

# Correlation



Image source: <http://commons.wikimedia.org/wiki/File:Gnome-power-statistics.svg>, GPL

## Lecture 4

Survey Research & Design in Psychology  
James Neill, 2016  
Creative Commons Attribution 4.0

## Overview



1. Covariation
2. Purpose of correlation
3. Linear correlation
4. Types of correlation
5. Interpreting correlation
6. Assumptions / limitations
7. Dealing with several correlations

2

## Readings

### Howitt & Cramer (2011/2014)

- Ch 6/7: Relationships between two or more variables: Diagrams and tables
- Ch 7/8: Correlation coefficients: Pearson correlation and Spearman's rho
- Ch 10/11: Statistical significance for the correlation coefficient: A practical introduction to statistical inference
- Ch 14/15: Chi-square: Differences between samples of frequency data
- **Note:** Howitt and Cramer doesn't cover point bi-serial correlation

## Covariation

4

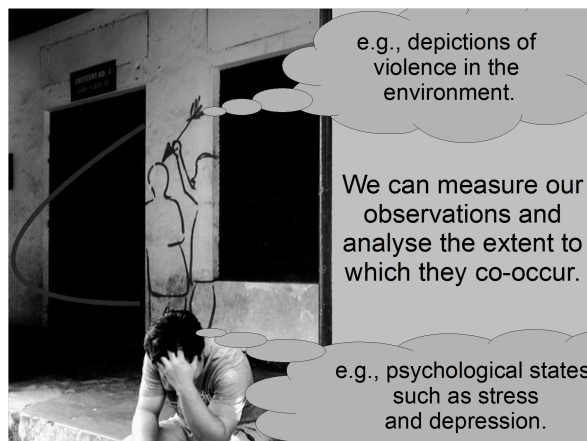
e.g., pollen and bees

e.g., study and grades

e.g., nutrients and growth

The world is made of  
co-variations

5



Co-variations are the basis of more complex models.

## Purpose of correlation

8

## Purpose of correlation

The underlying purpose of correlation is to help address the question:

What is the

- relationship or
- association or
- shared variance or
- co-relation

between **two variables**?

9

## Purpose of correlation

Other ways of expressing the underlying correlational question include:

To what extent do variables

- covary?
- depend on one another?
- explain one another?

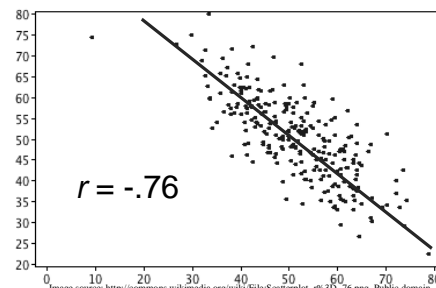
10

## Linear correlation

11

## Linear correlation

The extent to which two variables have a simple **linear** (straight-line) relationship.



12

## Linear correlation

Linear relations between variables are indicated by correlations':

- **Direction:** Sign (+ / -) indicates direction of relationship (+ve or -ve slope)
- **Strength:** Size indicates strength (values closer to -1 or +1 indicate greater strength)
- **Statistical significance:**  $p$  indicates likelihood that the observed relationship could have occurred by chance

13

## Types of relationships

- No relationship ( $r = 0$ ) (X and Y are independent)
- Linear relationship (X and Y are dependent)
  - As X ↑s, so does Y ( $r > 0$ )
  - As X ↑s, Y ↓s ( $r < 0$ )
- Non-linear relationship

14

## Types of correlation

To decide which type of correlation to use, consider the **levels of measurement** for each variable.

15

## Types of correlation

- Nominal by nominal: Phi ( $\Phi$ ) / Cramer's  $V$ , Chi-square
- Ordinal by ordinal: Spearman's rank / Kendall's Tau  $b$
- Dichotomous by interval/ratio: Point bi-serial  $r_{pb}$
- Interval/ratio by interval/ratio: Product-moment or Pearson's  $r$

16

## Types of correlation and LOM

	Nominal	Ordinal	Int/Ratio
Nominal	Clustered bar-chart Chi-square, Phi ( $\Phi$ ) or Cramer's $V$	← Recode	Clustered bar chart or scatterplot Point bi-serial correlation ( $r_{pb}$ )
Ordinal		Clustered bar chart or scatterplot Spearman's Rho or Kendall's Tau	← ↑ Recode
Int/Ratio			Scatterplot Product-moment correlation (17)

## Nominal by nominal

18

## Nominal by nominal correlational approaches

- Contingency (or cross-tab) tables
  - Observed
  - Expected
  - Row and/or column %s
  - Marginal totals
- Clustered bar chart
- Chi-square
- Phi ( $\phi$ ) / Cramer's V

19

## Contingency tables

- Bivariate frequency tables
- Marginal totals (blue)
- Cell frequencies (red)

	Disease		
	Diseased	Free	
Exposed	a	b	$n_1$
Not Exposed	c	d	$n_2$
	$m_1$	$m_2$	n

## Contingency table: Example

b2 Do you snore? \* b3r Smoker Crosstabulation

Count		b3r Smoker		Total
		0 No	1 Yes	
b2 Do you snore?	0 yes	50	16	66
	1 no	111	9	120
Total		161	25	186

BLUE = Marginal totals  
RED = Cell frequencies

21

## Contingency table: Example

b2 Do you snore? \* b3r Smoker Crosstabulation

		b3r Smoker		Total
		0 No	1 Yes	
b2 Do you snore?	0 yes	Count 50	16	66
		Expected Count 57.1	8.9	66.0
1 no	Count	111	9	120
	Expected Count	103.9	16.1	120.0
Total		Count 161	25	186
		Expected Count 161.0	25.0	186.0

- Expected counts are the cell frequencies for independent variables.
- Chi-square is based on the differences between the actual and expected cell counts.

22

b2 Do you snore? \* b3r Smoker Crosstabulation

% within b2. Do you snore?

		b3r Smoker		Total
		0 No	1 Yes	
b2 Do you snore?	0 yes	75.8%	24.2%	100.0%
	1 no	92.5%	7.5%	100.0%
Total		86.6%	13.4%	100.0%

Row and/or column cell percentages may also aid interpretation e.g., ~2/3rds of smokers snore, whereas only ~1/3<sup>rd</sup> of non-smokers snore.

b2 Do you snore? \* b3r Smoker Crosstabulation

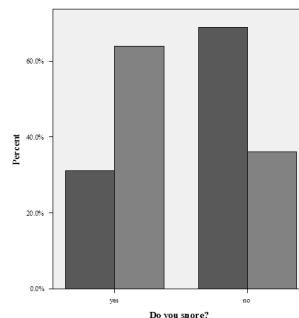
% within b3r Smoker

		b3r Smoker		Total
		0 No	1 Yes	
b2 Do you snore?	0 yes	31.1%	64.0%	35.5%
	1 no	68.9%	36.0%	64.5%
Total		100.0%	100.0%	100.0%

23

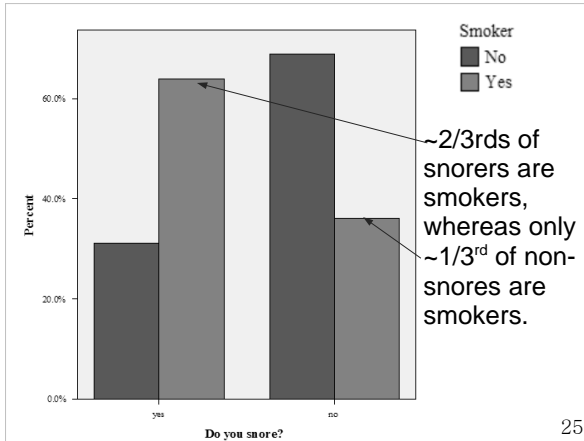
## Clustered bar graph

Bivariate bar graph of frequencies or percentages.



The category axis bars are clustered (by colour or fill pattern) to indicate the the second variable's categories.

24



25

## Pearson chi-square test

The value of the test-statistic is

$$X^2 = \sum \frac{(O - E)^2}{E}$$

where

$X^2$  = the test statistic that approaches a  $\chi^2$  distribution.

$O$  = frequencies observed;

$E$  = frequencies expected (asserted by the null hypothesis).

26

## Pearson chi-square test: Example

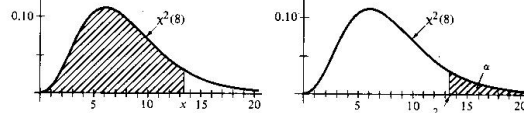
	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	10.259 <sup>a</sup>	1	.001
Continuity Correction <sup>a</sup>	8.870	1	.003
Likelihood Ratio	9.780	1	.002
Fisher's Exact Test			
Linear-by-Linear Association	10.204	1	.001
N of Valid Cases	186		

Write-up:  $\chi^2(1, 186) = 10.26, p = .001$

## Chi-square distribution: Example

The Chi-Square Distribution

The critical value for chi-square with 1 df and a critical alpha of .05 is 3.84



$$P(X \leq x) = \int_0^x \frac{1}{\Gamma(r/2)2^{r/2}} w^{r/2-1} e^{-w/2} dw$$

r	P(X ≤ x)							
	0.010	0.025	0.050	0.100	0.900	0.950	0.975	0.990
1	0.000	0.001	0.004	0.016	2.706	3.841	5.024	6.635
2	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210
3	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.34
4	0.297	0.484	0.711	1.064	7.779	9.488	11.14	13.28
5	0.554	0.831	1.145	1.610	9.236	11.07	12.83	15.09

## Phi (φ) & Cramer's V

(non-parametric measures of correlation)

### Phi (φ)

- Use for 2 x 2, 2 x 3, 3 x 2 analyses e.g., Gender (2) & Pass/Fail (2)

### Cramer's V

- Use for 3 x 3 or greater analyses e.g., Favourite Season (4) x Favourite Sense (5)

29

## Phi (φ) & Cramer's V: Example

### Symmetric Measures

	Value	Approx. Sig.
Nominal by Nominal	Phi	.235
	Cramer's V	.001
N of Valid Cases	186	

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

$\chi^2(1, 186) = 10.26, p = .001, \phi = .24$

Note that the sign is ignored here (because nominal coding is arbitrary, the researcher should explain the direction of the relationship)

## Ordinal by ordinal

31

## Ordinal by ordinal correlational approaches

- Spearman's rho ( $r_s$ )
- Kendall tau ( $\tau$ )
- Alternatively, use nominal by nominal techniques (i.e., recode the variables or treat them as having a lower level of measurement)

32

## Graphing ordinal by ordinal data

- Ordinal by ordinal data is difficult to visualise because its non-parametric, with many points.
- Consider using:
  - Non-parametric approaches (e.g., clustered bar chart)
  - Parametric approaches (e.g., scatterplot with line of best fit)

33

## Spearman's rho ( $r_s$ ) or Spearman's rank order correlation

- For ranked (ordinal) data
  - e.g., Olympic Placing correlated with World Ranking
- Uses product-moment correlation formula
- Interpretation is adjusted to consider the underlying ranked scales

34

## Kendall's Tau ( $\tau$ )

- Tau a
  - Does not take joint ranks into account
- Tau b
  - Takes joint ranks into account
  - For square tables
- Tau c
  - Takes joint ranks into account
  - For rectangular tables

35

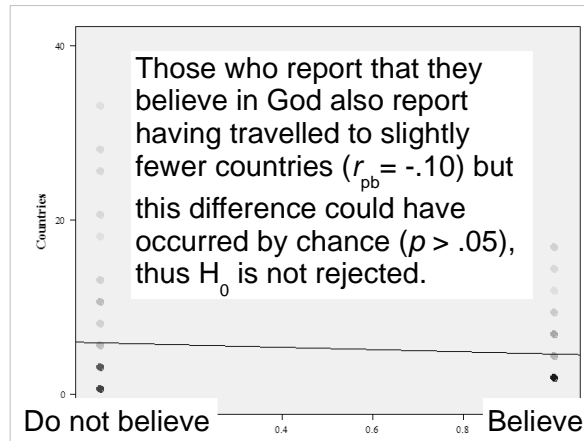
## Dichotomous by scale (interval/ratio)

36

## Point-biserial correlation ( $r_{pb}$ )

- One dichotomous & one continuous variable
  - e.g., belief in god (yes/no) and number of countries visited
- Calculate as for Pearson's product-moment  $r$
- Adjust interpretation to consider the direction of the dichotomous scales

37



## Point-biserial correlation ( $r_{pb}$ ): Example

Correlations

		b4r God	b8 Countries
b4r God 0 = No 1 = Yes	Pearson Correlation	1	-.095
	Sig. (2-tailed)		.288
	N	127	127
b8 Countries	Pearson Correlation	-.095	1
	Sig. (2-tailed)	.288	
	N	127	190

39

## Scale (interval/ratio) by Scale (interval/ratio)

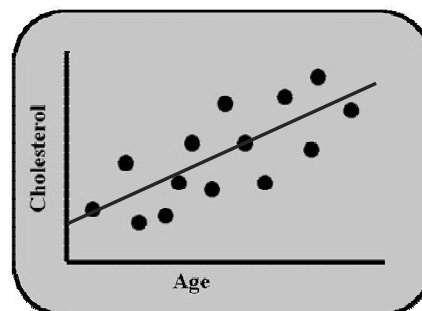
40

## Scatterplot

- Plot each pair of observations (X, Y)
  - x = predictor variable (independent; IV)
  - y = criterion variable (dependent; DV)
- By convention:
  - IV on the x (horizontal) axis
  - DV on the y (vertical) axis
- Direction of relationship:
  - +ve = trend from bottom left to top right
  - ve = trend from top left to bottom right

41

Scatterplot showing relationship between age & cholesterol with line of best fit



42

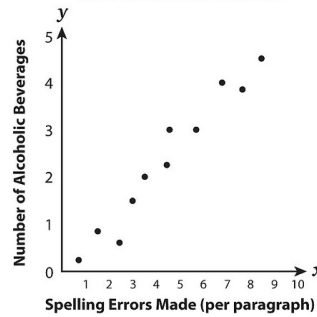
### Line of best fit

- The correlation between 2 variables is a measure of the degree to which pairs of numbers (points) cluster together around a best-fitting straight line
- Line of best fit:  $y = a + bx$
- Check for:
  - outliers
  - linearity

43

### What's wrong with this scatterplot?

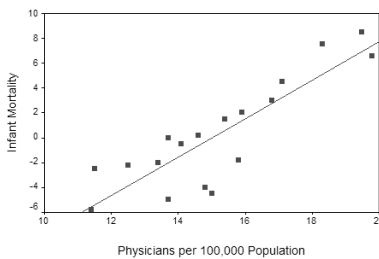
CORRELATION BETWEEN DRINKING AND SPELLING ERRORS



IV should be treated as X and DV as Y, although this is not always distinct.

44

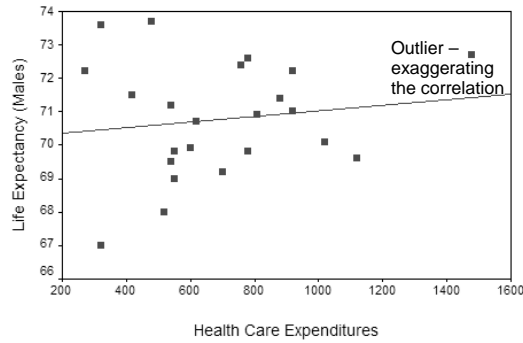
### Scatterplot example: Strong positive (.81)



Q: Why is infant mortality positively linearly associated with the number of physicians (with the effects of GDP removed)?

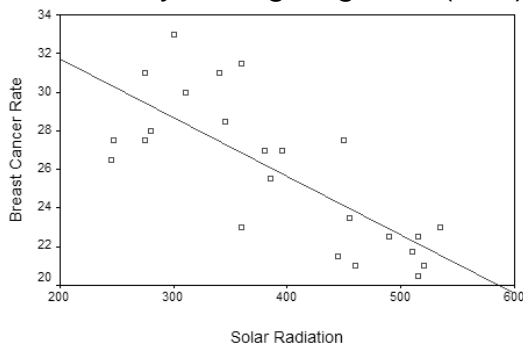
A: Because more doctors tend to be deployed to areas with infant mortality (socio-economic status aside).

### Scatterplot example: Weak positive (.14)



46

### Scatterplot example: Moderately strong negative (-.76)



47

### Pearson product-moment correlation ( $r$ )

- The product-moment correlation is the **standardised covariance**.

$$r_{X,Y} = \frac{\text{cov}(X, Y)}{S_X S_Y}$$

48



## Covariance

- Variance shared by 2 variables

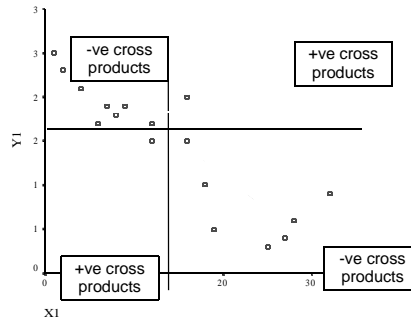
$$Cov_{XY} = \frac{\Sigma(X - \bar{X})(Y - \bar{Y})}{N - 1}$$

Cross products  
 $N - 1$  for the sample;  $N$  for the population

- Covariance reflects the direction of the relationship:
  - +ve cov indicates +ve relationship
  - ve cov indicates -ve relationship

49

## Covariance: Cross-products



50

## Covariance

- Size depends on the measurement scale → Can't compare covariance across different scales of measurement (e.g., age by weight in kilos versus age by weight in grams).
- Therefore, **standardise** covariance (divide by the cross-product of the SDs) → correlation
- Correlation is an effect size - i.e., standardised measure of strength of linear relationship

51

## Covariance, SD, and correlation: Example quiz question

The covariance between  $X$  and  $Y$  is 1.20. The  $SD$  of  $X$  is 2 and the  $SD$  of  $Y$  is 3. The correlation is:

- .20
- .30
- .40
- 1.20

Answer:  
 $1.20 / 2 \times 3 = .20$

52

## Hypothesis testing

Almost all correlations are not 0, therefore the question is:

“What is the **likelihood** that a relationship between variables is a ‘true’ relationship - or could it simply be a result of random sampling variability or ‘chance’?”

53

## Significance of correlation

- **Null hypothesis ( $H_0$ ):**  $\rho = 0$ : assumes that there is no ‘true’ relationship (in the population)
- **Alternative hypothesis ( $H_1$ ):**  $\rho \neq 0$ : assumes that the relationship is real (in the population)
- Initially assume  $H_0$  is true, and evaluate whether the data support  $H_1$ .
- $\rho$  (**rho**) = population product-moment correlation coefficient

54

## How to test the null hypothesis

- Select a critical value (alpha ( $\alpha$ )); commonly .05
- Can use a 1- or 2-tailed test
- Calculate correlation and its  $p$  value. Compare this to the critical value.
- If  $p <$  critical value, the correlation is statistically significant, i.e., that there is less than a  $x\%$  chance that the relationship being tested is due to random sampling variability.

55

## Correlation – SPSS output

Correlations		Cigarette Consumption per Adult per Day	CHD Mortality per 10,000
Cigarette Consumption per Adult per Day	Pearson Correlation		
	Sig. (2-tailed)		
	N		
CHD Mortality per 10,000	Pearson Correlation	.713*	
	Sig. (2-tailed)	.000	
	N	21	

\*\* . Correlation is significant at the 0.01 level (2-tailed).

56

## Imprecision in hypothesis testing

- **Type I error:** rejects  $H_0$  when  $H_0$  is true
- **Type II error:** accepts  $H_0$  when  $H_0$  is false
- A significance test result depends on the power of study, which is a function of:
  - Effect size ( $r$ )
  - Sample size ( $N$ )
  - Critical alpha level ( $\alpha_{crit}$ )

57

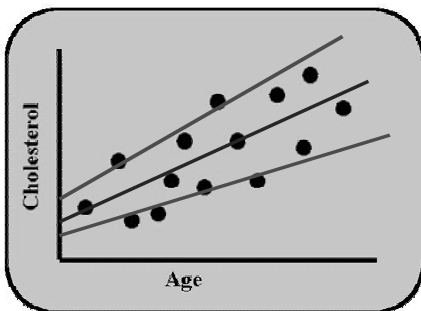
## Significance of correlation

$df$ ( $N - 2$ )	critical $p = .05$
5	.67
10	.50
15	.41
20	.36
25	.32
30	.30
50	.23
200	.11
500	.07
1000	.05

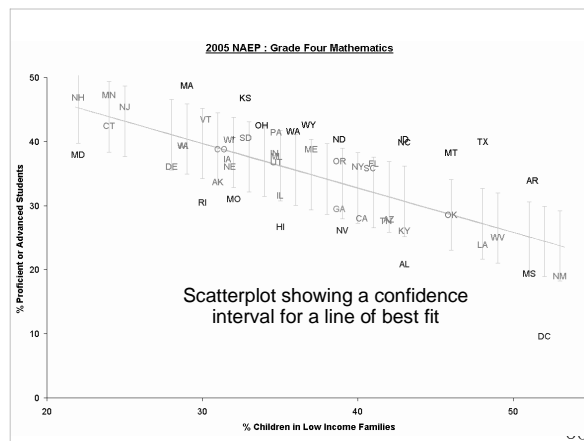
The size of correlation required to be significant decreases as  $N$  increases – why?

58

## Scatterplot showing a confidence interval for a line of best fit



59



### Practice quiz question: Significance of correlation

If the correlation between Age and test Performance is statistically significant, it means that:

- there is an important relationship between the variables
- the true correlation between the variables in the population is equal to 0
- the true correlation between the variables in the population is not equal to 0
- getting older causes you to do poorly on tests

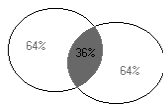
61

### Interpreting correlation

62

### Coefficient of Determination ( $r^2$ )

- CoD = The proportion of variance in one variable that can be accounted for by another variable.
- e.g.,  $r = .60$ ,  $r^2 = .36$



63

### Interpreting correlation (Cohen, 1988)

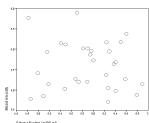
- A correlation is an **effect size**
- Rule of thumb:

<u>Strength</u>	<u>r</u>	<u>r<sup>2</sup></u>
Weak:	.1 - .3	1 - 10%
Moderate:	.3 - .5	10 - 25%
Strong:	>.5	> 25%

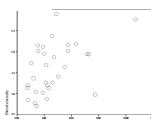
64

### Size of correlation (Cohen, 1988)

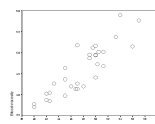
WEAK (.1 - .3)



MODERATE (.3 - .5)



STRONG (> .5)

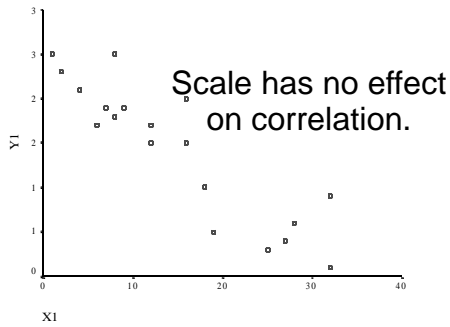


### Interpreting correlation (Evans, 1996)

<u>Strength</u>	<u>r</u>	<u>r<sup>2</sup></u>
very weak	0 - .19	(0 to 4%)
weak	.20 - .39	(4 to 16%)
moderate	.40 - .59	(16 to 36%)
strong	.60 - .79	(36% to 64%)
very strong	.80 - 1.00	(64% to 100%)

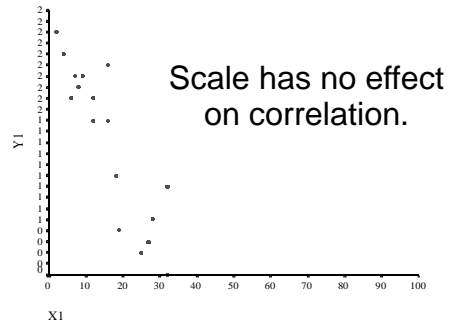
66

Correlation of this scatterplot =  $-.9$



67

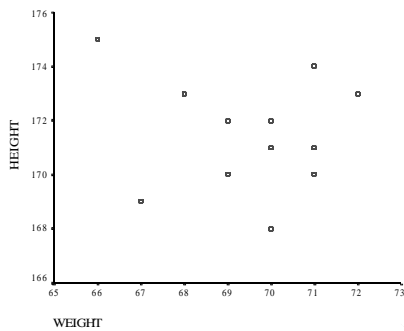
Correlation of this scatterplot =  $-.9$



68

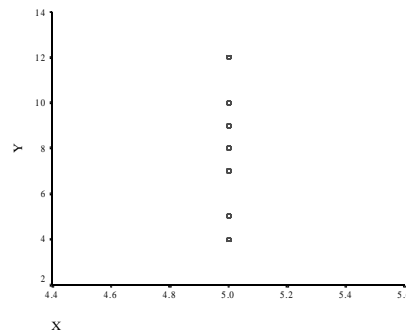
What do you estimate the correlation of this scatterplot of height and weight to be?

- a.  $-.5$
- b.  $-1$
- c.  $0$
- d.  $.5$
- e.  $1$



What do you estimate the correlation of this scatterplot to be?

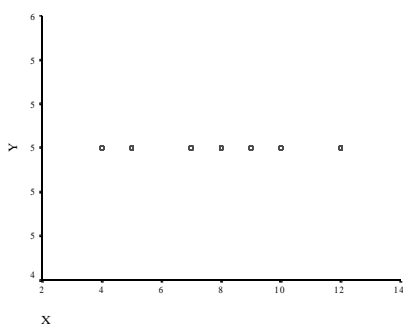
- a.  $-.5$
- b.  $-1$
- c.  $0$
- d.  $.5$
- e.  $1$



70

What do you estimate the correlation of this scatterplot to be?

- a.  $-.5$
- b.  $-1$
- c.  $0$
- d.  $.5$
- e.  $1$



1

### Write-up: Example

“Number of children and marital satisfaction were inversely related ( $r(48) = -.35, p < .05$ ), such that contentment in marriage tended to be lower for couples with more children. Number of children explained approximately 10% of the variance in marital satisfaction, a small-moderate effect.”

72

## Assumptions and limitations

(Pearson product-moment linear correlation)

73

## Assumptions and limitations

1. Levels of measurement
2. Normality
3. Linearity
  1. Effects of outliers
  2. Non-linearity
4. Homoscedasticity
5. No range restriction
6. Homogenous samples
7. Correlation is not causation

74

## Normality

- The X and Y data should be sampled from populations with normal distributions
- Do not overly rely on a single indicator of normality; use histograms, skewness and kurtosis (within -1 and +1)
- Inferential tests of normality (e.g., Shapiro-Wilks) are overly sensitive when sample is large

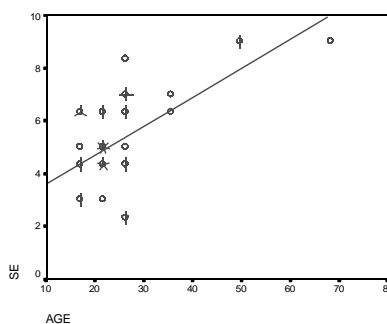
75

## Effect of outliers

- Outliers can disproportionately increase or decrease  $r$ .
- Options
  - compute  $r$  with & without outliers
  - get more data for outlying values
  - recode outliers as having more conservative scores
  - transformation
  - recode variable into lower level of measurement

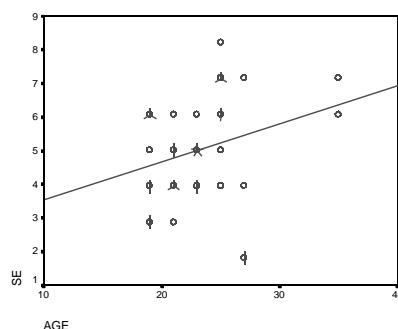
76

## Age & self-esteem ( $r = .63$ )



77

## Age & self-esteem (outliers removed) $r = .23$



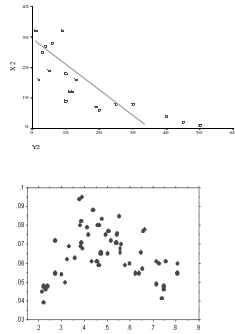
78

## Non-linear relationships

Check scatterplot

Can a linear relationship 'capture' the lion's share of the variance?

If so, use  $r$ .



79

## Non-linear relationships

If non-linear, consider

- Does a linear relation help?
- Transforming variables to 'create' linear relationship
- Use a non-linear mathematical function to describe the relationship between the variables

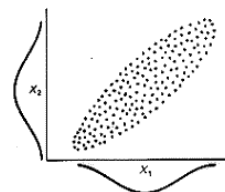
80

## Scedasticity

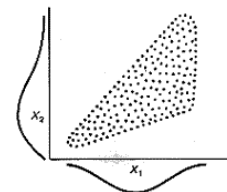
- **Homoscedasticity** refers to even spread about a line of best fit
- **Heteroscedasticity** refers to uneven spread about a line of best fit
- Assess visually and with Levene's test

81

## Scedasticity



Homoscedasticity with both variables normally distributed



Heteroscedasticity with skewness on one variable

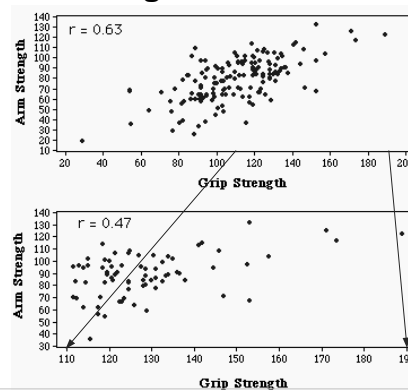
82

## Range restriction

- Range restriction is when the sample contains restricted (or truncated) range of scores
  - e.g., level of hormone X and age < 18 might have linear relationship
- If range restriction, be cautious in generalising beyond the range for which data is available
  - e.g., level of hormone X may not continue to increase linearly with age after age 18

83

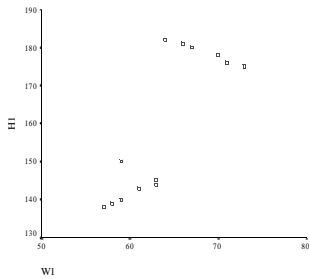
## Range restriction



84

## Heterogenous samples

- Sub-samples (e.g., males & females) may artificially increase or decrease overall  $r$ .
- Solution - calculate separately for sub-samples & overall; look for differences



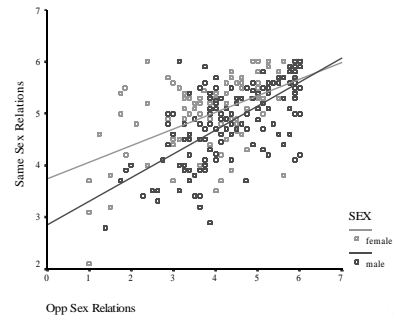
85

## Scatterplot of Same-sex & Opposite-sex Relations by Gender



$$r = .67$$

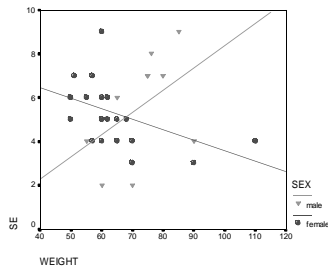
$$r = .52$$



## Scatterplot of Weight and Self-esteem by Gender

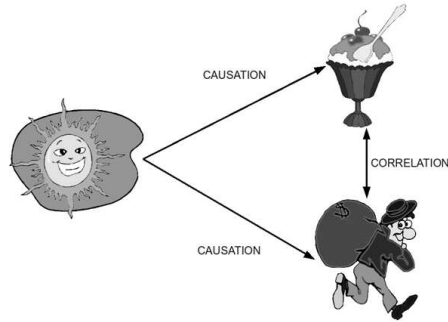
$$r = .50$$

$$r = -.48$$



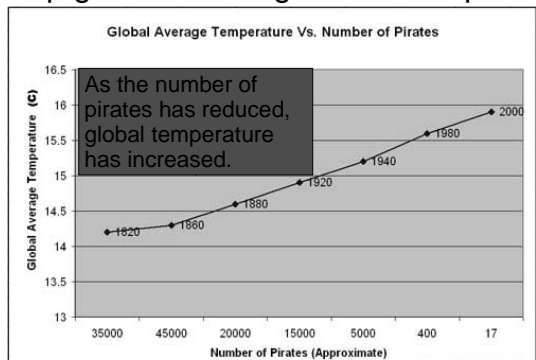
87

Correlation is not causation e.g.: correlation between ice cream consumption and crime, but actual cause is temperature



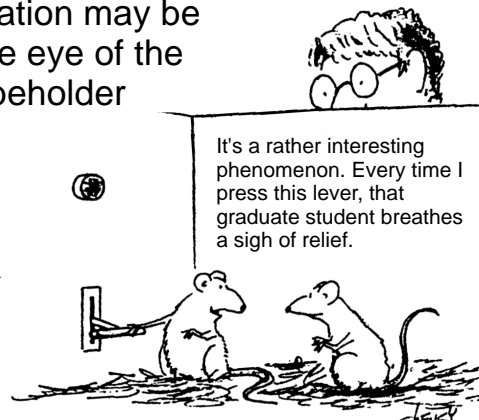
88

Correlation is not causation e.g.: Stop global warming: Become a pirate



89

Causation may be in the eye of the beholder



## Dealing with several correlations

91

## Dealing with several correlations

Scatterplot matrices organise scatterplots and correlations amongst several variables at once.

However, they are not sufficiently detailed for more than about five variables at a time.

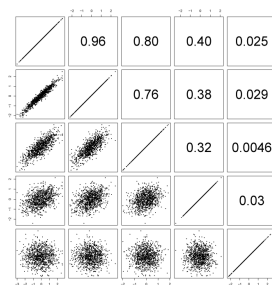


Image source: Unknown

92

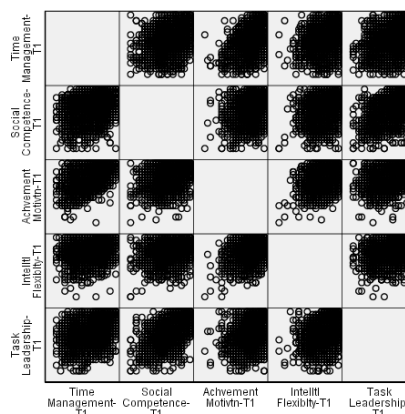
## Correlation matrix: Example of an APA Style Correlation Table

Table 1.  
*Correlations Between Five Life Effectiveness Factors for Adolescents and Adults (N = 3640)*

	Time Management	Social Competence	Achievement Motivation	Intellectual Flexibility	Task Leadership
Time Management		.36	.53	.31	.42
Social Competence			.37	.32	.57
Achievement Motivation				.42	.41
Intellectual Flexibility					.37
Task Leadership					

93

## Scatterplot matrix



94

## Summary

95

## Summary: Covariation

1. The world is made of covariations.
2. Covariations are the building blocks of more complex relationships which can be analysed through the use of:
  - factor analysis
  - reliability analysis
  - multiple regression

96



**Summary:  
Purpose of correlation**

1. Correlation is a standardised measure of the covariance (extent to which two phenomenon co-relate).
2. Correlation does not prove causation - may be opposite causality, bi-directional, or due to other variables.

97

**Summary:  
Types of correlation**

- Nominal by nominal:  
Phi ( $\Phi$ ) / Cramer's  $V$ , Chi-squared
- Ordinal by ordinal:  
Spearman's rank / Kendall's Tau  $b$
- Dichotomous by interval/ratio:  
Point bi-serial  $r_{pb}$
- Interval/ratio by interval/ratio:  
Product-moment or Pearson's  $r$

98

**Summary:  
Correlation steps**

1. Choose measure of correlation and graphs based on levels of measurement.
2. Check graphs (e.g., scatterplot):
  - Linear or non-linear?
  - Outliers?
  - Homoscedasticity?
  - Range restriction?
  - Sub-samples to consider?

99

**Summary:  
Correlation steps**

3. Consider
  - Effect size (e.g.,  $\Phi$ , Cramer's  $V$ ,  $r$ ,  $r^2$ )
  - Direction
  - Inferential test ( $p$ )
4. Interpret/Discuss
  - Relate back to hypothesis
  - Size, direction, significance
  - Limitations e.g.,
    - Heterogeneity (sub-samples)
    - Range restriction
    - Causality?

100

**Summary:  
Interpreting correlation**

- Coefficient of determination
  - Correlation squared
  - Indicates % of shared variance

<b>Strength</b>	<b><math>r</math></b>	<b><math>r^2</math></b>
Weak:	.1 - .3	1 - 10%
Moderate:	.3 - .5	10 - 25%
Strong:	> .5	> 25%

101

**Summary:  
Assumptions & limitations**

1. Levels of measurement
2. Normality
3. Linearity
4. Homoscedasticity
5. No range restriction
6. Homogenous samples
7. Correlation is not causation

102

## Summary: Dealing with several correlations

- Correlation matrix
- Scatterplot matrix

103

## References

Evans, J. D. (1996). *Straightforward statistics for the behavioral sciences*. Pacific Grove, CA: Brooks/Cole Publishing.

Howell, D. C. (2007). *Fundamental statistics for the behavioral sciences*. Belmont, CA: Wadsworth.

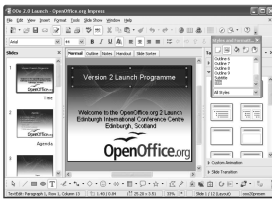
Howell, D. C. (2010). *Statistical methods for psychology* (7th ed.). Belmont, CA: Wadsworth.

Howitt, D. & Cramer, D. (2011). *Introduction to statistics in psychology* (5th ed.). Harlow, UK: Pearson.

104

## Open Office Impress

- This presentation was made using Open Office Impress.
- Free and open source software.
- <http://www.openoffice.org/product/impress.html>



105