

Introduction to Statistics with SPSS (15.0)

**Table of contents**

| | |
|---|-----------|
| Introduction | 3 |
| Chapter 1: Opening SPSS for the first time..... | 5 |
| An Excel file | 5 |
| A Text or a Data file | 6 |
| Step one | 6 |
| Step two | 7 |
| Step three | 7 |
| Step four | 8 |
| Step five | 8 |
| Step six | 9 |
| Reading data from a database | 9 |
| Typing all your data in the data editor | 9 |
| <i>Exercise</i> | 9 |
| Chapter 2: Basic structure of an SPSS data file | 10 |
| Data view | 10 |
| Variable view | 10 |
| <i>Exercise</i> | 11 |
| Chapter 3: SPSS Data Editor Menu | 12 |
| File | 12 |
| Edit and View | 12 |
| Data | 12 |
| Transform | 13 |
| Analyse and Graphs | 13 |
| Chapter 4: Qualitative data | 14 |
| Graph | 14 |
| <i>Exercise</i> | 15 |
| A bit of theory: the χ^2 test | 17 |
| A bit of theory: the null hypothesis and the error types | 21 |
| Chapter 5: Quantitative data | 23 |
| 5-1 A bit of theory: Assumptions of parametric data | 23 |
| How can you check that your data are parametric/normal? | 24 |
| <i>Example</i> | 24 |
| 5-2 A bit of theory: descriptive stats | 27 |
| The mean | 27 |
| The variance | 28 |
| The Standard Deviation | 28 |
| Standard Deviation vs. Standard Error | 29 |
| Confidence interval | 29 |
| <i>Quantitative data representation</i> | 30 |
| 5-3 A bit of theory: the t-test | 31 |
| Independent t-test | 33 |
| Paired t-test | 34 |
| <i>Exercise</i> | 34 |
| <i>Exercise</i> | 35 |
| 5-4 Comparison of more than 2 means: Analysis of variance | 36 |
| A bit of theory | 36 |
| <i>Exercise</i> | 39 |
| 5-5 Correlation | 45 |
| <i>Example</i> | 45 |
| A bit of theory: Correlation coefficient | 46 |
| EXERCISES | 50 |



Licence

This manual is © 2007-8, Anne Segonds-Pichon.

This manual is distributed under the creative commons Attribution-Non-Commercial-Share Alike 2.0 licence. This means that you are free:

- to copy, distribute, display, and perform the work
- to make derivative works

Under the following conditions:

- Attribution. You must give the original author credit.
- Non-Commercial. You may not use this work for commercial purposes.
- Share Alike. If you alter, transform, or build upon this work, you may distribute the resulting work only under a licence identical to this one.

Please note that:

- For any reuse or distribution, you must make clear to others the licence terms of this work.
- Any of these conditions can be waived if you get permission from the copyright holder.
- Nothing in this license impairs or restricts the author's moral rights.

Full details of this licence can be found at

<http://creativecommons.org/licenses/by-nc-sa/2.0/uk/legalcode>



Introduction

SPSS is the officially supported statistical package at Babraham. SPSS stands for “Statistical Package for the Social Sciences” as it was first designed by a psychologist. It has evolved a lot since then and is now widely used in many areas though a lot of the literature you can find on internet is still more related to psychology or social epidemiology than other areas.

It is a straight forward package with a friendly environment. There is a lot of easy to access documentation and the tutorials are very good.

However, unlike some other statistical packages, SPSS does not hold your hand all the way through your analysis. You have to make your own decisions and for that you need to have a basic knowledge of stats. The down side of this is that you can make mistakes but the up side is that you actually understand what you are doing. You are not just answering questions by clicking on window after window, you are doing your analysis for real, which means that you understand (well, more or less!) the analytical process but also when it comes to writing down your results, you will know exactly what to say. And, don't worry, if you are unsure about which test to choose or if you can apply the one you have chosen, you can always come to us.

Don't forget: you use stats to present your data in a comprehensible way and to make your point; this is just a tool, so don't hate it, use it!



Chapter 1: Opening SPSS for the first time

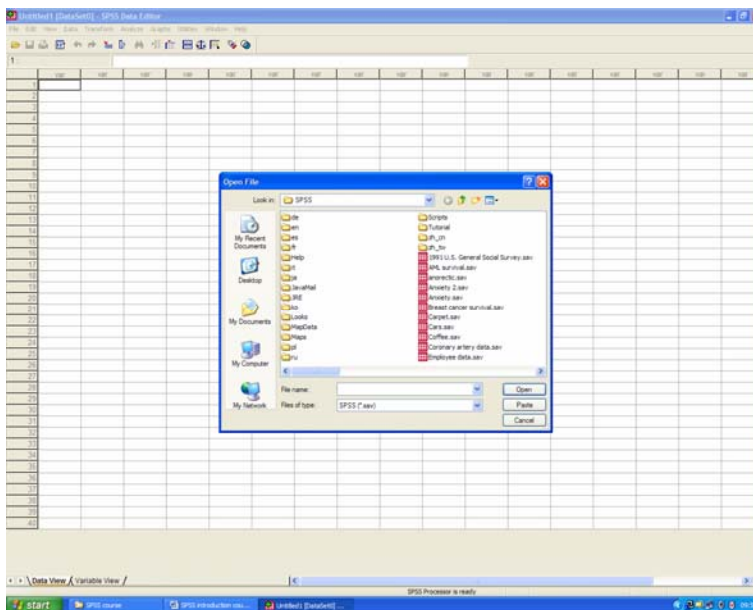
Click on the SPSS icon:

- a small window opens, giving you several choices: Run a tutorial, Type in data ...or opening existing SPSS files. If it is the first time you have run SPSS, it is likely you are not working on an SPSS file yet (!), it is then easier to close the window and to go to the file menu (top left of the screen) to look the file you want to work on :

file

open

data



By default it will look into the SPSS folder so unless you want to look at one of the example files, you want to go somewhere else. If you have never used SPSS before, you are likely to have your data stored as Excel, Text or Data files, so you have to select the format from the Type of files dropdown list (or select All Files if you are unsure).

An Excel file

If you are opening an Excel file, a window will appear and you will have to specify which worksheet your data are on and, if you don't want to import all of them, the range. By default SPSS will read variable names from the first row of data.

Tips: Make sure the work sheet you are opening only contains data (and graphs, or summary stats ...) and that the variable names are in the first row.

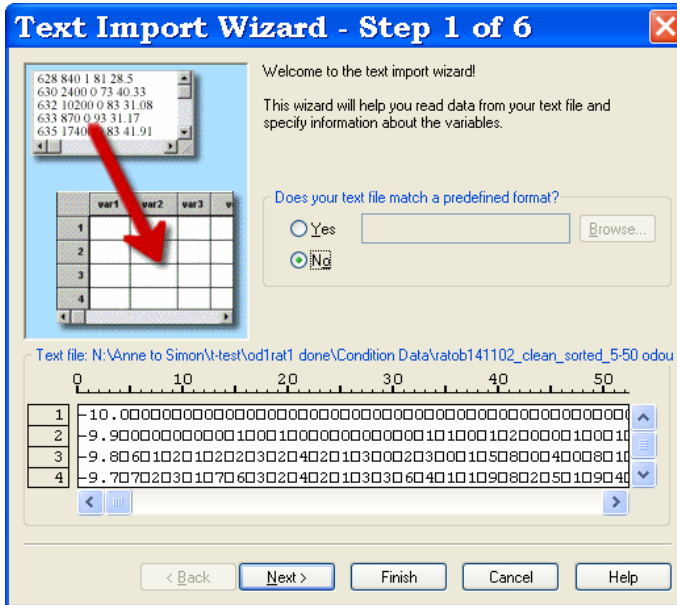
If you have formulas instead of values in some cells, SPSS will accept it but the variable(s) may be considered as string and not numerical data so you may have to change it before you start your analysis.

Finally, do not forget to close your Excel file before opening it through SPSS. It does not like to share!



A Text or a Data file

You open it the same way as an Excel file but instead of opening straight away, you will have to go through a Text Import Wizard:



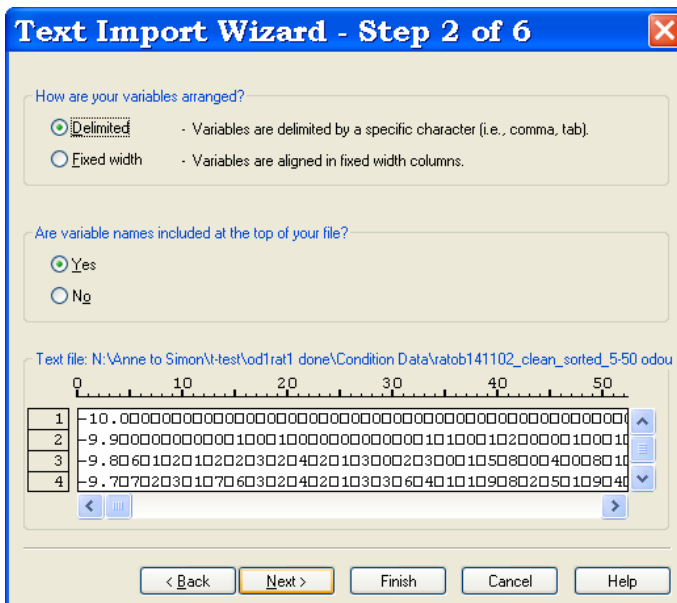
Step one

Just to check you are opening the right file.

It also asks if your text file matches a predefined format. If you know that you will be generating a lot of text files that will need to be imported into SPSS **in the same format**, then it is worth saving the processing of the file so that you only need to go through all the steps once. Now, it is the first time, so let's go through the steps.

At the bottom, is a data preview window showing you how your data look like at each step, so it should be ugly at the beginning and exactly as you want it at the end.

Click next.



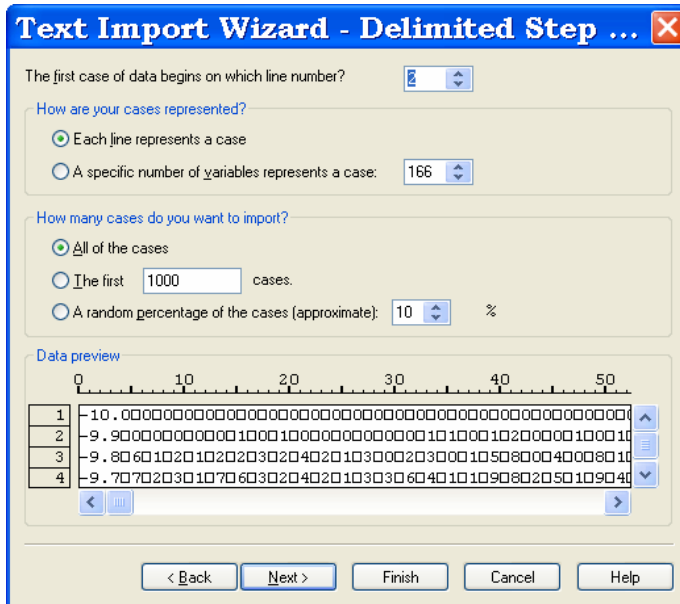


Step two

By default SPSS will consider your variables to be delimited by a specific character, which is usually the case.

Then it will ask if the variable names are included at the top of the file. By default SPSS says “no” but usually they are so you can change it to “yes”.

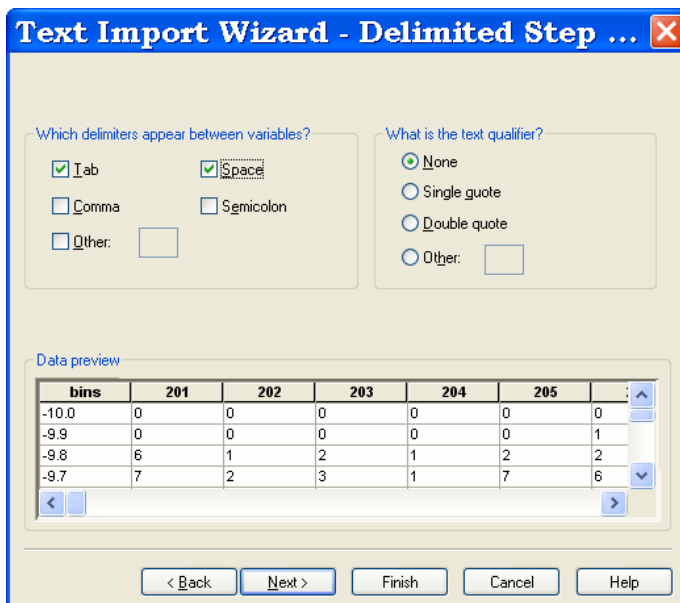
Click next.



Step three

The default settings on this window are usually the one you need: each line represents a case (i.e. all the data on one line correspond to a condition or an experiment or an animal) and you want to import all the cases. If not you can choose otherwise.

Go next.



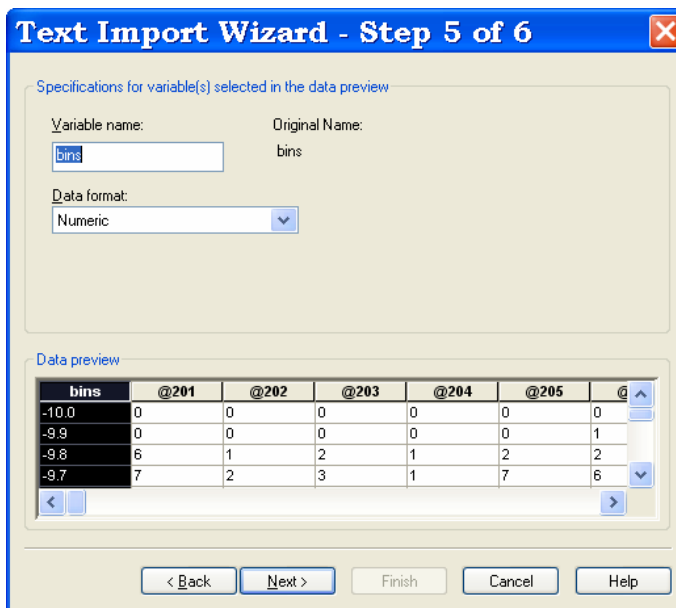


Step four

Your data should start to look better.

Which delimiters appear between variables? The default setting says “Tab” and “Space”, which will usually work on most data. If you are unsure, “play” with it by changing the settings (with and without the “Space” for instance) and see how your data look.

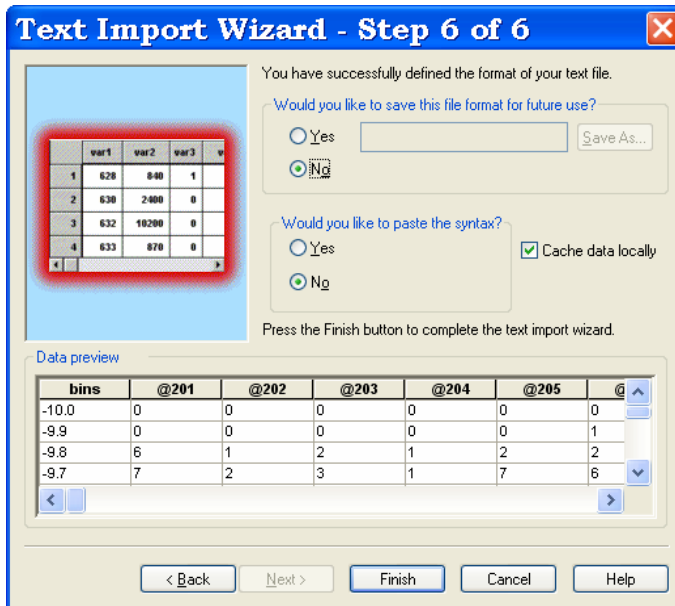
When you click on next, SPSS may tell you that some of the variable names are invalid. This can happen if, for instance, the variable name is numerical (see example above). If you click on OK, SPSS will transform the variable name(s) into valid one(s) (for instance by adding @ before the variable names it does not like) and the original column headings will be saved as variable labels (we will go back to this later). If you are not happy with the new variable names, you will be able to change them later.



Step five

You need to specify the format of your data (numerical, string ...). A numerical data is a number whereas a string data is any string of characters (e.g. for a variable mouse type, the data values could be either “wild type” or “mutant”).

Go next.



Step six

Your data should look perfect!

You can save the file format if you think you'll need it in the future. It will be saved as a normal file. So the next time you need to do the same file processing, in the first window, you answer yes to the question: does your text file match a predefined format?, then you browse, you select your format and click straight on Finish.

Reading data from a database

Alternatively, you can import your data from a database such as Access, using the Open Database command in the file menu. We will not go into any details since a previous knowledge of database system is needed.

Typing all your data in the data editor

Finally, you can type in your data directly into the SPSS data editor. There will be no problem afterwards to export it into Excel if you want to share your data with someone who does not have SPSS on his computer.

Exercise

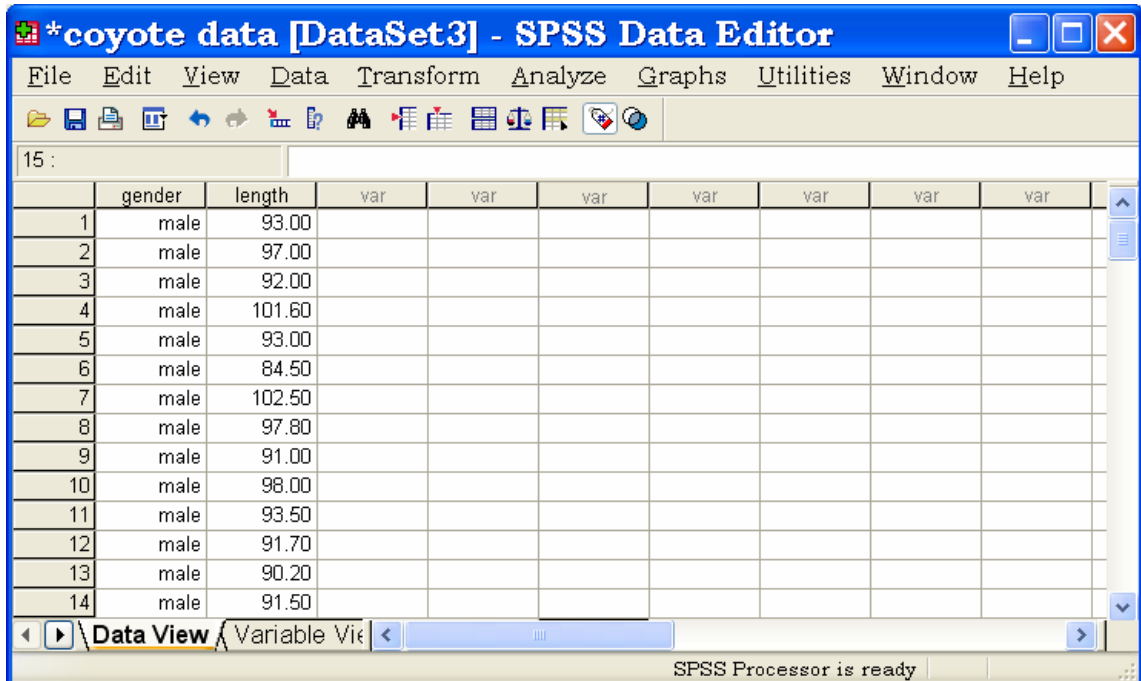
Import data from an Excel file: [cats and dogs.xls](#) and from a Text file: [coyote.txt](#)



Chapter 2: Basic structure of an SPSS data file

Unlike in Excel, SPSS files have 2 “sides”: the **Data view** which looks very much like an Excel file and a **Variable view** which is a kind of “behind the scenes” thing.

Data view



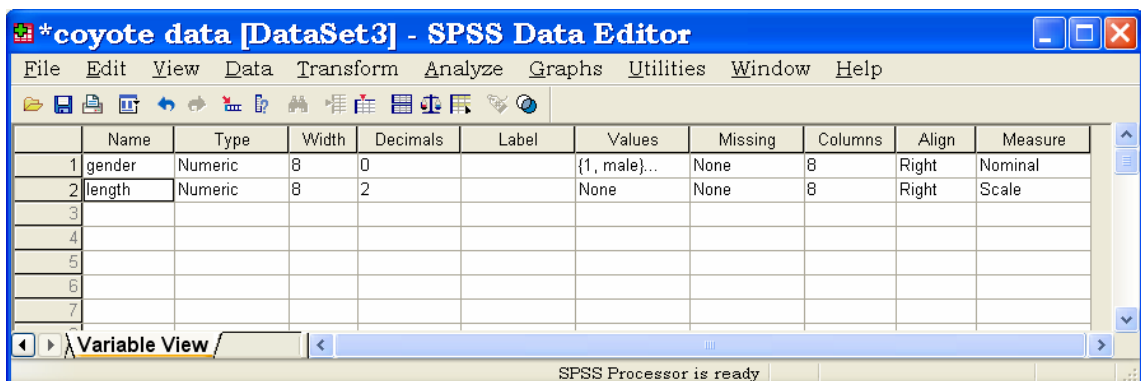
The screenshot shows the SPSS Data Editor window titled '*coyote data [DataSet3] - SPSS Data Editor'. The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Window, and Help. The toolbar contains various icons for file operations and data manipulation. The main area displays a data grid with 15 rows and 10 columns. The first three columns are labeled 'gender', 'length', and 'var'. The data for the first 14 rows is as follows:

| | gender | length | var | var | var | var | var | var | var |
|----|--------|--------|-----|-----|-----|-----|-----|-----|-----|
| 1 | male | 93.00 | | | | | | | |
| 2 | male | 97.00 | | | | | | | |
| 3 | male | 92.00 | | | | | | | |
| 4 | male | 101.60 | | | | | | | |
| 5 | male | 93.00 | | | | | | | |
| 6 | male | 84.50 | | | | | | | |
| 7 | male | 102.50 | | | | | | | |
| 8 | male | 97.80 | | | | | | | |
| 9 | male | 91.00 | | | | | | | |
| 10 | male | 98.00 | | | | | | | |
| 11 | male | 93.50 | | | | | | | |
| 12 | male | 91.70 | | | | | | | |
| 13 | male | 90.20 | | | | | | | |
| 14 | male | 91.50 | | | | | | | |

The status bar at the bottom indicates 'SPSS Processor is ready'.

In Data View, columns represent variables (e.g. gender, length), and rows represent cases (observations such as the sex and the length of the third coyote).

Variable view



The screenshot shows the SPSS Data Editor window titled '*coyote data [DataSet3] - SPSS Data Editor' in Variable View. The menu bar and toolbar are the same as in the Data View screenshot. The main area displays a table with 11 columns: Name, Type, Width, Decimals, Label, Values, Missing, Columns, Align, and Measure. The data for the first two variables is as follows:

| | Name | Type | Width | Decimals | Label | Values | Missing | Columns | Align | Measure |
|---|--------|---------|-------|----------|-------|--------------|---------|---------|-------|---------|
| 1 | gender | Numeric | 8 | 0 | | {1, male}... | None | 8 | Right | Nominal |
| 2 | length | Numeric | 8 | 2 | | None | None | 8 | Right | Scale |
| 3 | | | | | | | | | | |
| 4 | | | | | | | | | | |
| 5 | | | | | | | | | | |
| 6 | | | | | | | | | | |
| 7 | | | | | | | | | | |

The status bar at the bottom indicates 'SPSS Processor is ready'.

This is where you define the variables you will be using: to define/modify a property of a given variable, you click on the cell containing the property you want to define/modify.

You can modify:

- the name and the type of your variable,



- the width, which corresponds to the number of characters you can have in a cell,
 - the decimals, which corresponds to the number of decimals recorded,
- Tip: when importing data from Excel, SPSS would sometimes give extravagant number of decimals, like 12. Don't forget to check that before you start drawing graphs or analysing your data, otherwise you will be unable to read some of analysis outputs and you will get ugly graphs.
- the label is used when you want to define a variable more accurately or to describe it. In the example above, the label "length" could be "length of the body".
 - the values: useful for categorical data (e.g. gender: male=1 and female=2). This is quite an important characteristic:
 - o some analyses will not accept a string variable as a factor,
 - o when you draw a graph from your data, if you have not defined any values, you will only see numerical values on the x-axis. For example, you measure the level of a substance in 5 types of cell and you plot it. If you have not specified any values you'll get a x-axis with numbers from 1 to 5 instead of having the names of the types of cell.
 - o you will need to remember that you decided that male=1 and female=2!
 - missing: useful for epidemiological questionnaires,
 - column (see width),
 - align: like Excel: right, left or centre,
 - measure: scale (e.g. weight: quantitative variable), ordinal (e.g. no, a little, a lot) or nominal (e.g. male or female: qualitative variable).

Exercise (File: **cats and dogs.sav**)

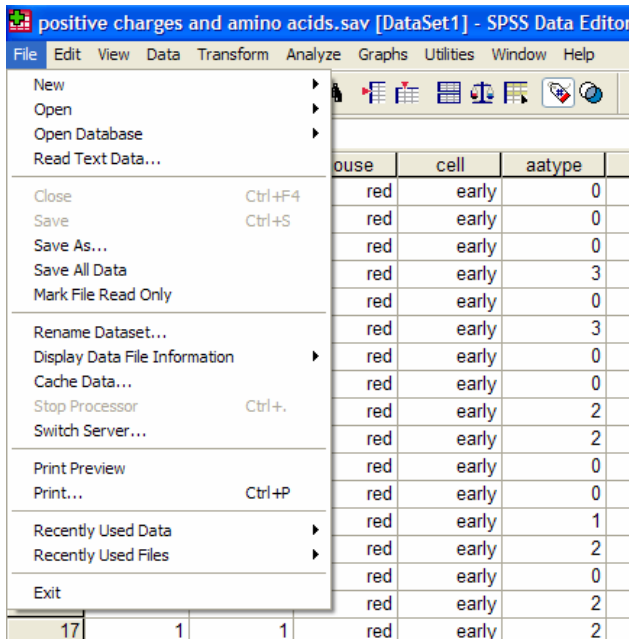
Recode the variables so that animal:1=cat and 2=dog, dance: 1=yes and 2=no and training: 1=food and 2=affection.

Label training as "Type of training" and dance as "Did they dance?".

Make sure the each variable in your file corresponds to the correct measure.



Chapter 3: SPSS Data Editor Menu



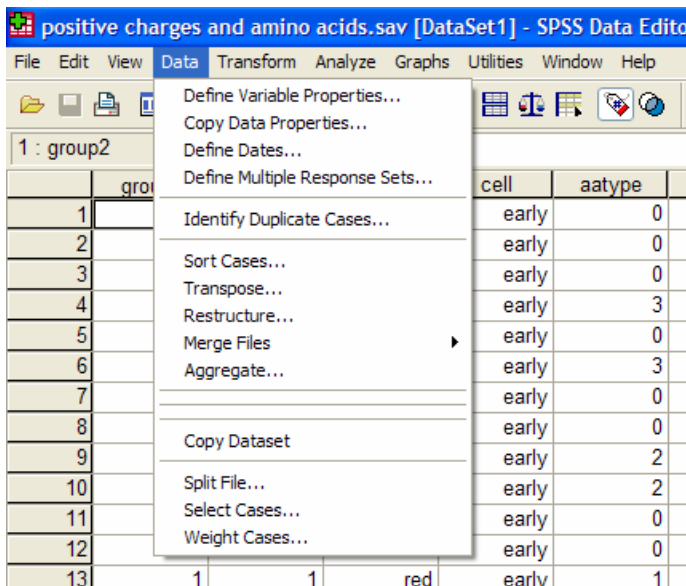
File

Same type of file menu as in Excel: you open and close files, save them, print data and have a look at the recently used files.

Edit and View

Very much like any Edit or View menu in a Window environment.

Data



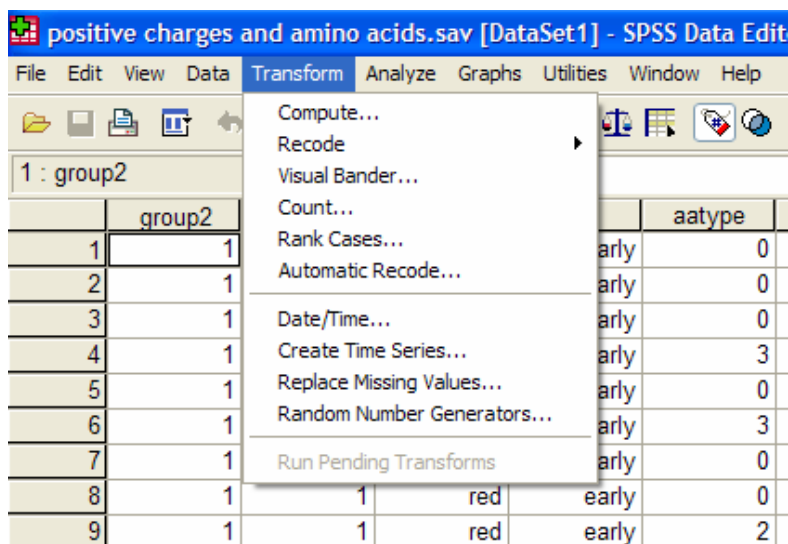


This is the menu which will allow you to tailor your data before the analysis.

The functions you will be likely to use the most:

- Sort cases: can also be accessed by right-clicking on the variable name,
- Transpose and Restructure: you can either restructure selected variables into cases or restructure selected cases into variables or transpose all data. You go through a Restructure Data Wizard. Tip: be careful with this function: instead of creating a new file, SPSS modifies your working file! So if you want to keep your original structure make sure you save the new one onto another name,
- Merge files: you can either add variables or cases. Tip: Make sure for the latest that the files have the exact same structure, including the variable properties: if a variable is a string in one file and a numeric one in the other file, they will be considered as 2 separate variables.
- Split File: could be very useful when you want to do several time the same analysis, like for each gender or for each cell types,
- Select cases: you can select the range of data that you want to look at.

Transform



- Compute variables: use the Compute dialog box to compute values for a variable based on numeric transformations of other variables (e.g. if you need to work out the log function of an existing variable).
- Recode into same variable or into differ rent variable: allows you to reassign the values of existing variables (categorical variables) or collapse ranges of existing values into new values (quantitative variables).

Analyse and Graphs

We will go through these menus in the following chapters.



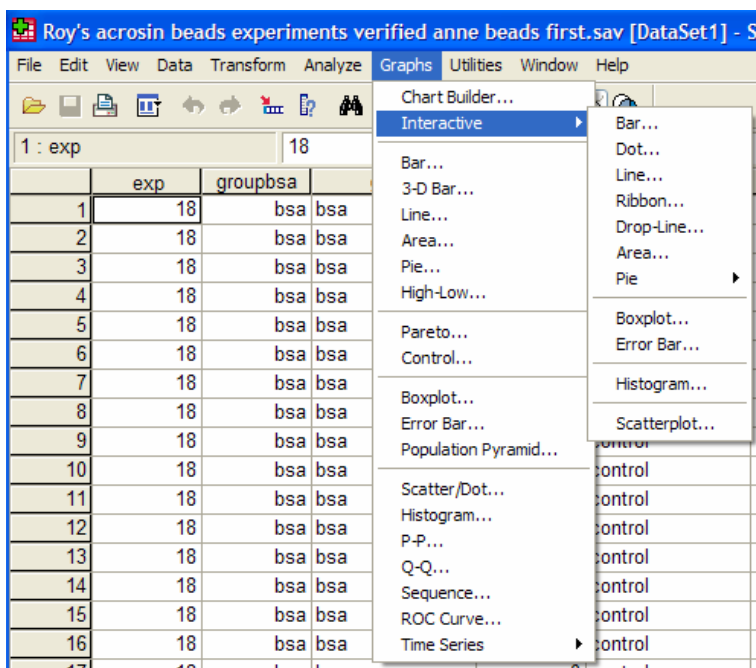
Chapter 4: Qualitative data

Now you know how to import data into SPSS and how to look at your data file. So it is time to talk about the data themselves. The first thing you need to do good stats is to know your data inside out. They are generally organised into variables, which can be divided into 2 categories: *qualitative* and *quantitative*.

Qualitative data are non numerical data and the values taken are usually names (also *nominal* data) (e.g. variable sex: male or female). The values can be numbers but not numerical (e.g. an experiment number is a numerical label but not a unit of measurement). A qualitative variable with intrinsic order in their categories is *ordinal*. Finally, there is the particular case of qualitative variable with only 2 categories, it is then said to be *binary* or *dichotomous* (e.g. alive/dead or male/female).

OK, so let's say you have collected your data and entered/imported them into SPSS. The first thing to do is to see how they look like. In order to do that, you have to go into the Graph menu.

Graph



This menu allows you to build different types of graphs from your data. What I tend to use the most is the interactive function: if you click on it, you get a sub menu from which you can choose the type of graph you want to build. It is very easy to use and very quick to play with if you want to look at your data through "different angles".

**Exercise** (File: **cats and dogs.sav**)

A researcher is interested in whether animals could be trained to line dance. He took some cats and dogs (**animal**) and tried to train them to dance by giving them either food or affection as a reward (**training**) for dance-like behaviour. At the end of the week a note was made of which animal could line dance and which could not (**dance**). All the variables are dummy variables (categorical).

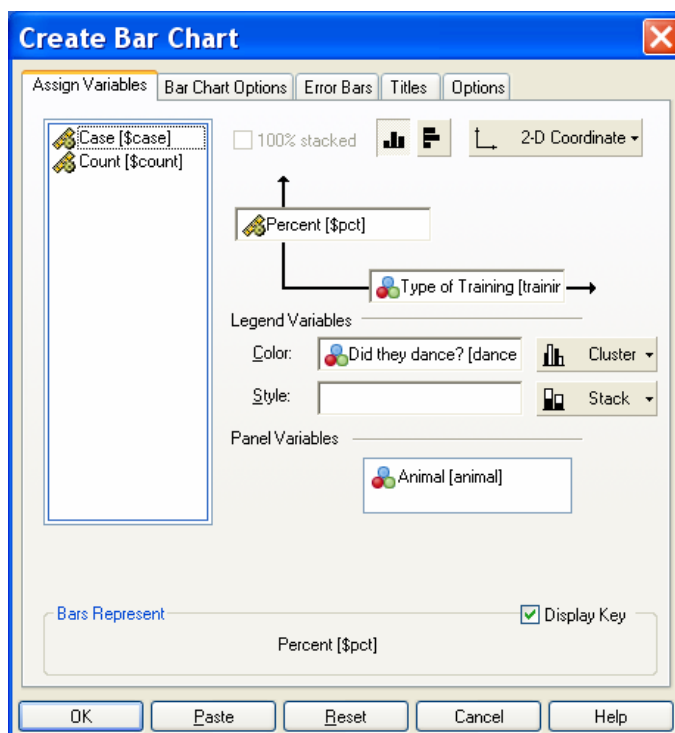
Is there an effect of training on dogs and cats' ability to learn to line dance?

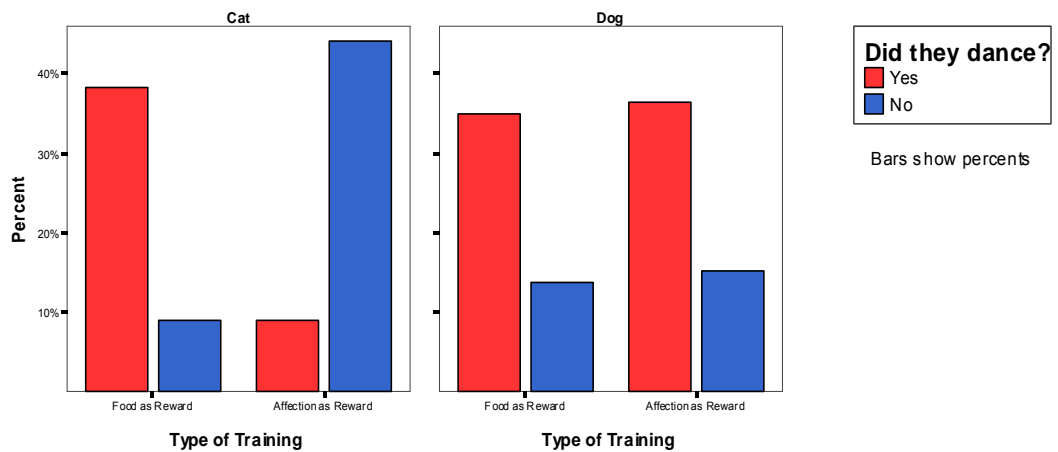
Plot the data so that you have one graph per species.

First, the bar chart: you go into Graph > Interactive > Bar. All you have to do is drag the variables from the list to the appropriate space.

A useful tool is the "Panel Variables" thing as it allows you to build several graphs in one go. It can be useful if you have made 3 or 4 times the same experiment, for example, and you want to have a quick look at the consistence of your results across your experiments.

Tip: you can put several variables in the panel variables window but with more than 2 it starts getting messy.

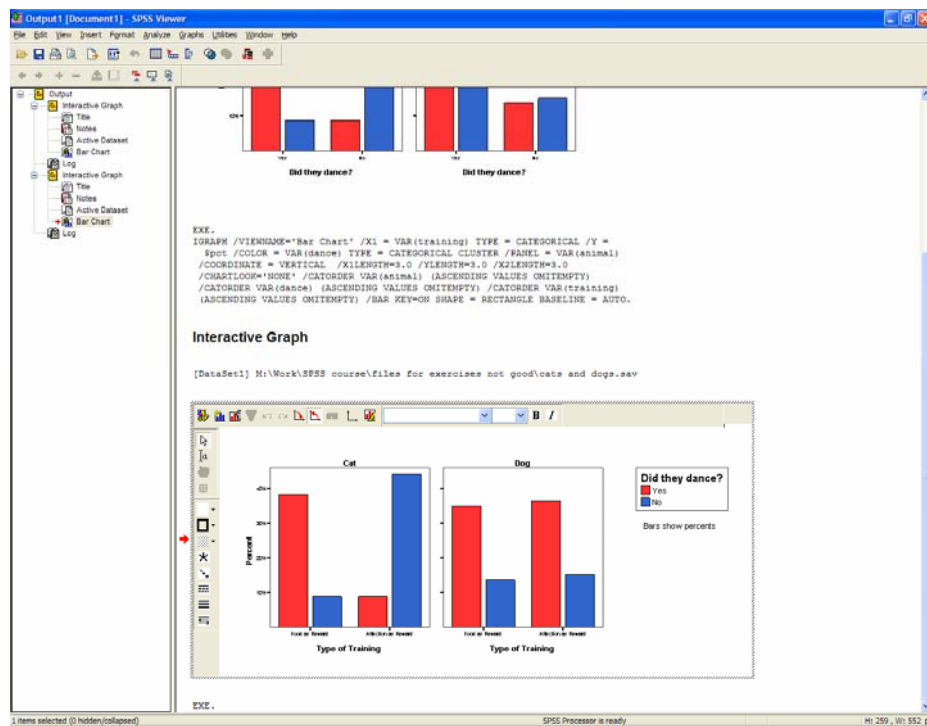





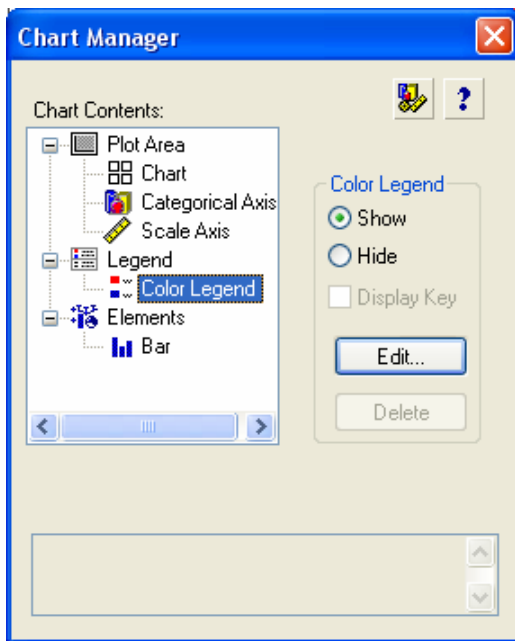
So clearly, from the graphs, you can say that there is an effect of training on cats but not on dogs. Now, you want to know if this effect is significant and to do so you need a χ^2 test.

About SPSS output:

The viewer window is divided into 2 panes. The outline pane (on the left) contains an outline of all of the information stored in the Viewer. If you have done several graphs/analyses on SPSS, you can scroll up and down and select the graph or the table you want to see on the contents pane (on the right), from which you can scroll up and down as well.



You can modify a graph by double-clicking on it. When the graph is "activated" you can either click on the bit you want to change (e.g. the y-axis) or choose the chart manager (top left corner ) from which you can choose any part of the graph, select it and go to Edit to make the changes.



A bit of theory: the χ^2 test

It could be either:

- a one-way χ^2 test, which is basically a test that compares the observed frequency of a variable in a single group with what would be the expected by chance.
- a two-way χ^2 test, the most widely used, in which the observed frequencies for two or more groups are compared with expected frequencies by chance. In other words, in this case, the χ^2 tells you whether or not there is an association between 2 categorical variables.

If you run a χ^2 on SPSS, it will do it in one step and will give you the level of significance of your test right away. But for you to understand what it is about, let's do it step by step.

Step 1: the contingency table

Some packages work out the χ^2 from such a table but SPSS will do it from the raw data.

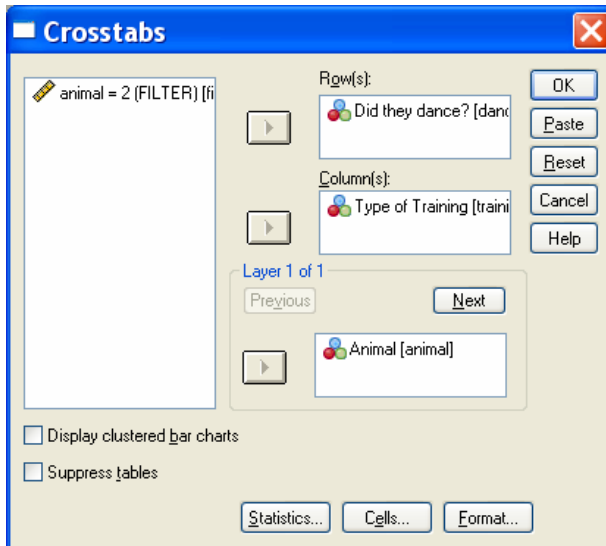
To obtain a contingency table with SPSS, you go: Analyse > Descriptive Statistics > Crosstabs.

An important thing to know about the χ^2 is that it does not tell you anything about causality; it is simply measuring the strength of the association between 2 variables and it is your knowledge of the biological system you are studying which will help you to interpret the result. Hence, you generally have an idea of which variable is acting the other.

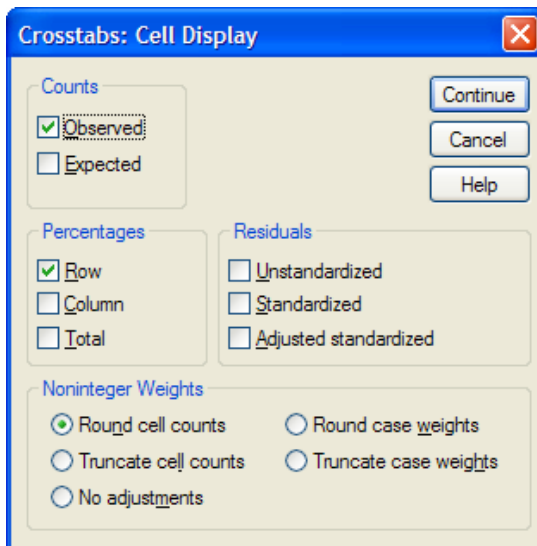
Traditionally in SPSS, the variable which you think is going to act on the other is put in rows. This variable is called the independent variable or the predictor as, in your hypothesis, its values will predict some of the variations of the other variable. The latter, also called the outcome or the dependent variable, as it depends on the values of the predictor, is in column.

The layer function allows you to run several tests at the same time.

So in our particular case (cats and dogs experiment), we should get the window below by simply dragging the variables.



It is likely that you want to express your results in percentages. To do so, you click on Cells (at the bottom of the Crosstabs window) and you get the following menu:



In this particular example, the comparison that makes more sense is the one between type of reward so, you choose the percentages in row and you get the table below.

Did they dance? * Type of Training * Animal Crosstabulation

| | | | | Type of Training | | Total |
|-------|-----------------|--------------------------|--------------------------|------------------|---------------------|--------|
| | | | | Food as Reward | Affection as Reward | |
| Cat | Did they dance? | Yes | Count | 26 | 6 | 32 |
| | | | % within Did they dance? | 81.3% | 18.8% | 100.0% |
| | No | Count | 6 | 30 | 36 | |
| | | % within Did they dance? | 16.7% | 83.3% | 100.0% | |
| Total | | Count | 32 | 36 | 68 | |
| | | % within Did they dance? | 47.1% | 52.9% | 100.0% | |
| Dog | Did they dance? | Yes | Count | 23 | 24 | 47 |
| | | | % within Did they dance? | 48.9% | 51.1% | 100.0% |
| | No | Count | 9 | 10 | 19 | |
| | | % within Did they dance? | 47.4% | 52.6% | 100.0% | |
| Total | | Count | 32 | 34 | 66 | |
| | | % within Did they dance? | 48.5% | 51.5% | 100.0% | |



You are going to use the values in this table to work out the χ^2 value:

$$\chi^2 = \sum \frac{(\text{Observed Frequency} - \text{Expected Frequency})^2}{\text{Expected Frequency}}$$

The observed frequencies are to one you measured, the values that are in your table. Now, you need to calculate the expected ones, which is done this way:

Expected frequency = (row total)*(column total)/grand total

So, for the cat, for example: the expected frequency of cat that would line dance after having received food as reward is:

- probability of line dancing: 32/68
- probability of receiving food: 32/68

So the expected frequency: (32/68)*(32/68) = 15.1

Did they dance? * Type of Training * Animal Crosstabulation

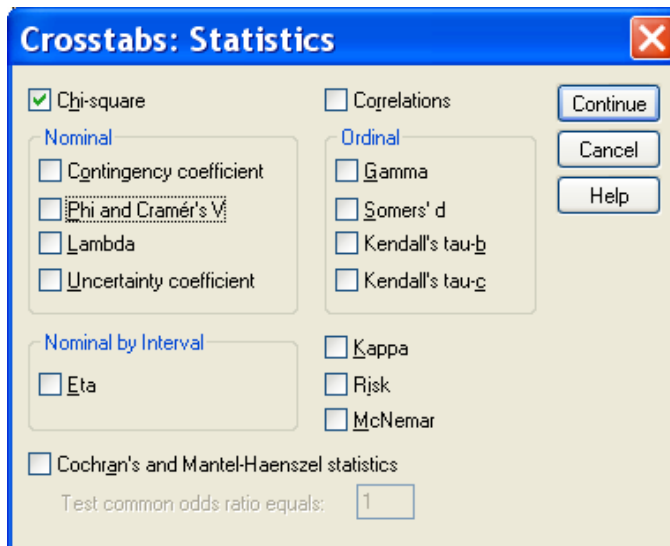
| Animal | | | | Type of Training | | Total |
|--------|-----------------|----------------|----------------|------------------|---------------------|-------|
| | | | | Food as Reward | Affection as Reward | |
| Cat | Did they dance? | Yes | Count | 26 | 6 | 32 |
| | | | Expected Count | 15.1 | 16.9 | 32.0 |
| | No | Count | 6 | 30 | 36 | |
| | | Expected Count | 16.9 | 19.1 | 36.0 | |
| | Total | Count | 32 | 36 | 68 | |
| | | Expected Count | 32.0 | 36.0 | 68.0 | |
| Dog | Did they dance? | Yes | Count | 23 | 24 | 47 |
| | | | Expected Count | 22.8 | 24.2 | 47.0 |
| | No | Count | 9 | 10 | 19 | |
| | | Expected Count | 9.2 | 9.8 | 19.0 | |
| | Total | Count | 32 | 34 | 66 | |
| | | Expected Count | 32.0 | 34.0 | 66.0 | |

Intuitively, one can see that we are kind of averaging things here, we try to find out the values we should have got by chance. If you work out the values for all the cells, you get:

So for the cat, the χ^2 value is:

$$(26-15.1)^2/15.1 + (6-16.9)^2/16.9 + (6-16.9)^2/16.9 + (30-19.1)^2/19.1 = 28.4$$

If you want SPSS to calculate the χ^2 , you click on Statistics at the bottom of the Crosstabs window and select Chi-square. The other options can be ignored today.



Then you get the following output.

Chi-Square Tests

| Animal | | Value | df | Asymp. Sig. (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
|--------|------------------------------------|---------------------|----|-----------------------|----------------------|----------------------|
| Cat | Pearson Chi-Square | 28.363 ^b | 1 | .000 | | |
| | Continuity Correction ^a | 25.830 | 1 | .000 | | |
| | Likelihood Ratio | 30.707 | 1 | .000 | | |
| | Fisher's Exact Test | | | | .000 | .000 |
| | Linear-by-Linear Association | 27.946 | 1 | .000 | | |
| | N of Valid Cases | 68 | | | | |
| Dog | Pearson Chi-Square | .013 ^c | 1 | .908 | | |
| | Continuity Correction ^a | .000 | 1 | 1.000 | | |
| | Likelihood Ratio | .013 | 1 | .908 | | |
| | Fisher's Exact Test | | | | 1.000 | .563 |
| | Linear-by-Linear Association | .013 | 1 | .909 | | |
| | N of Valid Cases | 66 | | | | |

a. Computed only for a 2x2 table

b. 0 cells (.0%) have expected count less than 5. The minimum expected count is 15.06.

c. 0 cells (.0%) have expected count less than 5. The minimum expected count is 9.21.

The line you are interested in is the first one, it gives you the value of the Pearson Chi-square and its level of significance.

Footnote b and c: it relates to the only assumption you have to be careful about when you run a χ^2 : with 2x2 contingency tables you should not have cells with an expected count below 5 as if it is the case it is likely that the test is not accurate (for larger table, all expected counts should be greater than 1 and no more than 20% of expected counts should be less than 5). If you have a high proportion of cells with a small value in it, there are 2 solutions to solve the problem: the first one is to collect more data or, if we have more than 2 categories, to group them to boost the proportions.

If you remember the χ^2 's formula, the calculation gives you an estimation (the Value) of the difference between your data and what you would have obtained if there was no association between your variables. Clearly, the bigger the value of the χ^2 , the bigger the difference between observed and expected frequencies and the more likely to be significant the difference is.



A bit of theory: the null hypothesis and the error types.

The null hypothesis (H_0) corresponds to the absence of effect (e.g.: the animals rewarded by food are as likely to line dance as the ones rewarded by affection) and the aim of a statistical test is to accept or to reject H_0 . Traditionally, a test or a difference are said to be “significant” if the probability of type I error is: $\alpha \leq 0.05$. It means that the level of uncertainty of a test usually accepted is 5%. It also means that there is a probability of 5% that you may be wrong when you say that your 2 means are different, for instance, or you can say that when you see an effect you want to be at least 95% sure that something is significantly happening.

| Statistical decision | True state of H_0 | |
|----------------------|---------------------|---------------|
| | H_0 True | H_0 False |
| Reject H_0 | Type I error | Correct |
| Do not reject H_0 | Correct | Type II error |

Tip: if your p-value is between 5% and 10% (0.05 and 0.10), I would not reject it too fast if I were you. It is often worth putting this result into perspective and asks yourself a few questions like:

- what the literature says about what am I looking at?
- what if I had a bigger sample?
- have I run other tests on similar data and were they significant or not?

The interpretation of a border line result can be difficult as it could be important in the whole picture.

So, for our “cats and dogs experiment”, you are more than 99% sure ($p < 0.0001$) that there is a significant effect of the reward in the ability of cats to learn to line dance.

About SPSS output:

The tables contain many statistical terms for which you can get definition directly from the viewer. To do so, you double-click on the table and then right-click on the word for which you want an explanation (e.g. Fisher’s Exact test). If you click on “What’s this?”, the definition will appear.



| | Value | df | Asymp. Sig. (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
|------------------------------------|---------------------|----|-----------------------|----------------------|----------------------|
| Pearson Chi-Square | 15.696 ^b | 1 | .000 | | |
| Continuity Correction ^a | 14.334 | 1 | .000 | | |
| Likelihood Ratio | 16.137 | 1 | .000 | | |
| Fisher's Exact Test | | | | .000 | .000 |
| Linear-by-Linear Association | | | .000 | | |

a. Computed only for a 2x2 table.
b. 0 cells (.0%) have expected count less than 5. The minimum expected count is 26.

What's This?

- Cut Ctrl+X
- Copy Ctrl+C
- Paste Ctrl+V
- Clear del
- Select
- Show Dimension Label
- Hide Category
- Ungroup
- Group
- Create Graph
- Table Properties...
- Cell Properties...
- TableLooks...
- Insert Footnote
- Delete Footnotes
- Hide Footnotes
- Pivoting Trays
- Toolbar
- Results Coach
- Case Studies

| | Value | df | Asymp. Sig. (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
|------------------------------------|---------------------|----|-----------------------|----------------------|----------------------|
| Pearson Chi-Square | 15.696 ^b | 1 | .000 | | |
| Continuity Correction ^a | 14.334 | 1 | .000 | | |
| Likelihood Ratio | 16.137 | 1 | .000 | | |
| Fisher's Exact Test | | | | .000 | .000 |
| Linear-by-Linear Association | | | .000 | | |

A test for independence in a 2x2 table. It is most useful when the total sample size and the expected values are small.

a. Computed only for a 2x2 table
b. 0 cells (.0%) have expected count less than 5. The minimum expected count is 26.

You can also play around with the tables. To do so you double-click on it and, if the Pivoting Trays window is not visible, you can get it from the menus.

Output1 - SPSS Viewer

File Edit View Insert Pivot Format Analyze Graphs Utilit

Output

- Interactive G
 - Title
 - Notes
 - Active D
 - Bar Char
- Crosstabs
 - Title
 - Notes

Bookmarks

- Transpose Rows and Columns
- Move Layers to Rows
- Move Layers to Columns
- Reset Pivots to Defaults
- Pivoting Trays
- Go to Layer...



Chapter 5: Quantitative data

When it comes to quantitative data, more tests are available but assumptions must be met to apply most of these tests. There are 2 types of stats tests: parametric and non-parametric ones. Parametric tests have 4 assumptions that must be met for the test to be accurate. Non-parametric tests are designed to be used with nominal or ordinal data (e.g. χ^2 test) and they make few or no assumptions about populations parameters (e.g. Mann-Whitney test).

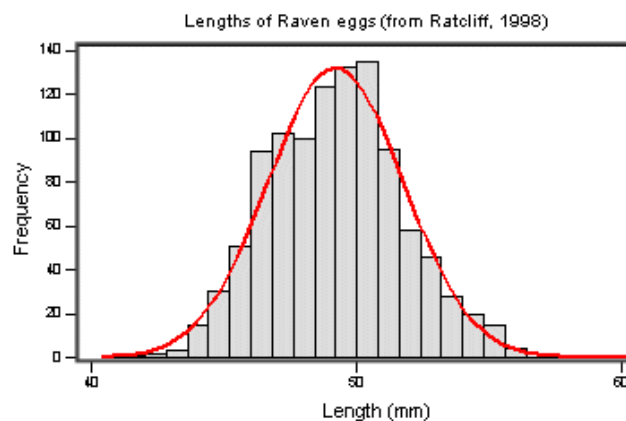
5-1 A bit of theory: Assumptions of parametric data

When you are dealing with quantitative data, the first thing you should look at is how they are distributed, how they look like. The distribution of your data will tell you if there is something wrong in the way you collected them or enter them and it will also tell you what kind of test you can apply to make them say something.

T-test, analysis of variance and correlation tests belong to the family of parametric tests and to be able to use them your data must comply with 4 assumptions.

1) The data have to be normally distributed (normal shape, bell shape, Gaussian shape). Departure from normality can be tested with SPSS. If the test tells you that your data are not normal, transformations can be made to make them suitable for parametric analysis.

Example of normally distributed data:



2) Homogeneity in variance: The variance should not change systematically throughout the data.

3) Interval data: The distance between points of the scale should be equal at all parts along the scale

4) Independence: Data from different subjects are independent so that values corresponding to one subject do not influence the values corresponding to another subject. There are specific designs for repeated measures experiments.

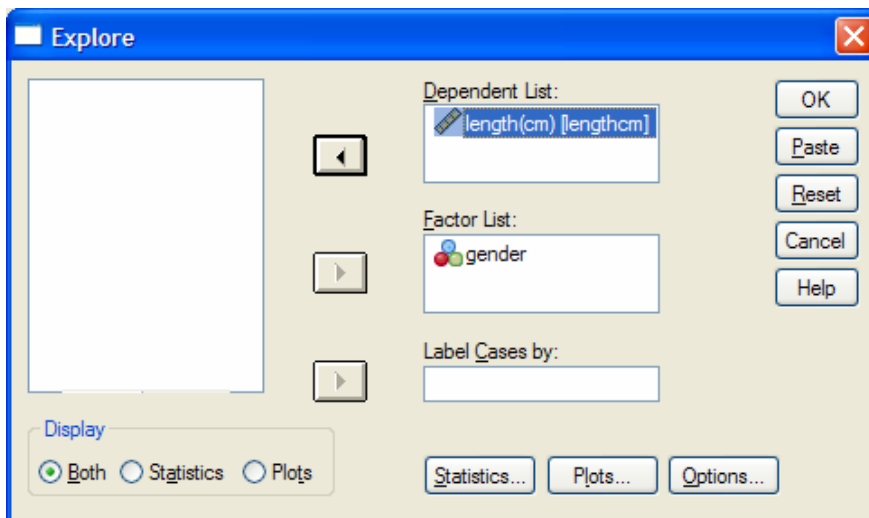


How can you check that your data are parametric/normal?

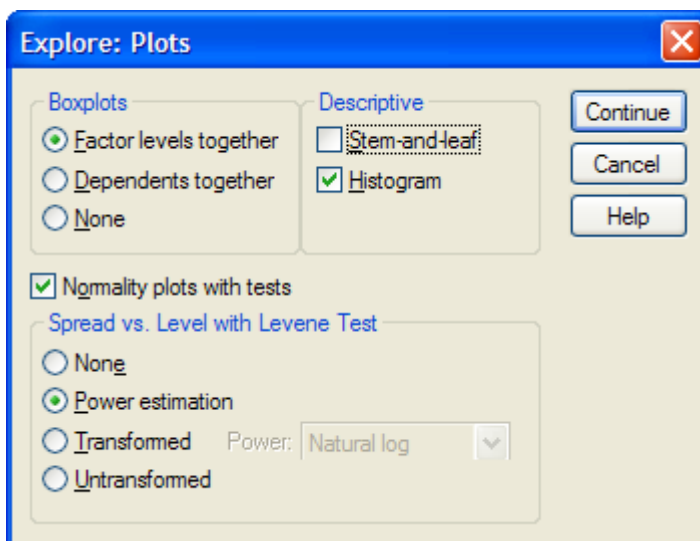
You can use the explore Menu. To do so, you go: Analyse>Descriptive Statistics>Explore. As SPSS says in its help menu: “Exploring data can help to determine whether the statistical techniques that you are considering for data analysis are appropriate. The Explore procedure provides a variety of visual and numerical summaries of the data, either for all cases or separately for groups of cases”. It can be useful to screen data and identify outliers. You can also check assumptions. Basically, it does all what you can find in Frequencies and Descriptives, only better as it is more complete.

Let's try it through an example.

Example (File: *coyote.sav*)



If you want to look at the distribution of your data, you click on Plots and you select Histogram and Normality plots and power estimation.

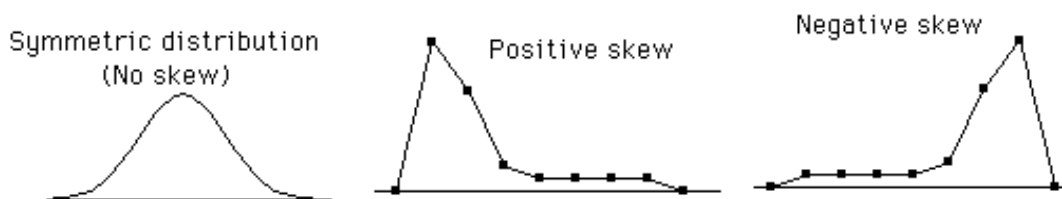


The first output you get is a summary of the descriptive stats. We will go through it in more details later on.



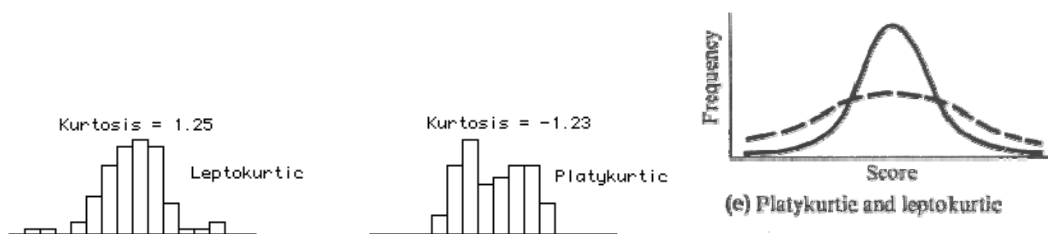
| Descriptives | | | | Statistic | Std. Error |
|---------------------|---------------------|----------------------------------|-------------|-----------|------------|
| length(cm) | Male | Mean | | 92.06 | 1.021 |
| | | 95% Confidence Interval for Mean | Lower Bound | 90.00 | |
| | | | Upper Bound | 94.12 | |
| | | 5% Trimmed Mean | | 92.09 | |
| | | Median | | 92.00 | |
| | | Variance | | 44.836 | |
| | | Std. Deviation | | 6.696 | |
| | | Minimum | | 78 | |
| | | Maximum | | 105 | |
| | | Range | | 27 | |
| | Interquartile Range | | 9 | | |
| | Skewness | | -.091 | .361 | |
| | Kurtosis | | -.484 | .709 | |
| | Female | Mean | | 89.71 | .999 |
| | | 95% Confidence Interval for Mean | Lower Bound | 87.70 | |
| | | | Upper Bound | 91.73 | |
| | | 5% Trimmed Mean | | 89.98 | |
| | | Median | | 90.00 | |
| | | Variance | | 42.900 | |
| | | Std. Deviation | | 6.550 | |
| Minimum | | | 71 | | |
| Maximum | | | 103 | | |
| Range | | | 32 | | |
| Interquartile Range | | 8 | | | |
| Skewness | | -.568 | .361 | | |
| Kurtosis | | .911 | .709 | | |

Skewness: lack of symmetry of a distribution



Kurtosis: measure of the degree of peakedness in the distribution

- The two distributions below have the same variance approximately the same skew, but differ markedly in kurtosis.



Then you get the results of the tests of normality (only relevant if you have around 20 data or more). If the tests are significant it means that there is departure from normality and you should not apply parametric test, unless you transform your data. So in the case of our coyotes, our data seem to be OK.



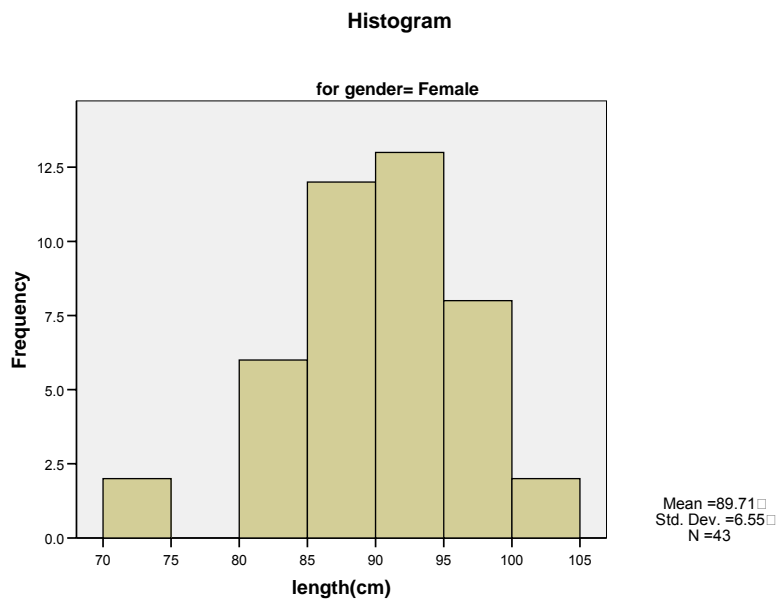
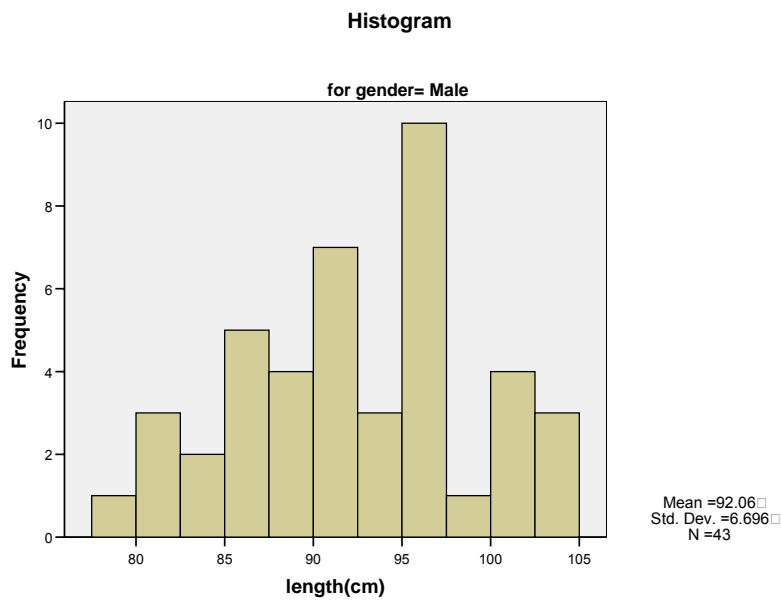
Tests of Normality

| gender | Kolmogorov-Smirnov ^a | | | Shapiro-Wilk | | |
|-------------------|---------------------------------|----|-------|--------------|----|------|
| | Statistic | df | Sig. | Statistic | df | Sig. |
| length(cm) Male | .089 | 43 | .200* | .984 | 43 | .819 |
| length(cm) Female | .078 | 43 | .200* | .970 | 43 | .316 |

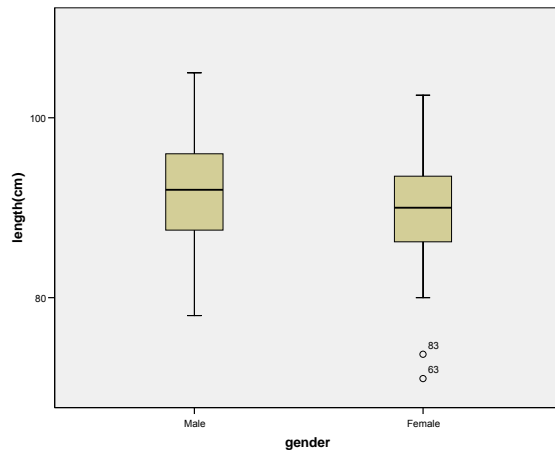
*. This is a lower bound of the true significance.

a. Lilliefors Significance Correction

To make sure, you can have a look at the histograms.



Though it is not the perfect bell shaped we all dream of, it looks OK. Finally, you can have a look at the boxplots.



No need for you to know every thing about boxplots. All you need to know, is that they should look about symmetrical and, when comparing different groups, about the same size as it is useful information for the interpretation of the tests. The other important information is given by the dots away from the boxplots: they are outliers and it is worth having a look at them (typo ...).

Finally, you can check the second assumption (equality of variances).

Test of Homogeneity of Variance

| | Levene Statistic | df1 | df2 | Sig. |
|--------------------------------------|------------------|-----|--------|------|
| length(cm) Based on Mean | .219 | 1 | 81 | .641 |
| Based on Median | .229 | 1 | 81 | .634 |
| Based on Median and with adjusted df | .229 | 1 | 80.423 | .634 |
| Based on trimmed mean | .231 | 1 | 81 | .632 |

In our case, the Levene's test is not significant so the variances are not significantly different from each other.

5-2 A bit of theory: descriptive stats

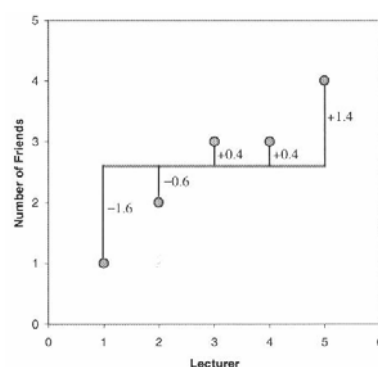
The mean (or average) μ = average of all values in a column

It can be considered as a model because it summaries the data.

- Example: number of friends of each members of a group of 5 lecturers: 1, 2, 3, 3 and 4

Mean: $(1+2+3+3+4)/5 = 2.6$ friends per lecturer: clearly an hypothetical value !

But if the values were: 0, 0, 1, 5 and 7, the mean would also be 2.6 but clearly it would not give an accurate picture of the data. So, how can you know that it is an accurate model? You look at the difference between the real data and your model. To do so, you calculate the difference between the real data and the model created and you make the sum so that you get the total error (or sum of differences).





$\sum(x_i - \mu) = (-1.6) + (-0.6) + (0.4) + (0.4) + (1.4) = 0$ And you get no errors !

Of course: positive and negative differences cancel each other out. So to avoid the problem of the direction of the error, you can square the differences and instead of sum of errors, you get the Sum of Squared errors (SS).

- In our example: $SS = (-1.6)^2 + (-0.6)^2 + (0.4)^2 + (0.4)^2 + (1.4)^2 = 5.20$

The variance

This SS gives a good measure of the accuracy of the model but it is dependent upon the amount of data: the more data, the higher the SS. The solution is to divide the SS by the number of observations (N). As we are interested in measuring the error in the sample to estimate the one in the population, we divide the SS by N-1 instead of N and we get the *variance* (S^2) = SS/N-1

- In our example: Variance (S^2) = 5.20 / 4 = 1.3

Why N-1 instead N?

If we take a sample of 4 scores in a population they are free to vary but if we use this sample to calculate the variance, we have to use the mean of the sample as an estimate of the mean of the population. To do that we have to hold one parameter constant.

- Example: mean of a sample is 10

We assume that the mean of the population from which the sample has been collected is also 10. If we want to calculate the variance, we must keep this value constant which means that the 4 scores cannot vary freely:

- If the values are 9, 8, 11 and 12 (mean = 10) and if we change 3 of these values to 7, 15 and 8 then the final value must be 10 to keep the mean constant.
- If we hold 1 parameter constant, we have to use N-1 instead of N.
- It is the idea behind the *degree of freedom*: one less than the sample size.

The Standard Deviation

The problem with the variance is that it is measured in squared units which is not very nice to manipulate. So for more convenience, the square root of the variance is taken to obtain a measure in the same unit as the original measure: the *standard deviation*.

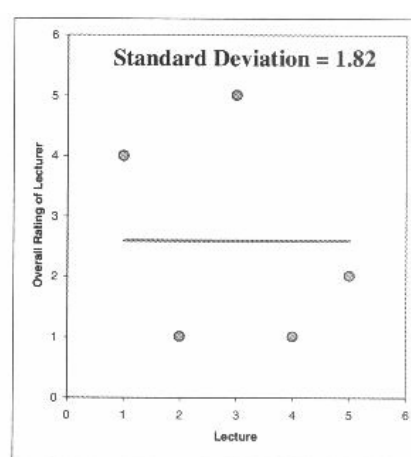
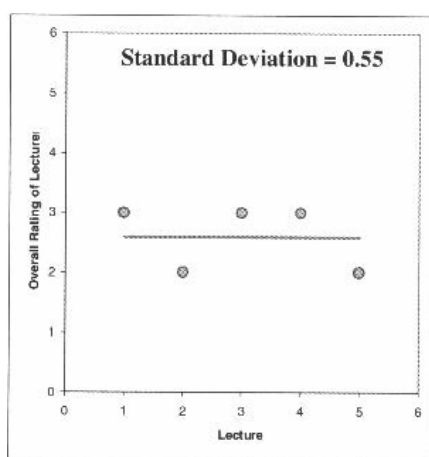
- S.D. = $\sqrt{SS/N-1} = \sqrt{S^2}$, in our example: S.D. = $\sqrt{1.3} = 1.14$

- So you would present your mean as follows: $\mu = 2.6 \pm 1.14$ friends

The standard deviation is a measure of how well the mean represents the data or how much your data are squattered around the mean.:

- small S.D.: data close to the mean: mean is a good fit of the data (graph on the left)

- large S.D.: data distant from the mean: mean is not an accurate representation (graph on the right)

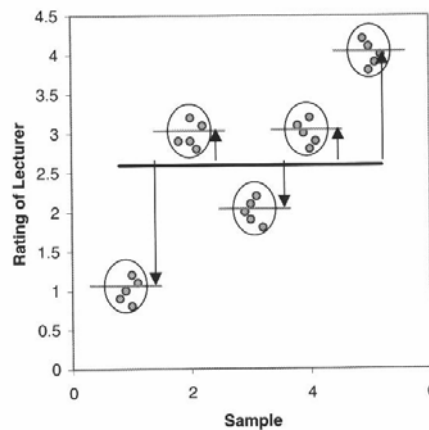
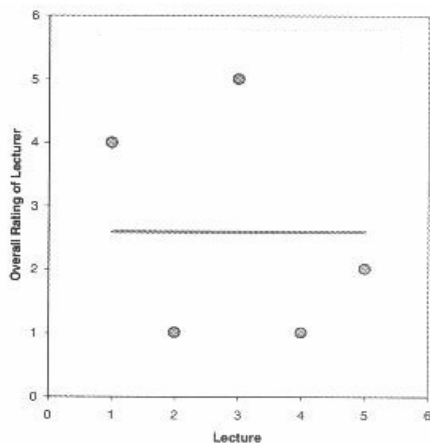




Standard Deviation vs. Standard Error

Many scientists are confused about the difference between the standard deviation (S.D.) and the *standard error of the mean* (S.E.M. = $S.D. / \sqrt{N}$).

- The S.D. (graph on the left) quantifies the scatter of the data and increasing the size of the sample does not increase the scatter (above a certain threshold).
- The S.E.M. (graph on the right) quantifies how accurately you know the true population mean, it's a measure of how much you expect sample means to vary. So the S.E.M. gets smaller as your samples get larger: the mean of a large sample is likely to be closer to the true mean than is the mean of a small sample.



A big S.E.M. means that there is a lot of variability between the means of different samples and that your sample might not be representative of the population.

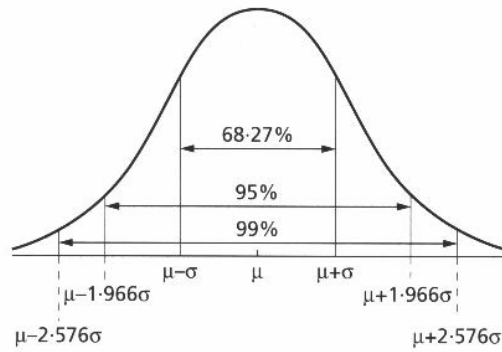
A small S.E.M. means that most samples means are similar to the population mean and so your sample is likely to be an accurate representation of the population.

Which one to choose?

- If the scatter is caused by biological variability, it is important to show the variation. So it is more appropriate to report the S.D. rather than the S.E.M. Even better, you can show in a graph all data points, or perhaps report the largest and smallest value.
- If you are using an in vitro system with no biological variability, the scatter can only result from experimental imprecision (no biological meaning). It is more sensible then to report the S.E.M. since the S.D. is less useful here. The S.E.M. gives your readers a sense of how well you have determined the mean.

Confidence interval

- The confidence interval quantifies the uncertainty in measurement. The mean you calculate from your sample of data points depends on which values you happened to sample. Therefore, the mean you calculate is unlikely to equal the true population mean exactly. The size of the likely discrepancy depends on the variability of the values (expressed as the S.D. or the S.E.M.) and the sample size. If you combine those together, you can calculate a 95% confidence interval (95% CI), which is a range of values. If the population is normal (or nearly so), you can be 95% sure that this interval contains the true population mean.

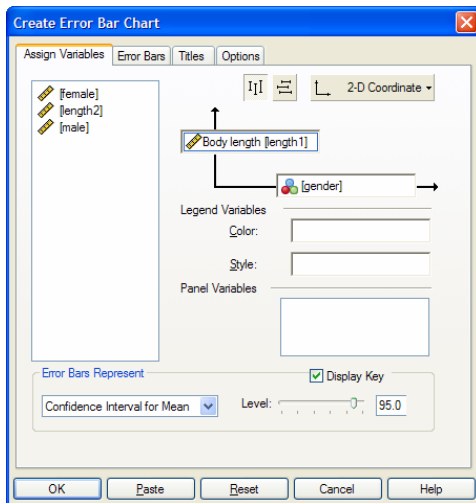


95% of observations in a normal distribution lie within +/- 1,96*SE

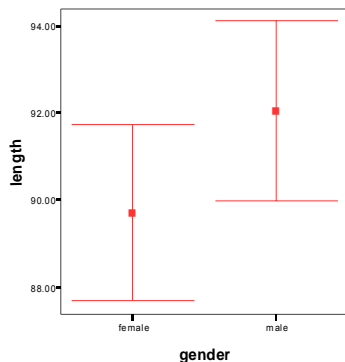
Quantitative data representation

OK so now, you have checked that you data were normally distributed and you know everything about descriptives stats. The next step is to plot your data.

Let's go back to our coyotes. What you want from your graph is to see if there is difference between males and females and possibly, have an idea of the significance of the difference. The best way to do it is to plot the error bars. To do so, you go Graphs>Interactive>Errors bar.



By default SPSS will go for the confidence interval and you will get the following graph.



Error Bars show 95.0% CI of Mean



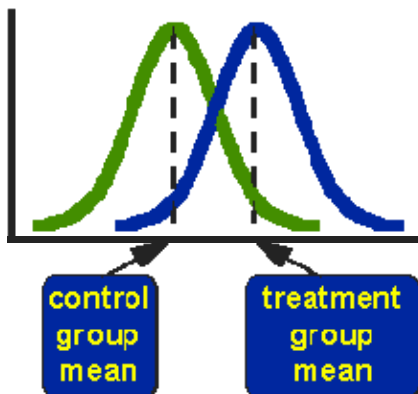
This is a very informative graph as you can spot the 2 means together with the confidence interval. We saw before that the 95% CI of the mean gives you the boundaries between which you 95% sure to find the true population mean. It is always better when you want to compare visually 2 or more groups to use the CI than the SD or the SEM. It gives you a better idea of the dispersion of your sample and it allows you to have an idea, before doing any stats, of the likelihood of a significant difference between your groups. Since your true group means have 95% chances of lying within their respective CI, an overlap between the CI tells you that the difference is probably not significant.

In our particular example, from the graph we can say that the average body length of female coyotes, for instance, is a little bit more than 92 cm and that 95 out of 100 samples from the same population would have means between about 90 and 94 cm. We can also say that despite the fact that the females appear longer than the males, this difference is probably not significant as the errors bars overlap considerably.

To check that, we can run a t-test.

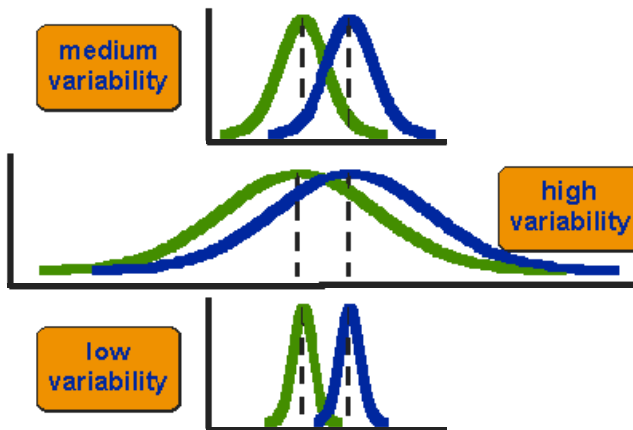
5-3 A bit of theory: the t-test

The t-test assesses whether the means of two groups are *statistically* different from each other. This analysis is appropriate whenever you want to compare the means of two groups.




The figure above shows the distributions for the treated (blue) and control (green) groups in a study. Actually, the figure shows the idealized distribution. The figure indicates where the control and treatment group means are located. The question the t-test addresses is whether the means are statistically different.

What does it mean to say that the averages for two groups are statistically different? Consider the three situations shown in the figure below. The first thing to notice about the three situations is that **the difference between the means is the same in all three**. But, you should also notice that the three situations don't look the same -- they tell very different stories. The top example shows a case with moderate variability of scores within each group. The second situation shows the high variability case. The third shows the case with low variability. Clearly, we would conclude that the two groups appear most different or distinct in the bottom or low-variability case. Why? Because there is relatively little overlap between the two bell-shaped curves. In the high variability case, the group difference appears least striking because the two bell-shaped distributions overlap so much.



This leads us to a very important conclusion: when we are looking at the differences between scores for two groups, we have to judge the difference between their means relative to the spread or variability of their scores. The t-test does just this.

The formula for the t-test is a ratio. The top part of the ratio is just the difference between the two means or averages. The bottom part is a measure of the variability or dispersion of the scores. Figure 3 shows the formula for the t-test and how the numerator and denominator are related to the distributions.

$$\begin{aligned} \frac{\text{signal}}{\text{noise}} &= \frac{\text{difference between group means}}{\text{variability of groups}} \\ &= \frac{\bar{X}_T - \bar{X}_C}{SE(\bar{X}_T - \bar{X}_C)} \\ &= \text{t-value} \end{aligned}$$


The diagram shows a normal distribution curve with two overlapping distributions. A red box highlights the difference between the two means, which corresponds to the numerator of the t-value formula. Another red box highlights the spread of the distributions, which corresponds to the denominator of the t-value formula. Arrows point from these boxes to the corresponding parts of the formula above.

The t-value will be positive if the first mean is larger than the second and negative if it is smaller.

To run a t-test on SPSS, you go: Analysis> Compare means and then you have to choose between different types of t-tests.

You can run a one-sample t-test which is when you want to compare a series of values (from one sample) to 0 for instance.

Then you have Independent-samples t-test and Paired-Samples t-test. The choice between the 2 is very intuitive. If you measure a variable in 2 **different populations**, you choose the independent t-test as the 2 populations are independent from each other. If you measure a variable 2 times in the **same population**, you go for the paired t-test.

So say, you want to compare the level of haemoglobin in the types of mouse (e.g. 2 breeds of sheep in terms of weight). To do so, you take a sample of each breed (the 2 samples have to be comparable) and you weigh each animal. You then run a Independent-samples t-test on your data to find out a difference.



If you want to compare 2 types of sheep food (A and B): you define 2 samples of sheep comparable in any other ways and you weigh them at day 1 and say at day 30. This time you apply a Paired-Samples t-test as you are interested in each individual difference in weight between day 1 and day 30.

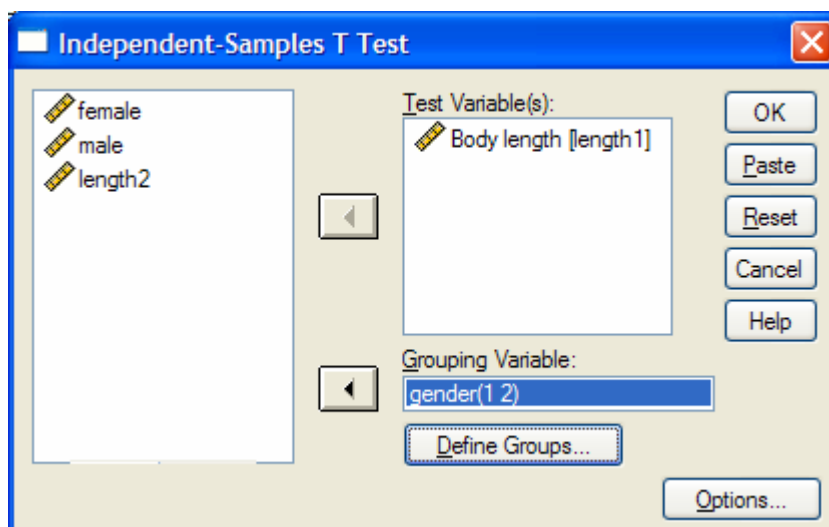
One last thing about the type of t-tests in SPSS: the structure of your data file will depend on your choice.

- If you run an independent t-test, you will need to organise your data in 2 columns as well but one will be a grouping variable and the other one will contain the data. In the sheep example, the grouping variable will be the breed and the data will be entered under the variable weight.
- If you run a paired t-test, you need 2 variables. To go back to the sheep example, you will have your data organised in 2 column: one for day 1 and the other for day 30.

Independent t-test

Let's go back to our example. You go Analysis>Compare means>Independent-samples t-test.

You define the grouping variable by entering the corresponding category: in our example, simply 1 (male) and 2 (female).



When you run the test in SPSS, you get the following out put.

The first table gives you the descriptive stats and the second one the results of the test.

Group Statistics

| | gender | N | Mean | Std. Deviation | Std. Error Mean |
|-------------|--------|----|-------|----------------|-----------------|
| Body length | male | 43 | 89.71 | 6.550 | .999 |
| | female | 43 | 92.06 | 6.696 | 1.021 |



Independent Samples Test

| | | Levene's Test for Equality of Variances | | t-test for Equality of Means | | | | | | |
|-------------|-----------------------------|---|------|------------------------------|--------|-----------------|-----------------|-----------------------|---|-------|
| | | F | Sig. | t | df | Sig. (2-tailed) | Mean Difference | Std. Error Difference | 95% Confidence Interval of the Difference | |
| | | | | | | | | | Lower | Upper |
| Body length | Equal variances assumed | .152 | .698 | -1.641 | 84 | .105 | -2.344 | 1.428 | -5.185 | .496 |
| | Equal variances not assumed | | | -1.641 | 83.959 | .105 | -2.344 | 1.428 | -5.185 | .496 |

A few words of explanation about this second table:

- Levene's test for equality of variance: we have seen before that the t-test compares the 2 means taking into account the variability within the groups. We have also seen that parametric test assume that the variances in experimental groups are roughly equals. Intuitively, one can see that if there is much more variability in 1 group than the other, the comparison between the means will be trickier. Fortunately there are adjustments that can be made in situations in which the variances are not equal. The Levene's test tells you if the variances are significantly different or not. In our case, the variances are considered as equal ($p=0.698$) so we can read the results of the t-test in the row "Equal variances assumed". Otherwise we would have looked at the results in the row below.

- $t = -1.641$ is the value of your t-test with 84 degrees of freedom and a p-value of 0.105 which tells you that the difference between males and females is not significant.

- Sig. (2-tailed) gives you the p-value of the test and 2-tailed means that you are looking at a difference either way.

NB: 1-tailed tests are mostly used in medical studies where researchers want to know if a treatment improves or not the condition of a patient.

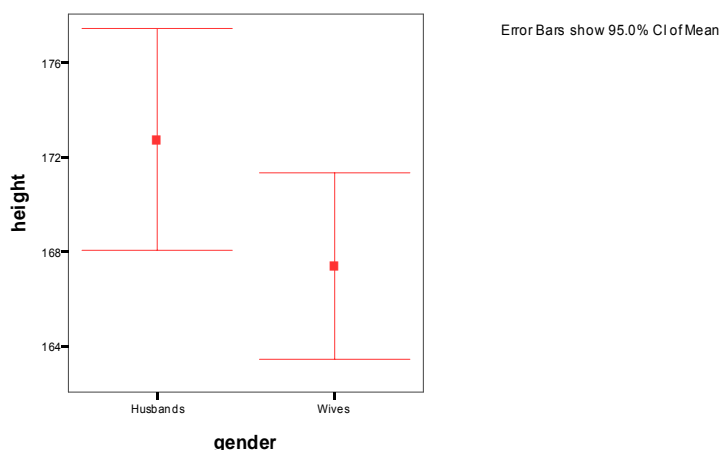
Paired t-test

Now let's try a Paired t-test. As we mentioned before, the idea behind the paired t-test is to look at a difference between 2 paired individuals or 2 measures for a same individual. For the test to be significant, the difference must be different from 0.

Exercise (File: [height husband wife.xls](#))

Import the data and make sure that the variables have the right measures and that 1=husband and 2=wife. Then plot the data.

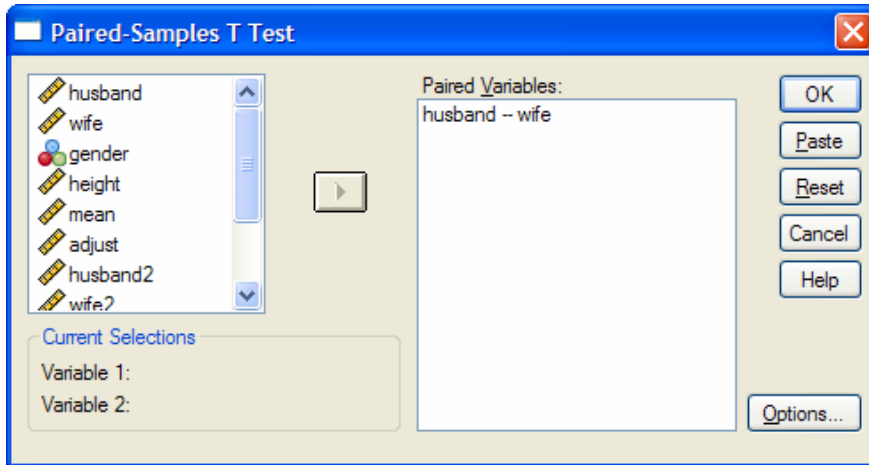
If everything goes right, you get the graph below.





From this graph, we can conclude that if husbands are taller than wives, this difference is not significant. So let's run a paired t-test to get a p-value.

To be able to run the test from the file, you are going to have a bit of copy and paste, so that you have 1 column with the values for the husbands and 1 with the values of the wives.



Paired Samples Statistics

| | | Mean | N | Std. Deviation | Std. Error Mean |
|--------|---------|--------|----|----------------|-----------------|
| Pair 1 | husband | 172.75 | 20 | 10.057 | 2.249 |
| | wife | 167.40 | 20 | 8.401 | 1.878 |

Paired Samples Test

| | | Paired Differences | | | | | t | df | Sig. (2-tailed) |
|--------|----------------|--------------------|----------------|-----------------|---|-------|-------|----|-----------------|
| | | Mean | Std. Deviation | Std. Error Mean | 95% Confidence Interval of the Difference | | | | |
| | | | | | Lower | Upper | | | |
| Pair 1 | husband - wife | 5.350 | 4.580 | 1.024 | 3.206 | 7.494 | 5.224 | 19 | .000 |

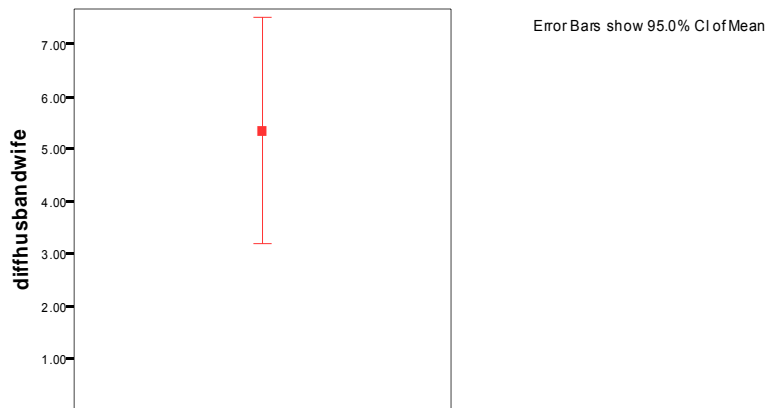
SPSS output for the paired t-test gives you the mean difference between husbands and wives pair wise. So you can say that on average husbands are 5.350 cm taller than their wives and that 95 out of 100 samples of the same population would have this mean difference between 3.206 cm and 7.494 cm. This interval does not include 0 which means that we can be pretty sure that the difference between the 2 groups is significant. This is confirmed by the p-value ($p < 0.0001$) which says that the test is highly significant.

So, how come the graph and the test tell us different things?

The problem is that we don't really want to compare the mean size of the wives to the mean size of the husband, we want to look at the difference pair-wise, in other words we want to know if, on average, a given wife is taller or smaller than her husband. So we are interested in the mean difference between husband and wife.

Exercise

Build a variable which contains the difference between the size of a husband and the one of his wife. Plot the difference on a graph and try to guess the result of the paired t-test.



With only 2 groups, you do not get a very nice graph but it is informative enough for you to see that the confidence interval does not include 0, so you are almost certain the result of the t-test is going to be significant.

Try to run a One Sample t-test.

5-4 Comparison of more than 2 means: Analysis of variance

A bit of theory

When we want to compare more than 2 means (e.g. more than 2 groups), we cannot run several t-test because it increases the **familywise error rate** which is the error rate across tests conducted on the same experimental data.

Example: if you want to compare 3 groups (1, 2 and 3) and you carry out 3 t-tests (groups 1-2, 1-3 and 2-3), each with an arbitrary 5% level of significance, the probability of not making the type I error is 95% (= 1 - 0.05). The 3 tests being independent, you can multiply the probabilities, so the overall probability of no type I errors is: $0.95 * 0.95 * 0.95 = 0.857$. Which means that the probability of making at least one type I error (to say that there is a difference whereas there is not) is $1 - 0.857 = 0.143$ or 14.3%. So the probability has increased from 5% to 14.3%. If you compare 5 groups instead of 3, the familywise error rate is 40% (= $1 - (0.95)^n$)

To overcome the problem of multiple comparisons, you need to run an **Analysis of variance (ANOVA)**, which is an extension of the 2 group comparison of a t-test but with a slightly different logic. If you want to compare 5 means, for example, you can compare each mean with another, which gives you 10 possible 2-group comparisons, which is quite complicated ! So, the logic of the t-test cannot be directly transferred to the analysis of variance. Instead the ANOVA compares variances: if the variance amongst the 5 means is greater than the random error variance (due to individual variability for instance), then the means must be more spread out than we would have explained by chance.

The statistic for ANOVA is the F ratio:

$$F = \frac{\text{variance among sample means}}{\text{variance within samples (=random. Individual variability)}}$$



also:

$$F = \frac{\text{variation explained by the model (systematic)}}{\text{variation explained by unsystematic factors}}$$

If the variance amongst sample mean is greater than the error variance, then $F > 1$. In an ANOVA, you test whether F is significantly higher than 1 or not.

Imagine you have a dataset of 50 data points, you make the hypothesis that these points in fact belong to 5 different groups (this is your hypothetical model). So you arrange your data into 5 groups and you run an ANOVA.

| Group | Data |
|-------|------|
| 1 A | 64 |
| 2 A | 1.52 |
| 3 A | 99 |
| 4 A | 17 |
| 5 A | 75 |
| 6 A | 1.48 |
| 7 A | 1.18 |
| 8 A | 33 |
| 9 A | 1.42 |
| 10 A | 2.09 |
| 11 B | 1.37 |
| 12 B | 1.23 |
| 13 B | 26 |
| 14 B | 47 |
| 15 B | 42 |
| 16 B | 64 |
| 17 B | 32 |
| 18 B | 65 |
| 19 B | 43 |
| 20 B | 67 |
| 21 C | 43 |
| 22 C | 70 |
| 23 C | 79 |
| 24 C | 89 |
| 25 C | 24 |
| 26 C | 25 |
| 27 C | 1.01 |
| 28 C | 77 |
| 29 C | 47 |
| 30 C | 47 |
| 31 D | 46 |
| 32 D | 65 |
| 33 D | 41 |
| 34 D | 81 |
| 35 D | 1.20 |
| 36 D | 1.08 |
| 37 D | 34 |
| 38 D | 1.98 |
| 39 D | 1.39 |

You get the table below.

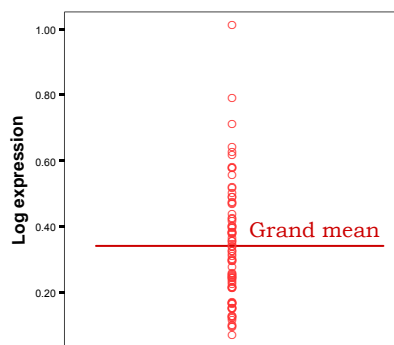
| Source of variation | Sum of Squares | df | Mean Square | F | Sig. |
|---------------------|----------------|----|-------------|-------|------|
| Between Groups | 351.520 | 4 | 87.880 | 9.085 | .000 |
| Within Groups | 435.300 | 45 | 9.673 | | |
| Total | 786.820 | 49 | | | |

Typical example of analyse of variance table

Let's go through the figures in the table. First the bottom row of the table:

$$\text{Total sum of squares} = \sum (x_i - \text{Grand mean})^2$$

In our case, Total SS = 786.820. If you were to plot your data to represent the total SS, you would produce the graph below. So the total SS is the squared sum of all the differences between each data point and the grand mean.



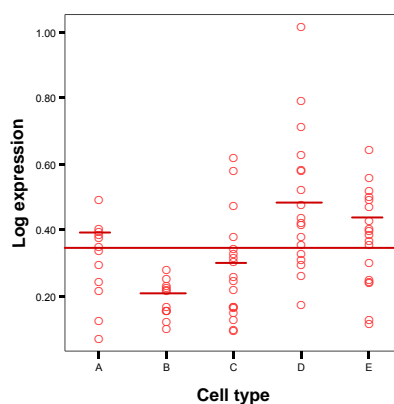
Now, you have an hypothesis to explain the variability, or at least you hope most of it: you think that your data can be split into 5 groups (e.g. 5 cell types), like in the graph below.

So you work out the mean for each cell type and you work out the squared differences between each of the mean and the grand mean, which gives you (second row of the table):

Between groups sum of squares = $\sum n_i (\text{Mean}_i - \text{Grand mean})^2$ where n is the number of data points in each of the i groups (see graph below).

In our example: Between groups SS = 351.520 and, since we have 5 groups, there are $5 - 1 = 4$ df, the mean SS = $351.520/4 = 87.880$.

If you remember the formula of the variance (= $SS / N-1$, with $df=N-1$), you can see that this value quantifies the variability between the groups' means, it is the between group variance.



There is one row left in the table, the within groups variability. It is the variability within each of the five groups, so it corresponds to the difference between each data point and its respective group mean:

Within groups sum of squares = $\sum (x_i - \text{Mean}_i)^2$ which in our case is equal to 435.300.

This value can also be obtained by doing $786.820 - 351.520 = 435.300$, which is logical since it is the amount variability left from the total variability after the variability explained by your model has been removed.

As there are 5 groups of $n=10$ values, $df = 5 \times (n - 1) = 5 \times (10 - 1) = 45$.

So the mean Within groups SS = $435.300/45 = 9.673$. This quantifies the remaining variability, the one not explained by the model, the individual variability between each value and the mean of the group to which it belongs according to your hypothesis.

At this point, you can see that the amount of variability explained by your model (87.880) is far higher than the remaining one (9.673).

So, you can work out the F-ratio: $F = 87.880 / 9.673 = 9.085$



SPSS calculates the level of significance of the test by taking into account the F ratio and the number of df for the numerator and the denominator. In our example, $p < 0.0001$, so the test is highly significant and you are more than 99% confident when you say that there is a difference between the groups' means.

Exercise (File: [protein expression.sav](#)): Find out if there is a significant difference in terms of protein expression between 5 cell types.

First, let's check for normality (remember: Analyze>Descriptive Statistics>Explore).

Tests of Normality

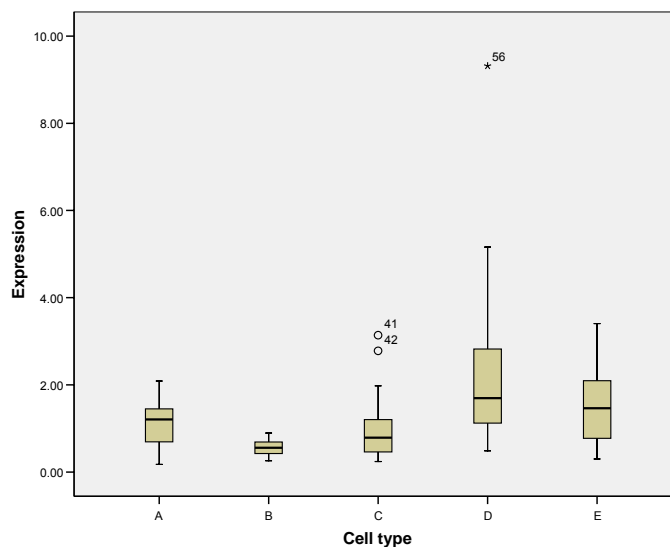
| Cell type | Kolmogorov-Smirnov(a) | | | Shapiro-Wilk | | | |
|------------|-----------------------|------|------|--------------|------|------|------|
| | Statistic | Df | Sig. | Statistic | df | Sig. | |
| Expression | A | .143 | 12 | .200(*) | .966 | 12 | .870 |
| | B | .170 | 12 | .200(*) | .954 | 12 | .700 |
| | C | .197 | 18 | .064 | .819 | 18 | .003 |
| | D | .206 | 18 | .042 | .753 | 18 | .000 |
| | E | .106 | 18 | .200(*) | .967 | 18 | .742 |

* This is a lower bound of the true significance.

a Lilliefors Significance Correction

Test of Homogeneity of Variance

| | | Levene Statistic | df1 | df2 | Sig. |
|------------|--------------------------------------|------------------|-----|--------|------|
| Expression | Based on Mean | 5.212 | 4 | 73 | .001 |
| | Based on Median | 2.888 | 4 | 73 | .028 |
| | Based on Median and with adjusted df | 2.888 | 4 | 24.977 | .043 |
| | Based on trimmed mean | 4.082 | 4 | 73 | .005 |



It does not look good: 2 out of 5 groups (C and D) show a significant departure from normality and there is no homogeneity of the variances ($p = 0.01$). The data from groups C and D are quite skewed and a look at the raw data shows more than a 10-fold jump between values of the same group (e.g. in



group A, value line 4 is 0.17 and value line 10 is 2.09). A good idea would be log-transform the data so that the spread is more balanced and to check again on the assumptions.

Tests of Normality

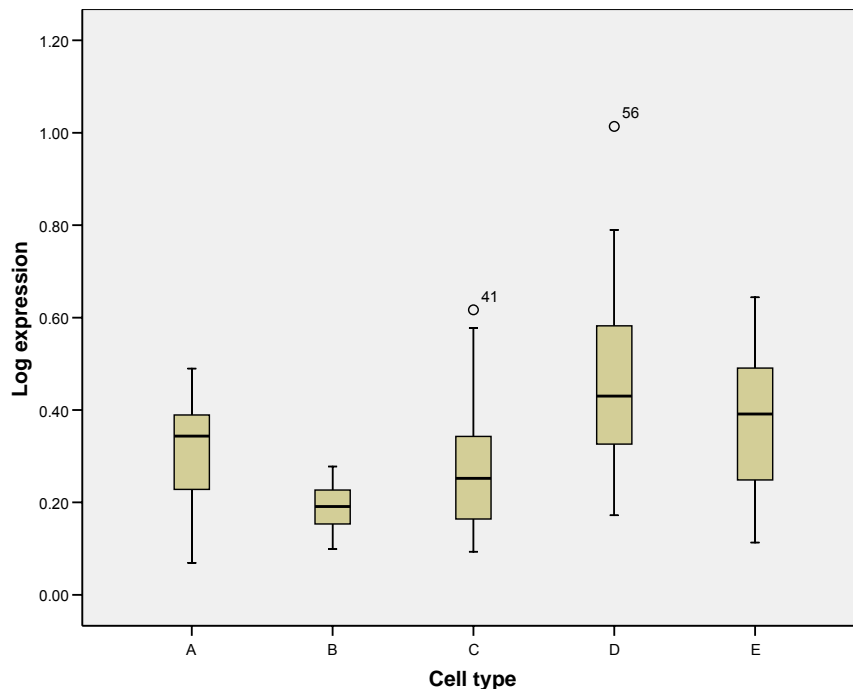
| | Cell type | Kolmogorov-Smirnov(a) | | | Shapiro-Wilk | | |
|----------------|-----------|-----------------------|----|---------|--------------|----|------|
| | | Statistic | df | Sig. | Statistic | df | Sig. |
| Log expression | A | .185 | 12 | .200(*) | .938 | 12 | .476 |
| | B | .182 | 12 | .200(*) | .955 | 12 | .713 |
| | C | .154 | 18 | .200(*) | .911 | 18 | .088 |
| | D | .142 | 18 | .200(*) | .942 | 18 | .309 |
| | E | .107 | 18 | .200(*) | .976 | 18 | .904 |

* This is a lower bound of the true significance.

a Lilliefors Significance Correction

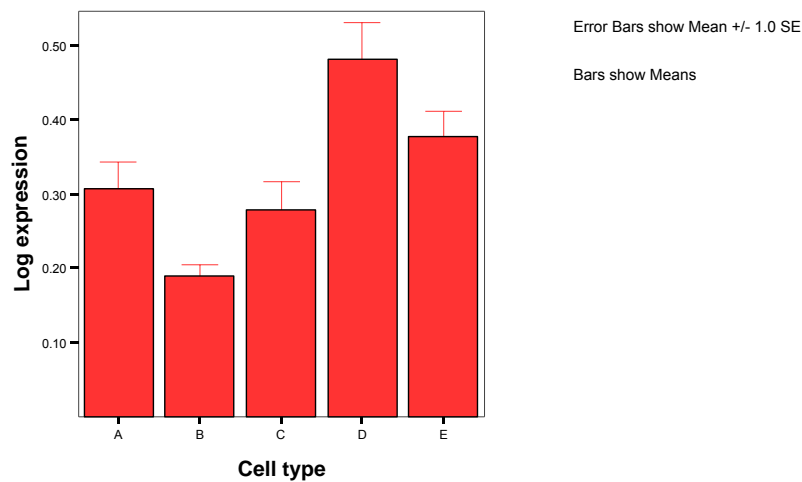
Test of Homogeneity of Variance

| | | Levene Statistic | df1 | df2 | Sig. |
|----------------|--------------------------------------|------------------|-----|--------|------|
| Log expression | Based on Mean | 3.008 | 4 | 73 | .024 |
| | Based on Median | 2.232 | 4 | 73 | .074 |
| | Based on Median and with adjusted df | 2.232 | 4 | 51.056 | .078 |
| | Based on trimmed mean | 2.793 | 4 | 73 | .032 |

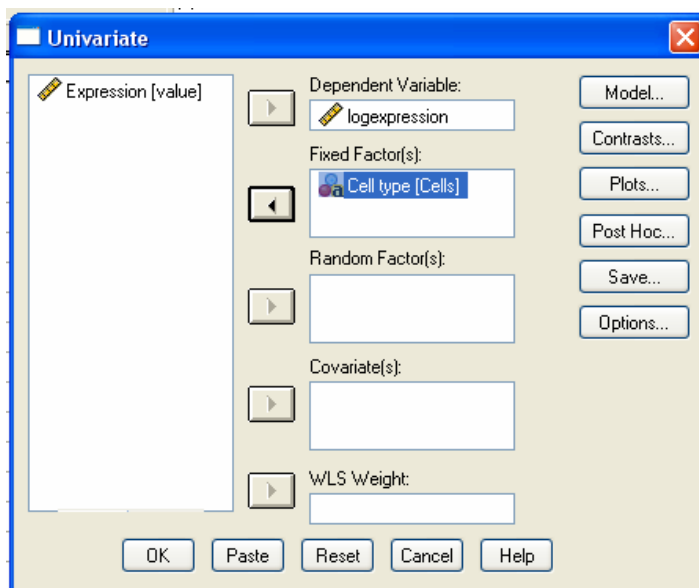


OK, the situation is getting better: data are (more) normal but the homogeneity of variance is not met though it has improved. Since the analysis of variance is a robust test (meaning that it behaves fairly well in front of moderate departure from both normality and equality of variance) and the variances are not “too” different, you can go ahead with the analysis.

First, we can plot the data and the graph below gives us hope in terms of significant difference between group means.



Then we run the ANOVA: to do so you go Analyze >General Linear Model >Univariate. Don't choose the One-Way ANOVA from "Compare Means" unless your samples are of the exact same size. All you have to do is dragging the variables you are interested in in the appropriate place.



You have several choices from this window:

- **Model:** you can include or not interactions when you have more than one factor.
- **Contrasts:** you can plan contrasts between groups before starting the analysis but often post hoc tests are easier to manipulate.
- **Plots:** you can plot the model, which is always good. By default, SPSS display line graphs but you can change it by activating the graph and then double-clicking on the lines to change it to bars, for instance.
- **Post Hoc:** when you have run your analysis of variance, saw that there is a significant difference between your groups and you want to know which group is actually different from which one, you run Post hoc tests.



- Save: you won't need this one.
- Options: allows you to run more tests and get a more detailed output.

Between-Subjects Factors

| | | N |
|-----------|---|----|
| Cell type | A | 12 |
| | B | 12 |
| | C | 18 |
| | D | 18 |
| | E | 18 |

Tests of Between-Subjects Effects

Dependent Variable: Log expression

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
|-----------------|-------------------------|----|-------------|---------|------|
| Corrected Model | .740 ^a | 4 | .185 | 7.906 | .000 |
| Intercept | 8.001 | 1 | 8.001 | 341.683 | .000 |
| Cells | .740 | 4 | .185 | 7.906 | .000 |
| Error | 1.709 | 73 | .023 | | |
| Total | 11.401 | 78 | | | |
| Corrected Total | 2.450 | 77 | | | |

a. R Squared = .302 (Adjusted R Squared = .264)

In the SPSS output for an ANOVA, the row showing the between group variation corresponds to the one with the group variable name (here Cells) and the one for the within groups variation is called "Error". The total variation is "Corrected total", so $0.740 + 1.709 = 2.450$. The rest of the table you can ignore, SPSS tending to produce "very talkative" outputs !

There is a significant difference between the means ($p < 0.0001$), but even if you have an indication from the graph, you cannot tell which mean is different which one. This is because the ANOVA is an "omnibus" test: it tells you that there is (or not) a difference between your means but not exactly which means are significantly different from which other ones. To find out, you need to apply **post hoc** tests. SPSS offers you several types of post hoc tests which you can choose depending on the difference in sample size and variance between your groups. These post hoc tests should only be used when the ANOVA finds a significant effect.

| Variance | Sample size | Post hoc test |
|---|------------------|---------------------|
| equal | equal | Tukey or Bonferroni |
| equal | Small difference | Gabriel |
| equal | Big difference | Hochberg's GT2 |
| different | - | Games-Howell |
| Comparisons of group means against control mean | | Dunnett |

In our example, since the sample sizes are different and the homogeneity of variance is not assumed we should run at least Gabriel and Games-Howell's tests. Usually, I recommend running them all.



Multiple Comparisons

Dependent Variable: Log expression
Gabriel

| (I) Cell type | (J) Cell type | Mean Difference (I-J) | Std. Error | Sig. | 95% Confidence Interval | |
|---------------|---------------|-----------------------|------------|-------|-------------------------|-------------|
| | | | | | Lower Bound | Upper Bound |
| A | B | .1176 | .06247 | .470 | -.0624 | .2977 |
| | C | .0274 | .05703 | 1.000 | -.1361 | .1909 |
| | D | -.1753* | .05703 | .028 | -.3388 | -.0118 |
| | E | -.0702 | .05703 | .908 | -.2337 | .0933 |
| B | A | -.1176 | .06247 | .470 | -.2977 | .0624 |
| | C | -.0902 | .05703 | .694 | -.2537 | .0733 |
| | D | -.2930* | .05703 | .000 | -.4565 | -.1294 |
| | E | -.1878* | .05703 | .014 | -.3513 | -.0243 |
| C | A | -.0274 | .05703 | 1.000 | -.1909 | .1361 |
| | B | .0902 | .05703 | .694 | -.0733 | .2537 |
| | D | -.2027* | .05101 | .002 | -.3497 | -.0557 |
| | E | -.0976 | .05101 | .447 | -.2446 | .0494 |
| D | A | .1753* | .05703 | .028 | .0118 | .3388 |
| | B | .2930* | .05703 | .000 | .1294 | .4565 |
| | C | .2027* | .05101 | .002 | .0557 | .3497 |
| | E | .1051 | .05101 | .345 | -.0419 | .2521 |
| E | A | .0702 | .05703 | .908 | -.0933 | .2337 |
| | B | .1878* | .05703 | .014 | .0243 | .3513 |
| | C | .0976 | .05101 | .447 | -.0494 | .2446 |
| | D | -.1051 | .05101 | .345 | -.2521 | .0419 |

Based on observed means.

*. The mean difference is significant at the .05 level.

Multiple Comparisons

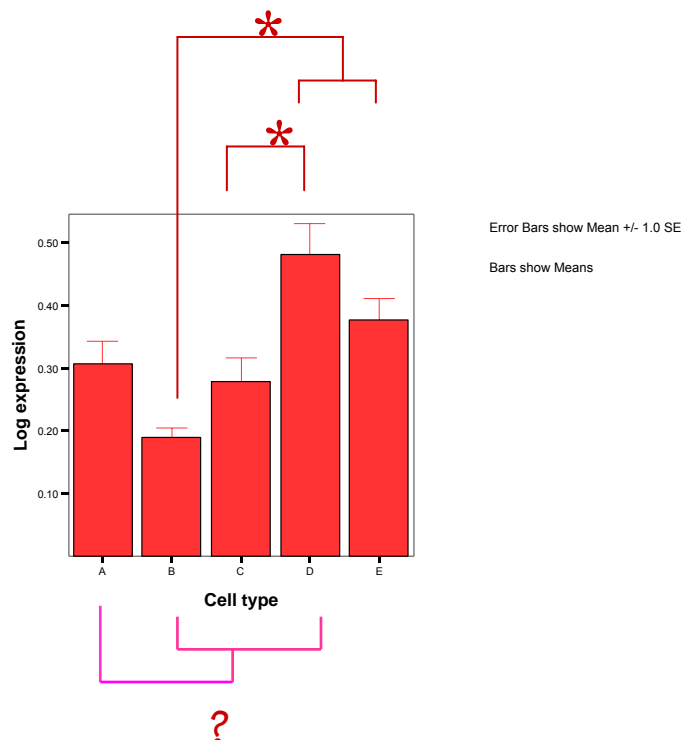
Dependent Variable: Log expression
Games-Howell

| (I) Cell type | (J) Cell type | Mean Difference (I-J) | Std. Error | Sig. | 95% Confidence Interval | |
|---------------|---------------|-----------------------|------------|------|-------------------------|-------------|
| | | | | | Lower Bound | Upper Bound |
| A | B | .1176 | .03882 | .055 | -.0020 | .2373 |
| | C | .0274 | .05095 | .983 | -.1214 | .1762 |
| | D | -.1753 | .06070 | .053 | -.3523 | .0017 |
| | E | -.0702 | .04921 | .617 | -.2142 | .0738 |
| B | A | -.1176 | .03882 | .055 | -.2373 | .0020 |
| | C | -.0902 | .03990 | .194 | -.2083 | .0278 |
| | D | -.2930* | .05178 | .000 | -.4476 | -.1383 |
| | E | -.1878* | .03766 | .000 | -.2990 | -.0767 |
| C | A | -.0274 | .05095 | .983 | -.1762 | .1214 |
| | B | .0902 | .03990 | .194 | -.0278 | .2083 |
| | D | -.2027* | .06140 | .019 | -.3803 | -.0251 |
| | E | -.0976 | .05007 | .312 | -.2418 | .0466 |
| D | A | .1753 | .06070 | .053 | -.0017 | .3523 |
| | B | .2930* | .05178 | .000 | .1383 | .4476 |
| | C | .2027* | .06140 | .019 | .0251 | .3803 |
| | E | .1051 | .05997 | .418 | -.0687 | .2790 |
| E | A | .0702 | .04921 | .617 | -.0738 | .2142 |
| | B | .1878* | .03766 | .000 | .0767 | .2990 |
| | C | .0976 | .05007 | .312 | -.0466 | .2418 |
| | D | -.1051 | .05997 | .418 | -.2790 | .0687 |

Based on observed means.

*. The mean difference is significant at the .05 level.

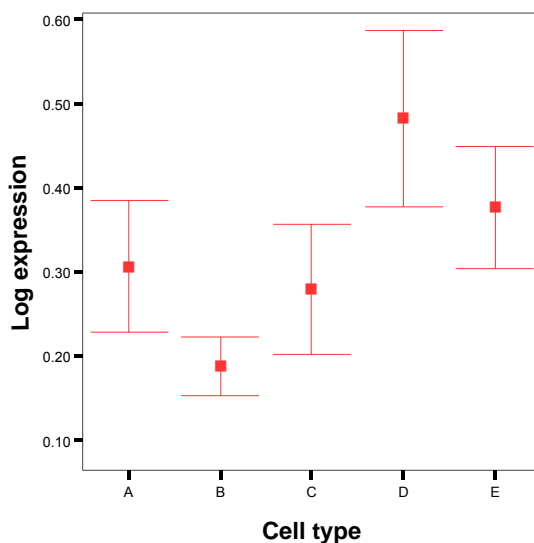
What the tests tell you is summarised in the graph below. Now, 2 things are puzzling: the first one is that the tests disagree about the difference between groups A and B. Gabriel says no ($p=0.470$) whereas Games-Howell says yes (well almost with $p=0.055$). The second one is about A and D: Games-Howell is border-line ($p=0.053$) whereas Gabriel is positive about the significance of the difference ($p=0.028$).



These problems can be solved (most of the time) by plotting the confidence intervals of the groups' means instead of the Standard Error (see graph below).

For the A-B difference there is no overlap but a difference in variance therefore you should trust the result of the Games-Howell test.

For group A and group D, there is a small overlap and may be more variability in group D than group A. Because of these reasons and because it is more convenient to report one test than 2, I would also go for Games-Howell this time. Remember, 5% is an arbitrary threshold meaning you cannot say that nothing is happening when you get a p-value of 0.053.





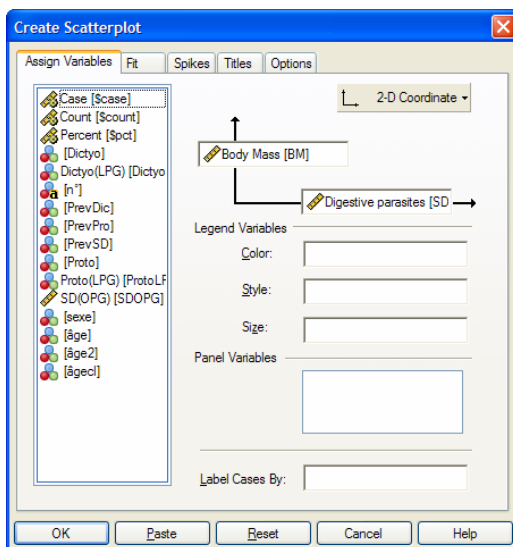
5-5 Correlation

If you want to find out about the relationship between 2 variables, you can run a correlation.

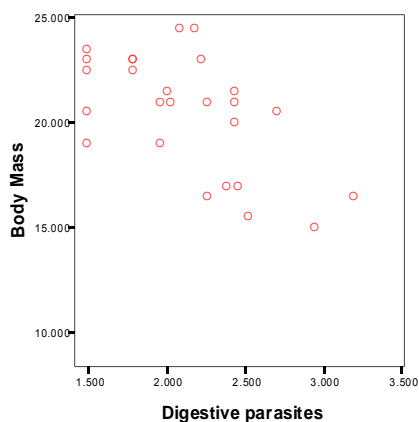
Example (File: roe deer.sav).

When you want to plot data from 2 quantitative variables between which you suspect (hope?) that there is a relationship, the best choice to have a first look at you data is the scatter plot. So on SPSS, you go: Graphs>Interactive>Scatterplot.

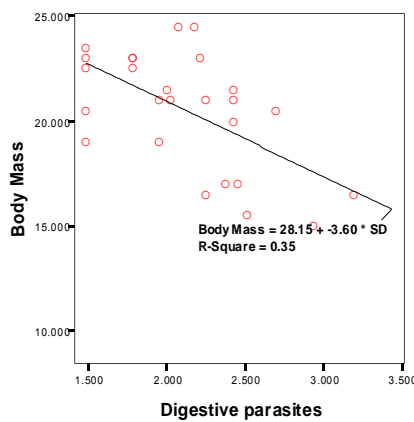
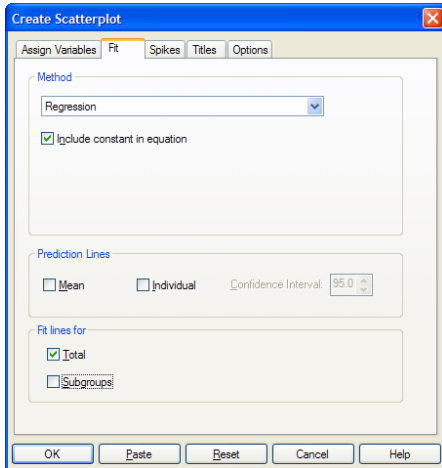
In our case we want to know if there is a relationship between the body mass and the parasite burden.



You have to choose between the x- and the y-axis for your 2 variables. It is usually considered that “x” predicts “y” ($y=f(x)$) so when looking at the relationship between 2 variables, you must have an idea of which one is likely to predict the other one. In our particular case, we want to know how an increase in parasite burden affects the body mass of the host.



By looking at the graph, one can think that something is happening here. To have a better idea, you can plot the regression line on the data.



Linear Regression

Now, the questions are: is the relationship significant? and what do these numbers on the graph mean?

To answer these questions, you need to run a correlation test.

A bit of theory: Correlation coefficient

A correlation is a measure of a linear relationship (can be expressed as straight-line graphs) between variables. The simplest way to find out whether 2 variables are associated is to look at whether they covary. To do so, you combine the variance of one variable with the variance of the other.

$$\text{cov}(X, Y) = \sum_{i=1}^N \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}.$$

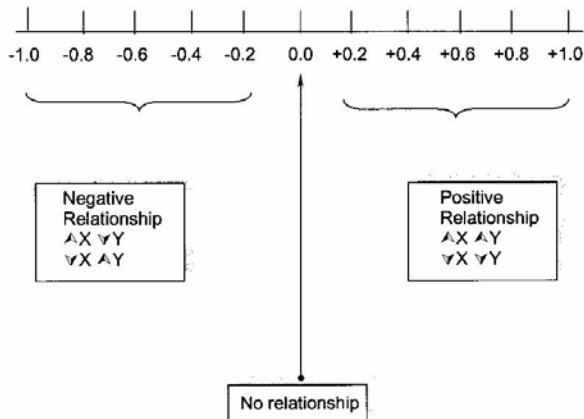
A positive covariance indicates that as one variable deviates from the mean, the other one deviates in the same direction, in other words if one variable goes up the other one goes up as well.

The problem with the covariance is that its value depends upon the scale of measurement used, so you won't be able to compare covariance between datasets unless both data are measured in the same units. To standardise the covariance, it is divided by the SD of the 2 variables. It gives you the most widely-used correlation coefficient: the Pearson product-moment correlation coefficient "r".



$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

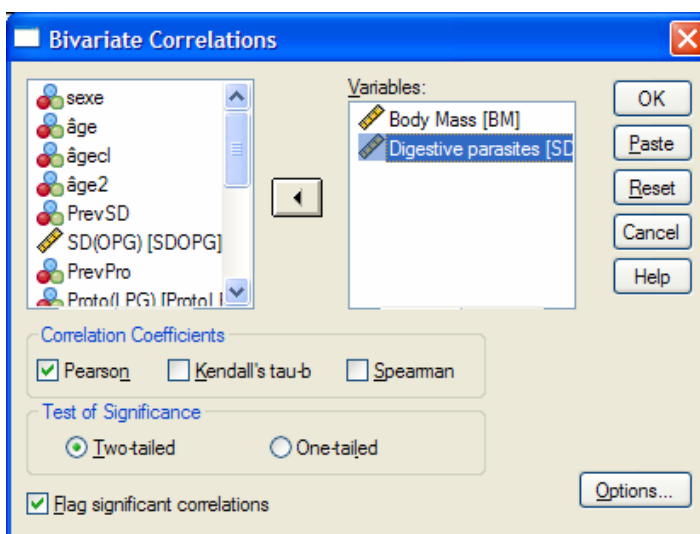
Of course, you don't need to remember that formula but it is important that you understand what the correlation coefficient does: it measures the magnitude and the direction of the relationship between two variables. It is designed to range in value between 0.0 and 1.0.



The 2 variables do not have to be measured in the same units but they have to be proportional (meaning linearly related)

One last thing before we go back to our example: the coefficient of determination r^2 : it gives you the proportion of variance in Y that can be explained by X, in percentage.

To run a correlation on SPSS, you go: Analysis>Correlate>Bivariate.





Correlations

| | | Body Mass | Digestive parasites |
|---------------------|---------------------|-----------|---------------------|
| Body Mass | Pearson Correlation | 1 | -.592** |
| | Sig. (2-tailed) | | .001 |
| | N | 26 | 26 |
| Digestive parasites | Pearson Correlation | -.592** | 1 |
| | Sig. (2-tailed) | .001 | |
| | N | 26 | 26 |

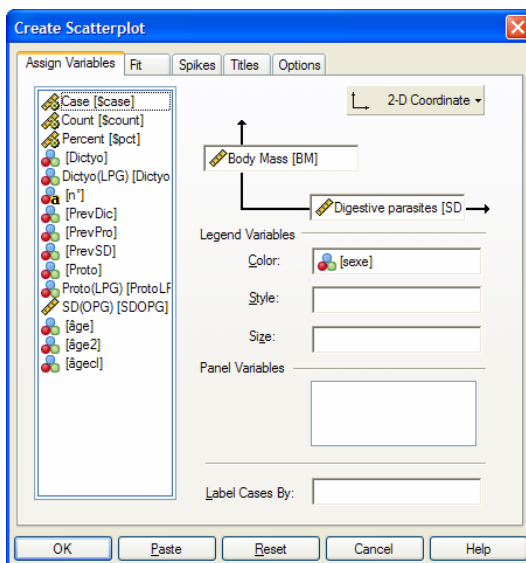
** . Correlation is significant at the 0.01 level (2-tailed).

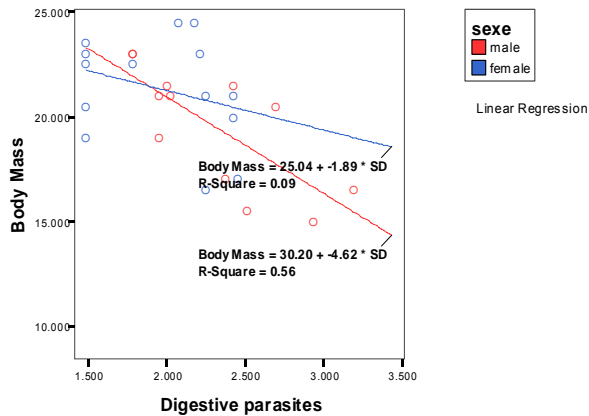
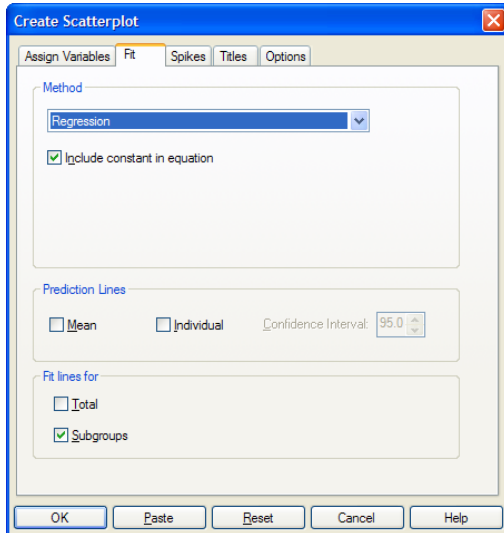
The SPSS output gives you a symmetrical matrix. So, this table tells us that there is a strong ($p=0.001$) negative ($r = -0.592$) relationship between the 2 variables, the body mass decreasing when the parasite burden increases.

If you square the correlation coefficient, you get: $r^2 = 0.3504$, which is the value you saw on the graph. It means the 35% of the variance in body mass is explained by the parasite burden.

The equation on the graph (Body mass = 28.15 – 3.6*Digestive parasites) tells you that for each increase of parasite burden of 1 unit, the animals loose 3.6 units of body mass and that the average body mass of the roe deers in that group is 28.15 kg.

Now, you may want to know if this relationship is the same for both sexes for instance. To do so, you go back to the scatterplot window and you add sex as, what SPSS called a legend variable.





Now you can see that you get 2 very different pictures according to the gender you are looking at: the effect of parasite burden is much stronger for males as it explains 56% of the variability in body mass whereas it only explains 9% of it in females.

If you run again the correlation, taking into account the sex, you get:

Correlations^a

| | | Body Mass | Digestive parasites |
|---------------------|---------------------|-----------|---------------------|
| Body Mass | Pearson Correlation | 1 | -.750** |
| | Sig. (2-tailed) | | .005 |
| | N | 12 | 12 |
| Digestive parasites | Pearson Correlation | -.750** | 1 |
| | Sig. (2-tailed) | .005 | |
| | N | 12 | 12 |

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

a. sexe = male

Correlations^a

| | | Body Mass | Digestive parasites |
|---------------------|---------------------|-----------|---------------------|
| Body Mass | Pearson Correlation | 1 | -.302 |
| | Sig. (2-tailed) | | .294 |
| | N | 14 | 14 |
| Digestive parasites | Pearson Correlation | -.302 | 1 |
| | Sig. (2-tailed) | .294 | |
| | N | 14 | 14 |

a. sexe = female

From the result of the tests, you can see that the correlation is only significant for the males and not for the females.

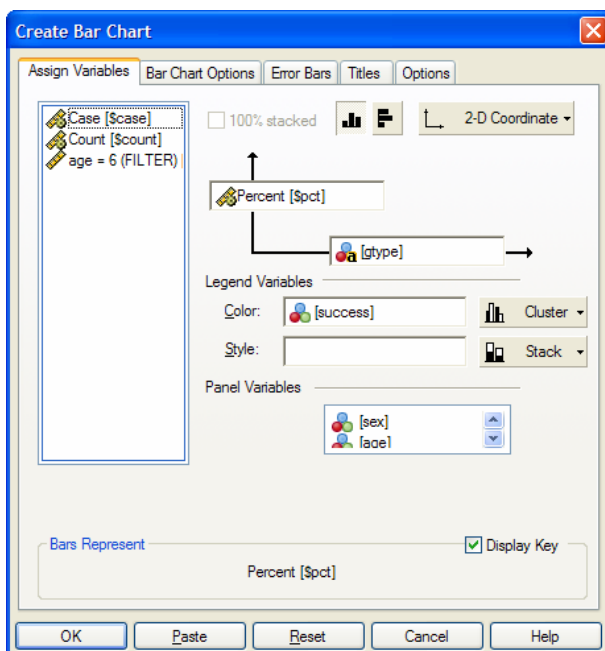
A key thing to remember when working with correlations is never to assume a correlation means that a change in one variable *causes* a change in another. Sales of personal computers and athletic shoes have both risen strongly in the last several years and there is a high correlation between them, but you cannot assume that buying computers causes people to buy athletic shoes (or vice versa).

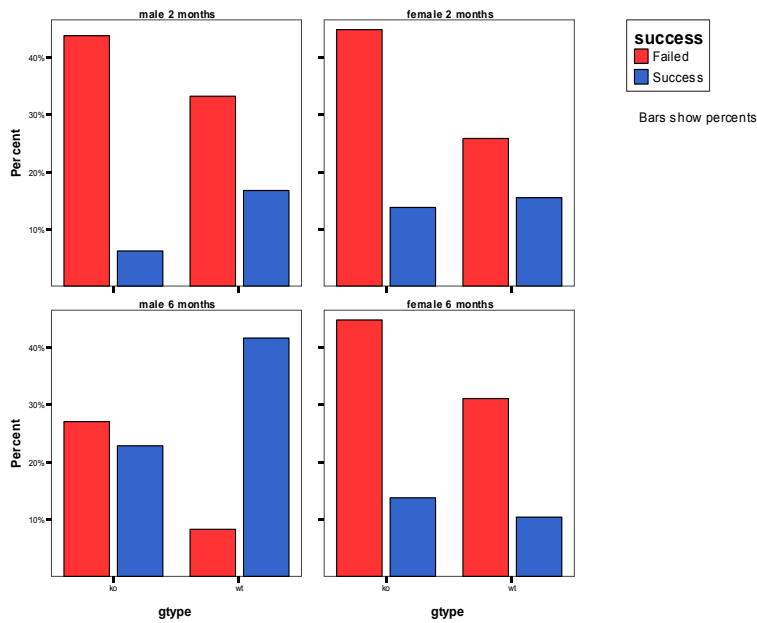
EXERCISES

File: [behavioural exp.xls](#)

A researcher wants to know if there is a difference between 2 types of mouse (wt and ko) in their ability to achieve a task in a behavioural experiment (failed=0 or success=1), taking into account the gender (1=male and 2=female) and the age (2 and 6 months-old).

Prepare the file and plot the data so that you get 4 graphs with males and females at 2 months-old on the top and males and females at 6 months-old at the bottom.





Find out if there is a difference in term of success between wt and ko 6 months-old mice. Do it separately for each gender.

success * gtype * sex Crosstabulation

| sex | | | | gtype | | Total |
|--------|---------|--------|----------------|--------|--------|--------|
| | | | | ko | wt | |
| male | success | Failed | Count | 13 | 4 | 17 |
| | | | % within gtype | 54.2% | 16.7% | 35.4% |
| | Success | | Count | 11 | 20 | 31 |
| | | | % within gtype | 45.8% | 83.3% | 64.6% |
| | Total | | Count | 24 | 24 | 48 |
| | | | % within gtype | 100.0% | 100.0% | 100.0% |
| female | success | Failed | Count | 26 | 18 | 44 |
| | | | % within gtype | 76.5% | 75.0% | 75.9% |
| | Success | | Count | 8 | 6 | 14 |
| | | | % within gtype | 23.5% | 25.0% | 24.1% |
| | Total | | Count | 34 | 24 | 58 |
| | | | % within gtype | 100.0% | 100.0% | 100.0% |



Chi-Square Tests

| sex | | Value | df | Asymp. Sig. (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
|--------|------------------------------------|--------------------|----|--------------------------|-------------------------|-------------------------|
| male | Pearson Chi-Square | 7.378 ^b | 1 | .007 | | |
| | Continuity Correction ^a | 5.829 | 1 | .016 | | |
| | Likelihood Ratio | 7.668 | 1 | .006 | | |
| | Fisher's Exact Test | | | | .015 | .007 |
| | N of Valid Cases | 48 | | | | |
| female | Pearson Chi-Square | .017 ^c | 1 | .897 | | |
| | Continuity Correction ^a | .000 | 1 | 1.000 | | |
| | Likelihood Ratio | .017 | 1 | .898 | | |
| | Fisher's Exact Test | | | | 1.000 | .568 |
| | N of Valid Cases | 58 | | | | |

a. Computed only for a 2x2 table

b. 0 cells (.0%) have expected count less than 5. The minimum expected count is 8.50.

c. 0 cells (.0%) have expected count less than 5. The minimum expected count is 5.79.

File: [bacteria count.xls](#)

Import the file, check for normality and plot the data so that you can see the difference in number of bacteria between the wt and the ko mice and have an idea of the significance of that difference. Run a t-test to check.



Descriptives

| type | | | Statistic | Std. Error | | |
|---------------------|---------------------|----------------------------------|----------------------------------|------------|-------------|--------|
| bact2 | ko | Mean | 256.83 | 26.280 | | |
| | | 95% Confidence Interval for Mean | Lower Bound | | 203.08 | |
| | | | Upper Bound | | 310.58 | |
| | 5% Trimmed Mean | 255.09 | | | | |
| | Median | 256.00 | | | | |
| | Variance | 20719.109 | | | | |
| | Std. Deviation | 143.941 | | | | |
| | Minimum | 22 | | | | |
| | Maximum | 541 | | | | |
| | Range | 519 | | | | |
| | Interquartile Range | 251 | | | | |
| | Skewness | .058 | .427 | | | |
| | Kurtosis | -.969 | .833 | | | |
| | wt | ko | Mean | | 365.03 | 17.643 |
| | | | 95% Confidence Interval for Mean | | Lower Bound | |
| Upper Bound | | | | 401.12 | | |
| 5% Trimmed Mean | | 367.94 | | | | |
| Median | | 374.00 | | | | |
| Variance | | 9338.033 | | | | |
| Std. Deviation | | 96.634 | | | | |
| Minimum | | 118 | | | | |
| Maximum | | 530 | | | | |
| Range | | 412 | | | | |
| Interquartile Range | | 135 | | | | |
| Skewness | | -.440 | .427 | | | |
| Kurtosis | | .031 | .833 | | | |

Tests of Normality

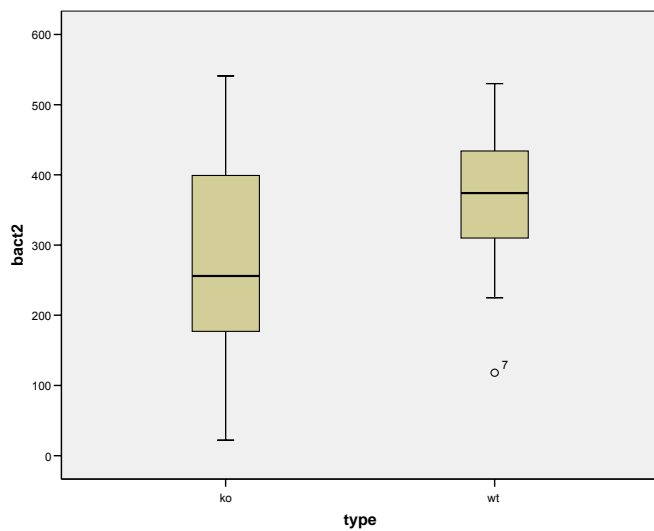
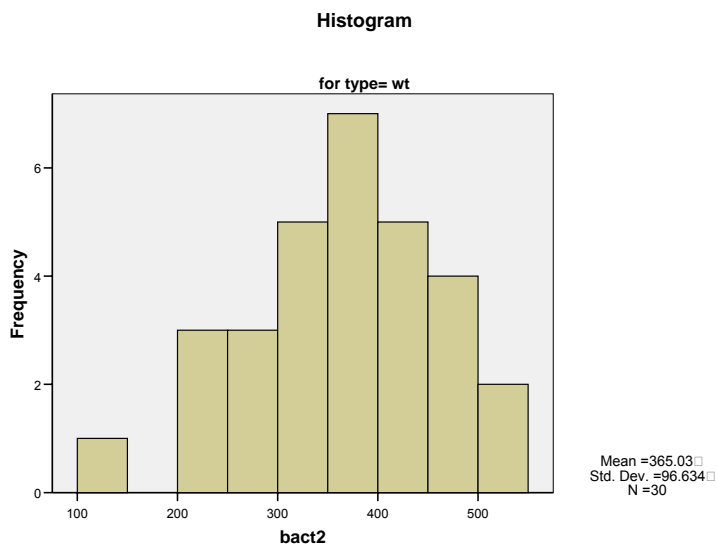
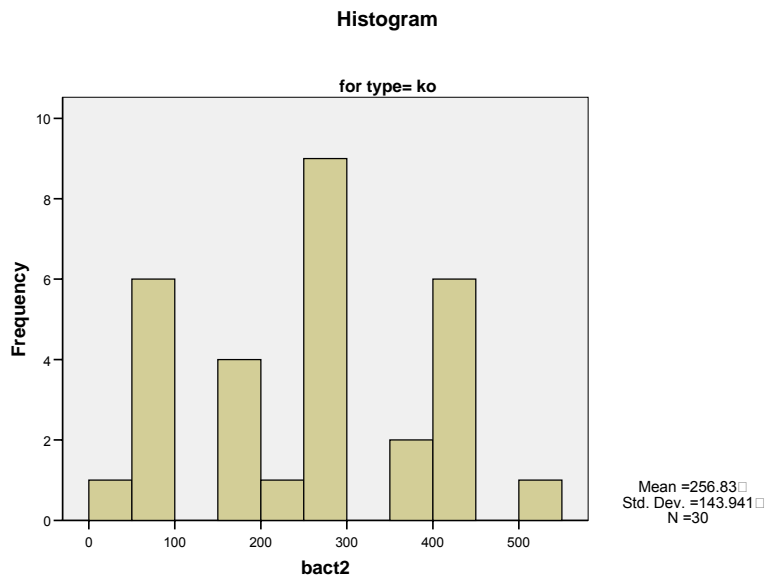
| type | Kolmogorov-Smirnov ^a | | | Shapiro-Wilk | | |
|----------|---------------------------------|----|-------|--------------|----|------|
| | Statistic | df | Sig. | Statistic | df | Sig. |
| bact2 ko | .134 | 30 | .180 | .943 | 30 | .108 |
| wt | .088 | 30 | .200* | .978 | 30 | .772 |

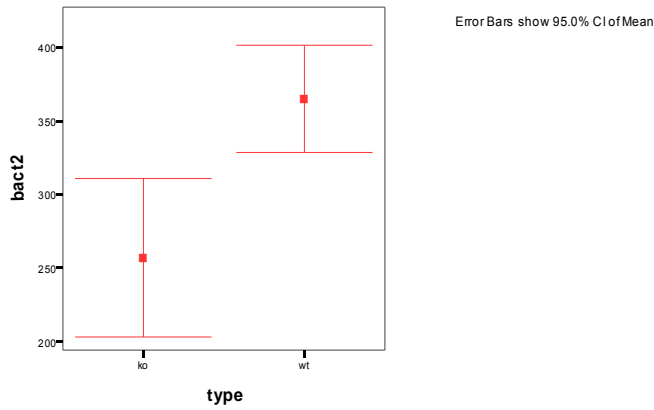
*. This is a lower bound of the true significance.

a. Lilliefors Significance Correction

Test of Homogeneity of Variance

| | | Levene Statistic | df1 | df2 | Sig. |
|------|--------------------------------------|------------------|-----|--------|------|
| bact | Based on Mean | 4.396 | 1 | 58 | .040 |
| | Based on Median | 4.413 | 1 | 58 | .040 |
| | Based on Median and with adjusted df | 4.413 | 1 | 51.843 | .041 |
| | Based on trimmed mean | 4.422 | 1 | 58 | .040 |





Group Statistics

| type | | N | Mean | Std. Deviation | Std. Error Mean |
|-------|----|----|--------|----------------|-----------------|
| bact2 | wt | 30 | 365.03 | 96.634 | 17.643 |
| | ko | 30 | 256.83 | 143.941 | 26.280 |

Independent Samples Test

| | | Levene's Test for Equality of Variances | | t-test for Equality of Means | | | | | | |
|-------|-----------------------------|---|------|------------------------------|--------|-----------------|-----------------|-----------------------|---|---------|
| | | F | Sig. | t | df | Sig. (2-tailed) | Mean Difference | Std. Error Difference | 95% Confidence Interval of the Difference | |
| | | | | | | | | | Lower | Upper |
| bact2 | Equal variances assumed | 4.396 | .040 | 3.418 | 58 | .001 | 108.200 | 31.653 | 44.840 | 171.560 |
| | Equal variances not assumed | | | 3.418 | 50.727 | .001 | 108.200 | 31.653 | 44.646 | 171.754 |