

# BC COMS 1016: Intro to Comp Thinking & Data Science

## Lecture 20 – Linear Regression, Least Squares, & Residuals

# Announcements



- Lab 8 – Regression
  - Due Friday 12/04
- Homework 7 - Confidence Intervals, Resampling, the Bootstrap, and the Central Limit Theorem
  - Due Thursday 12/03
- Homework 8 - Linear Regression
  - Due Monday 12/07



# Correlation



# Prediction

---

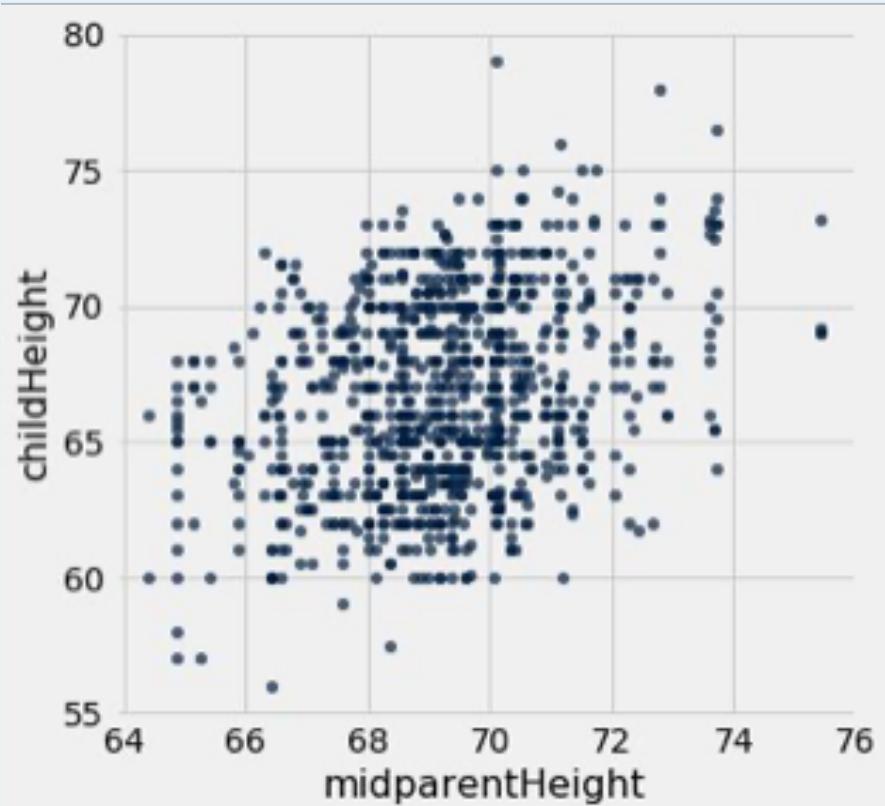
---

# Guess the future



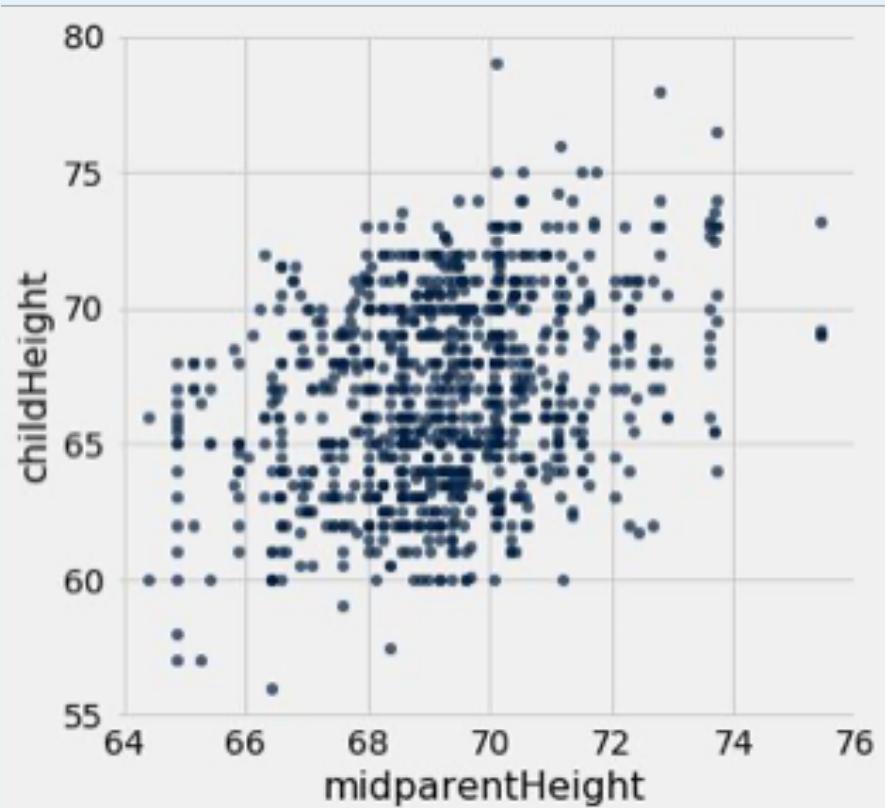
- Based on incomplete information
- One way of making predictions:
  - To predict an outcome for an individual,
  - find others who are like that individual
  - and whose outcomes you know.
  - Use those outcomes as the basis of your prediction.

# Galton's Heights



**Goal:** Predict the height of a new child, based on that child's midparent height

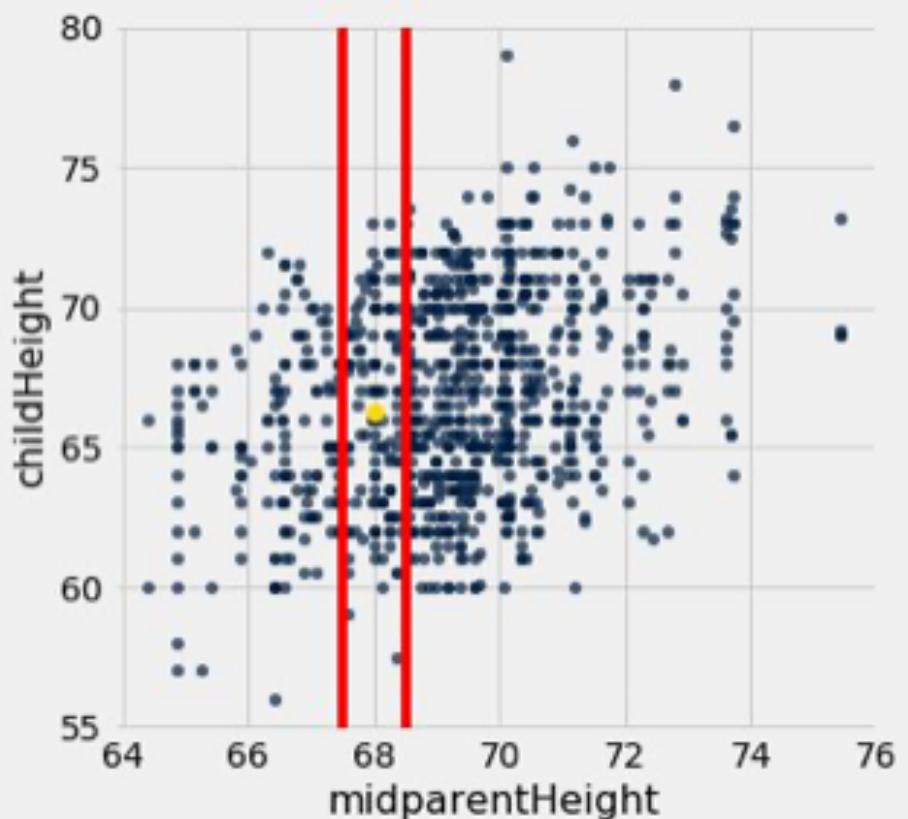
# Galton's Heights



How can we predict a child's height given a midparent height of 68 inches?



# Galton's Heights

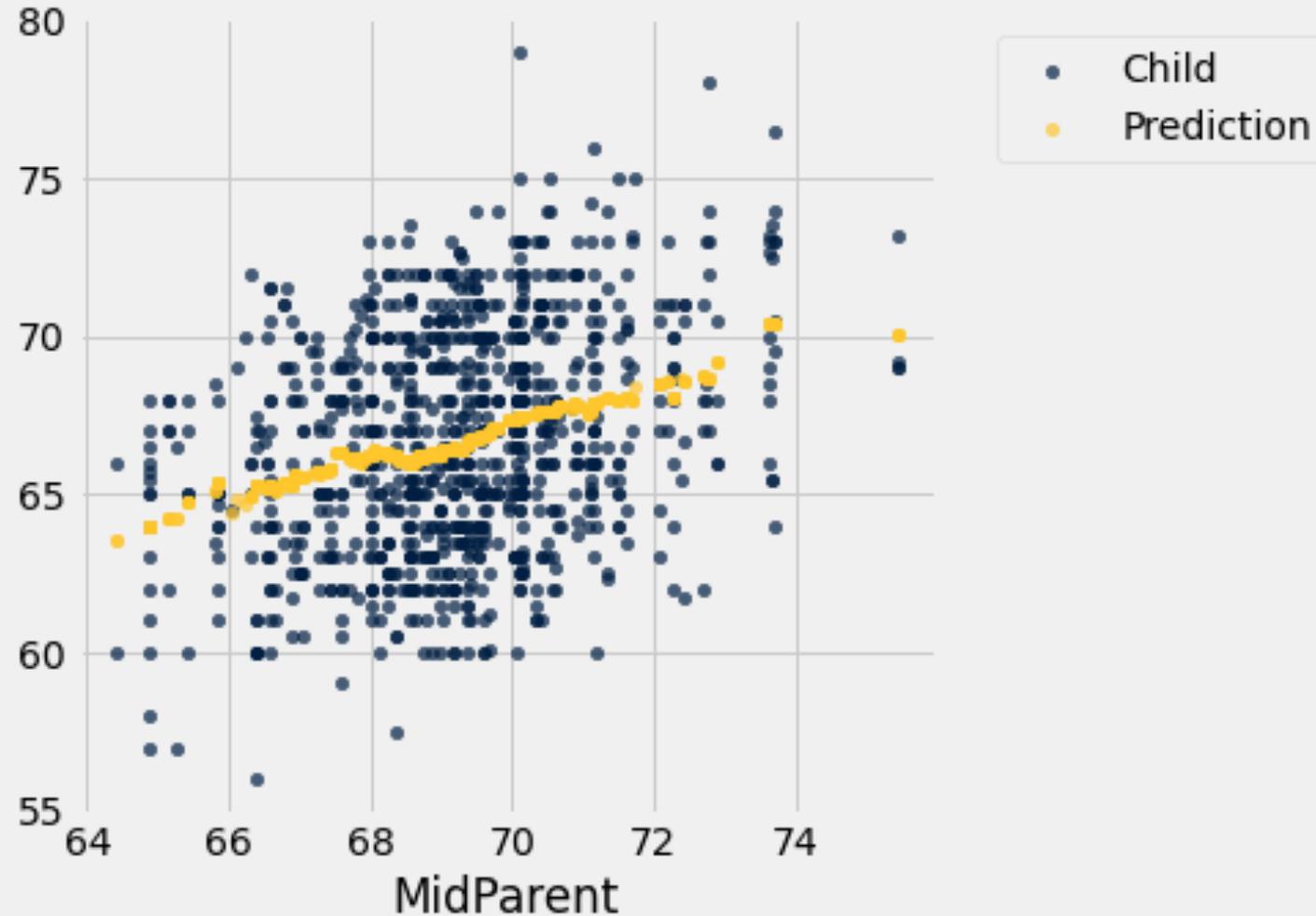


How can we predict a child's height given a midparent height of 68 inches?

**Idea:** Use the average height of the children of all families where the midparent Height is close to 68 inches



# Predicted Heights



# Graph of Average



For each  $x$  value, the prediction is the average of the  $y$  values in its nearby group.

The graph of these predictions is the  
**graph of averages**

If the association between  $x$  and  $y$  is linear, then points in the graph of averages tend to fall on a line. The line is called the **regression line**

# Nearest Neighbor Regression



A method for predicting a numerical  $y$ , given a value of  $x$ :

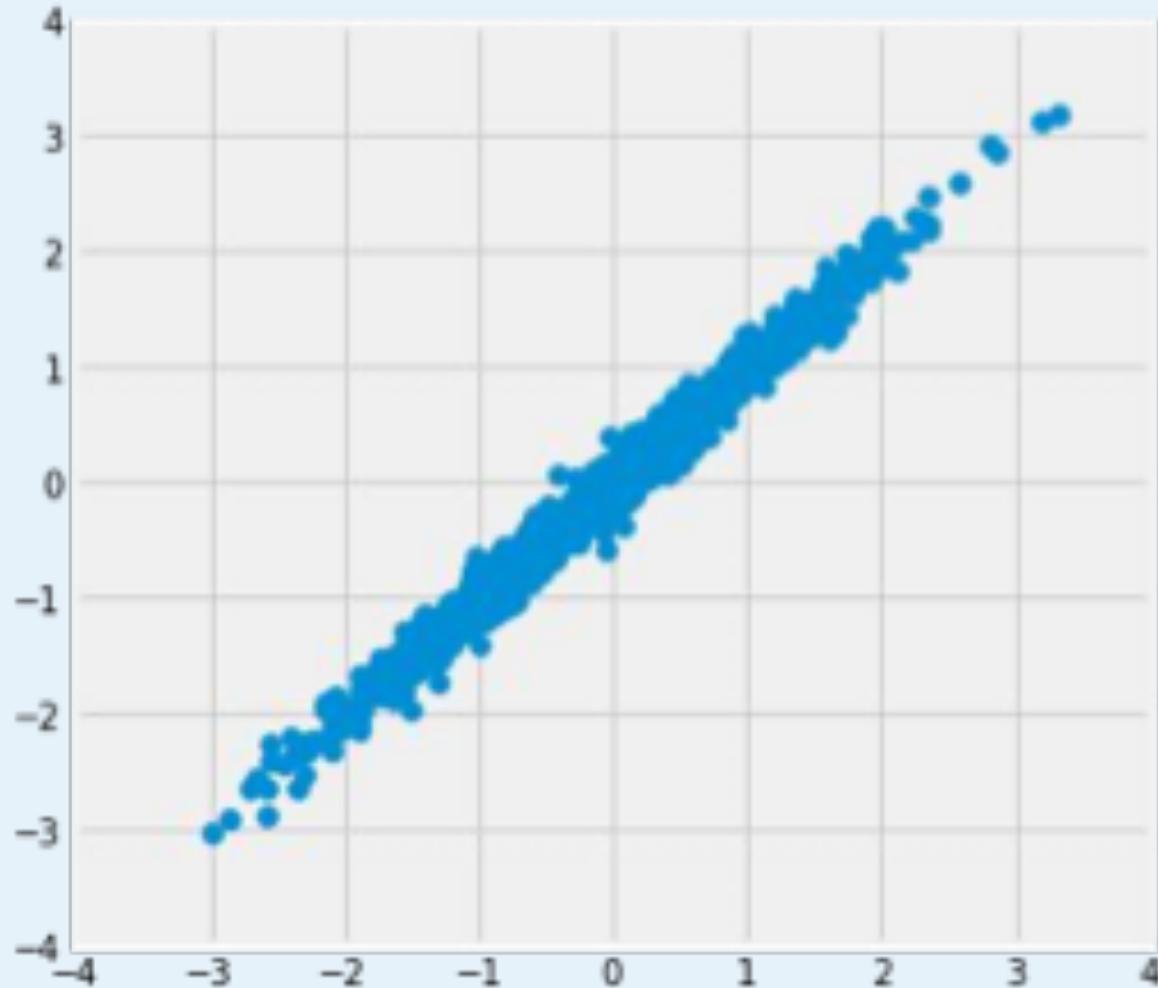
- Identify the group of points where the values of  $x$  are close to the given value
- The prediction is the average of the  $y$  values for the group



# Linear Regression



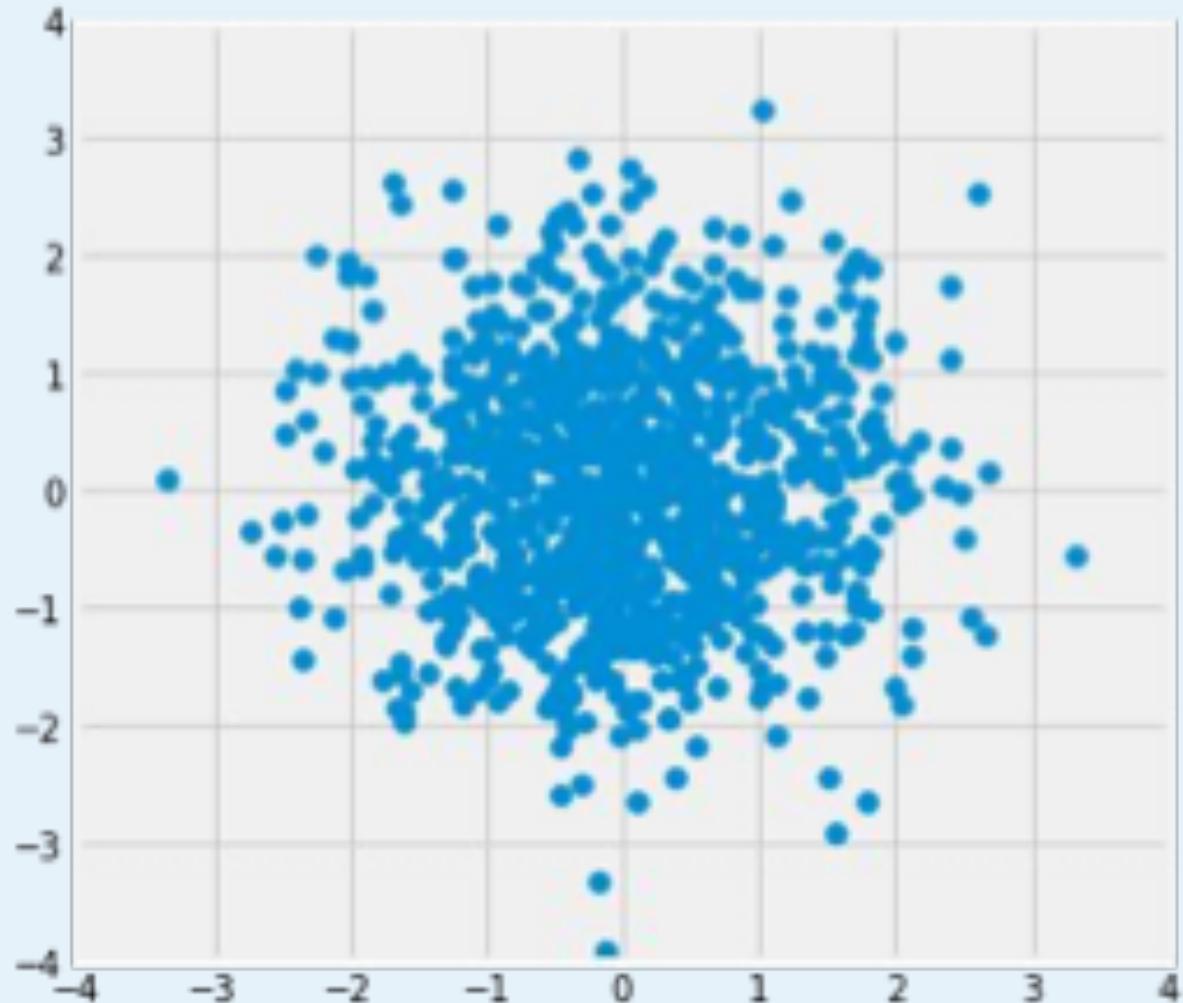
# Where is the prediction line?



$$r = 0.99$$



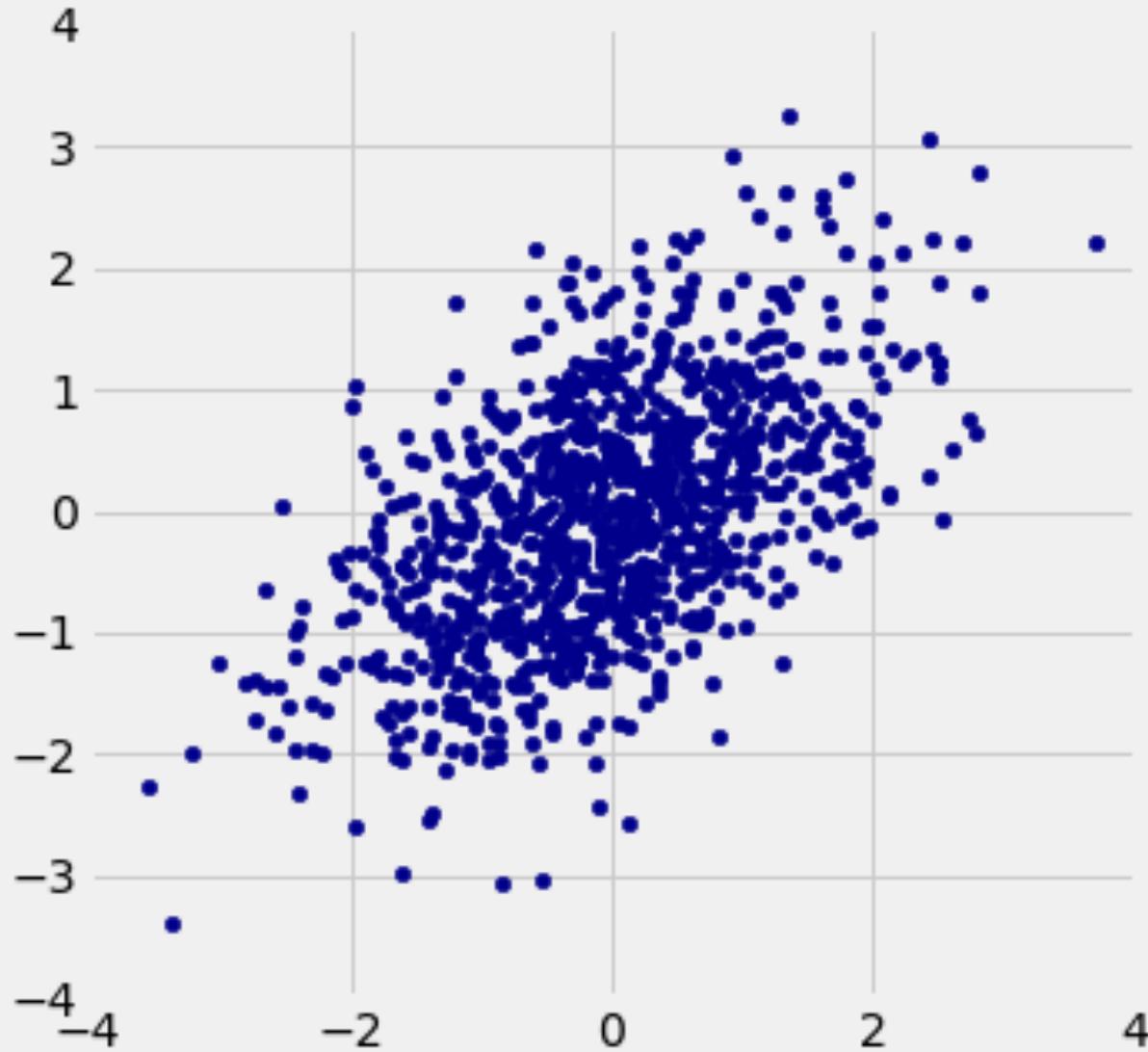
# Where is the prediction line?



$$r = 0.0$$



# Where is the prediction line?



$$r = 0.5$$

# Linear Regression



A statement about x and y pairs

- Measured in *standard units*
- Describing the deviation of x from 0
  - (the average of x's)
- And the deviation of y from 0
  - (the average of y's)

*On average,*

y deviates from 0 less than x deviates from 0

$$y_{su} = r \times x_{su}$$



# Slope and Intercept

# Regression Line Equation



In original units, the regression line has this equation:

$$y_{su} = r \times x_{su}$$

$$\frac{\text{estimate of } y - \text{mean}(y)}{\text{SD of } y} = r \times \frac{\text{given } x - \text{mean}(x)}{\text{SD of } x}$$

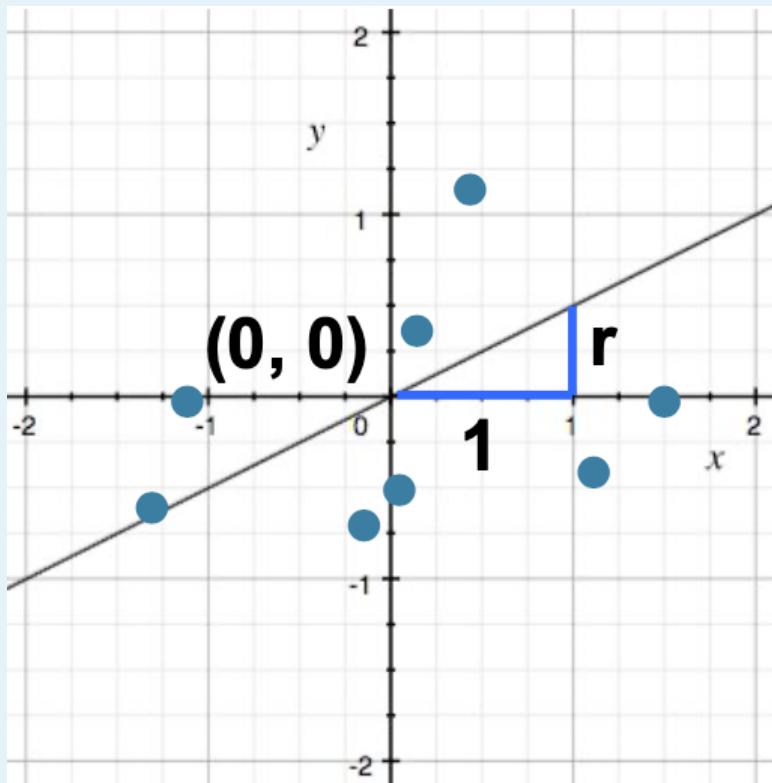
Lines can be expressed by *slope* & *intercept*

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

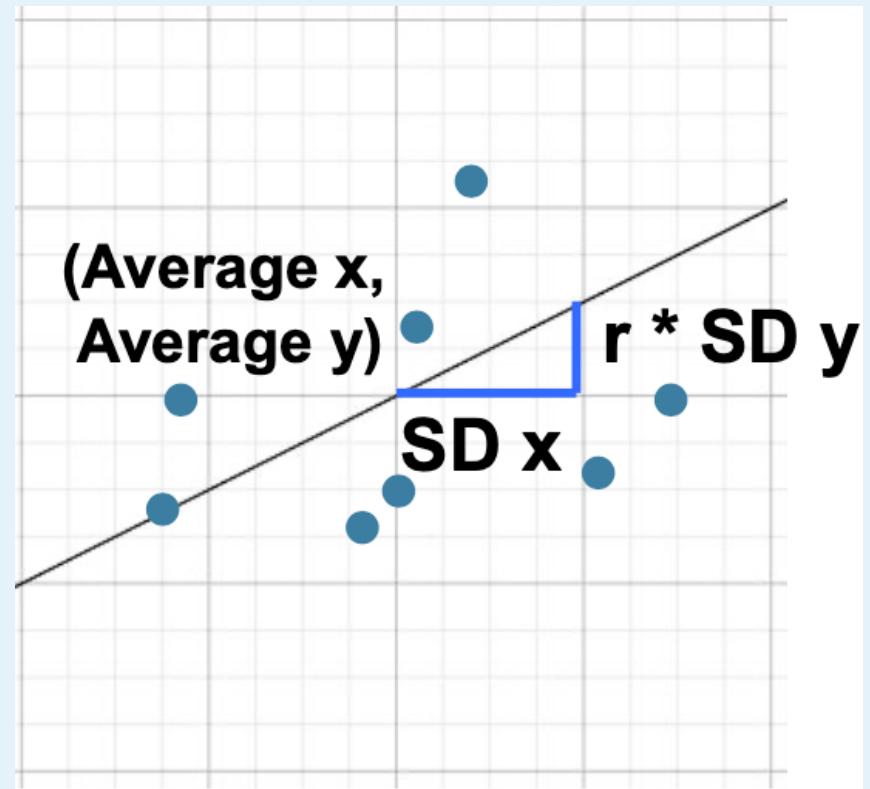


# Regression Line

## Standard Units



## Original Units





# Slope and Intercept

*estimate of  $y = slope * x + intercept$*

*slope of the regression line*

$$r * \frac{SD \text{ of } y}{SD \text{ of } x}$$

*intercept of the regression line*

$$\text{mean}(y) - \text{slope} \times \text{mean}(x)$$

# Prediction with Linear Regression



**Goal:** Predict  $y$  using  $x$

Examples:

- Predict # *hospital beds available* using *air pollution*
- Predict *house prices* using *house size*
- Predict # *app users* using # *app downloads*

# Regression Estimate



**Goal:** Predict  $y$  using  $x$

To find the regression estimate of  $y$ :

- Convert the given  $x$  to standard units
- Multiply by  $r$
- That's the regression estimate of  $y$ , but:
  - It's in standard units
  - So convert it back to the original units of  $y$



# Regression Line Equation

In original units, the regression line has this equation:

$$y_{su} = r \times x_{su}$$

$$\frac{\text{estimate of } y - \text{mean}(y)}{\text{SD of } y} = r \times \frac{\text{given } x - \text{mean}(x)}{\text{SD of } x}$$

Lines can be expressed by *slope & intercept*

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

What we want

What we observe

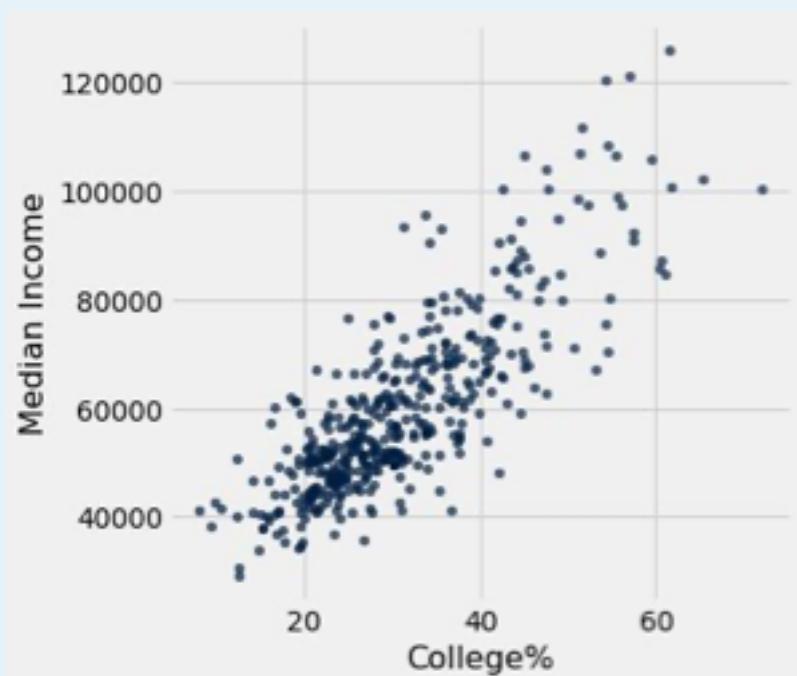


# Discussion Question

Based only on the graph,  
which must be true?

1. Going to college causes people to earn more.
2. For any district, having more college-educated people live there causes median incomes to rise.
3. For any district, having a higher median income causes more college-educated people to move there.

## USA Congressional Districts 2016



A woman with long, dark hair is walking across a city street. She is wearing a dark coat with a fur-trimmed hood and light-colored pants. A yellow taxi cab is visible in the background on the left. The scene has a blue-tinted, slightly blurred effect.

# Least Squares

---



# Error in Estimation

- **error = actual value – estimate**
- Typically, some errors are positive and some are negative
- To measure the rough size of the errors
  - **square the errors** to eliminate cancellation
  - Take the **mean** of the squared errors
  - Take the square **root** to fix the units
- **Root mean square error (rmse)**

# Least Squares Line



- Minimized the root mean squared error among all lines
- Equivalently, minimizes the mean squared error among all lines
- Names:
  - “Best fit” line
  - Least squares line
  - Regression line

# Numerical Optimization



- Numerical minimization is approximate but effective
- Lots of machine learning uses numerical minimization (demo)
- If the function **`mse(a, b)`** returns the mse of estimation using the line “estimate =  $ax + b$ ”,
  - then **`minimize(mse)`** returns array  $[a_0, b_0]$
  - $a_0$  is the slope and  $b_0$  the intercept of the line that *minimizes* the mse among lines with arbitrary slope  $a$  and arbitrary intercept  $b$  (that is, among all lines)

# Residuals



- Error in regression estimate
- One residual corresponding to each point  $(x, y)$
- **residual**
  - = **observed  $y$  - regression estimate of  $y$**
  - = observed  $y$  - height of regression line at  $x$
  - = vertical distance between the point and the best line



## A scatter diagram of residuals

- Should look like an unassociated blob for linear relations
- But will show patterns for non-linear relations
- Used to check whether linear regression is appropriate
- Look for curves, trends, changes in spread, outliers, or any other patterns



# Properties of residuals

- Residuals from a linear regression **always** have
  - Zero mean
    - (so rmse = SD of residuals)
  - **Zero** correlation with x
  - **Zero** correlation with the fitted values
- These are all true **no matter what the data look like**
  - Just like deviations from mean are zero on average