CSCI 3022 Intro to Data Science

Notebook day:
$$1819$$
 (maybe 20)

Prove that the OLS estimators: e_{S} in the continuous f_{S} in the continuous f_{S} in the continuous f_{S} is a law to f_{S} in the continuous f_{S} in the continuous f_{S} is a law to f_{S} in the continuous f_{S} in the continuous f_{S} is a law to f_{S} in the continuous f_{S} in the continuous f_{S} in the continuous f_{S} is a law to f_{S} in the continuous f_{S} in the continuous f_{S} in the continuous f_{S} is a law to f_{S} in the continuous f_{S} in the continuous

Mullen: Regression

satisfy
$$\widehat{\beta_0}, \widehat{\beta_1} = \operatornamewithlimits{argmin}_{\beta_0,\beta_1} \sum_{i=1}^n (Y_i - \overline{\beta_0} - \beta_1 X_i)^2$$

Hint: what are $\frac{d}{du}$ of

 $(2. \widehat{\hat{\beta}_1}) = \frac{Cov[X,Y]}{\sum_{i=1}^{n} (X_i - \bar{X})^2} = \frac{\sum_{i=1}^{n} (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$

+1 11-51= 152-65 (alculus/

SLR Overview

Problem: use predictor x to describe response y, using a line. We're subject to error, noise, or unexplained variability ε .

- 1. Model: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- 2. Terms: β_0 : the intercept of the goal line; β_1 : its slope
- 3. Assumptions on ε : Independence, Homoskedasticity, Normality

Goal:

Given sample data, which consists of n observed pairs, $(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)$, construct an estimated "line of best fit":

$$y = \hat{\beta}_0 + \hat{\beta}_1 x$$

This line can then be used to make predictions or provide explanations for unobserved phenomena.

Estimating SLR Parameters

A line provides the **best fit** to the data if the sum of the squared vertical distances (deviations) from the observed points to that line is as small as it can be.

The sum of squared vertical deviations from the data points to the line $y = \beta_0 + \beta_1 x$ is then

$$\sum_{i=1}^{n} \left(Y_i - \beta_0 - \beta_1 X_i \right)^2$$

The point estimates of β_0 and β_1 , denoted $\underline{\hat{\beta}_0}; \underline{\hat{\beta}_1}$ are called the *least squares estimates*. They are those values that minimize SSE or sum of squared errors.

$$\hat{\beta}_0, \hat{\beta}_1 = \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0) (X_i)^2$. Sounds like a Calculus problem!

$$\frac{1}{4} \frac{1}{4} \frac{1$$

$$\int_{\overline{B_i}} \sqrt{J_y} \geq (\# - X_i, y)^2 = Z(\# - X_i, y)$$

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$. Sounds like a Calculus problem!

$$\frac{df}{d\beta_0} = \frac{d}{d\beta_0} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

$$\frac{df}{d\beta_1} = \frac{d}{d\beta_1} \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2$$

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$. Sounds like a Calculus problem!

$$\frac{df}{d\beta_0} = \frac{d}{d\beta_0} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

$$\frac{df}{d\beta_1} = \frac{d}{d\beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

For finding the joint maximum/minimum of multiple inputs, we end up with a system of equations: set both equal to zero and find the values that make both equal to zero.

$$f(\beta_0, \beta_1) = \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2, \text{ so } :$$

$$\frac{\partial f}{\partial \beta_0} = \underbrace{\frac{\partial}{\partial \beta_0}}_{i=1} \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2$$

$$\frac{\partial f}{\partial \beta_1} = \underbrace{\frac{d}{\partial \beta_0}}_{i=1} \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2$$

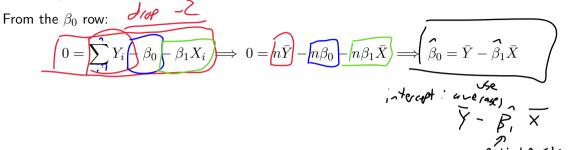
$$f(\beta_0,\beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2, \text{ so }:$$

$$\frac{\partial f}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 = \sum_{i=1}^n -2 (Y_i - \beta_0 - \beta_1 X_i)$$

$$\frac{\partial f}{\partial \beta_1} = \frac{d}{d\beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 = \sum_{i=1}^n -2 X_i (Y_i - \beta_0 - \beta_1 X_i)$$

$$f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2, \text{ so }:$$

$$\frac{\partial f}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 = \sum_{i=1}^n \frac{\partial}{\partial \beta_0} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} + \frac{\partial f}{\partial \beta_1} \right] = \frac{\partial}{\partial \beta_1} \left[\frac{\partial f}{\partial \beta_1} + \frac$$



From the β_0 row:

$$0 = \sum Y_i - \beta_0 - \beta_1 X_i \implies 0 = n\bar{Y} - n\beta_0 - n\beta_1 \bar{X} = \beta_0 = \bar{Y} - \beta_1 \bar{X}$$

Plugging that into the other row gives

$$0 = \sum X_i (Y_i - (\bar{Y} - \beta_1 \bar{X}) - \beta_1 X_i) = \sum X_i (Y_i - \bar{Y} + \beta_1 (\bar{X} - X_i))$$

$$0 = \sum_{X_i (Y_i - \bar{Y})} X_i (\bar{X} - \bar{X}) \implies \beta_1 = \frac{\sum_{X_i (Y_i - \bar{Y})}}{\sum_{X_i (X_i)} \bar{X}} \implies \beta_1 = \frac{\sum_{X_i (Y_i - \bar{Y})}}{\sum_{X_i (X_i)} \bar{X}}$$

From the β_0 row:

SLR Day 1 Review

$$0 = \sum Y_i - \beta_0 - \beta_1 X_i \implies 0 = n\bar{Y} - n\beta_0 - n\beta_1 \bar{X} \implies \beta_0 = \bar{Y} - \beta_1 \bar{X}$$

(2.3)

Plugging that into the other row gives

$$0 = \sum X_i (Y_i - (\bar{Y} - \beta_1 \bar{X}) - \beta_1 X_i) = \sum X_i (Y_i - \bar{Y} + \beta_1 (\bar{X} - X_i))$$

$$0 = \sum X_i (Y_i - \bar{Y}) + \beta_1 \sum X_i (\bar{X} - X_i) \implies \beta_1 = \frac{\sum X_i (Y_i - \bar{Y})}{\sum X_i (X_i - \bar{X})}$$

... are we done?

Note: $\sum (X_i - \bar{X}) = 0$, so $\bar{X} \sum (X_i - \bar{X}) = 0$. That's the difference between our current solution and the version on the prior slide.

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

1.
$$\hat{\beta_0} =$$

2.
$$\hat{\beta_1} =$$

What happens if $\beta_0 \approx 0$? If $\beta_1 \approx 0$?

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

1.
$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$
 integrate vsels
$$2. \hat{\beta}_1 = \frac{Cov[X,Y]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$
 What happens if $\beta_0 \approx 0$? If $\beta_1 \approx 0$?
$$5 = \frac{(\chi_i - \bar{X})^2}{(\chi_i - \bar{X})^2}$$

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

1.
$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

2.
$$\hat{\beta}_1 = \frac{Cov[X,Y]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

What happens if $\beta_0 \approx 0$? If $\beta_1 \approx 0$?

One result: the regression line goes through $(0,\beta_0)$. It also goes through $(\bar{X},\bar{Y})!$

Daily Recap

Today we learned

1. Regression!

Moving forward:

- nb day Friday

Next time in lecture:

- More Regression!