



statstutor community project
encouraging academics to share statistics support resources
All stcp resources are released under a Creative Commons licence

Stcp-marshallowen-7

The Statistics Tutor's



www.statstutor.ac.uk

© Ellen Marshall, University of Sheffield

Reviewer: Jean Russell University of Sheffield

Quick Guide to Commonly Used Statistical Tests

Contents

CONTENTS.....	2
INTRODUCTION.....	3
TIPS FOR TUTORING.....	4
SECTION 1	
GENERAL INFORMATION.....	5
DATA TYPES.....	6
SUMMARY STATISTICS.....	6
<i>Summary of descriptive and graphical statistics.....</i>	<i>6</i>
DECIDING ON APPROPRIATE STATISTICAL METHODS FOR RESEARCH.....	7
ORDINAL DATA.....	7
WHICH TEST SHOULD I USE?.....	7
<i>Common Single Comparison Tests.....</i>	<i>7</i>
<i>Tests of association.....</i>	<i>7</i>
<i>One scale dependent and several independent variables.....</i>	<i>9</i>
ASSUMPTION OF NORMALITY.....	10
<i>Statistical tests for normality.....</i>	<i>10</i>
<i>Non-parametric tests.....</i>	<i>11</i>
OTHER COMMON ASSUMPTIONS.....	12
<i>For independent t-tests and ANOVA.....</i>	<i>12</i>
<i>For repeated measures ANOVA.....</i>	<i>12</i>
<i>Independent observations.....</i>	<i>12</i>
CONFIDENCE INTERVALS.....	13
HYPOTHESIS TESTING.....	14
MULTIPLE TESTING.....	14
SAMPLE SIZE AND HYPOTHESIS TESTS.....	15
EFFECT SIZE.....	15
SECTION 2	
THE MOST COMMON STATISTICAL TECHNIQUES USED.....	16
INDEPENDENT T-TEST.....	17
MANN-WHITNEY TEST	18
PAIRED T-TEST.....	19
WILCOXON SIGNED RANK TEST	20
ONE-WAY ANOVA	21
KRUSKAL-WALLIS TEST	22
ONE-WAY ANOVA WITH REPEATED MEASURES (WITHIN SUBJECTS).....	23
FRIEDMAN TEST	24
TWO-WAY ANOVA.....	25
CHI-SQUARED TEST	26
ODDS AND RELATIVE RISK.....	27
<i>Odds.....</i>	<i>27</i>
<i>Odds Ratio.....</i>	<i>27</i>
<i>Relative Risk (RR).....</i>	<i>27</i>
CORRELATION.....	28
PEARSON'S CORRELATION COEFFICIENT.....	28
RANKED CORRELATION COEFFICIENTS.....	30
<i>Spearman's Rank Correlation Coefficient</i>	<i>30</i>
<i>Kendall's Tau Rank Correlation Coefficient.....</i>	<i>30</i>
PARTIAL CORRELATION.....	30
REGRESSION.....	31
LINEAR REGRESSION.....	31
LOGISTIC REGRESSION.....	33
SECTION 3	
OTHER STATISTICAL TESTS AND TECHNIQUES.....	34

PROPORTIONS TEST (Z-TEST)	35
RELIABILITY.....	37
<i>Interrater reliability</i>	37
<i>Cohen's Kappa</i>	37
<i>Intraclass Correlation Coefficient</i>	38
<i>Cronbach's alpha (reliability of scales)</i>	39
PRINCIPAL COMPONENT ANALYSIS (PCA).....	41
CLUSTER ANALYSIS.....	44
HIERARCHICAL CLUSTERING.....	45
K-MEANS CLUSTERING.....	47
 <u>SECTION 4</u>	
<u>SUGGESTED RESOURCES</u>	48
BOOKS.....	49
WEBSITES.....	50
<u>CONTRIBUTERS TO THE GUIDE</u>	51

Introduction

This guide is designed to help you quickly find the information you need about a particular statistical test.

Section 1

Section 1 contains general information about statistics including key definitions and which summary statistics and tests to choose. Use the “Which test should I use?” table to allow the student to choose the test they think is most appropriate, talking them through any assumptions or vocabulary they are unfamiliar with.

Section 2

Section 2 takes you through the most common tests used and those that are usually as complex as the students require. As a statistics tutor, you should be familiar with all these techniques.

Section 3

Section 3 contains tests and techniques that are more complex or are used less frequently. This section is aimed at tutors who have studied statistics in detail before.

Tips for tutoring

We are here to help the students to learn and not to do their work for them or to tell them exactly how to do their work (although this is very tricky sometimes!). Facilitating the understanding and ability to choose an appropriate statistical test is a success, even if the analysis is not as thorough as if we had done it ourselves.

Avoid maths! Most students carrying out project analysis do not need to know the maths behind the technique they are using. You may love maths but a lot of students are maths phobic – even an x can frighten them!

You cannot know everything! Statistics is a vast subject so don't be afraid to say that you don't know. Ask others for help or look up information on the internet to help.

Consider the students ability when advising on the best technique. They have to write up the analysis and therefore need to understand what has been done. Carrying out simple analysis or even just a graph to summarise their results may be enough for their project.

Don't assume that the student knows anything about the technique they are suggesting! Students with no statistical knowledge at all can come in saying their supervisor wants them to carry out multivariate analysis. Get them to explain their project and why they think the technique is suitable if you think that they know very little. In general, start with descriptive statistics and simpler analysis if they have not done any analysis yet. Some students do know exactly why they are doing something and have investigated the topic fully so just help them with the complex technique if you can.

Section 1

General information

Data types

In order to choose suitable summary statistics and analysis for the data, it is also important for students to distinguish between continuous (numerical/ scale) measurements and categorical variables.

Summary Statistics

Students often go straight to the hypothesis test rather than investigating the data with summary statistics and charts first. Encourage them to summarise their data first. As well as summarising their results, charts especially can show outliers and patterns.

For continuous normally distributed data, summarise using means and standard deviations. If the data is skewed or there are influential outliers, the median (middle value) and interquartile range (Upper quartile – lower quartile) are more appropriate.

Asking for the mean, median, minimum, maximum and standard deviation along with producing an appropriate chart will identify outliers and skewed data. A big difference between the mean and median indicates skewed data or influential outliers.

Summary of descriptive and graphical statistics

Chart	Variable type	Purpose	Summary Statistics
Pie Chart or bar chart	One Categorical	Shows frequencies/ proportions/percentages	Class percentages
Stacked / multiple bar	Two categorical	Compares proportions within groups	Percentages within groups
Histogram	One scale	Shows distribution of results	Mean and Standard deviation
Scatter graph	Two scale	Shows relationship between two variables and helps detect outliers	Correlation co-efficient
Boxplot	One scale/ one categorical	Compares spread of values	Median and IQR
Line Chart	Scale by time	Displays changes over time Comparison of groups	Means by time point

Means plot	One scale/ 2 categorical	Looks at combined effect of two categorical variables on the mean of one scale variable	Means
-------------------	--------------------------	---	-------

Deciding on appropriate statistical methods for research

This is the information you need from a student to help them decide on the most appropriate statistical techniques for their project

What is the main research question? This needs to be able to be defined with specific variables in mind. Which variables (types of measurement) will help answer the research question?

Which is the dependent (outcome) variable and what type of variable is it?

Which are the independent (explanatory) variables, how many are there and what data types are they?

Are relationships or differences between means of interest?

Are there repeated measurements of the same variable for each subject?

Ordinal data

Some departments routinely use parametric tests to analyse ordinal data. As a general rule of thumb, ordinal variables with seven or more categories can be analysed with parametric tests if the data is approximately normally distributed. However, if the students department/supervisor expect scales of five to be analysed as continuous data, warn them why you think this is not appropriate but let them do it. We are here to advice and help them make decisions. Sometimes, questionnaires have sets of questions trying to measure an underlying latent variable. In that situation, summing or averaging the scores gives a variable which could be considered as scale so parametric tests can be carried out.

Which test should I use?

Common Single Comparison Tests

Comparing:	Dependent (outcome) variable	Independent (explanatory) variable	Parametric test (data is normally distributed)	Non-parametric test (ordinal/ skewed data)
The averages of two	Scale	Nominal	Independent t-	Mann-Whitney test/

INDEPENDENT groups		(Binary)	test	Wilcoxon rank sum
The averages of 3+ independent groups	Scale	Nominal	One-way ANOVA	Kruskal-Wallis test
The average difference between paired (matched) samples e.g. weight before and after a diet	Scale	Time/ Condition variable	Paired t-test	Wilcoxon signed rank test
The 3+ measurements on the same subject	Scale	Time/ condition variable	Repeated measures ANOVA	Friedman test

Tests of association

Relationship between 2 continuous variables	Scale	Scale	Pearson's Correlation Coefficient	Spearman's Correlation Coefficient
Predicting the value of one variable from the value of a predictor variable or looking for significant relationships	Scale	Any	Simple Linear Regression	Transform the data
	Nominal (Binary)	Any	Logistic regression	
Assessing the relationship between two categorical variables	Categorical	Categorical		Chi-squared test

One scale dependent and several independent variables

1 st independent	2 nd independent	Test
Scale	Scale/ binary	Multiple regression
Nominal (Independent groups)	Nominal (Independent groups)	2 way ANOVA
Nominal (repeated measures)	Nominal (repeated measures)	2 way repeated measures ANOVA
Nominal (Independent groups)	Nominal (repeated measures)	Mixed ANOVA
Nominal	Scale	ANCOVA

Regression or ANOVA? Use regression if you have only scale or binary independent variables. Categorical variables can be recoded to dummy binary variables but if there are a lot of categories, ANOVA is preferable.

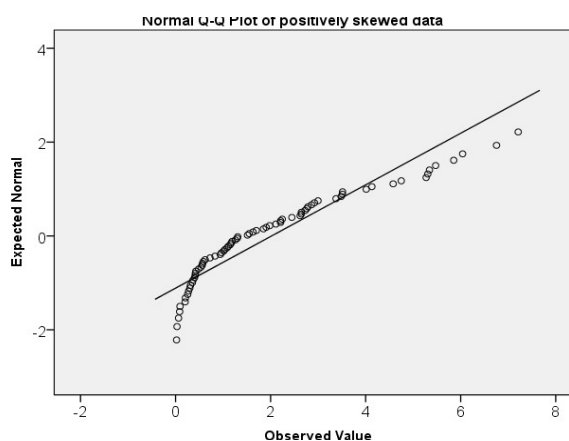
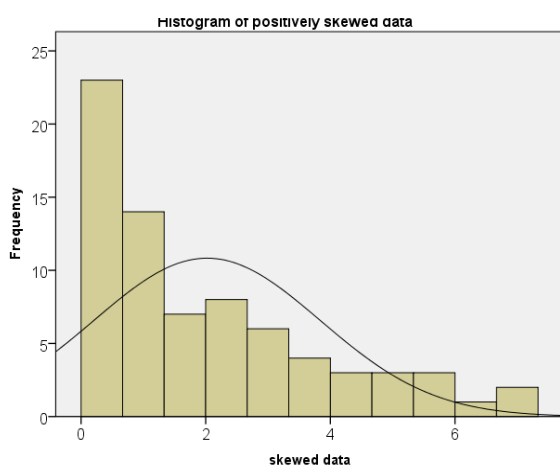
Assumption of normality

Parametric tests assume that the data follows a particular distribution e.g for t-tests, ANOVA and regression, the data needs to be normally distributed. Parametric tests are more powerful than non-parametric tests, when the assumptions about the distribution of the data are true. This means that they are more likely to detect true differences or relationships that exist.

The tests are quite robust to departures of non-normality so the data only needs to be approximately normally distributed.

Plotting a **histogram** or QQ plot of the variable of interest will give an indication of the shape of the distribution. Histograms should peak in the middle and be approximately symmetrical about the mean. If data is normally distributed, the points in QQ plots will be close to the line.

Below are some examples of very skewed data (i.e. non-normal).



Statistical tests for normality

There are statistical tests for normality such as the *Shapiro-Wilk* and *Kolmogorov-Smirnoff* tests but for small sample sizes ($n < 20$), the tests are unlikely to detect non-normality and for larger sample sizes ($n > 50$), the tests can be too sensitive. They are also sensitive to outliers so use histograms (large samples) or QQ plots (small samples).

Parametric test	What to check for normality
Independent t-test	Dependent variable by group
Paired t-test	Paired differences
One-way ANOVA	Residuals
Repeated measures ANOVA	Residuals at each time point
Pearson's correlation coefficient	Both variables are normally distributed

Non-parametric tests

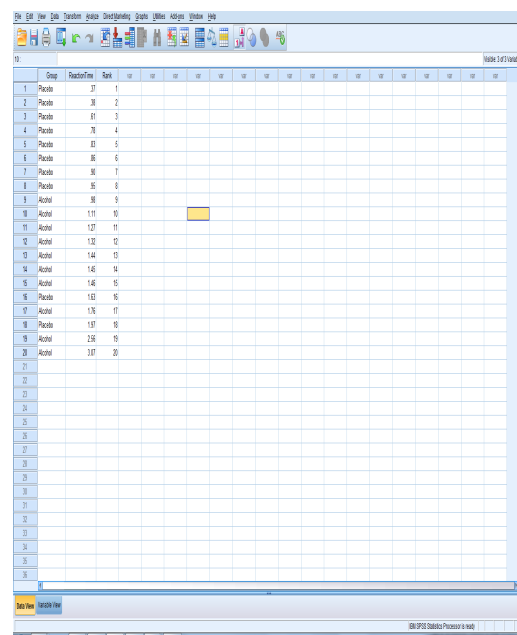
Non-parametric tests make no assumptions about the distribution of the data.

Nonparametric techniques are usually based on ranks or signs rather than the actual data and are usually less powerful than parametric tests.

The example to the right is data on reaction times after drinking either water or alcohol. The reaction times by group were not normally distributed so an independent t-test could not be used. The Mann-Whitney test is more appropriate. It tests the hypothesis that the two distributions are the same. All the data is ordered and ranked from the fastest to the slowest irrelevant of group. The sum of the ranks for each group is used to calculate a test statistic. If there is no difference between the groups the sum of the ranks will be similar. SPSS does all the ranking for you so you don't need to worry about that.

Non-parametric tests can also be used when other assumptions are not met e.g. equality of variance.

Some people also advise using non-parametric tests for small samples as it is difficult to assess normality.



	Group	ReactionTime	Rank
1	Pilsch	27	1
2	Pilsch	28	2
3	Pilsch	31	3
4	Pilsch	78	4
5	Pilsch	83	5
6	Pilsch	86	6
7	Pilsch	86	7
8	Pilsch	95	8
9	Pilsch	96	9
10	Alcohol	111	10
11	Alcohol	127	11
12	Alcohol	132	12
13	Alcohol	144	13
14	Alcohol	146	14
15	Alcohol	146	15
16	Pilsch	153	16
17	Alcohol	176	17
18	Pilsch	187	18
19	Alcohol	216	19
20	Alcohol	217	20
21			
22			
23			
24			
25			
26			
27			
28			
29			
30			
31			
32			
33			
34			
35			

Other common assumptions

For independent t-tests and ANOVA

Homogeneity of variances: Levene's test

Use: Used to test the equality of variances when comparing the means of independent groups e.g. Independent t-tests and ANOVA.

Note: The violation of this assumption is more serious than violation of the assumption of normality but both t-tests and ANOVA are fairly robust to deviations from this assumption. There are alternative tests within the t-test and ANOVA menus to deal with violations of this assumption.

Interpretation:

If the p-value is less than 0.05 reject H_0 and conclude that the assumption of equal variances has not been met.

For repeated measures ANOVA

Sphericity: Mauchly's test

Use: Tests for sphericity - a measure of whether variances of the differences between all repeated measures are all equal. If the assumption is not met, the F-statistic is positively biased leading to an increased risk of a type 1 error.

Interpretation:

Significant when p-value < 0.05 meaning there are significant differences between the variance of differences, i.e. condition of sphericity is not met. If the assumption is not met, use the Greenhouse-Geisser correction to the degrees of freedom which appears in the standard output.

Independent observations

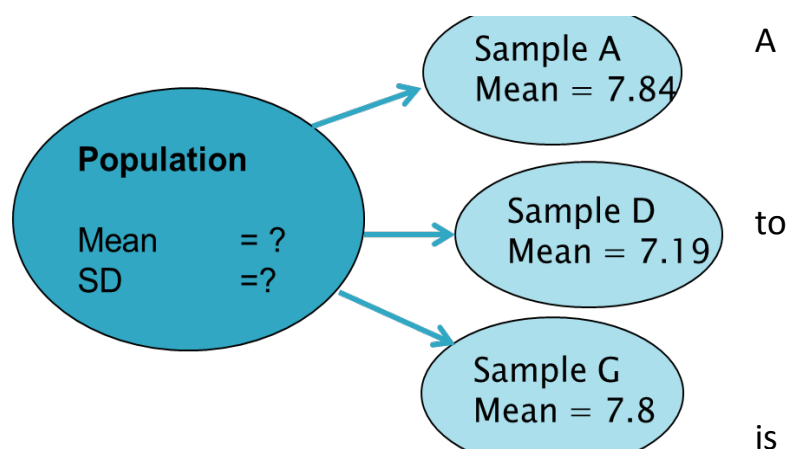
For most tests, it is assumed that the observations are independent. That is the results for one subject* are not affected by another. Examples of data which is not independent are repeated measures on the same subject (use the specific tests for this type of experiment) and observations over time (check the Durbin Watson test for regression). Another situation where observations are not independent is when subjects are nested within groups with a common influence e.g. children within classes who may be influenced by the teacher (use multilevel modelling to include class as an extra RANDOM factor). Time series analysis (which allows for non-independent measurements over time) and multilevel modelling are beyond the scope of most students.

*The subject is the unit of interest which could be a person, an observation, a day etc.

Confidence intervals

Most research uses sample data to make inferences about the wider population. population is the group of individuals you are interested in e.g. a study into the weight of babies born in Sheffield would use a sample but the results apply to the whole population.

Every sample taken from a population will contain different babies so the mean value varies especially if the sample size is small.



Confidence Intervals describe the variability surrounding the sample point estimate (the wider the interval, the less confident we can be about the estimate of the population mean). In general, all things being equal, the larger the sample size the better (more precise) the estimate is, as less variation between sample means is expected.

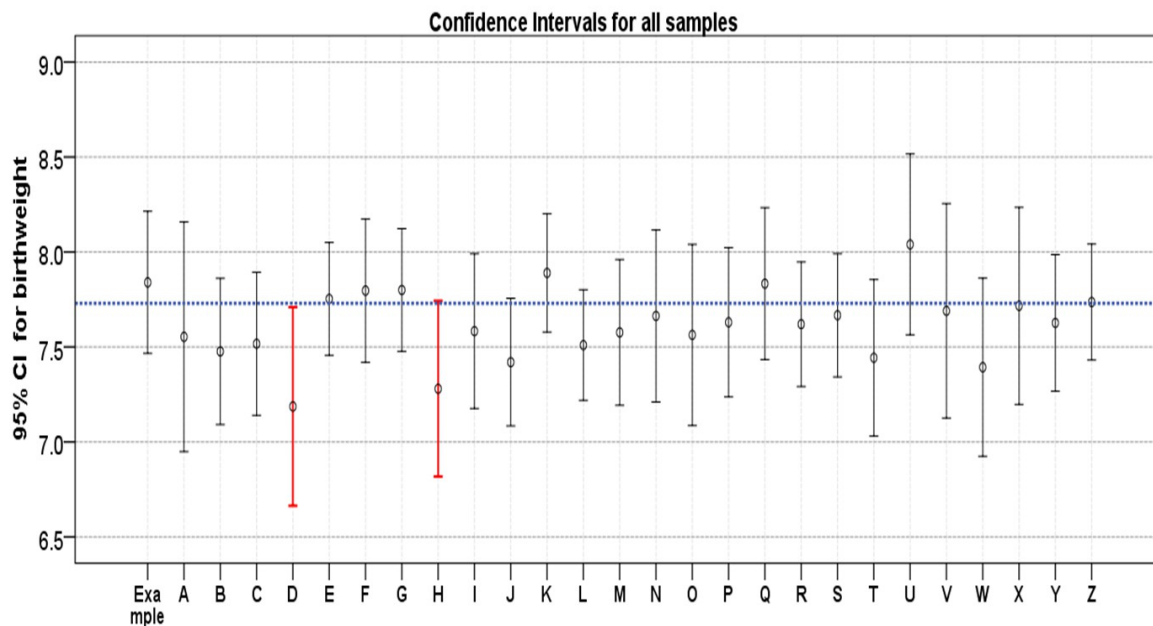
The equation for a 95% Confidence Interval for the population mean when the population standard deviation is unknown and the sample size is large (over 30) is

For example, sample D of 30 babies born in 2013 had a mean weight of 7.19lbs with a standard deviation of 1.4, the 95% Confidence Interval for the population mean of all babies is:

We would expect the population mean to be between 6.7 lbs and 7.7 lbs.

Confidence intervals give a range of values within which we are confident (in terms of probability) that the true value of a population parameter lies. A 95% CI is interpreted as 95% of the time the CI would contain the true value of the population parameter.

The diagram below shows the confidence intervals for 27 samples of babies taken from the same population. The actual population mean (which is not normally known) is 7.73 lbs. Two of the confidence intervals do not contain the population mean (don't overlap 7.73 lbs) including the one previously calculated.



The website CAST has some great applets for demonstrating concepts to students. This ebook contains core material including an applet for demonstrating confidence intervals
http://cast.massey.ac.nz/collection_public.html

There is a strong relationship between hypothesis testing and confidence intervals. For example, when carrying out a paired t-test, if the $p\text{-value} < 0.05$, the 95% confidence interval for the paired differences will not contain 0. However, a $p\text{-value}$ just concludes whether there is significant evidence of a difference or not. The confidence interval of the difference gives an indication of the size of the difference.

For more information on the use of confidence intervals see
<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1339793/>.

Hypothesis testing

Hypothesis testing is an **objective** method of making decisions or **inferences** from sample data (evidence). Sample data is used to choose between two choices i.e. **hypotheses** or statements about a population. Typically this is carried out by comparing what we have observed to what we expected if one of the statements (**Null Hypothesis**) was true.

Key terms:

NULL HYPOTHESIS (H_0) is a statement about the population & sample data used to decide whether to reject that statement or not. Typically the statement is that there is no difference between groups or association between variables.

ALTERNATIVE HYPOTHESIS (H_1) is often the research question and varies depending on whether the test is one or two tailed.

SIGNIFICANCE LEVEL: The probability of rejecting the null hypothesis when it is true, (also known as a type 1 error). This is decided by the individual but is normally set at 5% (0.05) which means that there is a 1 in 20 chance of rejecting the null hypothesis when it is true.

TEST STATISTIC is a value calculated from a sample to decide whether to accept or reject the null (H_0) and varies between tests. The test statistic compares differences between the samples or between observed and expected values when the null hypothesis is true.

P-VALUE: the probability of obtaining a test statistic at least as extreme as ours if the null is true and there really is no difference or association in the population of interest. P-values are calculated using different probability distributions depending on the test. A significant result is when the p-value is less than the chosen level of significance (usually 0.05).

The court case



Hypothesis testing can be thought of as a court case

Members of a jury have to decide whether a person is guilty or innocent based on evidence presented to them.

Null: The person is innocent

Alternative: The person is not innocent.

The null can only be rejected if there is enough evidence to disprove it and the jury do not know whether the person is really guilty or innocent so they may make a mistake.

If a court case was a hypothesis test, the jury consider the likelihood of innocence given the evidence and if there's less than a 5% chance that the person is innocent they reject the statement of innocence.

In reality, the person is actually Guilty (null false) **or** Innocent (null true) but we can only conclude that there is evidence to suggest that the null is false or not enough evidence to suggest it is false.

	The null hypothesis is actually:	
	False (i.e. there actually is a difference in the population)	True (i.e. there actually is no difference in the population)
You decide to:		
Reject the null hypothesis (i.e. conclude it is false and that there is a difference)	Correct ✓ POWER	False positive / type I error / α ✗
Not reject the null hypothesis (i.e. conclude it is not false and that there is no difference)	False negative / type II error / β ✗	Correct ✓

A type I error is equivalent to convicting an innocent person and is usually set at 5% (the magic 0.05!).

Multiple testing

Some students will try to perform a large number of tests on their data. The chance of a type I error increases with the number of tests. Adjustments to keep the type I error low for a larger number of tests are included as post hoc tests in ANOVA. This will mean less of the results are statistically significant. The most commonly used post hoc tests are Tukey and Sidak although Scheffe's is often used in medicine.

Suggest that the student looks in their notes or papers in their field when choosing. If adjustments need to be made by hand, the Bonferroni adjustment is the easiest to explain although it is the most conservative (least likely to lead to rejection of the null). Either divide the significance level initially used (probably 0.05) by the number of tests being carried out and compare the p-value with the new, smaller significance level. Alternatively, multiply the p-value by the number of tests being carried out and compare to 0.05. This is the standard adjustment made after the Kruskal-Wallis and has a maximum limit of 1 (as it's a probability!).

Sample size and hypothesis tests

The larger the sample size, the more likely a significant result is so for small sample sizes a huge difference is needed to conclude a significant difference. For large sample sizes, small differences may be significant but check if the difference is meaningful.

Effect size

An effect size is a measure of the strength or magnitude of the effect of an independent variable on a dependent variable which helps assess whether a statistically significant result is meaningful.

For example, **for a t-test**, the absolute effect size is just the difference between the two groups. A standardised effect size involves variability and can then be compared to industry standards.

Cohen gives the following guidance for the effect size d , although it is not always meaningful: 0.2 to 0.3 might be a "small" effect, around 0.5 a "medium" effect and 0.8 to infinity, a "large" effect.

Partial eta-squared

Partial eta-squared is a measure of variance. It represents the proportion of variance in the dependent variable that is explained by the independent variable. It also represents the effect size statistic. The effects sizes given in Cohen (1988) for the interpretation of the absolute effect sizes are:

$\eta^2 = 0.010$ is a small association.

$\eta^2 = 0.059$ is a medium association.

$\eta^2 = 0.138$ or larger is a large association.

Section 2

The most common statistical techniques used

Independent t-test

Dependent variable: Continuous

Independent variable: Binary (Group)

Use: A t-test is used to compare the means of two independent groups. Independent groups means that different people are in each group.

Plot: Box-plots (exploratory) or Confidence Interval plots with results

Assumptions	How to check	What to do if assumption is not met
Normality: dependent variables should be normally distributed within each group	Histograms of dependent variables per group / Shapiro Wilk	Mann-Whitney / Wilcoxon rank sum
Homogeneity of variance	Levene's test * (part of standard SPSS output)	Use bottom row of t-test output in SPSS

*Levenes test: If the assumption of homogeneity is not met, correct for this violation by not using the pooled estimate for the error term for the t-statistic and also by making adjustments to the degrees of freedom using the Welch-Satterthwaite method. SPSS does the automatically in the "Equal variances not assumed" row. Alternatively, a Mann-Whitney test can be carried out.

Interpretation:

If the p-value < 0.05, there is a significant difference between the means of the two groups. Report the means of the two groups or the mean difference and confidence interval from the SPSS output to describe the difference.

SPSS: Analyse > Compare means > Independent-samples T-test

Mann-Whitney test

(non-parametric equivalent to the independent t-test)

Dependent variable: Ordinal/ Continuous

Independent variable: Binary(Group)

The Mann-Whitney test is also known as the Mann–Whitney–Wilcoxon (MWW), Wilcoxon rank-sum test, or Wilcoxon–Mann–Whitney test.

Use: It is used to compare whether two groups containing different people are the same or not. The Mann-Whitney test ranks all of the data and then compares the sum of the ranks for each group to determine whether the groups are the same or not. There are two types of Mann-Whitney U tests. If the distribution of scores for both groups have the same shape, the medians can be compared. If not, use the default test which compares the mean ranks.

Plot: Histograms of the two groups

SPSS: Analyse > Nonparametric Tests > Independent Samples

1:
ts
g Values

☐ Automatically choose the tests based on the data
☒ Customize tests

Compare Distributions across Groups

☒ Mann-Whitney U (2 samples)
☐ Kruskal-Wallis 1-way ANOVA (k samples)
Multiple comparisons: All pairwise

☐ Kolmogorov-Smirnov (2 samples)
☐ Test for ordered alternatives (Jonckheere-Terpstra for k samples)
Hypothesis order: Smallest to largest

☐ Test sequence for randomness (Wald-Wolfowitz for 2 samples)
Multiple comparisons: All pairwise

Compare Ranges across Groups

☒ Moses extreme reaction (2 samples)
☒ Compute outliers from sample
☐ Custom number of outliers
Outliers: 1

Compare Medians across Groups

☒ Median test (k samples)
☒ Pooled sample median
☐ Custom
Median: 0
Multiple comparisons: All pairwise

Estimate Confidence Interval across Groups

☒ Hodges-Lehman estimate (2 samples)

Paired t-test

Dependent variable: Continuous (at least interval)

Independent variable: Time point 1 or 2/ condition

Use: A paired samples t-test can only be used when the data is paired or matched. Either there are before/after measurements of the same variable or the t-test can be used to compare how a group of subjects perform under two different test conditions. The test assesses whether the mean of the paired differences is zero.

Plot: Histogram of differences

Assumptions	How to check	What to do if assumption is not met
Normality: paired differences* should be normally distributed	Histogram of differences / Shapiro Wilk	Wilcoxon signed rank

Interpretation:

If the p-values < 0.05 then there is a statistically significant difference between the two time points/experiments. Report the mean difference.

SPSS: Analyse \hat{c} Compare means \hat{c} Paired-samples T-test

*Paired differences can be calculated using *Transform \hat{c} Compute variable*

Wilcoxon Signed Rank test

(non-parametric equivalent to the paired t-test)

Dependent variable: Ordinal/ Continuous

Independent variable: Time/ Condition (binary)

Use: The Wilcoxon signed rank test is used to compare two related samples, matched samples or repeated measurements on a single sample to assess whether their population mean ranks differ. It is a paired difference test and is the non-parametric alternative to the paired t-test. The absolute differences are ranked then the signs of the actual differences used to add the negative and positive ranks.

Plot: Histogram of differences

Interpretation:

If $p\text{-value} < 0.05$ then there is evidence that the population mean ranks differ. Report the medians of the two sets of measurements

SPSS: Analyse > Nonparametric tests > Legacy dialogs > 2 related samples

Other related tests:

Sign: Compares the number of negative and positive differences

McNemar: Can be used for binary nominal variables when changes in a subjects score are of interest. It compares the number of subjects who have changed their score in a positive direction with those changing their score in a negative direction.

One-way ANOVA

Dependent variable: Continuous

Independent variable: Categorical (at least 3 categories)

Use: Used to detect the difference in means of 3 or more independent groups. It can be thought of as an extension of the t-test for 3 or more independent groups. ANOVA uses the ratio of the between group variance to the within group variance to decide whether there are statistically significant differences between the groups or not.

Plot: Box-plots or confidence interval plots

Assumptions	How to check	What to do if assumption is not met
Residuals should be normally distributed	Histogram/ QQ plot of residuals / SW	Kruskall-Wallis test (non-parametric)
Homogeneity of variance	Levene's test / Bartlett's test	Welch test instead of ANOVA (adjusted for the differences in variance) and Games-Howell post hoc or Kruskall-Wallis

Interpretation:

ANOVA tests " H_0 : all group means are equal" using an F-test. The p-value concludes whether or not there is at least one pairwise difference. If the p-value < 0.05 , we reject H_0 and conclude that there is a significant difference between at least one pair of means. Post-hoc tests are used to test where the pairwise differences are. Report the significant pairwise differences and the means.

Post-hoc adjustments:

Tukey or Scheffe are most commonly used but check if their department uses something else. If treatments are being tested against a control use Dunnett. If group sample sizes vary use Hochberg's GT2. If there is a difference between the group variances (Levene's test gives a p-value < 0.05 so we reject H_0) use Games-Howell.

SPSS: Analyse \hat{c} General Linear model \hat{c} Univariate (Welch test unavailable)

Or Compare means \hat{c} ANOVA although dependent variable by group will have to be checked for normality as there's no option for calculating residuals.

Kruskal-Wallis test

(non-parametric equivalent to the one-way ANOVA)

Dependent variable: Ordinal/ Continuous

Independent variable: Categorical

Use: Kruskal-Wallis compares the medians of two or more samples to determine if the samples have come from different populations. It is an extension of the Mann–Whitney U test to 3 or more groups. The distributions do not have to be normal and the variances do not have to be equal.

Plot: Box-plots

Assumptions	How to check	What to do if assumption is not met
Independent observations	Check data	Friedman
Similar sample sizes	Check data	
>5 data points per sample (ideally)	Frequencies of group	

Interpretation:

When the Kruskal-Wallis test leads to significant results, then at least one of the samples is different from the other samples. The test does not identify where the differences occur or how many differences actually occur. The Mann-Whitney U test would help analyse the specific sample pairs for significant differences. Make sure a p-value correction for multiple testing is used in the post-hoc tests.

SPSS: Analyse > Nonparametric tests > Independent samples

Note: Double click on the output for the test and a second screen appears with a lot more information including the post-hoc tests with Bonferroni adjustments. Select 'Pairwise comparisons' from the list in the bottom right hand corner if the main test is significant. Use the 'adjusted sig.' column which is the Mann-Whitney p-value multiplied by the number of pairwise tests (Bonferroni correction).

One-way ANOVA with repeated measures (within subjects)

Dependent variable: Continuous

Independent variable: categorical with “levels” as the within subject factor

Use: Tests the equality of means in 3 or more groups. All sample members characteristics must be measured under multiple conditions i.e. the dependent variable is repeated. Standard ANOVA cannot be used as the assumption of independence has been violated. This is the equivalent of a one-way ANOVA but for repeated samples and is an extension of a paired-samples t-test. It is used to analyse (1) changes in mean score over 3 or more time points (2) differences in mean score under 3 or more conditions. It separates the variance from measures and from people, hence decreasing the mean squared error.

Plot: Means plot over time. For instructions see Brunel ASK video.

https://www.youtube.com/watch?v=duUEb_j9wfY&list=UUdb6U06idJlt7IWNr5YzAdA

Assumptions	How to check	What to do if assumption is not met
Check normality of residuals by time point	Histograms of residuals etc	Friedman test (non-parametric)
Sphericity: variances of the differences between each pair of repeated measures are all equal	Mauchly's test	a Greenhouse-Geisser correction or the Huynh-Feldt correction to the df

Interpretation:

If the main ANOVA is significant, there is a difference between at least two time points. The Bonferroni post hoc tests will conclude where those differences are. Report the significant post hoc results and means at each time point.

Post-hoc adjustments:

Of the three post hoc adjustments in 'Options', Bonferroni is most commonly used.

SPSS Analyse Ć General Linear Model Ć Repeated measures

Note: The factor e.g. time needs to be specified before the main analysis screen appears. The repeated measures are in different columns and are entered in the main screen.

Friedman test

(non-parametric equivalent to repeated measures ANOVA)

Dependent variable: Ordinal/ Continuous measured on at least 3 occasions or 3 measures under different conditions

Independent variable: Time/ Condition

Use: The Friedman test is used to detect differences in scores across multiple occasions or conditions. The scores for each subject are ranked and then the sums of the ranks for each condition are used to calculate a test statistic. The Friedman test can also be used when subjects have ranked a list e.g. rank these pictures in order of preference.

Interpretation:

If the Friedman test is significant ($p\text{-value} < 0.05$) then there are differences in the distributions across the time points/ conditions.

Post-hoc tests:

To examine where the differences actually occur, separate Wilcoxon signed-rank tests on the different combinations of related groups are run with the Bonferroni adjustment. The adjusted p-value column is the Wilcoxon p-value multiplied by the number of tests.

SPSS: Analyse > Non-parametric tests > Related samples

Note: Double click on the output for the test and a second screen appears with a lot more information including the post-hoc tests with Bonferroni adjustments. Select 'Pairwise comparisons' from the list in the bottom right hand corner.

Two-way ANOVA

Dependent variable: Continuous

Independent variables: Two categorical (2+ levels within each)

Use: Comparing means for combinations of two independent categorical variables (factors).

There are three sets of hypothesis with a two-way ANOVA. H_0 for each set is as follows:

- ⤵ The population means of the first factor are equal – equivalent to a one-way ANOVA for the row factor.
- ⤵ The population means of the second factor are equal – equivalent to a one-way ANOVA for the column factor.
- ⤵ There is no interaction between the two factors – equivalent to performing a test for independence with contingency tables (a chi-squared test for independence).

Plot: Means plot to look at interaction between the two independent variables. Use the lines plot option but ask for the mean rather than frequencies of the dependent variable.

Assumptions	How to check	What to do if assumption is not met
Residuals should be normally distributed	Histogram/ QQ plot of residuals / SW	Transform the data
Homogeneity of variance	Levene's test / Bartlett's test	Compare p-values with a smaller significance level e.g. 0.05

Interpretation:

When interpreting the results you need to return to the hypotheses and address each one in turn.

Post-hoc adjustments:

Tukey or Scheffe are generally used. For testing treatments against a control use Dunnett, if group sample sizes vary use Hochberg's GT2 and if there is a difference between the group variances (Levene's test gives a p-value < 0.05) use Games-Howell.

SPSS: Analyse > General Linear Model > Univariate

Chi-squared test

(non-parametric)

Dependent variable: Categorical

Independent variable: Categorical

Use: The null hypothesis is that there is no relationship/association between the two categorical variables. The chi-squared test compares expected frequencies, assuming the null is true, with the observed frequencies from the study. When obtaining a significant chi-squared result, calculate percentages in a table to summarise where the differences between the groups are.

Plot: Stacked/ multiple bar chart with percentages

Assumptions	How to check	What to do if assumption is not met
80% of expected cell counts >5	SPSS tells you under the test	Fisher's exact (usually for 2x2 tables, but can also be used for others) or merge categories where sensible
No cells with expected frequency below 1		

Interpretation:

If $p < 0.05$, there is significant evidence of relationship between the two variables. Use %'s to describe what the relationship is.

For 2 x 2 tables, use Yates continuity correction.

For comparing 2 paired proportions e.g. proportions of people changing their response given some information about the topic, use McNemar's test.

SPSS: Analyse > Descriptive > Crosstabs > Statistics

Odds and Relative Risk

Odds

The odds of an event happening is defined as follows where p_1 is the probability of an event happening:

Odds Ratio

The odds ratio is a measure of effect size, describing the strength of association or non-independence between two binary data values. If the probabilities of the event happening in each of the groups are (Group 1) and (Group 2), then the odds ratio is:

An odds ratio of 1 implies that the event is equally likely in both groups. An odds ratio greater than one implies that the event is more likely to occur in the first group, while less than one implies that it has more chance to happen in group 2.

Relative Risk (RR)

Relative risk is a similar and more direct measure of comparing the probabilities in two groups. As with the odds ratio, a relative risk equal to 1 means that the event is equally probable in both groups. A relative risk larger than 1 implies that the chance of the event occurring is higher in group 1 and smaller for group 2.

Summary

	Number of times the event happened	Number of times the event did not happen
Group 1	a	b
Group 2	c	d

Correlation

Use: Correlation (r) is used to measure the strength of association between two variables and ranges between -1 (perfect negative correlation) to 1 (perfect positive correlation). Cohen (1992) has the following interpretation of the absolute value of the correlation:

Correlation coefficient value	Association
-0.3 to +0.3	Weak
-0.5 to -0.3 or 0.3 to 0.5	Moderate
-0.9 to -0.5 or 0.5 to 0.9	Strong
-1.0 to -0.9 or 0.9 to 1.0	Very strong

Cohen, L. (1992). *Power Primer. Psychological Bulletin*, 112(1) 155-159

Plot: Scatterplot

SPSS: Analyse > Correlate > Bivariate Correlation

Pearson's correlation coefficient

Dependent variable: Continuous

Independent variable: Continuous

Pearson's correlation coefficient is the most common measure of correlation.

ρ (ρ) = population correlation and r = sample correlation

Assumptions	How to check	What to do if assumption is not met
Continuous data for each variable	Check data	If ordinal data use Spearman's or Kendall tau
Linearly related variables	Scatter plot	Transform data
Both variables are normally distributed	Histograms of variables/ Shapiro Wilk	Use rank correlation: Spearman's or Kendall tau

Ranked correlation coefficients

Dependent variable: Continuous/ Ordinal

Independent variable: Continuous/ Ordinal

Spearman's Rank Correlation Coefficient

Spearman's rank correlation coefficient is a non-parametric statistical measure of the strength of a monotonic relationship between paired data. The notation used for the sample correlation is r_s .

Assumptions	How to check	What to do if assumption is not met
Linearly related variables	Scatter plot	Transform data

Kendall's Tau Rank Correlation Coefficient

Use for small data sets with a large number of tied ranks

Use: Kendall's tau rank correlation coefficient is used to measure the association between two measured quantities. A tau test is a non-parametric hypothesis test for statistical dependence based on the tau coefficient. Specifically, it is a measure of rank correlation, i.e. the similarity of the orderings of the data when ranked by each of the quantities.

The Tau-b statistic makes adjustments for ties; values of Tau-b range from -1 (100% negative association or perfect inversion) to $+1$ (100% positive association or perfect agreement). A value of zero indicates the absence of association.

Partial correlation

Partial correlation allows for a third continuous or binary variable to be controlled for. The correlation then measures the association between the independent and dependent variables after removing the variation due to the control.

Regression

There are many different types of regression. The type of regression that is used depends on the type of dependent variable e.g. linear regression is used with a continuous dependent variable, logistic regression with a binary dependent variable and Poisson regression with a Poisson counts dependent variable.

Use: Regression can be used with many continuous and binary independent variables (x). It gives a numerical explanation of how variables relate, enables prediction of the dependent variable (y) given the independent variable (x) and can be used to control for confounding factors when describing a relationship between two variables. Regression produces a line of best fit by minimising the RSS (residual sum of squares) and tests that the slope is 0 for each independent variable. Note: A residual is the difference between an observed y and that predicted by the model.

Linear Regression

Dependent variable: Continuous

Independent variables: Any but categorical must be turned into binary dummy variables

Assumptions	How to check	What to do if assumption is not met
Independent observations (no correlation between successive values)	Durbin Watson = 1.5 – 2.5	Time series – beyond the scope of the tutor and the student!
Residuals should be normally distributed	Histogram/ QQ plot of residuals / SW	Transform the dependent variable
The relationship between the independent and dependent variables is linear	Scatterplot of independent and dependent variables	Transform either the independent or dependent variable
Homoscedasticity: The variance of the residuals about predicted responses should be the same for all predicted responses	Scatterplot of predicted values against residuals	Transform the dependent variable
No observations have a large overall influence (leverage)	Look at Cook's and Leverage distances	Remove observation with very high leverage

Interpretation

ANOVA table: The ANOVA table decides whether the model as a whole is significant. The model is compared to a 'null' model where every observation is predicted to be the same.

Coefficients table: The 'B' column in the coefficients table, gives us the values of the slope and intercept terms for the regression line. For multiple regression, (where there are several predictor variables), the coefficients table shows the significance of each variable individually after controlling for the other variables in the model.

Model summary: The R^2 value shows the proportion of the variation in the dependent variable which is explained by the model. It varies from 0 to 1 but is usually reported as a percentage. The better the model, the higher the R^2 value. The level for a 'good model' varies by discipline but above 70% is generally considered to be good for prediction.

TIP: The majority of students coming to the centre will just be using regression to look for significant relationships so don't confuse them with model selection or comparing models unless they ask about it.

Comparing regression models using R^2

R^2 increases for each additional variable added so the adjusted R^2 is better for comparing models where the dependent variable is the same. The adjusted R^2 takes into account the number of degrees of freedom of the model. When comparing regression models where the dependent variable was transformed or which used different sets of observations, R^2 is not a reliable guide to model quality. There are options for adding variables in blocks in SPSS which enables comparison of models in the output. The change in R^2 is tested formally.

Model selection

There are several methods for model selection (Forwards, Backwards and stepwise) available within SPSS which result in a model only containing significant variables.

Dummy variables

A dummy variable is a binary variable representing one category of a categorical variable e.g. for marital status, code as separate variables Married yes/ no, Divorced yes/ no etc.

Interactions

Interactions need to be created by multiplying the two variables of interest using *Transform*

→ *Compute variable.*

Logistic Regression

Dependent variable: Binary

Independent variables: Any (Use the 'Categorical' option to specify categorical variables and SPSS creates the dummy variables)

Use: One of the most commonly used tests for categorical variables is the Chi-squared test which looks at whether or not there is a relationship between two categorical variables but this doesn't make an allowance for the potential influence of other explanatory variables. For continuous outcome variables, multiple regression can be used for

- controlling for other explanatory variables when assessing relationships between a dependent variable and several independent variables

- predicting outcomes of a dependent variable using a linear combination of explanatory (independent) variables

Logistic regression does the same but the outcome variable is categorical. It leads to a model which can predict the probability of the event happening for an individual.

Note: This is a parametric test that does not require the data to be normally distributed.

Interpretation:

If probabilities (p) of the event of interest happening for individuals are needed, the logistic regression equation can be written as:

Initially look for those variables where the p -value is less than 0.05. These are the factors that have a statistically significant effect on the prediction of the dependent variable. When interpreting the differences, it is easier to look at the $\exp(\beta_i)$ which represents the odds ratio for the individual variable.

Section 3

Other statistical tests and techniques

Proportions test (z-test)

Use: The proportions test is used to test whether the proportions of two populations differ significantly with respect to one characteristic.

The test:

The null hypothesis is $p_1 - p_2 = 0$.

The Z-test statistic is calculated as follows:

Interpretation:

If the resulting p-value is significant (p-value < 0.05) then there is evidence of a statistically significant difference between the two proportions.

Assumptions	How to check	What to do if assumption is not met
At least 10 observations in each group	Frequencies	-

One-sample

SPSS> Analyze > Non Parametric > Chi-square (but weight cases)

Two-samples

SPSS> Analyze > Crosstabs > Options – Choose z-test

Reliability

Reliability can be divided into two main sections. The reliability between raters (interrater agreement) and the reliability of a set of questions when measuring an underlying variable (Cronbach's alpha).

Interrater reliability

The technique for assessing agreement between raters/ instruments depends on the type of variable being compared. For nominal data, Cohen's kappa should be used and for scale data, the Intraclass Correlation Coefficient (ICC) should be used. Data should be entered with one column for each rater and one row for each subject being rated. Intrarater reliability is when measurements from the same person are being compared.

Cohen's Kappa

Use: Assessing agreement between raters for categorical variables. Kappa compares the proportion of actual agreement between raters (P_A) to the proportion expected to agree by chance (P_C). The P_C values use expected values calculated as in the Chi-squared test for association (row total x column total/grand total).

Interpretation:

There is a test for agreement which tests the hypothesis that agreement is 0 but like correlation, the interpretation of the coefficient itself is more important. The following guidelines were devised by Landis and Koch (1977).

Cohen's kappa	Strength of agreement
< 0	Poor (Agreement worse than by chance)
0 - 0.2	Slight
0.21 - 0.4	Fair
0.41 - 0.6	Moderate
0.61 - 0.8	Good
0.81 - 1	Very good

SPSS: Analyze > Descriptive Statistics > Crosstabs

Select 'Kappa' from the statistics options

Note: Kappa looks at exact matches only so for ordinal data, weighted Kappa is preferable but this is not an option in SPSS. For ordinal data, if exact matches are required, use Kappa or consider using the ICC for scale data for close matches.

Intraclass Correlation Coefficient

Use: A measure of agreement of continuous measurements for two or more raters. There are several options for the ICC in SPSS. To choose the right combination of model, type and form, you need to ask the student the following questions:

Do all the raters rate all the subjects?

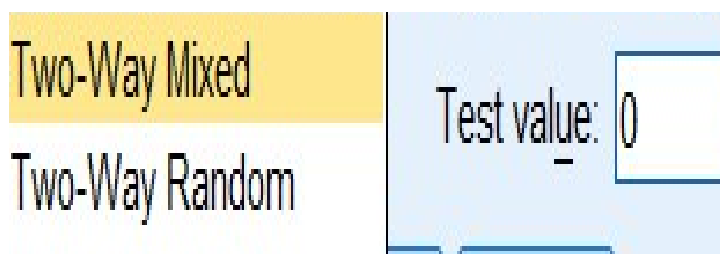
Are the raters the only possible raters (whole population) or a sample of raters?

Do raters need to match exactly or just be consistent (so raters may rank the subjects the same even if their scores don't match)?

Would you normally use just one rater or take an average of several raters?

SPSS: *Analyse* > *Scale* > *Reliability Analysis*

In the 'Statistics' options, select 'Intraclass correlation coefficient'



Models:

One way random	Not all raters have scored all the subjects
Two way random	Random selection of raters and subjects (default)
Two way mixed	Analysis contains whole population of raters

Type:

Absolute agreement	Exact matches on scores/ measurements are required
Consistency	Raters are consistent with their scoring e.g. rater A consistently scores lower than rater B

Form (appears in output): Single measures is used if when not testing for reliability, a subjects score is from one rater only. If a score is usually based on an average of several (k) raters, use average measures.

Interpretation:

Interpret in the same way as Cohen's kappa.

Cronbach's alpha (reliability of scales)

Researchers often use sets of likert style questions to measure an underlying latent variable that cannot be measured exactly. The set of questions is called a scale and the individual questions are called items. The scores for the items can be added or averaged to give an overall score for the scale. It is important that the items are all measuring the same underlying variable.

Use: Cronbach's alpha is a measure of internal consistency (how closely related the items are as a group). If the questions relate to the same issue, participants will be expected to get similar scores on each question.

This measure is not robust against missing data.

Interpretation:

Cronbach's alpha ranges from 0 to 1 and scores are expected to be between 0.7 and 0.9. Below is a commonly accepted rule of thumb for interpreting Cronbach's alpha.

Cronbach's alpha	Internal consistency
$\alpha \geq 0.9$	Very high consistency (the items are so similar that some may not be needed)
$0.8 \leq \alpha < 0.9$	Good
$0.7 \leq \alpha < 0.8$	Acceptable
$\alpha < 0.7$	Poor internal consistency

SPSS: Analyse > Scale > Reliability Analysis.

Make sure 'alpha' is selected as the Model in the dialogue box that appears.

Multivariate techniques

Multivariate techniques generally have more than one dependent variable but Multiple Regression and Discriminant Analysis are also referred to as multivariate despite only having one dependent variable. Multivariate techniques tend to involve classification or data reduction. The following table contains some of the more commonly used techniques.

Purpose	Dependent variable (type)	Independent variable (type)	Analysis
Data reduction	2+ (Scale/ binary although ordinal often used)		Principal Components Analysis
			Factor Analysis
	2+ (categorical)		Correspondence Analysis
Identify groups of similar subjects	2+ (Any)		Cluster Analysis
Compare groups	2+ (Scale)	1+ (Categorical)	MANOVA
		1+ (Categorical) & 1+ (Scale)	MANCOVA
Predict group membership	1 (Categorical)	Any	Linear Discriminant Analysis

Principal Component Analysis (PCA)

Use: Principal component analysis aims to reduce the number of inter-correlated variables to a smaller set which explains the overall variability almost as well. It produces new variables which are linear combinations of the original variables called Principal Components (PC's) or factors with the first PC explaining the most variation. These new variables can be used in further analysis e.g. regression. PCA can be based on either the correlation or covariance matrix. If variables are not on the same scale, the variable with the largest variance will dominate the first PC. If the variables are measured on a similar scale, use the covariance matrix, otherwise use the correlation matrix (default in SPSS).

Principal components analysis is not usually used to identify underlying latent variables but if interpretation of which variables contribute most to each PC is of interest, choose a method of rotation of the loadings (correlations between individual variables and the PC's) to identify clearer patterns. The varimax (variance maximising) rotation of the loadings maximises the variability of the new PC whilst minimising the variance around the new variable. It assumes that the PC's are not related.

Dependents: Scale/ binary (preferably mostly scale or mostly binary)

Numerical ordinal variables are often included but the choice of numbering e.g. 1, 2, 3 will impact on the results

Assumptions	How to check
Normality of variables	Histograms
Minimum sample size	Although PCA can be carried out on any number of cases, 5 – 10 cases per variable are often suggested for reliable results. Smaller sample sizes are reasonable if high loadings are achieved (0.8+) but the last few PC's are non-informative and shouldn't be used.
Variables should be adequately related	Correlation matrix shows coefficients above 0.3. Additional checks: KMO > 0.6, Bartlett's $p < 0.05$
No significant outliers	Component scores more than 3 SD's from mean

SPSS: Analyse \hat{c} Dimension Reduction \hat{c} Factor

Non-default options to select:

Menu	What to select
Descriptives	Coefficients, KMO and Bartlett's
Extraction	Scree plot (Note: Correlation matrix is the default so select 'Covariance matrix' if variables are measured on a similar scale)
Rotation	Varimax (most commonly used), Rotated solution, Loadings plot
Factor Scores	Save as variables (if PC scores for each case are required), factor score coefficient matrix
Options	Sorted by size, Suppress small coefficients (set min = 0.3)

Interpretation:

Are the variables adequately related?

Correlation matrix: If $r > 0.9$ for two variables, they are too related and only one is needed. If one variable consistently has correlations under 0.1, it is not related enough to the other variables and is likely to form a PC of its own.

KMO (Kaiser-Meyer-Olkin Measure of Sampling Adequacy): Varies from 0 to 1 but the closer to 1 the better. Should be above 0.6 to use PCA.

Bartlett's test: Bartlett's test for sphericity tests the hypothesis that correlation matrix is an identity matrix (correlations between variables are 0). The test should be significant for PCA to be appropriate.

How many PC's/ factors should be used?

By default SPSS only retains PC's with eigenvalues above 1 (Kaiser Criterion) but this method sometimes selects more PC's than needed. The *Screeplot* is an alternative method for choosing the optimal number of PC's. It plots the eigenvalues for each PC. When the line plateaus, no more PC's are needed. Parallel analysis is an alternative method which is becoming popular <http://pareonline.net/pdf/v12n2.pdf>.

Other output

Total variance explained: If a rotation has been requested, the percentage of the common variance explained by the retained PC's is in the 'Extraction sums of squared loadings' section whereas it would contain the percentage of total variance if no rotation has been requested. The 'Rotation Sums of Squared Loadings' section contains the distribution of the variation after the Varimax rotation. The varimax (variance maximising) rotation of the loadings maximises the variability of the new PC whilst minimising the variance around the new variable.

Factor scores: The standardised PC scores for each subject will be added to the data set for use in further analysis or to identify subjects scoring highly on certain PC's. Producing a correlation matrix of the scores checks that the PC's are not related. If they are, re-run using the 'Direct Oblimin' rotation instead of the Varimax rotation. It is also useful to examine scatter plots of successive PC's to look for unexpected structure/ subgroups.

Communalities table: The 'Extraction' column contains the proportion of each variables variance explained by the extracted PC's/ factors. The highest coefficients are the strongest variables. Ideally, coefficients should be above 0.5.

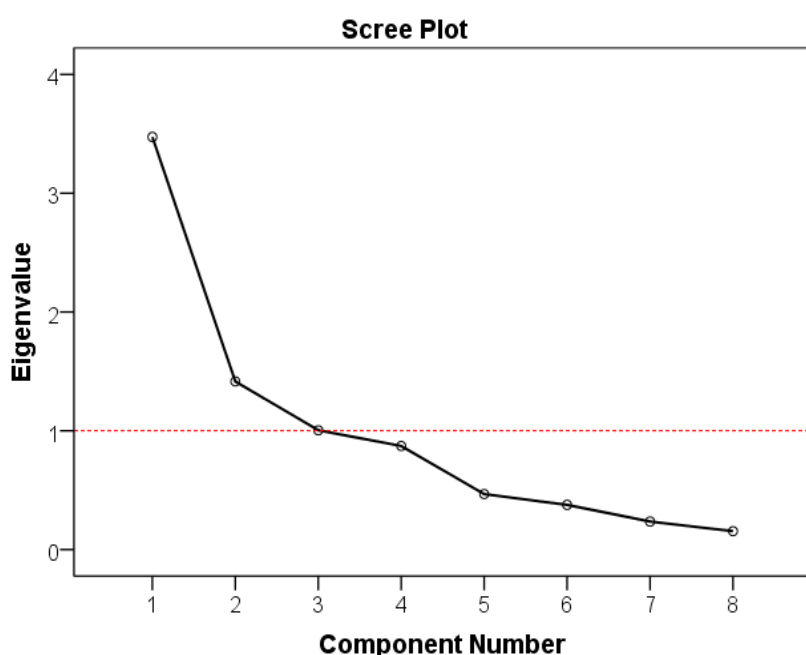
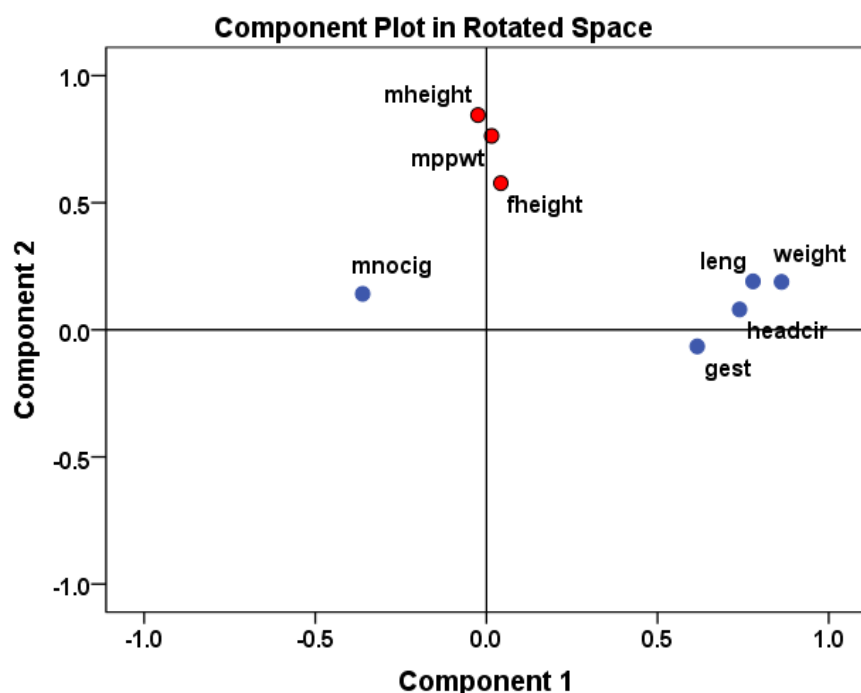
Component matrix: Shows the component loadings (correlations between individual variables and the PC's). Although the values can range between -1 and +1, we chose to suppress coefficients between -0.3 and +0.3 and concentrate on those with higher loadings. Look at the variables for each PC with the highest loadings. Do similar types of variable all have higher loadings on one PC? The component plot is helpful when grouping raw variables together based on their loadings on each PC.

Component score coefficient matrix: Displays the coefficients for the linear combination of variables for the calculation of individual PC scores.

Example data set: Birth weight data set

The results of PCA on a data set containing information on 680 newborn babies and their parents are displayed below. Variables include Baby's length (leng), weight (weight), headcircumference (headcir), gestational age at birth (gest), and mother's height (mheight), pre-pregnancy weight (mppwt) number of cigarettes smoked per day (mnocig) and father's height (fheight).

The component matrix and plot suggest that PC 1 relates mainly to the measurements of the babies (length, weight, gestational age at birth and head circumference) although the number of cigarettes smoked by the mother has a negative loading. PC 2 relates to the height and weight of the parents.



The scree plot suggests that three PC's would be better than the two chosen using an eigenvalue cut off of 1. Further investigation found that the 3rd PC contained only mncig (number of cigarettes smoked by the mother per day) which had very weak correlation with the other variables.

Main references for this section: <http://www.statisticshell.com/docs/factor.pdf>
http://statistics.ats.ucla.edu/stat/spss/output/principal_components.htm

Exploratory Factor Analysis (EFA)

Use: There are two types of Factor Analysis. Exploratory Factor Analysis (EFA) aims to group together and summarise variables which are correlated and can therefore identify possible underlying latent variables which cannot be measured directly whereas Confirmatory Factor Analysis (CFA) tests theories about latent factors. Confirmatory Factor Analysis is performed using additional SPSS software and is beyond the scope of stats support but EFA is commonly used in disciplines such as Psychology and can be found in standard textbooks.

Exploratory Factor Analysis and Principal Component Analysis are very similar. The main differences are:

PCA uses all the variance in the variables analysed whereas EFA uses only the common (shared) variance between the variables

EFA aims to identify underlying latent variables (factors) rather than just reduce the number of variables

Dependents: Scale/ binary(preferably mostly scale or mostly binary) Numerical ordinal variables are often included but the choice of numbering e.g. 1, 2, 3 will impact on the results

Factor analysis is commonly used on questionnaire data with likert style questions attempting to measure underlying latent variables such as depression which cannot be measured directly.

SPSS: Analyse > Dimension reduction > Factor

Run the analysis with two possible changes:

Extraction: Change the method to 'Principal axis factoring'

Rotation: Variamax assumes that the factors are not related. If it is likely that the underlying factors do correlate, use 'Direct Oblimin'

Interpretation:

Interpret the output in the same way as PCA although the Component matrix is called the Factor Matrix if the extraction method has changed.

Look at the Rotated Factor matrix to see which variables contribute most to each factor (PC). Variables measuring the same underlying latent variable should all have high loadings on a particular factor and by looking at the raw variables, a sensible name can be given to the factor. The next factor should be measuring another latent variable etc.

The factor plot is useful for assessing grouping of variables on more than one factor. If there are two factors, the variables appear on a scatterplot.

Main reference for this section: <http://www.statisticshell.com/docs/factor.pdf>

Cluster Analysis

Cluster analysis is a multivariate technique used to group individuals/ variables based on common characteristics. The groups are unknown.

There are three main types of cluster analysis:

Procedure	When to use
Hierarchical clustering	Small data sets of one data type (e.g. continuous) where different numbers of clusters are to be investigated. Both cases and variables can be clustered
K-means	Moderate sized data sets of continuous variables where the number of clusters (k) is specified. Several values of k can be run and optimal chosen.
Two-step	Large data sets or those with a mixture of data types

Hierarchical Clustering

Hierarchical clustering in SPSS is agglomerative (each subject starts in its' own cluster and then clusters are merged until all subjects are in one cluster). The way in which merging occurs depends on the way in which dissimilarity between individuals/ clusters is measured and the method for combining clusters. Once a subject has joined a cluster, it cannot leave. In SPSS, the variables need to be all of one data type.

SPSS: Analyse > Classify > Hierarchical Cluster

Select cases or variables to be clustered.

Statistics:

Proximity matrix (optional – a dissimilarity matrix of distances between subjects)

Range of solutions (a, b) – Output will be compared for a to b number of clusters

Plots: Dendrogram

Method:

Cluster method: Cluster method is how the clusters are merged. The single linkage method joins clusters with the smallest distance between two cases in different clusters (Nearest neighbour) whereas for complete linkage (Furthest neighbour), the distance between two clusters is defined as the distance between the two furthest points. The default method is Between-groups linkage (distance between clusters is the average distance of all data points within these clusters). Wards method (which can only be used with Euclidean measure), calculates the means of all variables within a cluster, then the squared Euclidean distance between cases and means are summed. Clusters are merged where the minimum increase in

sums of squares occurs. Wards is commonly used but susceptible to outliers. Single linkage will isolate outliers into groups which can be removed from the data set and the analysis run again using Wards.

Measure: this is how dissimilarities are measured using distances.

For continuous (interval data) and binary data use Squared Euclidean distance which is the sum of the squared differences between every pair of subjects. This is a dissimilarity measure. For ordinal (count) variables use either Chi-squared or Phi-squared.

Transform: If the variables are measured on different scales, those with larger values/ variances will dominate the distance calculations. Standardisation can be used although variables with more variation are often more likely to separate clusters.

Save: Here the cluster membership for each case is requested. This can be for a set number of clusters or several columns for several cluster numbers.

Interpretation

The main interpretation from the output regards choosing the right number of clusters. The **dendrogram** shows the cases joining (vertical lines) and the distances between clusters. The horizontal axis are the distances scaled to a range of 1 – 25 so are proportional to the actual distances. Long horizontal lines suggest that there has been a large change in the average distance within the cluster so choosing the clusters before the longest jumps is preferable.

The **Proximity matrix** contains the dissimilarity scores between each pair of subjects with small numbers indicating similarity between subjects. The **Agglomeration schedule** shows how the distance measure increases as additional cases are merged into clusters in the 'Coefficients' column. The first step joins the two subjects with the lowest dissimilarity then further steps either join two different subjects or merges a new subject into a formed cluster. The dissimilarity measure increases with the number of steps. Where there is a large increase in the coefficients value, the step merged clusters that were not very similar.

Finally, Hierarchical clustering does not distinguish between the groups but if group membership has been saved for the best number of clusters, descriptives can be calculated by cluster to see which characteristics differ.

K-means clustering

K-means clustering doesn't require a dissimilarity matrix for all pairs as each subject is assigned to the cluster with the mean closest to its value. The means are recalculated at each step (iteration) and subjects can be reassigned to another cluster at any step. The iterations stop when changes in cluster centres don't change much anymore or the maximum iterations are reached. K-means does require a set number of clusters to be specified in the beginning and initial cluster centres (means which are far apart) which are estimated by SPSS if not given. Some people carry out Hierarchical clustering on a sub set of the data to get an idea about the optimal number of clusters and starting cluster means and then carry out K-means. Only continuous variables can be included, Euclidean distances are used and K-means clustering is very sensitive to outliers which may form clusters of their own.

SPSS: Analyze → Classify → K-Means Cluster

In the K-Means menu, specify the number of clusters and request these extras:

Save: Cluster Membership

Options: ANOVA table

Interpretation: The *initial cluster centres* are given in the first table followed by the changes to cluster centres in the *iteration history*. The last row should show negligible change. The *final cluster centres* show how the variables differ in each cluster. It should be clear which variables are most different and therefore define each cluster but the *ANOVA table* shows which variables contribute most to the separation (highest F-statistics) and least. The F-tests should be used for descriptive purposes rather than formal tests as clustering is based on maximising the between cluster to within cluster variation.

Section 4

Suggested resources

Books

SPSS Survival guide. Julie Pallent.

A clear and fairly concise guide to performing the most common tests in SPSS including assumptions, steps in SPSS and interpreting the output.

SPSS for Psychologists. Brace, Kemp and Snelgar.

We like this book because: It offers students quick examples of using SPSS to undertake statistical analyses and interpret the results.

Discovering statistics using SPSS. Andy Field.

A favourite in Psychology with interesting examples but quite wordy. Good for extra information not included in more concise books but the additional detail can distract from the main points.

Oxford Handbook of Medical Statistics. Janet and Philip Peacock.

Great summary guide covering a wide range statistical techniques and definitions.

Problem Solving: A Statistician's Guide. Chris Chatfield.

We like this book because: It provides ideas and examples of how to go about undertaking a statistical analysis. It also provides a quick overview of many different statistical analyses so students can see if that might be useful.

Elementary Survey Sampling. Scheaffer, Mendenhall and Ott.

We like this book because: It includes simple formula to calculate margins of error (and sample sizes for a required margin of error) from sample surveys and is especially useful where the population being studied is not large.

Multivariate Statistical Methods: A Primer. Bryan Manly.

We like this book because: It gives a good overview of multivariate methods that allows a student to assess whether these are useful. It does include some mathematics but should be accessible to anyone having studied mathematics modules as part of their undergraduate degree, but this could be omitted anyway for others and still be useful.

100 Statistical Tests. Gopal Kanji.

We like this book because: A great resource if you can't remember the details of a particular test. Also useful to find a test for less common situations.

Websites

Statstutor: <http://www.statstutor.ac.uk/>

A growing collection of statistics teaching resources in different media formats.

STEPS glossary: <http://www.stats.gla.ac.uk/steps/glossary/index.html>

A useful resource for quick definitions.

BrunelASK videos: <https://www.youtube.com/user/BrunelASK>

These short videos, usually about using SPSS, were created by Christine Pereira at Brunel University. Once reviewed they will appear on statstutor.

CAST: http://cast.massey.ac.nz/collection_public.html

A collection of computer assisted statistics textbooks including core statistics ebooks and apps for lecturers to use in class.

Statistics Fun <http://www.youtube.com/user/statisticsfun>

A YouTube channel of statistics videos.

WhatTest: <http://whattest.lboro.ac.uk>

A website aimed at new researchers who need help deciding on an appropriate analysis plan for a study or experiment.

Online Statistics Education: A Multimedia Course of Study:
<http://onlinestatbook.com/>

A teaching resource with apps for use to demonstrate techniques.

Contributors to the guide

This guide was produced primarily by the Maths and Statistics Help centre (MASH) at Sheffield University as part of a statistics tutors training project funded by SIGMA.

Main authors

Ellen Marshall (University of Sheffield)

Elizabeth Boggis (University of Sheffield)

Other contributors

Chetna Patel (University of Sheffield)

Marta Emmett (University of Sheffield)

Alun Owen (University of Worcester)

Reviewers (University of Sheffield)

Jean Russell and Nick Fieller (Multivariate)

Statistics: 1.3 The Chi-squared test for two-way tables

Rosie Shier. 2004.

1 Introduction

If we have two categorical variables we can look at the relationship between these variables by putting the data in a two-way table.

Example

A sample of 200 components is selected from the output of a factory that uses three different machines to manufacture these components. Each component in the sample is inspected to determine whether or not it is defective. The machine that produced the component is also recorded. The results are as follows:

	Machine			
Outcome	A	B	C	Total
Defective	8 (12.9%)	6 (8.8%)	12 (17.1%)	26 (13.0%)
Non-defective	54	62	58	174
Total	62	68	70	200

The manager wishes to determine whether or not there is a relationship between the proportion of defectives and the machine used.

In general, the null hypothesis in a two-way table is that there is no association between the row variable and the column variable and the alternative hypothesis is that there is an association. In this case, this is equivalent to saying:

H_0 : There are no differences between machines in the percentage of defectives produced.

H_1 : There are differences ...

To test the null hypothesis we compare the observed cell counts with **expected** cell counts calculated under the assumption that the null hypothesis is true. If the null hypothesis were true the row (or column) percentages would all be the same. Therefore:

$$\begin{aligned}
 \text{Expected cell count} &= \frac{\text{Row percentage}}{100} \times \text{Column total} \\
 &= \frac{\text{Row total}}{\text{Overall total}} \times \text{Column Total} \\
 &= \frac{\text{Row total} \times \text{Column Total}}{n}
 \end{aligned}$$

where n is the overall total.

Looking at our example – expected cell counts:

	Machine		
Outcome	A	B	C
Defective	$\frac{26 \times 62}{200} = 8.06$	$\frac{26 \times 68}{200} = 8.84$	$\frac{26 \times 70}{200} = 9.10$
Non-defective	$\frac{174 \times 62}{200} = 53.94$	$\frac{174 \times 68}{200} = 59.16$	$\frac{174 \times 70}{200} = 60.90$

Notice that (as specified by the null hypothesis) the expected row percentages are all the same (i.e. $\frac{8.06}{62} = \frac{8.84}{68} = \frac{9.10}{70} = 13\% =$ overall percentage of defectives).

2 Carrying out the chi-squared test

To test the null hypothesis we now compute a statistic that compares the entire set of observed counts with the set of expected counts. This statistic is called the **chi-squared statistic** and is given by:

$$\chi^2 = \sum \frac{(O - E)^2}{E} \quad \text{where } O = \text{observed cell count} \ \& \ E = \text{expected cell count}$$

and the sum is over all $r \times c$ cells in the table, where r = number of rows, c = number of columns in the table.

In our example:

$$\chi^2 = \frac{(8 - 8.06)^2}{8.06} + \frac{(6 - 8.84)^2}{8.84} + \dots + \frac{(62 - 59.16)^2}{59.16} + \frac{(58 - 60.90)^2}{60.90} = 2.11 \quad (3 \text{ d.p.})$$

The number of degrees of freedom is simply $(r - 1) \times (c - 1)$.

In our case this is $(3 - 1) \times (2 - 1) = 2$. The p-value for the chi-squared test in this case is $P(\chi_2^2 > 2.11)$.

Looking this up in tables of the chi-squared distribution gives $0.25 < p < 0.5$.

In this case we have no real evidence that the percentage of defectives varies from machine to machine.

3 Validity of chi-squared tests for two-way tables

Chi-squared tests are only valid when you have a reasonable sample size. The following guidelines can be used:

1. For 2 x 2 tables:

- If the total sample size is greater than 40, χ^2 can be used.
- If the total sample size is between 20 and 40 and the smallest expected frequency is at least 5, χ^2 can be used.
- Otherwise Fisher's exact test must be used.

2. For other tables:

- χ^2 can be used if no more than 20% of the expected frequencies are less than 5 and none is less than 1.

4 Carrying out a chi-squared test in SPSS

Your data could be in one of two formats:

1. Individual data	Machine	Outcome
	A	0
	A	1
	B	1
	C	0
	A	1
	B	0
	\vdots	\vdots

Where 0 stands for defective and 1 for non-defective. In this case, you would have 200 lines of data, one for each component.

2. Grouped data	Machine	Outcome	Frequency
	A	0	8
	A	1	54
	B	0	6
	B	1	62
	C	0	12
	C	1	58

In this case, you only have $r \times c = 6$ lines of data, one for each cell in the table.

If you have grouped data, you need to do the following before carrying out the chi-squared test:

- **Data**
- **Weight Cases**
- Select **Weight cases by** and choose your frequency variable as the **Frequency Variable**.

You are now ready to carry out the chi-squared test (the procedure is now the same whether you have individual or grouped data):

- **Analyze**
- **Descriptive Statistics**
- **Crosstabs**
- Choose your outcome variable as the **Row Variable** and your explanatory variable as the **Column Variable** (it actually doesn't matter which way round they are, but the commands below are based on them being this way around).
- Click on the **Cells** button and select **Column** under **Percentages**. (You can also ask SPSS to display expected cell counts by clicking on the appropriate button here.) Now click on **Continue**.
- Click on the **Statistics** button and select **Chi-square** in the top left-hand corner. Now click on **Continue**.
- Click on **OK**

(See over for output)

Your output will look like this:

Case Processing Summary

	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
OUTCOME * MACHINE	200	100.0%	0	.0%	200	100.0%

OUTCOME * MACHINE Crosstabulation

OUTCOME		MACHINE			Total
		A	B	C	
Defective	Count	8	6	12	26
	% within MACHINE	12.9%	8.8%	17.1%	13.0%
Non-defective	Count	54	62	58	174
	% within MACHINE	87.1%	91.2%	82.9%	87.0%
Total	Count	62	68	70	200
	% within MACHINE	100.0%	100.0%	100.0%	100.0%

Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	2.112	2	.348
Likelihood Ratio	2.144	2	.342
Linear-by-Linear Association	.585	1	.444
N of Valid Cases	200		

- a. 0 cells (.0%) have expected count less than 5.
The minimum expected count is 8.06.

We are interested in the top row of the last table — i.e. the **Pearson Chi-Square**. The last column in the table gives the p-value — in this case we have $p = 0.348$. We would normally round this to $p = 0.3$.

1 Introduction

This handout is designed to provide only a brief introduction to cluster analysis and how it is done. Books giving further details are listed at the end.

Cluster analysis is a multivariate method which aims to classify a sample of subjects (or objects) on the basis of a set of measured variables into a number of different groups such that similar subjects are placed in the same group. An example where this might be used is in the field of psychiatry, where the characterisation of patients on the basis of clusters of symptoms can be useful in the identification of an appropriate form of therapy. In marketing, it may be useful to identify distinct groups of potential customers so that, for example, advertising can be appropriately targeted.

WARNING ABOUT CLUSTER ANALYSIS

Cluster analysis has no mechanism for differentiating between relevant and irrelevant variables. Therefore the choice of variables included in a cluster analysis must be underpinned by conceptual considerations. This is very important because the clusters formed can be very dependent on the variables included.

2 Approaches to cluster analysis

There are a number of different methods that can be used to carry out a cluster analysis; these methods can be classified as follows:

- Hierarchical methods
 - Agglomerative methods, in which subjects start in their own separate cluster. The two 'closest' (most similar) clusters are then combined and this is done repeatedly until all subjects are in one cluster. At the end, the optimum number of clusters is then chosen out of all cluster solutions.
 - Divisive methods, in which all subjects start in the same cluster and the above strategy is applied in reverse until every subject is in a separate cluster. Agglomerative methods are used more often than divisive methods, so this handout will concentrate on the former rather than the latter.
- Non-hierarchical methods (often known as **k-means clustering** methods)

1 Introduction

This handout is designed to provide only a brief introduction to factor analysis and how it is done. Books giving further details are listed at the end.

As for principal components analysis, factor analysis is a multivariate method used for data reduction purposes. Again, the basic idea is to represent a set of variables by a smaller number of variables. In this case they are called **factors**. These factors can be thought of as underlying constructs that cannot be measured by a single variable (e.g. happiness).

2 Assumptions

Factor analysis is designed for interval data, although it can also be used for ordinal data (e.g. scores assigned to Likert scales). The variables used in factor analysis should be linearly related to each other. This can be checked by looking at scatterplots of pairs of variables. Obviously the variables must also be at least moderately correlated to each other, otherwise the number of factors will be almost the same as the number of original variables, which means that carrying out a factor analysis would be pointless.

3 The steps in factor analysis

The factor analysis model can be written algebraically as follows. If you have p variables X_1, X_2, \dots, X_p measured on a sample of n subjects, then variable i can be written as a linear combination of m factors F_1, F_2, \dots, F_m where, as explained above $m < p$. Thus,

$$X_i = a_{i1}F_1 + a_{i2}F_2 + \dots + a_{im}F_m + e_i$$

where the a_i s are the factor loadings (or scores) for variable i and e_i is the part of variable X_i that cannot be 'explained' by the factors.

There are three main steps in a factor analysis:

1. Calculate initial factor loadings.

This can be done in a number of different ways; the two most common methods are described very briefly below:

- **Principal component method**

As the name suggests, this method uses the method used to carry out a principal

Statistics: 1.5 Oneway Analysis of Variance

Rosie Cornish. 2006.

1 Introduction

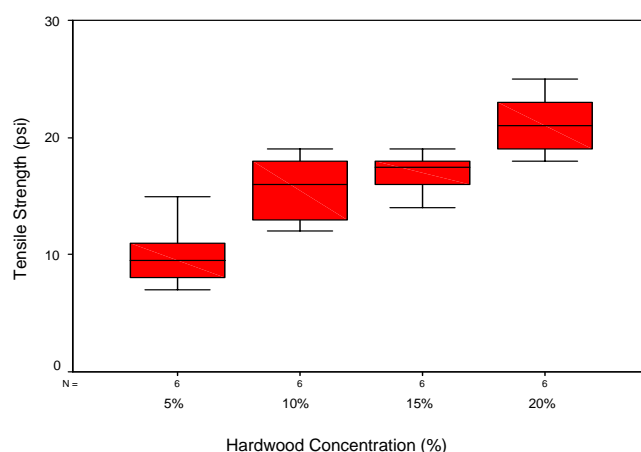
Oneway analysis of variance (ANOVA) is used to compare several means. This method is often used in scientific or medical experiments when treatments, processes, materials or products are being compared.

Example:

A paper manufacturer makes grocery bags. They are interested in increasing the tensile strength of their product. It is thought that strength is a function of the hardwood concentration in the pulp. An investigation is carried out to compare four levels of hardwood concentration: 5%, 10%, 15% and 20%. Six test specimens are made at each level and all 24 specimens are then tested in random order. The results are shown below:

Hardwood Concentration (%)	Tensile strength (psi)						Mean	Standard Deviation
5	7	8	15	11	9	10	10.00	2.83
10	12	17	13	18	19	15	15.67	2.81
15	14	18	19	17	16	18	17.00	1.79
20	19	25	22	23	18	20	21.17	2.64
All							15.96	4.72

Source: Applied Statistics and Probability for Engineers - Montgomery and Runger



As stated above, in ANOVA we are asking the question, "Do all our groups come from populations with the same mean?". To answer this we need to compare the sample means. However,

Multiple regression

Introduction

Multiple regression is a logical extension of the principles of simple linear regression to situations in which there are several predictor variables. For instance if we have two predictor variables, X_1 and X_2 , then the form of the model is given by:

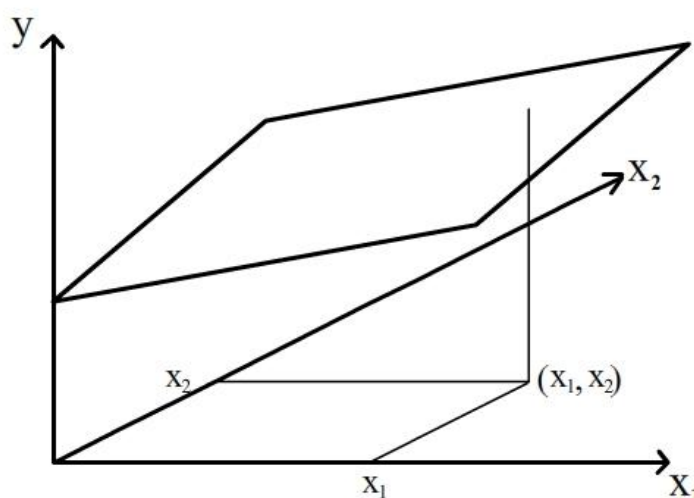
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$$

which comprises a deterministic component involving the three *regression coefficients* (β_0 , β_1 and β_2) and a random component involving the *residual* (error) term, e .

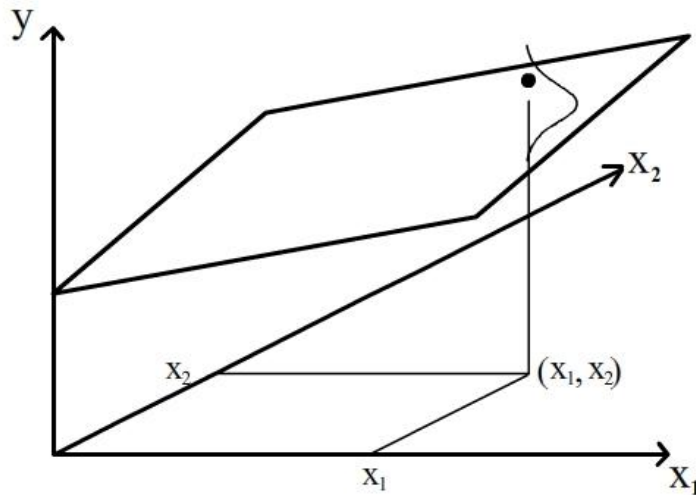
Note that the predictor variables can be either continuous or categorical. In the case of the latter these variables need to be coded as dummy variables (not considered in this tutorial). The response variable must be measured on a continuous scale.

The residual terms represent the difference between the predicted and observed values of individuals. They are assumed to be independently and identically distributed normally with zero mean and variance σ^2 , and account for natural variability as well as maybe measurement error.

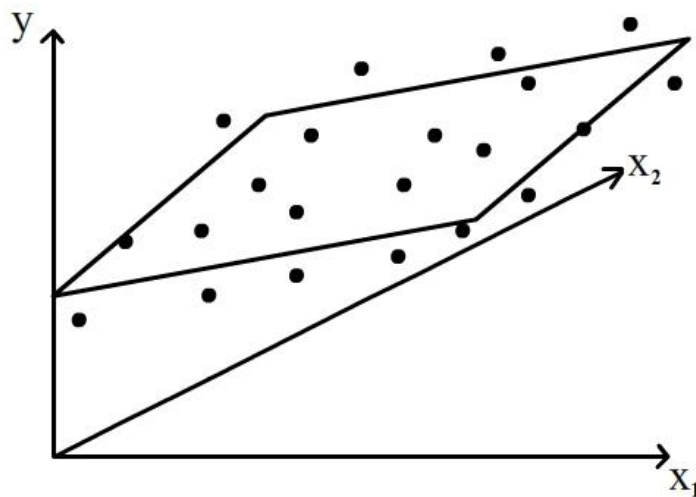
For the two (continuous) predictor example the deterministic component is in the form of a plane which provides the predicted (mean/expected) response for given predictor variable value combinations. Thus if we want the expected value for the specific values x_1 and x_2 , then this is obtained from the orthogonal projection from the point (x_1, x_2) in the $X_1 - X_2$ plane to the expected value plane in the 3D space. The resulting Y value is the expected value from this explanatory variable combination.



Observed values for this combination of explanatory variables are drawn from a normal distribution with variance σ^2 centred on the expected value point:



Our data should thus appear to be a collection of points that are randomly scattered with constant variability around the plane.



The multiple regression model fitting process takes such data and estimates the regression coefficients (β_0 , β_1 and β_2) that yield the plane that has best fit amongst all planes.

Model assumptions

The assumptions build on those of simple linear regression:

- *Ratio of cases to explanatory variables.* Invariably this relates to research design. The minimum requirement is to have at least five times more cases than explanatory variables. If the response variable is skewed then this number may be substantially more.
- *Outliers.* These can have considerable impact upon the regression solution and their inclusion needs to be carefully considered. Checking for extreme values should form part of the initial data screening process and should be performed on both the response and explanatory variables. Univariate outliers can simply be identified by considering the distributions of individual variables say by using boxplots. Multivariate outliers can be detected from residual scatterplots.
- *Multicollinearity and singularity.* Multicollinearity exists when there are high correlations among the explanatory variables. Singularity exists when there is perfect correlation between explanatory variables. The presence of either affect the interpretation of the explanatory variables effect on the response variable. Also it can lead to numerical problems in finding the regression solution. The presence of multicollinearity can be detected by examining the correlation matrix (say $r = \pm 0.9$ and above). If there is a pair of variables that appear to be highly multicollinear then only one should be used in the regression. Note; some context dependent thought has to be given as to which one to retain!
- *Normality, linearity, homoscedasticity and independence of residuals.* The first three of these assumptions are checked using residual diagnostic plots after having fit a multiple regression model. The independence of residuals is usually assumed to be true if we have indeed collected a random sample from the relevant population.

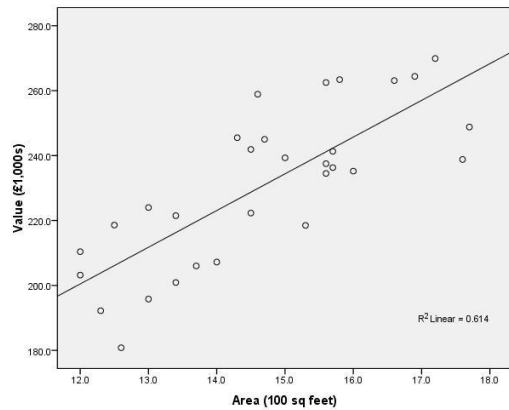
Example

Suppose we are interested in predicting the current market value of houses in a particular city. We have collected data that comprises a random sample of 30 house current values (£1,000s) together with their corresponding living area (100 ft²) and the distance in miles from the city centre.

Case Summaries			
	Value (£1,000s)	Area (100 sq feet)	City centre distance (miles)
1	210.4	12.0	1.2
2	262.5	15.6	1.5
3	258.9	14.6	1.6
4	245.0	14.7	2.5
5	239.3	15.0	2.7
6	263.1	16.6	2.6
7	203.2	12.0	3.2
8	221.5	13.4	3.3
9	207.2	14.0	4.1
10	234.5	15.6	4.2
11	195.8	13.0	4.4
12	222.3	14.5	4.7
13	192.2	12.3	5.1
14	248.8	17.7	5.3
15	218.5	15.3	5.5
16	224.0	13.0	1.3
17	241.9	14.5	1.6
18	245.5	14.3	1.9
19	263.4	15.8	2.4
20	264.4	16.9	2.6
21	269.9	17.2	4.0
22	236.3	15.7	3.2
23	235.2	16.0	4.2
24	218.6	12.5	3.9
25	241.3	15.7	3.8
26	237.5	15.6	4.5
27	180.8	12.6	4.8
28	200.9	13.4	5.0
29	206.0	13.7	5.1
30	238.8	17.6	6.3

Can we build a multiple regression model that can successfully predict house values using the living area and distance variables?

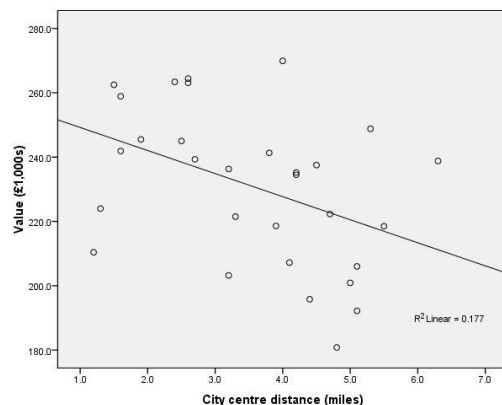
If we consider the relationship between value and area it appears that there is a very significant positive correlation between the two variables (i.e. value increases with area). Fitting a simple linear regression model indicates that 61.4% of the variability in the values is explained by the area.



Correlations			
		Value (£1,000 s)	Area (100 sq feet)
Value (£1,000s)	Pearson Correlation	1	.784**
	Sig. (2-tailed)		.000
	N	30	30
Area (100 sq feet)	Pearson Correlation	.784**	1
	Sig. (2-tailed)	.000	
	N	30	30

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

If we consider the relationship between value and distance it appears that there is a significant negative correlation between the two variables (i.e. value decreases with distance). Fitting a simple linear regression model indicates that 17.7% of the variability in the values is explained by the distance from the city centre.



Correlations			
		Value (£1,000 s)	City centre distance (miles)
Value (£1,000s)	Pearson Correlation	1	-.421*
	Sig. (2-tailed)		.020
	N	30	30
City centre distance (miles)	Pearson Correlation	-.421*	1
	Sig. (2-tailed)	.020	
	N	30	30

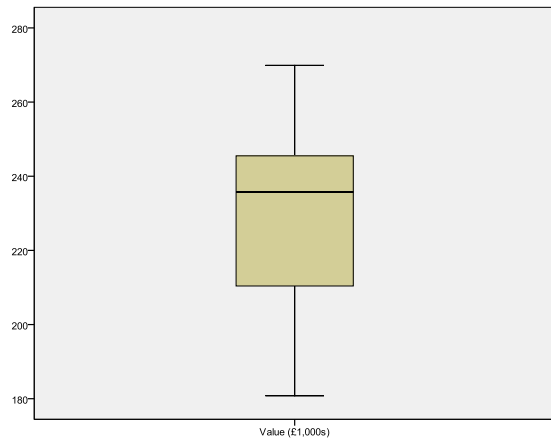
*. Correlation is significant at the 0.05 level (2-tailed).

Thus individually either variable is useful for predicting a house value. We shall now consider the fitting of a multiple regression model that uses both variables for predictions.

First of all we need to address the assumptions that we check before fitting a multiple regression model.

Ratio of cases to explanatory variables.

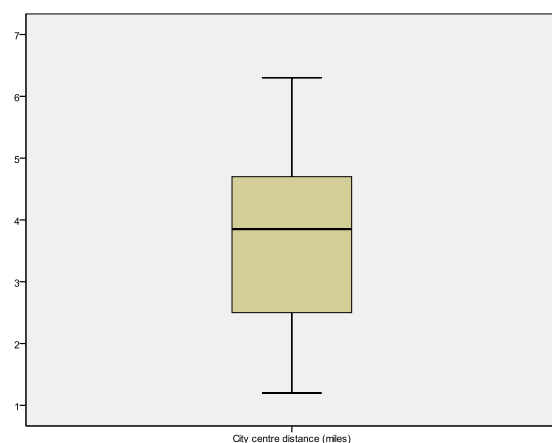
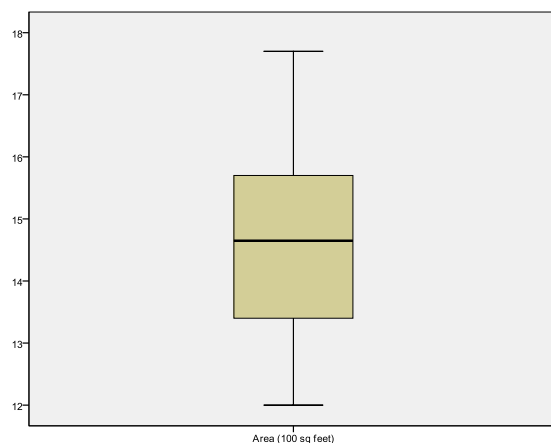
We have 30 cases and 2 explanatory variables. Looking at the boxplot of the response variable value does not overtly worry us that there is a skewness problem.



Thus as we have 15 times more cases than explanatory variables we should have an adequate number of cases.

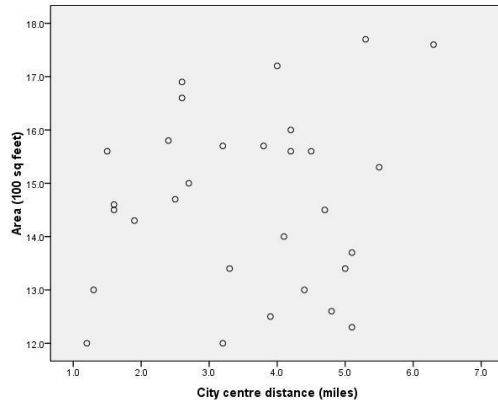
Outliers

The boxplot above of the response variable does not identify any outliers and neither do the two boxplots below of the explanatory variables:



Multicollinearity and singularity

Examining the correlation between the two explanatory variables reveals that there is not a significant correlation between them. Thus we have no concerns over multicollinearity.



Correlations			
		Area (100 sq feet)	City centre distance (miles)
Area (100 sq feet)	Pearson Correlation	1	.154
	Sig. (2-tailed)		.415
	N	30	30
City centre distance (miles)	Pearson Correlation	.154	1
	Sig. (2-tailed)	.415	
	N	30	30

Independence of residuals

Our data has come from a random sample and thus the observations should be independent and hence the residuals should be too.

It appears that our pre model fitting assumption checks are satisfactory, and so we can now consider the multiple regression output.

Coefficients ^a								
		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95.0% Confidence Interval for B	
Model		B	Std. Error	Beta			Lower Bound	Upper Bound
1	(Constant)	80.121	12.051		6.648	.000	55.393	104.848
	City centre distance (miles)	-9.456	.965	-.555	-9.799	.000	-11.436	-7.476
	Area (100 sq feet)	12.548	.818	.869	15.340	.000	10.870	14.226

a. Dependent Variable: Value (£1,000s)

The unstandardized coefficients are the coefficients of the estimated regression model. Thus the expected value of a house is given by:

$$value = 80.121 - 9.456 \times distance + 12.548 \times area.$$

Recalling that value is measured in £1,000s and area is in units of 100 ft², we can interpret the coefficients (and associated 95% confidence intervals) as follows.

- For each one mile increase in distance from the city centre, the expected change in house value is -£9,456 (-£11,456, -£7,476). Thus house values drop by £9,456 for each one mile from the city centre.
- For each 100 ft² increase in area, the expected house value is expected to increase by £12,548 (£10,870, £14,226).

The significance tests of the two explanatory variable coefficients indicate that both of the explanatory variables are significant ($p < .001$) for predicting house values. If however either had a p -value $> .05$, then we could infer that the offending variable(s) are not significant for predicting house values.

Note that the intercept here gives the expected value of £80,121 for what would be a house of no area in the exact middle of the city. It is debateable whether this makes any sense and can be dismissed by the fact that these values of the explanatory variables are an extrapolation from what we have observed.

The standardized coefficients are appropriate in multiple regression when we have explanatory variables that are measured on different units (which is the case here). These coefficients are obtained from regression after the explanatory variables are all standardized. The idea is that the coefficients of explanatory variables can be more easily compared with each other as they are then on the same scale. Here we see that the *area* standardised coefficient is larger in absolute value than that of *distance*: thus we can conclude that a change in area has a greater relative effect on house value than does a change in distance from the city centre.

Model Summary ^b				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.957 ^a	.915	.909	7.2459

a. Predictors: (Constant), Area (100 sq feet), City centre distance (miles)

b. Dependent Variable: Value (£1,000s)

Examining the model summary table:

- The multiple correlation coefficient, R, indicates that we have a very high correlation of .957 between our response variable and the two explanatory variables.
- From the R squared value (*coefficient of determination*) we can see that the model fits the data reasonably well; 91.5% of the variation in the house values can be explained by the fitted model together with the house area and distance from the city centre values.
- The adjusted R square value is attempts to correct for this. Here we can see it has slightly reduced the estimated proportion. If you have a small data set it may be worth reporting the adjusted R squared value.
- The standard error of the estimate is the estimate of the standard deviation of the error term of the model, σ . This gives us an idea of the expected variability of predictions and is used in calculation of confidence intervals and significance tests.

The remaining output is concerned with checking the model assumptions of normality, linearity and homoscedasticity of the residuals. Residuals are the differences between the observed and predicted responses. The residual scatterplots allow you to check:

- *Normality*: the residuals should be normally distributed about the predicted responses;
- *Linearity*: the residuals should have a straight line relationship with the predicted responses;
- *Homoscedasticity*: the variance of the residuals about predicted responses should be the same for all predicted responses.

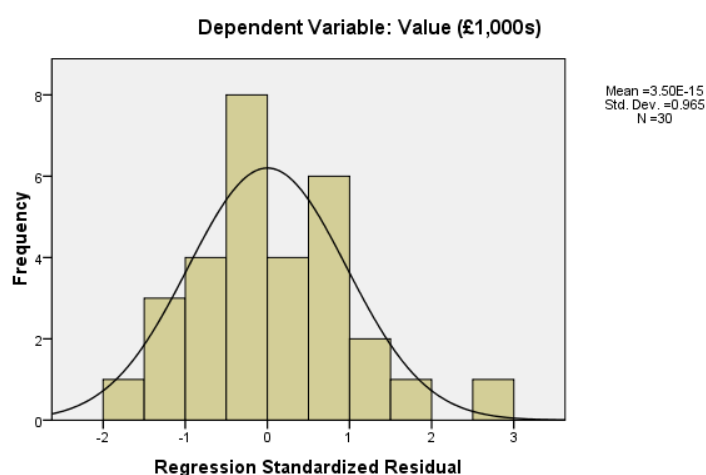
Residuals Statistics^a

	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	186.234	267.596	230.923	22.9872	30
Residual	-12.0357	18.5084	.0000	6.9916	30
Std. Predicted Value	-1.944	1.595	.000	1.000	30
Std. Residual	-1.661	2.554	.000	.965	30

a. Dependent Variable: Value (£1,000s)

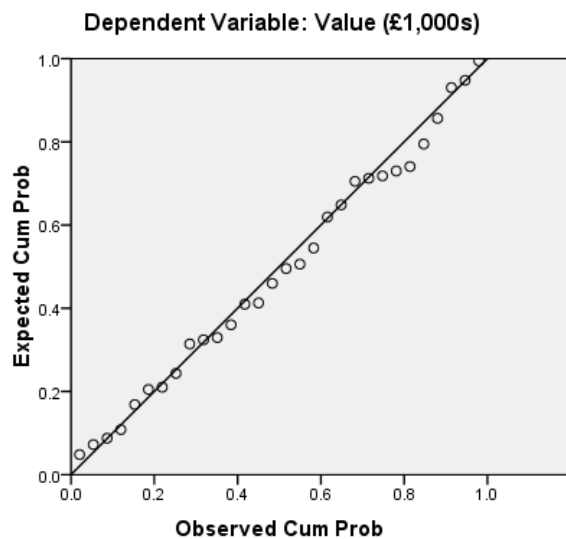
The above table summarises the predicted values and residuals in unstandardised and standardised forms. It is usual practice to consider standardised residuals due to their ease of interpretation. For instance outliers (observations that do not appear to fit the model that well) can be identified as those observations with standardised residual values above 3.3 (or less than -3.3). From the above we can see that we do not appear to have any outliers.

Histogram



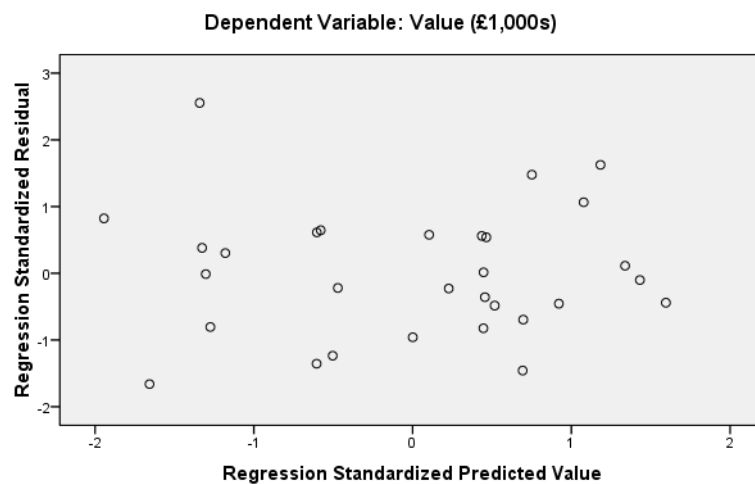
The above plot is a check on normality; the histogram should appear normal; a fitted normal distribution aids us in our consideration. Serious departures would suggest that normality assumption is not met. Here we have a histogram that does look reasonably normal given that we have only 30 data points and thus we have no real cause for concern.

Normal P-P Plot of Regression Standardized Residual



The above plot is a check on normality; the plotted points should follow the straight line. Serious departures would suggest that normality assumption is not met. Here we have no major cause for concern.

Scatterplot



The above scatterplot of standardised residuals against predicted values should be a random pattern centred around the line of zero standard residual value. The points should have the same dispersion about this line over the predicted value range. From the above we can see no clear relationship between the residuals and the predicted values which is consistent with the assumption of linearity.

Thus we are happy that the assumptions of the model have been met and thus would be confident about any inference/predictions that we gain from the model.

Predictions

In order to get an expected house *value* for particular *distance* and *area* values we can use the fitted equation. For example, for a house that is 5 miles from the city centre and is 1,400 ft²:

$$\begin{aligned} \text{value} &= 80.121 - 9.456 \times 5 + 12.548 \times 14 \\ &= 208.513 \end{aligned}$$

i.e. £208,513.

Alternatively, we could let a statistics program do the work and calculate confidence or prediction intervals at the same time. For instance, when requesting a predicted value in SPSS we can also obtain the following:

- the predicted values for the various explanatory variable combinations together with the associated standard errors of the predictions;
- 95% CI for the expected response;
- 95% CI for individual predicted responses.

Returning to our example we get the following:

- the expected house value is £208,512 (s.e. = 2,067.6);
- we are 95% certain that interval from £204,269 to £212,754 covers the unknown expected house value;
- we are 95% certain that interval from £193,051 to £223,972 covers the range of predicted individual house value observations.

even if all the population means were identical, we would not expect the sample means to be exactly equal — there will be always be some differences due to sampling variation. The question therefore becomes, “Are the observed differences between the sample means simply due to sampling variation or due to real differences in the population means?” This question cannot be answered just from the sample means — we also need to look at the variability of whatever we’re measuring. In analysis of variance we compare the variability between the groups (how far apart are the means?) to the variability within the groups (how much natural variation is there in our measurements?). This is why it is called analysis of variance.

ANOVA is based on two assumptions. Therefore, before we carry out ANOVA, we need to check that these are met:

- 1) The observations are random samples from normal distributions.
- 2) The populations have the same variance, σ^2 .

Fortunately, ANOVA procedures are not very sensitive to unequal variances — the following rule can be applied:

If the largest standard deviation (not variance) is less than twice the smallest standard deviation, we can use ANOVA and our results will still be valid.

So, before carrying out any tests we must first look at the data in more detail to determine whether these assumptions are satisfied:

- i) Normality: If you have very small samples, it can sometimes be quite difficult to determine whether they come from a normal distribution. However, we can assess whether the distributions are roughly symmetric by (a) comparing the group means to the medians — in a symmetric distribution these will be equal and (b) looking at boxplots or histograms of the data
- ii) Equal variances: We can simply compare the group standard deviations.

In our example the medians are very close to the means. Also the standard deviations in the four groups (see table on page 1) are quite similar. These results, together with the boxplots given above, indicate that the distribution of tensile strength at each hardwood concentration is reasonably symmetric and that its variability does not change markedly from one concentration to another. Therefore, we can proceed with the analysis of variance.

2 The ANOVA Model

2.1 Notation

In general, if we sample n observations from each of k populations (groups), the total number of observations is $N = nk$. The following notation is used:

- y_{ij} represents the j^{th} observation in group i (e.g. y_{13} is the 3rd observation in the first group, y_{31} is the first observation in the third group, and so on).
- \bar{y}_i represents the mean in group i
- \bar{y} represents the mean of all the observations

2.2 Sums of squares

The total variation of all the observations about the overall mean is measured by what is called the **Total sum of squares**, given by:

$$SS_T = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y})^2$$

This variation can be split into two components:

- 1) the variation of the group means about the overall mean (between-group variation)
- 2) the variation of the individual observations about their group mean (within-group variation)

It can be shown that:

$$\sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y})^2 = n \sum_{i=1}^k (\bar{y}_i - \bar{y})^2 + \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2$$

or

$$SS_T = SS_B + SS_W$$

In words this is written as:

Total sum of squares = Between groups sum of squares + Within groups sum of squares

Degrees of freedom and mean squares

Each sum of squares has a certain number of degrees of freedom:

SS_T compares N observations to the overall mean, so has $N - 1$ degrees of freedom.

SS_B compares k means to the overall mean, so has $k - 1$ degrees of freedom.

SS_W compares N observations to k sample means, so has $N - k$ degrees of freedom.

Notice that $N - 1 = (N - k) + (k - 1)$ (i.e. the degrees of freedom are related in the same way as the sums of squares: $df_T = df_B + df_W$).

Degrees of freedom:

The degrees of freedom basically indicates how many 'values' are free to vary. When we are considering variances or sums of squares, because the sum of the deviations is always zero, the last deviation can be found if we know all the others. So, if we have n deviations, only $n - 1$ are free to vary.

The **mean square** for each source of variation is defined as being the sum of squares divided by its degrees of freedom. Thus:

$$MS_B = SS_B / (k - 1) \quad \text{and} \quad MS_W = SS_W / (N - k)$$

3 The F Test in ANOVA

It can be shown that if the null hypothesis is true and there are no differences between the (unknown) population means, MS_B and MS_W will be very similar. On the other hand, if the (unknown) means are different, MS_B will be greater than MS_W (this makes sense intuitively — if the population means are very different, we would expect the sample means to be quite far apart and therefore the between group variability will be large). Therefore, the ratio

MS_B/MS_W is a statistic that is approximately equal to 1 if the null hypothesis is true but will be larger than 1 if there are differences between the population means.

The ratio MS_B/MS_W is a ratio of variances, and follows what is called an F distribution with $k - 1$ and $N - k$ degrees of freedom. In summary:

To test $H_0 : \mu_1 = \mu_2 = \dots = \mu_k = \mu$ use the statistic $F = \frac{MS_B}{MS_W}$ and compare this to the F distribution with $k - 1$ and $N - k$ degrees of freedom.

The F Distribution and F Tests

A statistical test called the **F test** is used to compare variances from two normal populations.

It is tested by the **F-statistic**, the ratio of the two sample variances: $F = \frac{s_1^2}{s_2^2}$.

Under the null hypothesis, this statistic follows an F distribution with $n_1 - 1$ and $n_2 - 1$ degrees of freedom, written $F(n_1 - 1, n_2 - 1)$.

The **F distributions** are a family of distributions which depend on two parameters: the degrees of freedom of the sample variances in the numerator and denominator of the F statistic. The degrees of freedom in the numerator are always given first. The F distributions are not symmetric and, since variances cannot be negative, cannot take on values below zero. The peak of any F-distribution is close to 1; values far from 1 provide evidence against the null hypothesis. Two examples of the F distribution with different degrees of freedom are shown in Figure 1.

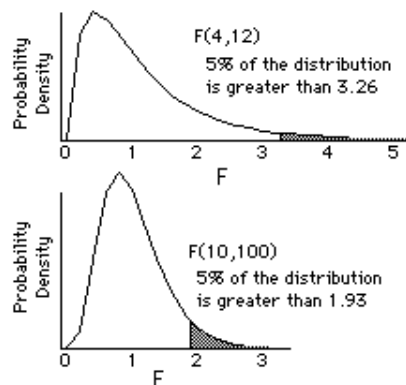


Figure 1: Probability density functions of two F distributions

4 The ANOVA Table

When you carry out an analysis of variance on a computer, you will get an **analysis of variance** (or **ANOVA**) table, as shown below.

The ANOVA table: Tensile strength of paper

Source of variation	Sum of squares	Degrees of freedom	Mean square	F	p
Between groups	382.79	3	127.60	$\frac{127.60}{6.51} = 19.61$	$p < 0.001$
Within groups	130.17	20	6.51		
Total	512.96	23			

From our results we can say that there is strong evidence that the mean tensile strength varies with hardwood concentration. Although the F -test does not specify the nature of these differences it is evident from our results that, as the hardwood concentration increases, so does the tensile strength of the paper. It is possible to test more specific hypotheses — for example, that there is an increasing or decreasing trend in the means — but these tests will not be covered in this leaflet.

What exactly is the p-value?

If the (true) mean tensile strength of paper made with different concentrations of hardwood were actually constant (i.e. if the hardwood concentration had no effect on tensile strength whatsoever), the probability of getting sample means as far apart as, or further apart than, we did (i.e. means of 10.0, 15.7, 17.0, and 21.1, or values further apart than this) is incredibly small — less than 0.001. The p-value represents this probability.

We turn this around and conclude that the true mean tensile strength is very unlikely to be constant (i.e. we conclude that the hardwood concentration does seem to have an effect on tensile strength).

Note: The ANOVA results given above are based on the assumption that the sample size in each group is equal. Usually this will be the case — most experiments are designed with equal sized samples in each experimental group. However, it is also possible to carry out an analysis of variance when the sample sizes are not equal. In this case, the formulae for the sums of squares etc. are modified to take account of the different sample sizes. If you use a statistical package to carry out your analysis, this is done automatically.

5 Carrying out oneway ANOVA in SPSS

- **Analyze**
- **Compare Means**
- **One Way ANOVA**
- Choose your outcome variable (in our case tensile strength) to go in **Dependent List**
- Choose the variable that defines the groups as the **Factor** then click on **OK**.

The output will look something like the ANOVA table given above.

components analysis. However, the factors obtained will not actually be the principal components (although the loadings for the k^{th} factor will be proportional to the coefficients of the k^{th} principal component).

- **Principal axis factoring**

This is a method which tries to find the lowest number of factors which can account for the variability in the original variables that is associated with these factors (this is in contrast to the principal components method which looks for a set of factors which can account for the total variability in the original variables).

These two methods will tend to give similar results if the variables are quite highly correlated and/or the number of original variables is quite high. Whichever method is used, the resulting factors at this stage will be uncorrelated.

2. Factor rotation

Once the initial factor loadings have been calculated, the factors are rotated. This is done to find factors that are easier to interpret. If there are 'clusters' (groups) of variables — i.e. subgroups of variables that are strongly inter-related — then the rotation is done to try to make variables within a subgroup score as highly (positively or negatively) as possible on one particular factor while, at the same time, ensuring that the loadings for these variables on the remaining factors are as low as possible. In other words, the object of the rotation is to try to ensure that all variables have high loadings only on one factor.

There are two types of rotation method, **orthogonal** and **oblique** rotation. In orthogonal rotation the rotated factors will remain uncorrelated whereas in oblique rotation the resulting factors will be correlated. There are a number of different methods of rotation of each type. The most common orthogonal method is called **varimax** rotation; this is the method that many books will recommend.

3. Calculation of factor scores

When calculating the final factor scores (the values of the m factors, F_1, F_2, \dots, F_m , for each observation), a decision needs to be made as to how many factors to include. This is usually done using one of the following methods:

- Choose m such that the factors account for a particular percentage (e.g. 75%) of the total variability in the original variables.
- Choose m to be equal to the number of eigenvalues over 1 (if using the correlation matrix). [A different criteria must be used if using the covariance matrix.]
- Use the scree plot of the eigenvalues. This will indicate whether there is an obvious cut-off between large and small eigenvalues.

In some statistical packages (e.g. SPSS) this choice is actually made at the outset. The second method, choosing eigenvalues over 1, is probably the most common one.

The final factor scores are usually calculated using a regression-based approach.

4 Carrying out factor analysis in SPSS

- **Analyze**
- **Data Reduction**
- **Factor**
- Select the variables you want the factor analysis to be based on and move them into the **Variable(s)** box.
- In the **Descriptives** window, you should select **KMO and Bartlett's test of sphericity**. KMO is a statistic which tells whether you have sufficient items for each factor. It should be over 0.7. Bartlett's test is used to check that the original variables are sufficiently correlated. This test should come out significant ($p < 0.05$) — if not, factor analysis will not be appropriate. Click on **Continue**.
- In the **Extraction** window, you can select the extraction method you want to use (e.g. principal components, etc.). Under **Analyze** ensure that **Correlation Matrix** is selected (this is the default). The default is also to extract eigenvalues over 1 but if you want to extract a specific number of factors you can specify this. Click on **Continue**.
- In the **Rotation** window you can select your rotation method (as mentioned above, **Varimax** is the most common). You can also ask SPSS to display the rotated solution. Once you have selected this click on **Continue**.
- In the **Scores** window you can specify whether you want SPSS to save factor scores for each observation (this will save them as new variables in the data set). Under **Method** choose **Regression**. You can also ask SPSS to display the factor score coefficients (the a_{is}). Click on **Continue**.
- **OK**

5 References

- Manly, B.F.J. (2005), **Multivariate Statistical Methods: A primer**, Third edition, Chapman and Hall.
- Rencher, A.C. (2002), **Methods of Multivariate Analysis**, Second edition, Wiley.

3 Types of data and measures of distance

The data used in cluster analysis can be interval, ordinal or categorical. However, having a mixture of different types of variable will make the analysis more complicated. This is because in cluster analysis you need to have some way of measuring the distance between observations and the type of measure used will depend on what type of data you have.

A number of different measures have been proposed to measure 'distance' for binary and categorical data. For details see the book by Everitt, Landau and Leese. Readers are also referred to this text for details of what to do if you have a mixture of different data types. For interval data the most common distance measure used is the **Euclidean distance**.

3.1 Euclidean distance

In general, if you have p variables X_1, X_2, \dots, X_p measured on a sample of n subjects, the observed data for subject i can be denoted by $x_{i1}, x_{i2}, \dots, x_{ip}$ and the observed data for subject j by $x_{j1}, x_{j2}, \dots, x_{jp}$. The Euclidean distance between these two subjects is given by

$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

When using a measure such as the Euclidean distance, the scale of measurement of the variables under consideration is an issue, as changing the scale will obviously effect the distance between subjects (e.g. a difference of 10cm could being a difference of 100mm). In addition, if one variable has a much wider range than others then this variable will tend to dominate. For example, if body measurements had been taken for a number of different people, the range (in mm) of heights would be much wider than the range in wrist circumference, say. To get around this problem each variable can be standardised (converted to z-scores). However, this in itself presents a problem as it tends to reduce the variability (distance) between clusters. This happens because if a particular variable separates observations well then, by definition, it will have a large variance (as the between cluster variability will be high). If this variable is standardised then the separation between clusters will become less. Despite this problem, many textbooks do recommend standardisation. If in doubt, one strategy would be to carry out the cluster analysis twice — once without standardising and once with — to see how much difference, if any, this makes to the resulting clusters.

4 Hierarchical agglomerative methods

Within this approach to cluster analysis there are a number of different methods used to determine which clusters should be joined at each stage. The main methods are summarised below.

- Nearest neighbour method (single linkage method)

In this method the distance between two clusters is defined to be the distance between the two closest members, or neighbours. This method is relatively simple but is often criticised because it doesn't take account of cluster structure and can result in a problem called **chaining** whereby clusters end up being long and straggly. However, it is better than the other methods when the natural clusters are not spherical or elliptical in shape.

- **Furthest neighbour method (complete linkage method)**
In this case the distance between two clusters is defined to be the maximum distance between members — i.e. the distance between the two subjects that are furthest apart. This method tends to produce compact clusters of similar size but, as for the nearest neighbour method, does not take account of cluster structure. It is also quite sensitive to outliers.
- **Average (between groups) linkage method (sometimes referred to as UPGMA)**
The distance between two clusters is calculated as the average distance between all pairs of subjects in the two clusters. This is considered to be a fairly robust method.
- **Centroid method**
Here the centroid (mean value for each variable) of each cluster is calculated and the distance between centroids is used. Clusters whose centroids are closest together are merged. This method is also fairly robust.
- **Ward's method**
In this method all possible pairs of clusters are combined and the sum of the squared distances within each cluster is calculated. This is then summed over all clusters. The combination that gives the lowest sum of squares is chosen. This method tends to produce clusters of approximately equal size, which is not always desirable. It is also quite sensitive to outliers. Despite this, it is one of the most popular methods, along with the average linkage method.

It is generally a good idea to try two or three of the above methods. If the methods agree reasonably well then the results will be that much more believable.

4.1 Selecting the optimum number of clusters

As stated above, once the cluster analysis has been carried out it is then necessary to select the 'best' cluster solution. There are a number of ways in which this can be done, some rather informal and subjective, and some more formal. The more formal methods will not be discussed in this handout. Below, one of the informal methods is briefly described.

When carrying out a hierarchical cluster analysis, the process can be represented on a diagram known as a **dendrogram**. This diagram illustrates which clusters have been joined at each stage of the analysis and the distance between clusters at the time of joining. If there is a large jump in the distance between clusters from one stage to another then this suggests that at one stage clusters that are relatively close together were joined whereas, at the following stage, the clusters that were joined were relatively far apart. This implies that the optimum number of clusters may be the number present just before that large jump in distance. This is easier to understand by actually looking at a dendrogram — see references for further information.

5 Non-hierarchical or k-means clustering methods

In these methods the desired number of clusters is specified in advance and the 'best' solution is chosen. The steps in such a method are as follows:

1. Choose initial cluster centres (essentially this is a set of observations that are far apart — each subject forms a cluster of one and its centre is the value of the variables for that subject).
2. Assign each subject to its 'nearest' cluster, defined in terms of the distance to the centroid.
3. Find the centroids of the clusters that have been formed
4. Re-calculate the distance from each subject to each centroid and move observations that are not in the cluster that they are closest to.
5. Continue until the centroids remain relatively stable.

Non-hierarchical cluster analysis tends to be used when large data sets are involved. It is sometimes preferred because it allows subjects to move from one cluster to another (this is not possible in hierarchical cluster analysis where a subject, once assigned, cannot move to a different cluster). Two disadvantages of non-hierarchical cluster analysis are: (1) it is often difficult to know how many clusters you are likely to have and therefore the analysis may have to be repeated several times and (2) it can be very sensitive to the choice of initial cluster centres. Again, it may be worth trying different ones to see what impact this has.

One possible strategy to adopt is to use a hierarchical approach initially to determine how many clusters there are in the data and then to use the cluster centres obtained from this as initial cluster centres in the non-hierarchical method.

6 Carrying out cluster analysis in SPSS

6.1 Hierarchical cluster analysis

- **Analyze**
- **Classify**
- **Hierarchical cluster**
- Select the variables you want the cluster analysis to be based on and move them into the **Variable(s)** box.
- In the **Method** window select the clustering method you want to use. Under **Measure** select the distance measure you want to use and, under **Transform values**, specify whether you want all variables to be standardised (e.g. to z-scores) or not.
- In the **Statistics** window you can specify whether you want to see the **Proximity Matrix** (this will give the distance between all observations in the data set — only really recommended for relatively small data sets!). You can also specify whether you want the output to include details of cluster membership — either for a fixed number of clusters or for a range of cluster solutions (e.g. 2 to 5 clusters).
- In the **Save** window you can specify whether you want SPSS to save details of cluster membership — again, either for a fixed number of clusters or for a range of cluster solutions (e.g. 2 to 5 clusters). If you ask it to do this, this information will be included as additional variables at the end of the data set.
- In the **Plots** window you can specify which plots you would like included in the output.
- **OK**

6.2 K-means cluster analysis

- **Analyze**
- **Classify**
- **K-means cluster**
- Select the variables you want the cluster analysis to be based on and move them into the **Variable(s)** box.
- Under **Method**, ensure that **Iterate and Classify** is selected (this is the default).
- In the **Iterate** window you can specify how many iterations you would like SPSS to perform before stopping. The default is ten. It might be worth leaving it as ten to start with and then increasing this if convergence doesn't occur (i.e. a stable cluster solution is not reached) within ten iterations.
- In the **Save** window you can specify whether you want SPSS to save details of cluster membership and distance to the cluster centre for each subject (observation).
- **OK**

7 References

- Everitt, B.S., Landau, S. and Leese, M. (2001), **Cluster Analysis**, Fourth edition, Arnold.
- Manly, B.F.J. (2005), **Multivariate Statistical Methods: A primer**, Third edition, Chapman and Hall.
- Rencher, A.C. (2002), **Methods of Multivariate Analysis**, Second edition, Wiley.

Statistics: 1.1 Paired t-tests

Rosie Shier. 2004.

1 Introduction

A paired t-test is used to compare two population means where you have two samples in which observations in one sample can be paired with observations in the other sample. Examples of where this might occur are:

- Before-and-after observations on the same subjects (e.g. students' diagnostic test results before and after a particular module or course).
- A comparison of two different methods of measurement or two different treatments where the measurements/treatments are applied to the same subjects (e.g. blood pressure measurements using a stethoscope and a dynamap).

2 Procedure for carrying out a paired t-test

Suppose a sample of n students were given a diagnostic test before studying a particular module and then again after completing the module. We want to find out if, in general, our teaching leads to improvements in students' knowledge/skills (i.e. test scores). We can use the results from our sample of students to draw conclusions about the impact of this module in general.

Let x = test score before the module, y = test score after the module

To test the null hypothesis that the true mean difference is zero, the procedure is as follows:

1. Calculate the difference ($d_i = y_i - x_i$) between the two observations on each pair, making sure you distinguish between positive and negative differences.
2. Calculate the mean difference, \bar{d} .
3. Calculate the standard deviation of the differences, s_d , and use this to calculate the standard error of the mean difference, $SE(\bar{d}) = \frac{s_d}{\sqrt{n}}$
4. Calculate the t-statistic, which is given by $T = \frac{\bar{d}}{SE(\bar{d})}$. Under the null hypothesis, this statistic follows a t-distribution with $n - 1$ degrees of freedom.
5. Use tables of the t-distribution to compare your value for T to the t_{n-1} distribution. This will give the p-value for the paired t-test.

NOTE:

For this test to be valid the differences only need to be approximately normally distributed. Therefore, it would not be advisable to use a paired t-test where there were any extreme outliers.

Example

Using the above example with $n = 20$ students, the following results were obtained:

Student	Pre-module score	Post-module score	Difference
1	18	22	+4
2	21	25	+4
3	16	17	+1
4	22	24	+2
5	19	16	-3
6	24	29	+5
7	17	20	+3
8	21	23	+2
9	23	19	-4
10	18	20	+2
11	14	15	+1
12	16	15	-1
13	16	18	+2
14	19	26	+7
15	18	18	0
16	20	24	+4
17	12	18	+6
18	22	25	+3
19	15	19	+4
20	17	16	-1

Calculating the mean and standard deviation of the differences gives:

$$\bar{d} = 2.05 \text{ and } s_d = 2.837. \text{ Therefore, } SE(\bar{d}) = \frac{s_d}{\sqrt{n}} = \frac{2.837}{\sqrt{20}} = 0.634$$

So, we have:

$$t = \frac{2.05}{0.634} = 3.231 \quad \text{on 19 df}$$

Looking this up in tables gives $p = 0.004$. Therefore, there is strong evidence that, on average, the module does lead to improvements.

3 Confidence interval for the true mean difference

The in above example the estimated average improvement is just over 2 points. Note that although this is statistically significant, it is actually quite a small increase. It would be useful to calculate a confidence interval for the mean difference to tell us within what limits the true difference is likely to lie. A 95% confidence interval for the true mean difference is:

$$\bar{d} \pm t^* \frac{s_d}{\sqrt{n}} \quad \text{or, equivalently} \quad \bar{d} \pm (t^* \times SE(\bar{d}))$$

where t^* is the 2.5% point of the t-distribution on $n - 1$ degrees of freedom.

Using our example:

We have a mean difference of 2.05. The 2.5% point of the t-distribution with 19 degrees of freedom is 2.093. The 95% confidence interval for the true mean difference is therefore:

$$2.05 \pm (2.093 \times 0.634) = 2.05 \pm 1.33 = (0.72, 3.38)$$

This confirms that, although the difference in scores is statistically significant, it is actually relatively small. We can be 95% sure that the true mean increase lies somewhere between just under one point and just over 3 points.

4 Carrying out a paired t-test in SPSS

The simplest way to carry out a paired t-test in SPSS is to compute the differences (using **Transform, Compute**) and then carrying out a one-sample t-test as follows:

- **Analyze**
- **Compare Means**
- **One-Sample T Test**
- Choose the difference variable as the **Test Variable** and click **OK**

The output will look like this:

One-Sample Statistics

	N	Mean	Std. Deviation	Std. Error Mean
Difference	20	2.0500	2.83725	.63443

One-Sample Test

Test Value = 0						
	t	df	Sig. (2-tailed)	Mean Difference	95% Confidence Interval	
					Lower	Upper
Difference	3.231	19	.004	2.0500	.7221	3.3779