# Math 183 Statistical Methods

Eddie Aamari S.E.W. Assistant Professor

eaamari@ucsd.edu math.ucsd.edu/~eaamari/ AP&M 5880A

Today: Chapter 7 (beginning)

- Learn to describe a scatterplot using FDSO
- Get a visual and intuitive sense of the correlation coefficient  ${\cal R}$
- Residuals
- Least square method for regression

## More Inference for Two Numerical Variables

Two categorical variables:

- Display via contingency table
- Inference via  $\chi^2$

Two quantitative variables:

- Display via scatterplot
- Inference depends if data are about 1 population or 2 and what you're exploring.



Each row holds two pieces of data about the **same** person.

### Response and Explanatory Variables



A Scatterplot shows the relationship between two <u>numeric</u> variables.

## How to Discus Scatterplots: FDSO

Form:



Direction:



## How to Discus Scatterplots: FDSO

Strength: Very strong Strong Weak

Outliers:



This scatterplot has a weak to moderate, positive, linear association. We have an outlier.

As with a single quantitative variable, the notion of outlier is vague. Just think: "A point that stands far apart from the overall trend in the scatterplot."

### Your Turn! Describing Real Data Using FDSO

Age of husband and Age of wife



- Form: Linear
- Direction: Positive
- Strenght: Strong
- Outliers: None

### Your Turn! Describing Real Data Using FDSO

Engine size of a car (liters) and Fuel economy (MPG)



- Form: Curved
- Direction: Negative
- Strenght: Moderate
- Outliers: Maybe the top red plus...?

### Moving Beyond Vague Language

**Correlation:** A statistic that measures the Srength and Direction of a linear association (Form) between two quantitative variables where <u>no</u> Outliers are present.

Reported using either an uppercase R or a lowercase r:

R = r = -0.3.

<u>Note</u>: The programming language R is NOT named after the R statistic, but for the first names of its inventers (Ross Ihaka and Robert Gentleman) and as a reference to a related language S.

### Deriving Knowledge about R Using Examples

Looking at these examples, make observations about the R statistic.



### Facts About R, The Correlation Statistic

- We always have  $-1 \le R \le 1$ .
- The Direction is encoded in the sign:
  - + : positive association
  - : negative association
- The Srength is encoded in the value:

- 1 and -1 are possible: the data must perfectly lie on a line
- Choice of predictor/response doesn't matter:

$$\operatorname{cor}(X,Y)=\operatorname{cor}(Y,X)$$

- Correlation is a unitless idea
- Correlation is unaffected by linear scale changes

cor(X,Y) = cor(X,3.14Y) = cor(X,Y+100) = cor(2X-6.8,95Y+7)

# Cautions About R

Never use R to measure correlation for associations that are non-linear.



Never use R to measure correlation if outliers are present (even one!).



Correlation doesn't measure strength in non-linear associations

The effect of a single outlier can be huge.

#### How is R Calculated?



Numerator: Factors that accumulate positive and negative directions based on the individual points

Denominator: Factors that eliminate scaling for each axis, and ensure that  $-1 \le R \le 1$ .

(R is almost always calculated in a statistical program, not by hand)

### Your Turn!

On your cell phone, go to

http://guess the correlation.com/

and play for one minute.



#### Moving Beyond Correlation to Prediction

If we had to draw a straight line through our scatterplot that "fit the data the best", it might look like the red line.



A **linear model** is an equation the "best fits" the data. It is also as the "line of best fit", the "regression line", and the "least squares line".

### Observed Value, Predicted Value, Residual

Given any point (x, y) in a scatter plot, we now have two key ideas: for x fixed,

- y: the value observed from actual data point
- $\hat{y}$ : the **value predicted** from model

From these values, we can calculate the **residual**  $e = y - \hat{y}$ , which measure how off the model is at the value x.



#### What is the "Best Fit"?

Ideally, all the residuals would be 0 (a perfect model!).

Decause data are rarely perfect, we could define

model error = 
$$\sum$$
 (residuals for each data point)?

Then, the line of best fit makes this expression minimal.

Two choices to fix this issue are:

model error = 
$$\sum |\text{residuals}|$$
 or model error =  $\sum (\text{residuals})^2$ .

The second choice is standard, penalizes large residuals, and is easily differentiable.

### Exploring Linear Regression Visually

Try out this interactive GeoGebra tool:

https://www.geogebra.org/m/dlsxY1uX



- Move  $p_1$  and  $p_2$  until you *feel* you have the line of best fit. Check your guess by showing the line of best fit.
- Refresh the page. Try checking "Show Residuals" and "Residuals<sup>2</sup>", and then minimize the sum of squared residuals (red number). Then check your guess.