

STAT 2593

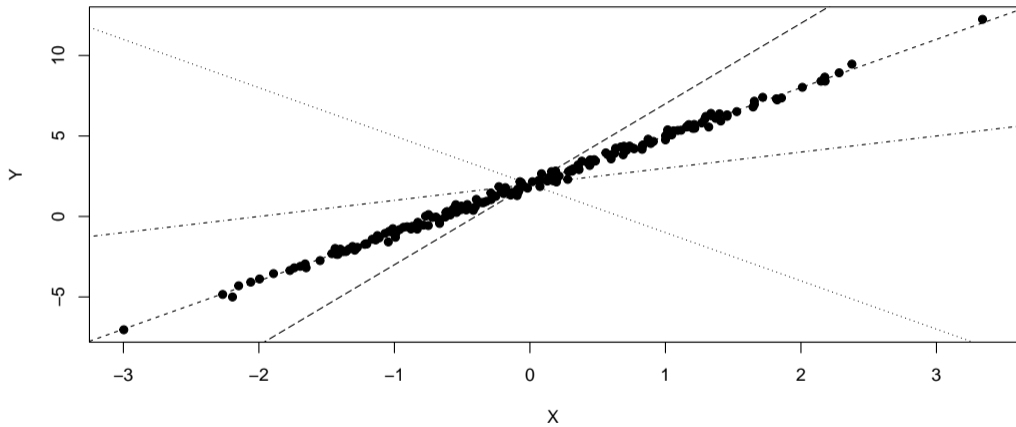
Lecture 039 - Estimating Model Parameters

Dylan Spicker

Estimating Model Parameters

Learning Objectives

1. Understand how the best linear regression line is solved for.
2. Interpret the regression coefficients and understand the limitations.
3. Understand the variance breakdown, what is explained by the model, and what is not.
4. Understand the coefficient of determination and its limitations.



Estimating Parameters for the Simple Linear Regression

- ▶ We want the *best* possible line, as defined by β_0 and β_1 .

Estimating Parameters for the Simple Linear Regression

- ▶ We want the *best* possible line, as defined by β_0 and β_1 .
- ▶ We will use estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, which gives $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.

Estimating Parameters for the Simple Linear Regression

- ▶ We want the *best* possible line, as defined by β_0 and β_1 .
- ▶ We will use estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, which gives $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.
 - ▶ The mistakes we make we call **residuals**, and denote this $e_i = y_i - \hat{y}_i$.

Estimating Parameters for the Simple Linear Regression

- ▶ We want the *best* possible line, as defined by β_0 and β_1 .
- ▶ We will use estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, which gives $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.
 - ▶ The mistakes we make we call **residuals**, and denote this $e_i = y_i - \hat{y}_i$.
 - ▶ The idea is to minimize the **squared residuals**, given by $\sum_{i=1}^n e_i^2$.

Estimating Parameters for the Simple Linear Regression

- ▶ We want the *best* possible line, as defined by β_0 and β_1 .
- ▶ We will use estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, which gives $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.
 - ▶ The mistakes we make we call **residuals**, and denote this $e_i = y_i - \hat{y}_i$.
 - ▶ The idea is to minimize the **squared residuals**, given by $\sum_{i=1}^n e_i^2$.
- ▶ Doing this results in

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.
 - ▶ Note, $X = 0$ does not always have a substantive meaning.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.
 - ▶ Note, $X = 0$ does not always have a substantive meaning.
- ▶ Recall that β_1 is the slope of the regression line.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.
 - ▶ Note, $X = 0$ does not always have a substantive meaning.
- ▶ Recall that β_1 is the slope of the regression line.
 - ▶ $\hat{\beta}_1$ is the *best guess* we have at the slope, given the data.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.
 - ▶ Note, $X = 0$ does not always have a substantive meaning.
- ▶ Recall that β_1 is the slope of the regression line.
 - ▶ $\hat{\beta}_1$ is the *best guess* we have at the slope, given the data.
 - ▶ For a 1 unit increase in X , we expect that Y will change by $\hat{\beta}_1$.

Regression Parameter Interpretation

- ▶ Recall that β_0 is the intercept of the regression line.
 - ▶ $\hat{\beta}_0$ is the *best guess* we have at the intercept, given the data.
 - ▶ If $X = 0$, this is the value that we would expect Y to take on.
 - ▶ Note, $X = 0$ does not always have a substantive meaning.
- ▶ Recall that β_1 is the slope of the regression line.
 - ▶ $\hat{\beta}_1$ is the *best guess* we have at the slope, given the data.
 - ▶ For a 1 unit increase in X , we expect that Y will change by $\hat{\beta}_1$.
 - ▶ Be careful for *extrapolation* and for *bad model fits*.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.
- ▶ We can estimate the variance, σ^2 , as $\hat{\sigma}^2 = \frac{SSE}{n-2}$.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.
- ▶ We can estimate the variance, σ^2 , as $\hat{\sigma}^2 = \frac{SSE}{n-2}$.
- ▶ We can also consider the **total variation** in Y , which is given by

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = S_y^2.$$

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.
- ▶ We can estimate the variance, σ^2 , as $\hat{\sigma}^2 = \frac{SSE}{n-2}$.
- ▶ We can also consider the **total variation** in Y , which is given by

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = S_y^2.$$

- ▶ We will always have that $SSE \leq SST$.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.
- ▶ We can estimate the variance, σ^2 , as $\hat{\sigma}^2 = \frac{SSE}{n-2}$.
- ▶ We can also consider the **total variation** in Y , which is given by

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = S_y^2.$$

- ▶ We will always have that $SSE \leq SST$.
- ▶ The difference, $SST - SSE$, is the variation **explained by the model**.

Error Sums of Squares

- ▶ We call the sum of squared residuals the **error sum of squares**, $SSE = \sum_{i=1}^n e_i^2$.
 - ▶ This is the variation left unexplained by the model.
- ▶ We can estimate the variance, σ^2 , as $\hat{\sigma}^2 = \frac{SSE}{n-2}$.
- ▶ We can also consider the **total variation** in Y , which is given by

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = S_y^2.$$

- ▶ We will always have that $SSE \leq SST$.
- ▶ The difference, $SST - SSE$, is the variation **explained by the model**.
- ▶ This is called the **regression sums of squares**, denoted SSR.

The Coefficient of Determination

- ▶ If we consider the ratio of the variance that is explained by the model we write

$$r^2 = \frac{SSR}{SST}.$$

The Coefficient of Determination

- ▶ If we consider the ratio of the variance that is explained by the model we write

$$r^2 = \frac{SSR}{SST}.$$

- ▶ This is the **r-squared** value, or the **coefficient of determination**.

The Coefficient of Determination

- ▶ If we consider the ratio of the variance that is explained by the model we write

$$r^2 = \frac{SSR}{SST}.$$

- ▶ This is the **r-squared** value, or the **coefficient of determination**.
- ▶ It gives the proportion of variance in Y which is captured by the model.

The Coefficient of Determination

- ▶ If we consider the ratio of the variance that is explained by the model we write

$$r^2 = \frac{SSR}{SST}.$$

- ▶ This is the **r-squared** value, or the **coefficient of determination**.
- ▶ It gives the proportion of variance in Y which is captured by the model.
- ▶ It is typically used to indicate the strength of the relationship, with values closer to 1 being preferable.

The Coefficient of Determination

- ▶ If we consider the ratio of the variance that is explained by the model we write

$$r^2 = \frac{SSR}{SST}.$$

- ▶ This is the **r-squared** value, or the **coefficient of determination**.
- ▶ It gives the proportion of variance in Y which is captured by the model.
- ▶ It is typically used to indicate the strength of the relationship, with values closer to 1 being preferable.
- ▶ Note: The coefficient of determination has received a lot of criticism. It is probably best to steer largely clear of it!

Summary

- ▶ Linear regression estimates are determined through the **least squares procedure**.
- ▶ There is a closed form expression for both the slope and the intercept.
- ▶ The intercept gives the value we expect to observe at $X = 0$ and the slope captures the expected change in outcome for a unit change in X .
- ▶ The total variance can be decomposed into the error sum of squares and the regression sum of squares.
- ▶ The proportion of variance which is explained by the model is called the coefficient of determination.