

DSFA
Spring 2021

Lecture 29

Least Squares

Announcements

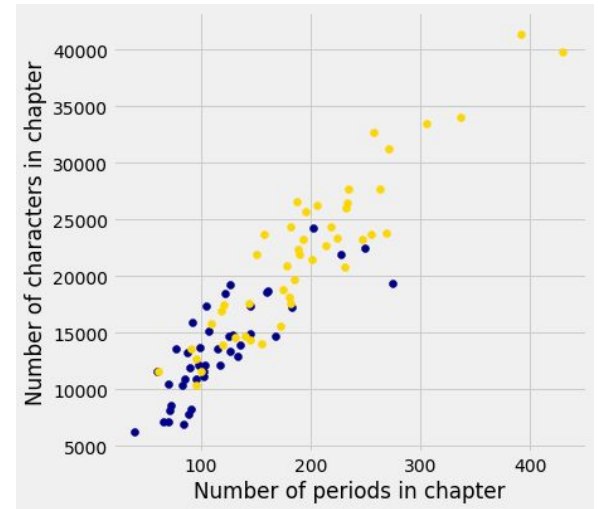
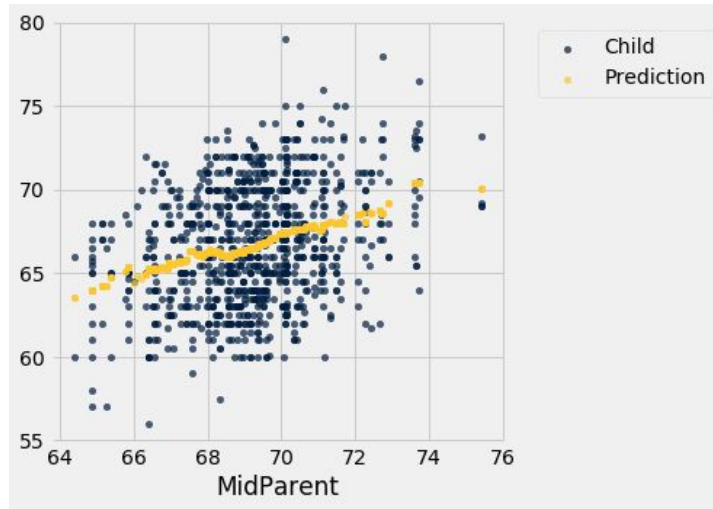
- Prelim 2, April 20, 8:30PM-10PM in Kennedy 116 (here) for Ithaca-resident students, assigned seating
 - Coverage from Lecture 12 - Lecture 26 (Monday)
 - Review sheet and sample exam posted on Canvas.
 - Table of functions included again, allowed a double-sided sheet of notes you make yourself
-

Announcements: Coming up

- HW 5 out, due 4/30
 - This Wednesday/Thursday: Lab 8
 - Friday 4/23, Monday 4/26: Wellness days, no lecture
 - Wednesday 4/28, Thursday 4/29: Lab 9
 - Friday 4/30: Project 3 out, Part 1 due 5/7, Part 2 due 5/14
 - Labs 5/5, 5/6; 5/12, 5/13 dedicated to working on Project 3.
 - Final: May 22.
-

Prediction

If we have a line describing the relation between two variables, we can make predictions

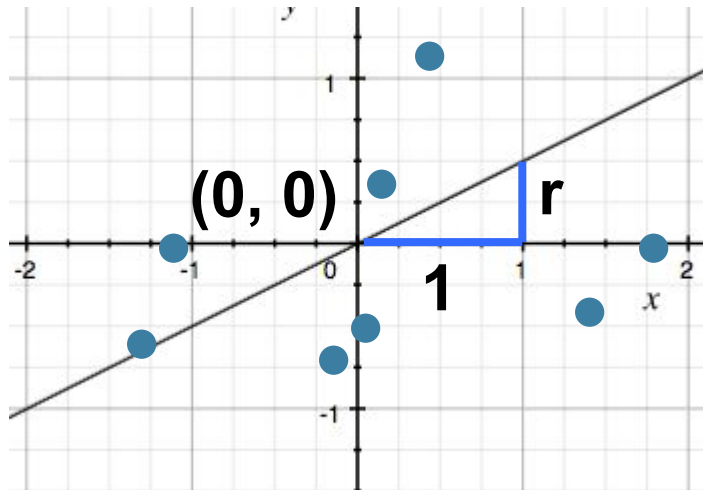


Prediction

- **Problem:** given a known x value, predict y , where both are in standard units
 - **Solution:**
 - Compute correlation coefficient r
 - Predict that $y = r * x$
 - Why is that a line? (slope = r , intercept = 0)
 - Why use that line?
 - It is a version of the graph of averages, smoothed to a line
-

Regression Line Equation

In standard units, the equation of the regression line is:



Fitted value

Observed value

$$y_{(\text{su})} = r \times x_{(\text{su})}$$

Correlation coefficient

Regression Line Equation

In original units, the regression line has this equation:

$$\frac{\text{estimate of } y - \text{average of } y}{\text{SD of } y} = r \times \frac{\text{the given } x - \text{average of } x}{\text{SD of } x}$$

y in standard units

x in standard units

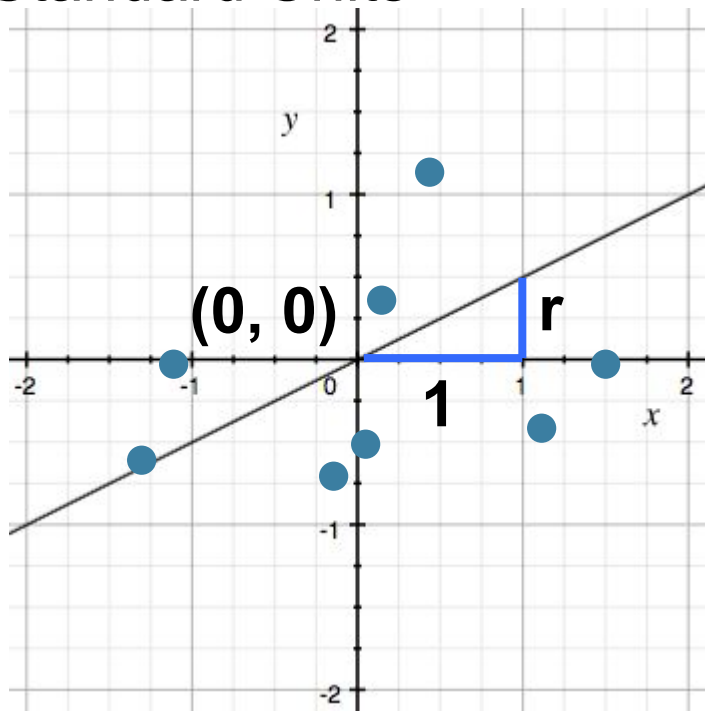
$$y = \text{slope} \times x + \text{intercept}$$

$$\text{slope of the regression line} = r \cdot \frac{\text{SD of } y}{\text{SD of } x}$$

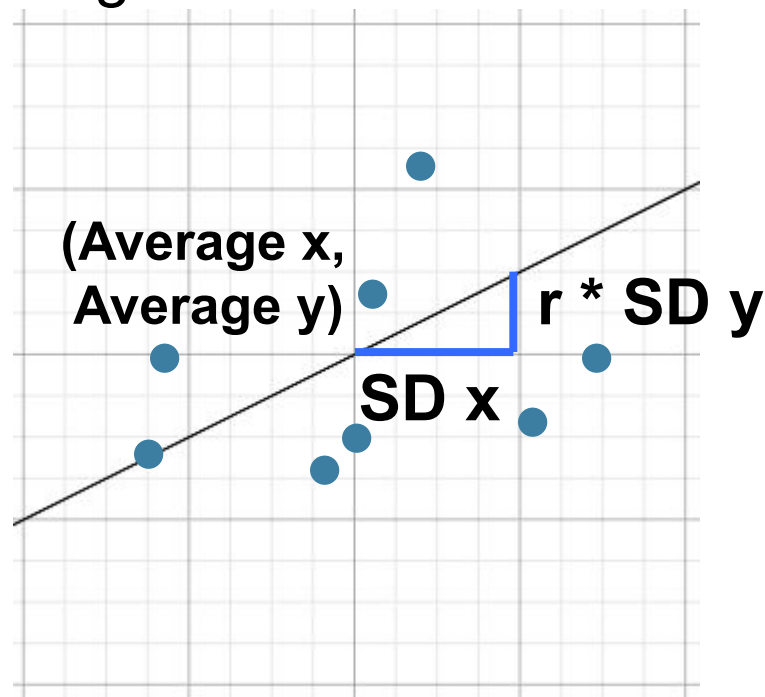
$$\text{intercept of the regression line} = \text{average of } y - \text{slope} \cdot \text{average of } x$$

Regression Line

Standard Units



Original Units



Which of these to predict a final exam score?

$$y = r * x$$

$$y = r * x + i$$

$$y = s * x + i$$

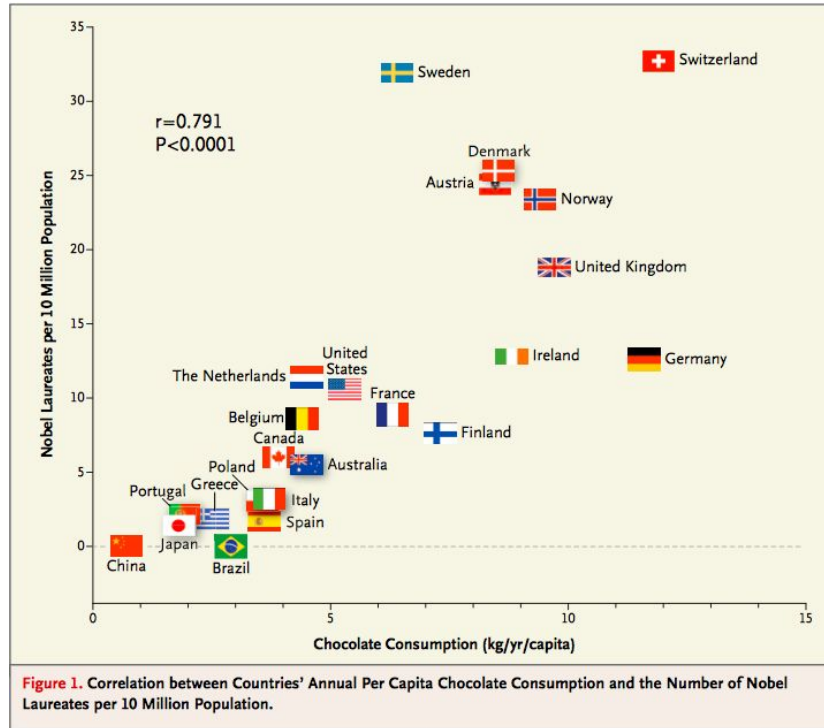
None of the above



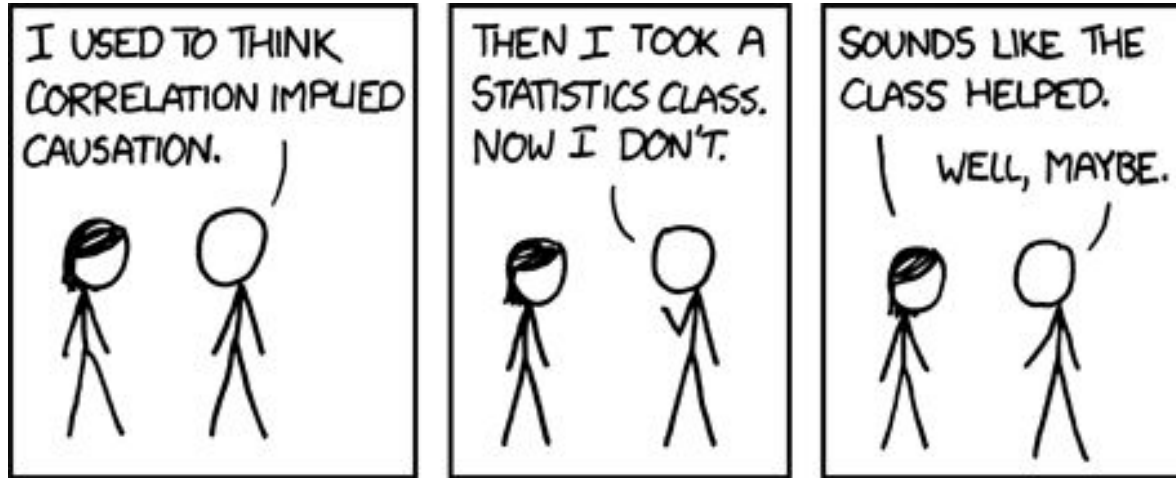
Abuses of r

- Summarizing non-linear data with r
 - Eliminating outliers to “improve” r
 - Drawing conclusions about individuals based on data about groups (*ecological* correlations)
 - Jumping to conclusions about causality
-

Correlation is not causation



Correlation is not causation



Quantifying Error

Error in Prediction

- How good is the regression line at making predictions?
 - Hard to say for unknown data
 - But easy for data we already have

- **error = actual value – prediction**

(Demo)

RMSE

RMSE = root mean square error

$$\text{RMSE} = \text{std}(y) * \text{sqrt}(1 - r^2)$$

- If $r = 1$, what is RMSE? 0
- If $r = 0$, what is RMSE? $\text{std}(y)$

Compare regression line to other lines using RMSE...

(Demo)

Line with smallest RMSE?

- SciPy function `minimize(f)` returns the value \mathbf{x} that produces the minimum output $\mathbf{f}(\mathbf{x})$ from \mathbf{f}
- Also works for functions that make multiple arguments
- How to use to find best line:
 - Write function `rmse(a, b)` that returns the RMSE for line with slope \mathbf{a} and intercept \mathbf{b}
 - Call `minimize(rmse)` and get output array $[\mathbf{a}_0, \mathbf{b}_0]$
 - \mathbf{a}_0 is slope and \mathbf{b}_0 intercept of line that minimizes RMSE

(Demo)

Regression line

- Regression line has the minimum RMSE of all lines
- Names:
 - Regression line
 - Least squares line
 - “Best fit” line