression 0000000000

Scatterplots and Linear Relationships I







Fitting a Line by Least-Squares Regression 0000000000

Scatterplots and Linear Relationships I





- Explanatory Variable:
 Poverty Rate (X)
- Response Variable:

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships I





- Explanatory Variable:
 Poverty Rate (X)
- Response Variable:
 - High School Graduation Rate (Y)

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships I





- Explanatory Variable:
 Poverty Rate (X)
- Response Variable:
 - High School Graduation Rate (Y)
- Relationship:

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships I





- Explanatory Variable:
 Poverty Rate (X)
- Response Variable:
 - High School Graduation Rate (Y)
- Relationship:
 - Linear, negative, moderately strong

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships II





• Model (hand-fitted):

$$\hat{Y} = \beta_0 + \beta_1 X = 90 - 0.4X$$

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships II

State-by-State Graduation and Poverty Rates





Model (hand-fitted):

$$\hat{Y} = \beta_0 + \beta_1 X = 90 - 0.4X$$

 Hat (Ŷ) indicates this is an estimate of Y

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships II

State-by-State Graduation and Poverty Rates





Model (hand-fitted):

$$\hat{Y} = \beta_0 + \beta_1 X = 90 - 0.4X$$

- Hat (\hat{Y}) indicates this is an estimate of γ
- Slope of β₁ = -0.4 means every 1 unit increase in Poverty corresponds to a 0.4 unit decrease on average in Graduation.

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships II

State-by-State Graduation and Poverty Rates





Model (hand-fitted):

$$\hat{Y} = \beta_0 + \beta_1 X = 90 - 0.4X$$

- Hat (\hat{Y}) indicates this is an estimate of Y
- Slope of β₁ = -0.4 means every 1 unit increase in Poverty corresponds to a 0.4 unit decrease on average in Graduation.
- Intercept of $\beta_0 = 90$ means model predicts graduation rate of 90% when poverty rate is 0%.

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships III

State-by-State Graduation and Poverty Rates





Model:

$$\hat{Y} = 90 - 0.4 \cdot X$$

• What does the model predict to be the graduation rate for a state with theoretical poverty rate 10%?

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships III

State-by-State Graduation and Poverty Rates





Model:

$$\hat{Y} = 90 - 0.4 \cdot X$$

• What does the model predict to be the graduation rate for a state with theoretical poverty rate 10%?

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships III





State-by-State Graduation and Poverty Rates

Model:

$$\hat{Y} = 90 - 0.4 \cdot X$$

• What does the model predict to be the graduation rate for a state with theoretical poverty rate 7%?

$$\hat{Y} = 90 - 0.4 \cdot 10 = 86$$

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships IV

State-by-State Graduation and Poverty Rates





Model:

$$\hat{Y} = 90 - 0.4 \cdot X$$

• Oregon has a poverty rate of 14. What does the model predict is Oregon's graduation rate?

$$\hat{Y} = 90 - 0.4 \cdot 14 = 84.4$$

Fitting a Line by Least-Squares Regression 000000000

Scatterplots and Linear Relationships IV

State-by-State Graduation and Poverty Rates





• Model:

$$\hat{Y} = 90 - 0.4 \cdot X$$

• Oregon has a poverty rate of 14. What does the model predict is Oregon's graduation rate?

$$\hat{Y} = 90 - 0.4 \cdot 14 = 84.4$$

But Oregon's actual graduation rate is 80

Fitting a Line by Least-Squares Regression 000000000

Residuals

- Residuals are the leftover variation in the data after accounting for model fit.
- Each observation (X_i, Y_i) has its own residual e_i, which is the difference between the observed (Y_i) and predicted (Ŷ_i) value:

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

Fitting a Line by Least-Squares Regression 000000000

Residuals

- Residuals are the leftover variation in the data after accounting for model fit.
- Each observation (X_i, Y_i) has its own residual e_i, which is the difference between the observed (Y_i) and predicted (Ŷ_i) value:

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

State-by-State Graduation and Poverty Rates, with Residual Heights



Fitting a Line by Least-Squares Regression 0000000000

Residuals

- Residuals are the leftover variation in the data after accounting for model fit.
- Each observation (X_i, Y_i) has its own residual e_i, which is the difference between the observed (Y_i) and predicted (Ŷ_i) value:

$$e_i = Y_i - \hat{Y}$$

State-by-State Graduation and Poverty Rates, with Residual Heights



• Oregon's residual is

 $e = Y - \hat{Y} = 80 - 84.4 = -4.4$

Residual Plot

• To visualize the degree of accuracy of a linear model, we use residual plots:



Residual Plot for Graduation and Poverty Rates

Residual Plot

• To visualize the degree of accuracy of a linear model, we use residual plots:



• Points preserve original *x*-position, but with *y*-position equal to residual.

Residual Plot

• To visualize the degree of accuracy of a linear model, we use residual plots:



• Points preserve original *x*-position, but with *y*-position equal to residual.
Section 2

Quantifying Goodness-of-Fit

Quantifying Goodness-of-Fit 0000 Fitting a Line by Least-Squares Regression 000000000

Correlation Coefficient

• The **Correlation Coefficient** R describes the strength of a *linear* relationship between two quantitative variables, and is always a number between -1 and 1.

Correlation Coefficient

- The **Correlation Coefficient** *R* describes the strength of a *linear* relationship between two quantitative variables, and is always a number between -1 and 1.
- The *sign* of *R* indicates the direction of relationship, while the *magnitude* of *R* indicates the strength.

Correlation Coefficient

- The **Correlation Coefficient** *R* describes the strength of a *linear* relationship between two quantitative variables, and is always a number between -1 and 1.
- The *sign* of *R* indicates the direction of relationship, while the *magnitude* of *R* indicates the strength.

If $ R $ is between	Then the linear relationship is
0.7 and 1	strong
0.3 and .07	moderate
0 and 0.3	weak

Correlation Coefficient

- The **Correlation Coefficient** *R* describes the strength of a *linear* relationship between two quantitative variables, and is always a number between -1 and 1.
- The *sign* of *R* indicates the direction of relationship, while the *magnitude* of *R* indicates the strength.



• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

1

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

• Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.
 - What will the scatterplot of X and Y look like?

1

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.
 - What will the scatterplot of X and Y look like?



• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.
 - What will the scatterplot of X and Y look like?



• Based on the formula, is correlation positive or negative?

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.
 - What will the scatterplot of X and Y look like?



• Based on the formula, is correlation positive or negative?

$$R = \sum$$
(Pos.)(Pos.) + \sum (Neg.)(Neg.)

• The Correlation Coefficient is defined via formula using the means (\bar{x}, \bar{y}) and standard deviations (s_x, s_y) of the variables X and Y:

$$R = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- Suppose that when X is above its mean, then Y also tends to be above its mean. And similarly, when X is below its mean, then Y is also.
 - What will the scatterplot of X and Y look like?



• Based on the formula, is correlation positive or negative?

$$R = \sum$$
(Pos.)(Pos.) + \sum (Neg.)(Neg.)

R = 0.6995848

Quantifying Goodness-of-Fit

Fitting a Line by Least-Squares Regression

Correlation is not Association

• Correlation measures strength of *LINEAR* relationship:

Correlation is not Association

- Correlation measures strength of *LINEAR* relationship:
- Which of the following has the strongest correlation (largest value of |R|)?



Correlation is not Association

- Correlation measures strength of *LINEAR* relationship:
- Which of the following has the strongest correlation (largest value of |R|)?



Answer: (b), not (a)

• Computing a correlation coefficient is no substitute for data visualization.

- Computing a correlation coefficient is no substitute for data visualization.
- All of the following have identical, strong positive correlation (R = 0.82):

- Computing a correlation coefficient is no substitute for data visualization.
- All of the following have identical, strong positive correlation (R = 0.82):



I II III IV ## Correlation 0.82 0.82 0.82 0.82

- Computing a correlation coefficient is no substitute for data visualization.
- All of the following have identical, strong positive correlation (R = 0.82):



I II III IV ## Correlation 0.82 0.82 0.82 0.82

• However, each graphic tells a radically different story about the relationship between the variables.

Section 3

Fitting a Line by Least-Squares Regression

Quantifying Goodness-of-Fit 00000 Fitting a Line by Least-Squares Regression

Measure for **BEST** Line

• The line of best fit to a scatterplot should minimize residuals, meaning:

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

• Option 2: Minimize the sum of squares

$$e_1^2 + e_2^2 + \dots + e_n^2$$

• Option 2 is usually preferred.

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- Option 2 is usually preferred.
 - Most commonly used.

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \cdots + e_n^2$$

- Option 2 is usually preferred.
 - Most commonly used.
 - **2** More computationally efficient.

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \cdots + e_n^2$$

- Option 2 is usually preferred.
 - 1 Most commonly used.
 - Ø More computationally efficient.
 - 8 Has theoretical advantages (by analogy with distance and pythagorean thm.)

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \cdots + e_n^2$$

- Option 2 is usually preferred.
 - Most commonly used.
 - **2** More computationally efficient.
 - 8 Has theoretical advantages (by analogy with distance and pythagorean thm.)
 - **()** Appropriately weights one large residuals as "worse" than many small ones.

- The line of best fit to a scatterplot should minimize residuals, meaning:
 - Option 1: Minimizing the sum of absolute values

 $|e_1|+|e_2|+\cdots+|e_n|$

$$e_1^2 + e_2^2 + \cdots + e_n^2$$

- Option 2 is usually preferred.
 - 1 Most commonly used.
 - **2** More computationally efficient.
 - 8 Has theoretical advantages (by analogy with distance and pythagorean thm.)
 - O Appropriately weights one large residuals as "worse" than many small ones.
 - **6** Has well-understood properties for inference

Quantifying Goodness-of-Fit 00000 Fitting a Line by Least-Squares Regression

Line of Best Fit

• For the three data points below, consider candidates for the line of best fit:



Goal: minimize $e_1^2 + e_2^2 + e_3^2$

Line of Best Fit

• For the three data points below, consider candidates for the line of best fit:



 $\mathrm{Purple\ line:}\quad e_1^2+e_2^2+e_3^2=1^2+0^2+1^2=2$

Line of Best Fit

• For the three data points below, consider candidates for the line of best fit:



 $\mathrm{Maroon\ line}: \quad e_1^2 + e_2^2 + e_3^2 = 0^2 + 1^2 + 1^2 = 2$

Quantifying Goodness-of-Fit 00000 Fitting a Line by Least-Squares Regression 0000000000

Line of Best Fit

• For the three data points below, consider candidates for the line of best fit:



 $\mathrm{Blue\ line:}\quad e_1^2+e_2^2+e_3^2=0.5^2+1^2+0.5^2=1.5$

A Formula for the Least Squares Regression Line

• Suppose *n* observations for variables *X* and *Y* are collected:

 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ with means \bar{x}, \bar{y} , standard deviations s_x, s_y , and correlation R.

A Formula for the Least Squares Regression Line

• Suppose *n* observations for variables *X* and *Y* are collected:

 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ with means \bar{x}, \bar{y} , standard deviations s_x, s_y , and correlation R.

• The Least Squares Regression Line modeling Y as a function of X is

$$\hat{Y} = \beta_0 + \beta_1 X$$
A Formula for the Least Squares Regression Line

• Suppose *n* observations for variables *X* and *Y* are collected:

 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ with means \bar{x}, \bar{y} , standard deviations s_x, s_y , and correlation R.

• The Least Squares Regression Line modeling Y as a function of X is

$$\hat{Y} = \beta_0 + \beta_1 X$$

where the slope β_1 is given by

$$\beta_1 = \frac{s_y}{s_x}R$$

A Formula for the Least Squares Regression Line

• Suppose *n* observations for variables *X* and *Y* are collected:

 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ with means \bar{x}, \bar{y} , standard deviations s_x, s_y , and correlation R.

• The Least Squares Regression Line modeling Y as a function of X is

$$\hat{Y} = \beta_0 + \beta_1 X$$

where the slope β_1 is given by

$$\beta_1 = \frac{s_y}{s_x}R$$

and where the intercept is given by

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

• The least squares line is

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

• The line *always* goes through the point (\bar{x}, \bar{y})

٠

Quantifying Goodness-of-Fit 00000 Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

• The least squares line is

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$



The line *always* goes through the point (\bar{x}, \bar{y})

##		mean_x	mean_y
##	1	2.108887	5.179967

Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

Fitting a Line by Least-Squares Regression 0000000000

Properties of the Least-Squares Regression line

• The least squares line is

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

• The slope has the same sign as the correlation coefficient

Fitting a Line by Least-Squares Regression 0000000000

Properties of the Least-Squares Regression line

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$





Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

• The least squares line is

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

• The slope is close to 0 when either $R \approx 0$ or when s_x is much bigger than s_y

Fitting a Line by Least-Squares Regression

Properties of the Least-Squares Regression line

• The least squares line is

$$\hat{Y} = \beta_0 + \beta_1 X$$
 $\beta_1 = \frac{s_y}{s_x} R$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

• The slope is close to 0 when either $R \approx 0$ or when s_x is much bigger than s_y

