

DSFA

Spring 2021

Lecture 28

Linear Regression

Announcements

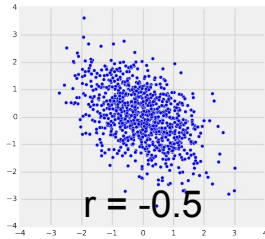
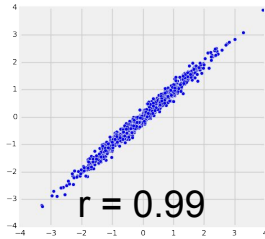
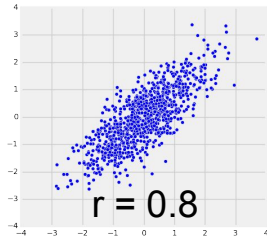
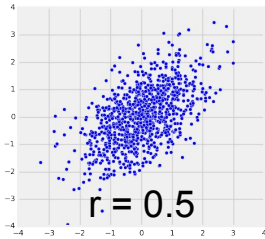
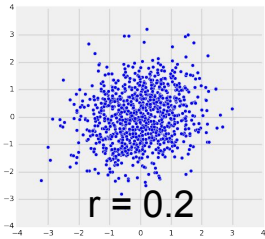
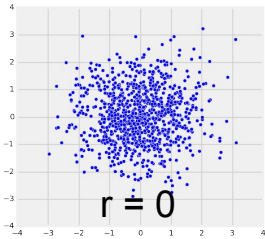
- Project 2, Part 2, due Friday 5:59PM
 - Prelim 2, April 20, 8:30PM-10PM in Kennedy 116 (here) for Ithaca-resident students, assigned seating
 - Coverage from Lecture 12 - Lecture 26 (Monday)
 - Review session on Saturday 3:30PM-5:30PM, room Uris G01
 - Review sheet and sample exam posted on Canvas.
 - NB: The sample exam is not one I wrote, and is likely to be somewhat different than what I will do.
 - Table of functions included again, allowed a double-sided sheet of notes you make yourself
-

Announcements

- HW 5 out this weekend, not due until Friday 4/30

The Correlation Coefficient r

- Measures linear association
- Based on standard units
- $-1 \leq r \leq 1$
 - $r = 1$: scatter is perfect straight line sloping up
 - $r = -1$: scatter is perfect straight line sloping down
- $r = 0$: No linear association; *uncorrelated*



Definition of r

Correlation Coefficient (r) =

average of	product of	x in standard units	and	y in standard units
---------------	------------	---------------------------	-----	---------------------------

Measures how clustered the scatter is around a straight line

(Demo)

Properties of r

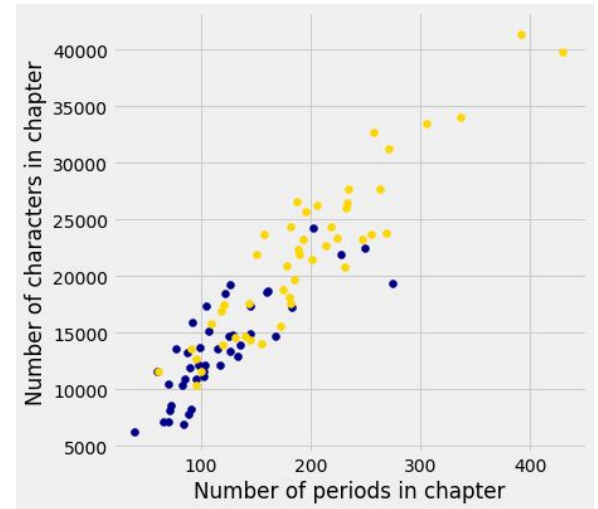
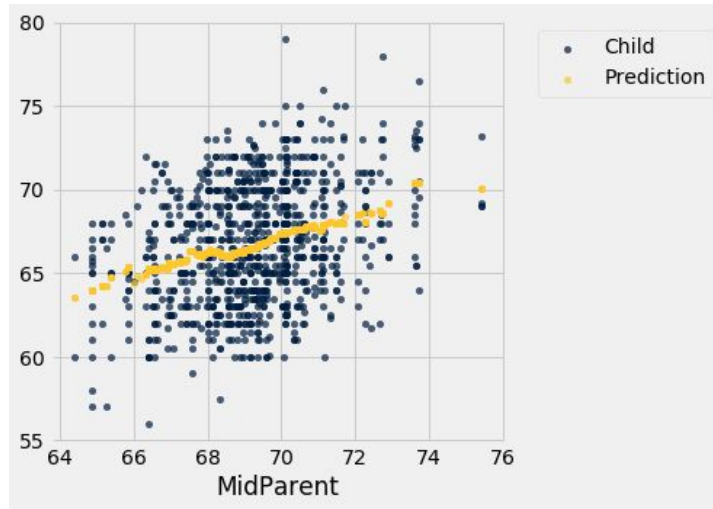
- r is a pure number, with no units
- r is not affected by changing units of measurement
- r is not affected by switching the horizontal and vertical axes

(Demo)

Prediction

Prediction

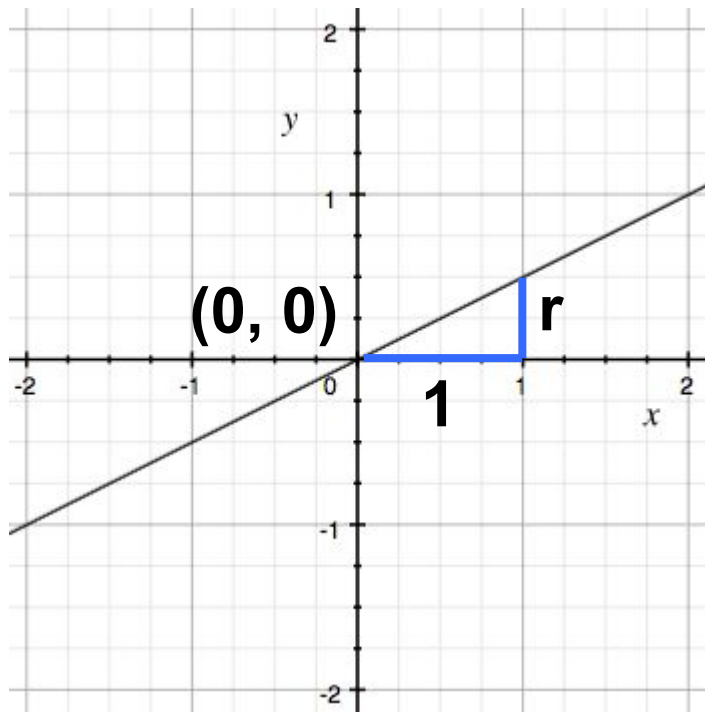
If we have a line describing the relation between two variables, we can make predictions



Prediction

- **Problem:** given a known x value, predict y , where both are in standard units
- **Solution:**
 - Compute r
 - Predict that $y = r * x$
- Why is that a line?

Equation of a Line



$$y = r * x$$

In general:

$$y = a * x + b$$

(a is slope, b is intercept)

Prediction

- **Problem:** given a known x value, predict y , where both are in standard units
- **Solution:**
 - Compute r
 - Predict that $y = r * x$
- Why is that a line?
- Why use *that* line?

(Demo)

Prediction

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

Prediction

- Predict $y = r * x$ (in standard units)
 - Example:
 - $x = 2$ (in standard units)
 - $r = .75$
 - What is the prediction for y (in standard units)?
 - A. 0.0
 - B. 0.75
 - C. 1.5
 - D. 2.0
-

Prediction

- **Predict** $y = r * x$ (in standard units)
 - Example:
 - A course has a typical prelim (mean=70, std=10), and a hard final (mean=50, std=12)
 - The scores on the exams look linearly related when visualized, with $r = .75$
 - **Predict** a student's final exam score, given that their prelim score was 90 (*go ahead and work on that*)
-

When poll is active, respond at pollev.com/dsfa

Text **DSFA** to **22333** once to join

Final score?

50
67.5
68
74

None of the above



Prediction

- Prelim: mean=70, std=10
 - $x = 90 = 70 + 2 * 10$ in original units = 2 standard units
 - Prediction:
 - $y = r * x = .75 * 2 = 1.5$ standard units
 - Final: mean=50, std=12
 - $y = 50 + 1.5 * 12 = \mathbf{68}$ in original units
-

Prediction

- Predict $y = r * x$ (in standard units)
 - If $r = .75$ and x is 2 std above mean, then prediction for y is 1.5 std above mean
 - So y predicted to be **closer to mean** than x

 - “Regression to the mean”
 - Children with exceptionally tall parents tend not to be as tall
 - Galton called it “regression to mediocrity” (Demo)
-

Linear Regression

(Demo)

Equation for regression line

$$(y \text{ in } su) = r * (x \text{ in } su)$$

Equation for regression line

$$(y \text{ in su}) = r * \frac{x - \text{mean}(\text{all } x)}{\text{std}(\text{all } x)}$$

Equation for regression line

$$\frac{y - \text{mean}(\text{all } y)}{\text{std}(\text{all } y)} = r * \frac{x - \text{mean}(\text{all } x)}{\text{std}(\text{all } x)}$$

Equation for regression line

$$\frac{y - \text{mean}(\text{all } y)}{\text{std}(\text{all } y)} = r * \frac{x - \text{mean}(\text{all } x)}{\text{std}(\text{all } x)}$$

Do some algebra to put that in the form $y = \text{slope} * x + \text{intercept}$...

Slope and Intercept

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

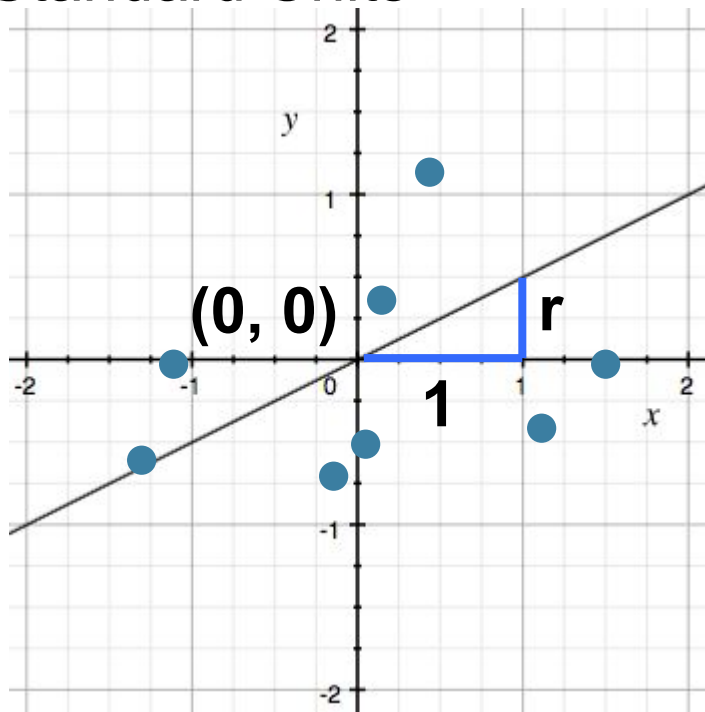
$$\text{slope of the regression line} = r \cdot \frac{\text{SD of } y}{\text{SD of } x}$$

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

(Demo)

Regression Line

Standard Units



Original Units

