

Stat 202-2015S W6 - Friday

(Pg 1)

## Review - Chapter 2 - Relationships between pairs of variables

Associated - two variables measured on same cases are associated if knowing values of one tells you something about values of other.

dep  
indep Response Variable - measures outcome of a study  
Explanatory variable - explains or causes change in response variable

## Scatterplot

Oasis - strong positive  
BioClocks - weak positive  
Fidget - moderate negative

Form - eg line, also could be something else. A curve for example, sine wave.  
We will consider only lines

Strength - How closely do data follow form?

Direction Pos/neg depending on whether above average values of one occur with above average values of other (pos) or below average values of other (neg)

# Correlation

A number which quantifies the direction and strength of the linear relationship ~~at~~ between two variables.



To the extent that the form is a line this is useful. ~~Not necessarily~~ Not necessarily as useful if form is something else.

~~Direction~~

Direction :  $\begin{cases} \text{positive association} \\ r \text{ is pos} \\ \text{negative association} \\ r \text{ is negative} \end{cases}$

Strength:  $\begin{cases} \text{strong association} \\ r \text{ is close to } \pm 1 \\ \text{weak association} \\ r \text{ is close to } 0 \end{cases}$

Defined

$$r = \frac{1}{n-1} \sum \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

$s_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$       Standardized  $x$       Standardized  $y$       Z-score

Properties

- \*  $r$  does not depend on units  
change data with a linear transformation  
to change unit cm to inches or lightyears  
doesn't change correlation
- \*  $r$  has no units
- \*  $r > 0$  means pos association
- \*  $r < 0$  means negatve association
- ~~with  $r = 0$  means no~~
- \* no distinction between explanatory / response vars
- \* Both variables must be quantitative
- \*  $-1 \leq r \leq 1$
- \* values near zero indicate a weak linear ~~relationship~~  
relationship.
- \* values close to  $\pm 1$  indicate data lie  
close to a straight line
- \*  $r = \pm 1$  exactly means points lie exactly on a line
- \* measures only strength of linear relationship  
not curved relationship no matter how strong.

Oasis  
Bioclocks  
Fidget } guess

Regression

Regression vs Correlation

\* Correlation is a measure which quantifies the strength of the linear relationship between two variables.  $r$

\* Regression finds the best fitting line

$y = mx + b$  Finds  $m$  and  $b$

Statisticians use other notation

Slope and intercept

$y = b_0 + b_1 x$

$m$  or  $b_1$

Slope =  $r \frac{S_y}{S_x}$

intercept =  $\bar{y} - b_1 \bar{x}$

Show StatCrunch

\* Regression line passes through  $(\bar{x}, \bar{y})$

$y = (\bar{y} - b_1 \bar{x}) + b_1 x$   
plug  $x = \bar{x}$  get  $y = \bar{y}$

Correlation line on scatter plot intercept / slope

\* Has different slope depending on what is explanatory / response even when plotted on same axes, \* must decide which is explanatory which is response.