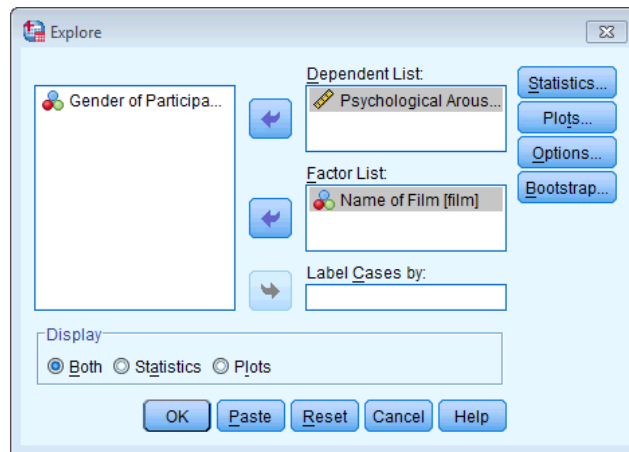# Chapter 5: The beast of bias

## Smart Alex's Solutions

## Task 1

*Using the **ChickFlick.sav** data from Chapter 4, check the assumptions of normality and homogeneity of variance for the two films (ignore gender): are the assumptions met?*

The output you should get is reproduced below (I used the *explore* function described in Chapter 5).



The skewness statistics gives rise to a *z*-score of −0.378/0.512 = −0.74 for *Bridget Jones's Diary*, and 0.04/0.512 = 0.08 for *Memento*. These show no significant skewness. For kurtosis these values are −0.254/0.992 = −0.26 for *Bridget Jones's Diary*, and −1.024/0.992 = −1.03, which again are both non-significant.

The Q-Q plots confirm these findings: for both films the expected quantile points are close to those that would be expected from a normal distribution (i.e. the dots fall close to the diagonal line).
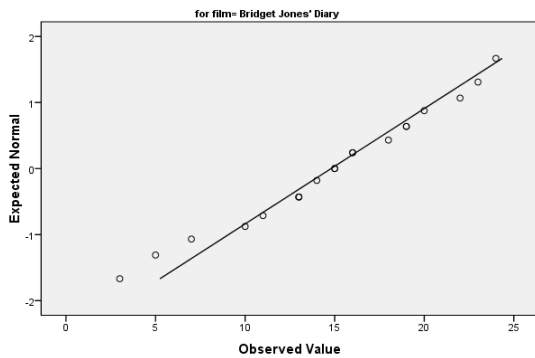
The K-S tests show no significant deviation from normality for both films. We could report that arousal scores for *Bridget Jones's Diary*, *D*(20) = 0.13, *ns*, and *Memento*, *D*(20) = 0.10, *ns*, were both not significantly different from a normal distribution. Therefore we can assume normality in the sample data.

In terms of homogeneity of variance, Levene's test shows that the variances of arousal for the two films were not significantly different, *F*(1, 38) = 1.90.
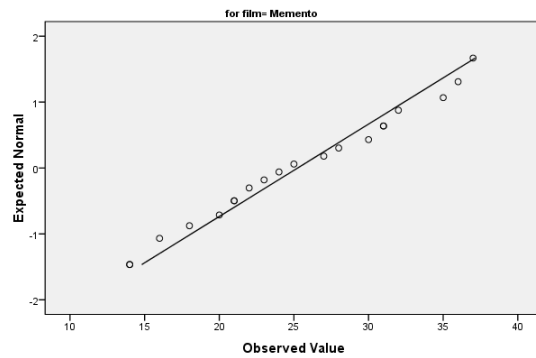
**Descriptives**

| | Film | | | Statistic | Std. Error |
|---|---|---|---|---|---|
| Arousal | Bridget Jones' Diary | Mean | | 14.8000 | 1.28062 |
| | | 95% Confidence Interval for Mean | Lower Bound | 12.1196 | |
| | | | Upper Bound | 17.4804 | |
| | | 5% Trimmed Mean | | 14.9444 | |
| | | Median | | 15.0000 | |
| | | Variance | | 32.800 | |
| | | Std. Deviation | | 5.72713 | |
| | | Minimum | | 3.00 | |
| | | Maximum | | 24.00 | |
| | | Range | | 21.00 | |
| | | Interquartile Range | | 7.50 | |
| | | Skewness | | -.378 | .512 |
| | | Kurtosis | | -.254 | .992 |
| | Memento | Mean | | 25.2500 | 1.59419 |
| | | 95% Confidence Interval for Mean | Lower Bound | 21.9133 | |
| | | | Upper Bound | 28.5867 | |
| | | 5% Trimmed Mean | | 25.2222 | |
| | | Median | | 24.5000 | |
| | | Variance | | 50.829 | |
| | | Std. Deviation | | 7.12944 | |
| | | Minimum | | 14.00 | |
| | | Maximum | | 37.00 | |
| | | Range | | 23.00 | |
| | | Interquartile Range | | 10.75 | |
| | | Skewness | | .040 | .512 |
| | | Kurtosis | | -1.024 | .992 |

Normal Q-Q Plot of Arousal
for film= Bridget Jones' Diary

Normal Q-Q Plot of Arousal
for film= Memento

**Tests of Normality**

| | Film | Kolmogorov-Smirnov<sup>a</sup> | | | Shapiro-Wilk | | |
|---|---|---|---|---|---|---|---|
| | | Statistic | df | Sig. | Statistic | df | Sig. |
| Arousal | Bridget Jones' Diary | .127 | 20 | .200<sup>*</sup> | .972 | 20 | .788 |
| | Memento | .097 | 20 | .200<sup>*</sup> | .960 | 20 | .552 |

a. Lilliefors Significance Correction

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

**Test of Homogeneity of Variance**

| | | Levene Statistic | df1 | df2 | Sig. |
|---|---|---|---|---|---|
| Arousal | Based on Mean | 1.898 | 1 | 38 | .176 |
| | Based on Median | 1.805 | 1 | 38 | .187 |
| | Based on Median and with adjusted df | 1.805 | 1 | 37.847 | .187 |
| | Based on trimmed mean | 1.912 | 1 | 38 | .175 |

# Task 2

*The file **SPSSExam.sav** contains data regarding students' performance on an SPSS exam. Four variables were measured: **exam** (first-year SPSS exam scores as a percentage), **computer** (measure of computer literacy in percent), **lecture** (percentage of SPSS lectures attended) and **numeracy** (a measure of numerical ability out of 15). There is a variable called **uni** indicating whether the student attended Sussex University (where I work) or Duncetown University. Compute and interpret descriptive statistics for **exam, computer, lecture**, and **numeracy** for the sample as a whole.*

To see the distribution of the variables, we can use the *frequencies* command, which we came across in the previous section (see Section 5.3.2). Use this dialog box and place all four variables (**exam**, **computer**, **lecture** and **numeracy**) in the *Variable(s)* box. Then click on <kbd>Statistics...</kbd> to select the *Statistics* dialog box and select some measures of central tendency (mean, mode, median), measures of variability (range, standard deviation, variance, quartile splits) and measures of shape (kurtosis and skewness). Also click on <kbd>Charts...</kbd> to access the *Charts* dialog box and select a frequency distribution of scores with a normal curve. Return to the main dialog box by clicking on <kbd>Continue</kbd> and once in the main dialog box, click on <kbd>OK</kbd> to run the analysis.

The output shows the table of descriptive statistics for the four variables in this example. From this table, we can see that, on average, students attended nearly 60% of lectures, obtained 58% in their SPSS exam, scored only 51% on the computer literacy test, and only 5 out of 15 on the numeracy test. In addition, the standard deviation for computer literacy was relatively small compared to that of the percentage of lectures attended and exam scores. These latter two variables had several modes (multimodal).

The output provides tabulated frequency distributions of each variable (not reproduced here). These tables list each score and the number of times that it is found within the data set. In addition, each frequency value is expressed as a percentage of the sample (in this case the frequencies and percentages are the same because the sample size was 100). Also, the cumulative percentage is given, which tells us how many cases (as a percentage) fell below a certain score. So, for example, we can

see that 66% of numeracy scores were 5 or less, 74% were 6 or less, and so on. Looking in the other direction, we can work out that only 8% (100 − 92%) got scores greater than 8.
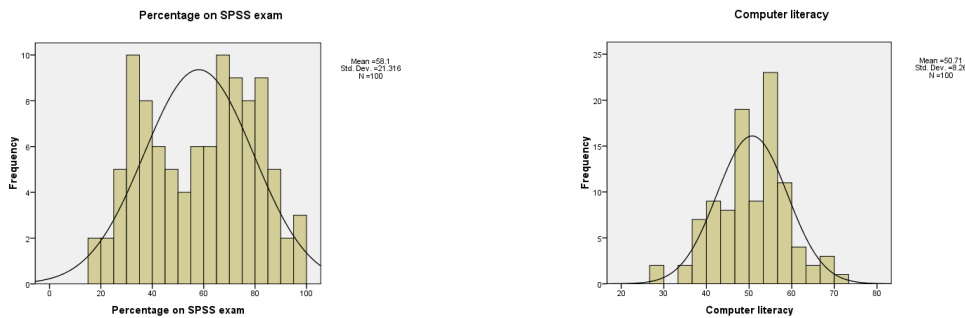
**Statistics**

| | | Percentage on SPSS exam | Computer literacy | Percentage of lectures attended | Numeracy |
|---|---|---|---|---|---|
| N | Valid | 100 | 100 | 100 | 100 |
| | Missing | 0 | 0 | 0 | 0 |
| Mean | | 58.10 | 50.71 | 59.765 | 4.85 |
| Std. Error of Mean | | 2.132 | .826 | 2.1685 | .271 |
| Median | | 60.00 | 51.50 | 62.000 | 4.00 |
| Mode | | 72ᵃ | 54 | 48.5ᵃ | 4 |
| Std. Deviation | | 21.316 | 8.260 | 21.6848 | 2.706 |
| Variance | | 454.354 | 68.228 | 470.230 | 7.321 |
| Skewness | | -.107 | -.174 | -.422 | .961 |
| Std. Error of Skewness | | .241 | .241 | .241 | .241 |
| Kurtosis | | -1.105 | .364 | -.179 | .946 |
| Std. Error of Kurtosis | | .478 | .478 | .478 | .478 |
| Range | | 84 | 46 | 92.0 | 13 |
| Minimum | | 15 | 27 | 8.0 | 1 |
| Maximum | | 99 | 73 | 100.0 | 14 |

a. Multiple modes exist. The smallest value is shown

**Output**

Finally, we are given histograms of each variable with the normal distribution overlaid. These graphs show us several things. The exam scores are very interesting because this distribution is quite clearly not normal; in fact, it looks suspiciously bimodal (there are two peaks, indicative of two modes). This observation corresponds with the earlier information from the table of descriptive statistics. It looks as though computer literacy is fairly normally distributed (a few people are very good with computers and a few are very bad, but the majority of people have a similar degree of knowledge) as is the lecture attendance. Finally, the numeracy test has produced very positively skewed data (the majority of people did very badly on this test and only a few did well). This corresponds to what the skewness statistic indicated.

Descriptive statistics and histograms are a good way of getting an instant picture of the distribution of your data. This snapshot can be very useful: for example, the bimodal distribution of SPSS exam scores instantly indicates a trend that students are typically either very good at statistics or struggle with it (there are relatively few who fall in between these extremes). Intuitively, this finding fits with the nature of the subject: statistics is very easy once everything falls into place, but before that enlightenment occurs it all seems hopelessly difficult!
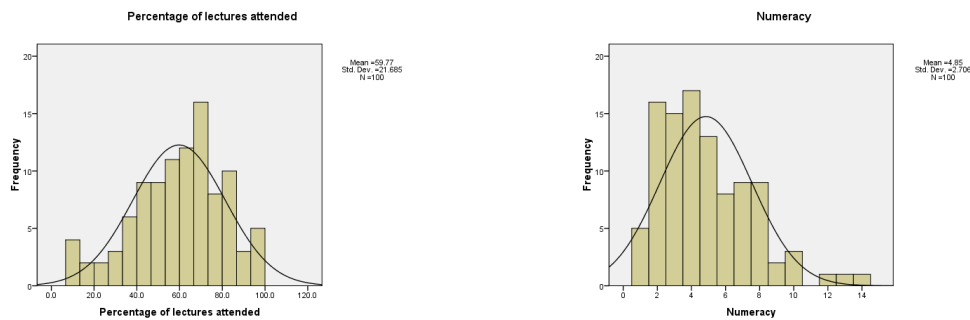
**Figure: Histograms of the SPSS exam data**

## Task 3

*Calculate and interpret the z-scores for skewness for all variables.*

For the SPSS exam scores, the *z*-score of skewness is −0.107/0.241 = −0.44. For numeracy, the *z*-score of skewness is 0.961/0.241 = 3.99. For computer literacy, the z-score of skewness is −0.174/0.241 = −0.72. For lectures attended, the z-score of skewness is −0.422/0.241 = −1.75. It is pretty clear then that the numeracy scores are significantly positively skewed ($p < .05$) because the *z*-score is greater than 1.96, indicating a pile-up of scores on the left of the distribution (so most students got low scores). For the other three variables, the skewness is non-significant, $p < .05$, because the values lie between −1.96 and 1.96.

## Task 4

*Calculate and interpret the z-scores for kurtosis for all variables.*

For SPSS exam scores, the *z*-score of kurtosis is −1.105/0.478 = −2.31, which is significant, $p < .05$, because it lies outside −1.96 and 1.96.
For computer literacy, the *z*-score of kurtosis is 0.364/0.478 = 0.76, which is non-significant, $p < .05$, because it lies between −1.96 and 1.96.
For lectures attended, the *z*-score of kurtosis is −0.179/0.478 = −0.37, which is non-significant, $p < .05$, because it lies between −1.96 and 1.96.
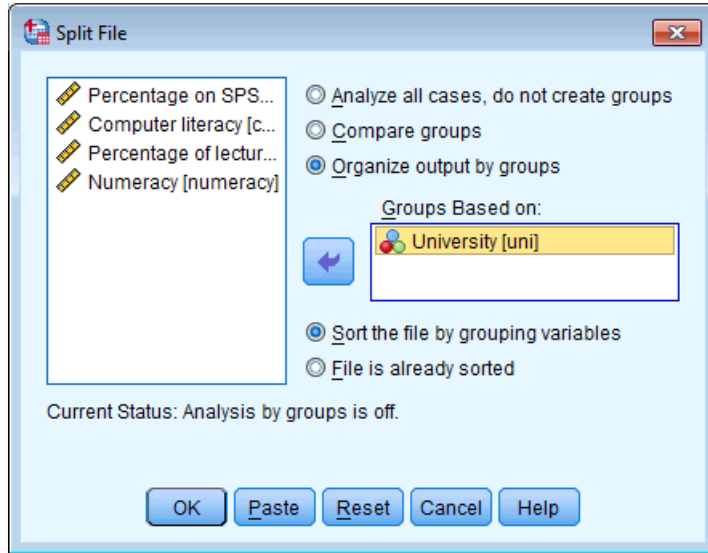For numeracy, the *z*-score of kurtosis is 0.946/0.478 = 1.98, which is significant, $p < .05$, because it lies outside −1.96 and 1.96.

## Task 5

*Use the* split file *command to look at and interpret the descriptive statistics for **numeracy** and **exam**.*

If we want to obtain separate descriptive statistics for each of the universities, we can split the file, and then proceed using the *frequencies* command. To split the file, select Data Split File... or click on . In the resulting dialog box select the option *Organize output by groups*. Once this option is selected, the *Groups Based on* box will activate. Select the variable containing the group codes by

which you wish to repeat the analysis (in this example select **Uni**), and drag it to the box or click on

. By default, SPSS will sort the file by these groups (i.e. it will list one category followed by the other in the data editor window). Once you have split the file, use the *frequencies* command. Let's request statistics for only **numeracy** and **exam** scores for the time being.



**Dialog box for the *split file* command**

|  | Duncetown University | | | | Sussex University | |
|---|---|---|---|---|---|---|

Duncetown University

**Statistics**[b]

|  |  | Percentage on SPSS exam | Numeracy |
|---|---|---|---|
| N | Valid | 50 | 50 |
|  | Missing | 0 | 0 |
| Mean |  | 40.18 | 4.12 |
| Std. Error of Mean |  | 1.780 | .292 |
| Median |  | 38.00 | 4.00 |
| Mode |  | 34[a] | 4 |
| Std. Deviation |  | 12.589 | 2.067 |
| Variance |  | 158.477 | 4.271 |
| Skewness |  | .309 | .512 |
| Std. Error of Skewness |  | .337 | .337 |
| Kurtosis |  | -.567 | -.484 |
| Std. Error of Kurtosis |  | .662 | .662 |
| Range |  | 51 | 8 |
| Minimum |  | 15 | 1 |
| Maximum |  | 66 | 9 |

a. Multiple modes exist. The smallest value is shown
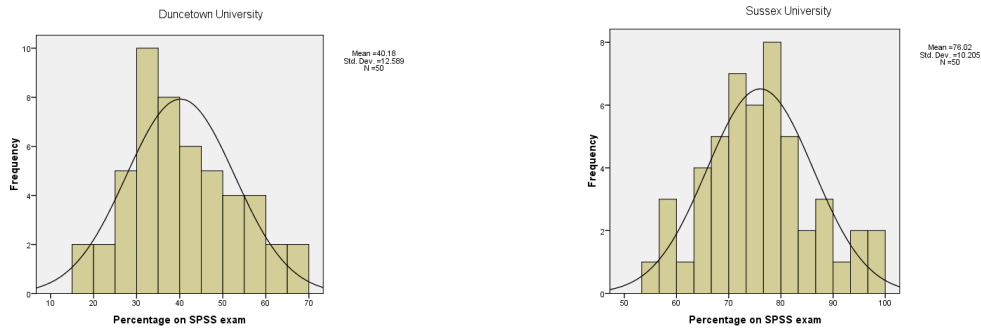b. University = Duncetown University

Sussex University

**Statistics**[b]

|  |  | Percentage on SPSS exam | Numeracy |
|---|---|---|---|
| N | Valid | 50 | 50 |
|  | Missing | 0 | 0 |
| Mean |  | 76.02 | 5.58 |
| Std. Error of Mean |  | 1.443 | .434 |
| Median |  | 75.00 | 5.00 |
| Mode |  | 72[a] | 5 |
| Std. Deviation |  | 10.205 | 3.071 |
| Variance |  | 104.142 | 9.432 |
| Skewness |  | .272 | .793 |
| Std. Error of Skewness |  | .337 | .337 |
| Kurtosis |  | -.264 | .260 |
| Std. Error of Kurtosis |  | .662 | .662 |
| Range |  | 43 | 13 |
| Minimum |  | 56 | 1 |
| Maximum |  | 99 | 14 |

a. Multiple modes exist. The smallest value is shown
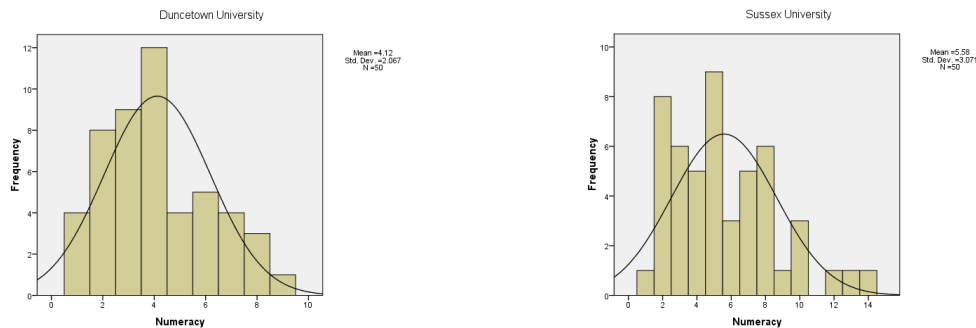b. University = Sussex University

**SPSS output**

The SPSS output is split into two sections: first the results for students at Duncetown University, then the results for those attending Sussex University. The output shows the two main summary tables. From these tables it is clear that Sussex students scored higher on both their SPSS exam and the numeracy test than their Duncetown counterparts. In fact, looking at the means reveals that, on average, Sussex students scored an amazing 36% more on the SPSS exam than Duncetown students, and had higher numeracy scores too (what can I say, my students are the best).

SPSS Exam Mark



Numeracy



**Distributions of exam and numeracy scores for Duncetown University and Sussex University students**

The figure shows the histograms of these variables split according to the university attended. The first interesting thing to note is that for exam marks, the distributions are both fairly normal. This seems odd because the overall distribution was bimodal. However, it starts to make sense when you consider that for Duncetown the distribution is centred around a mark of about 40%, but for Sussex the distribution is centred around a mark of about 76%. This illustrates how important it is to look at distributions within groups. If we were interested in comparing Duncetown to Sussex it wouldn't matter that overall the distribution of scores was bimodal; all that's important is that each group comes from a normal distribution, and in this case it appears to be true. When the two samples are combined, these two normal distributions create a bimodal one (one of the modes being around the centre of the Duncetown distribution, and the other being around the centre of the Sussex data!). For numeracy scores, the distribution is slightly positively skewed (there is a larger concentration at the lower end of scores) in both the Duncetown and Sussex groups. Therefore, the overall positive skew observed before is due to the mixture of universities. When you have finished with the *split file* command, remember to *switch it off* (otherwise SPSS will carry on doing every analysis on each group separately). To switch this function off, return to the *Split File* dialog box and select <u>A</u>*nalyze all cases, do not create groups*.

# Task 6

*Repeat Task 5 but for the computer literacy and percentage of lectures attended.*

The SPSS output is split into two sections: first, the results for students at Duncetown University, then the results for those attending Sussex University. From these tables it is clear that Sussex and Duncetown students scored similarly on computer literacy (both means are very similar). Sussex students attended slightly more lectures (63.27%) than their Duncetown counterparts (56.26%).

The histograms are also split according to the university attended. All of the distributions look fairly normal. The only exception is the computer literacy scores for the Sussex students. This is a fairly flat distribution apart from a huge peak between 50 and 60%. It's slightly heavy-tailed (right at the very ends of the curve the bars come above the line) and very pointy. This suggests positive kurtosis. If you examine the values of kurtosis you will find that there is significant ($p < .05$) positive kurtosis: 1.38/0.662 = 2.08, which falls outside of −1.96 and 1.96.

Duncetown University

Sussex University

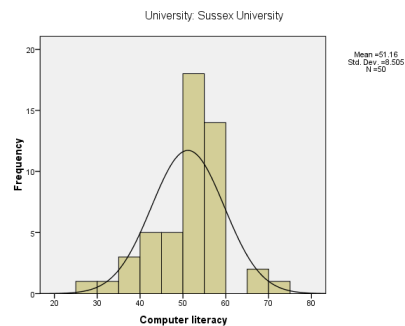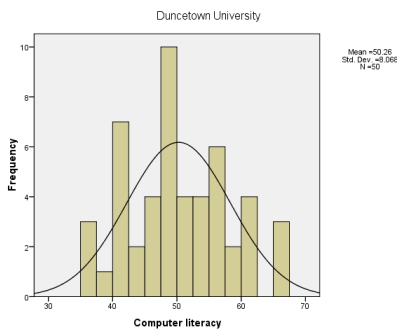**Statistics[b]**

| | | Computer literacy | Percentage of lectures attended |
|---|---|---|---|
| N | Valid | 50 | 50 |
| | Missing | 0 | 0 |
| Mean | | 50.26 | 56.260 |
| Std. Error of Mean | | 1.141 | 3.3619 |
| Median | | 49.00 | 60.500 |
| Mode | | 48[a] | 48.5[a] |
| Std. Deviation | | 8.068 | 23.7726 |
| Variance | | 65.094 | 565.135 |
| Skewness | | .225 | -.309 |
| Std. Error of Skewness | | .337 | .337 |
| Kurtosis | | -.515 | -.383 |
| Std. Error of Kurtosis | | .662 | .662 |
| Range | | 32 | 92.0 |
| Minimum | | 35 | 8.0 |
| Maximum | | 67 | 100.0 |

a. Multiple modes exist. The smallest value is shown

b. University = Duncetown University

**Statistics[b]**

| | | Computer literacy | Percentage of lectures attended |
|---|---|---|---|
| N | Valid | 50 | 50 |
| | Missing | 0 | 0 |
| Mean | | 51.16 | 63.270 |
| Std. Error of Mean | | 1.203 | 2.6827 |
| Median | | 54.00 | 65.750 |
| Mode | | 54 | 42.0[a] |
| Std. Deviation | | 8.505 | 18.9697 |
| Variance | | 72.341 | 359.849 |
| Skewness | | -.538 | -.365 |
| Std. Error of Skewness | | .337 | .337 |
| Kurtosis | | 1.379 | -.221 |
| Std. Error of Kurtosis | | .662 | .662 |
| Range | | 46 | 87.5 |
| Minimum | | 27 | 12.5 |
| Maximum | | 73 | 100.0 |

a. Multiple modes exist. The smallest value is shown

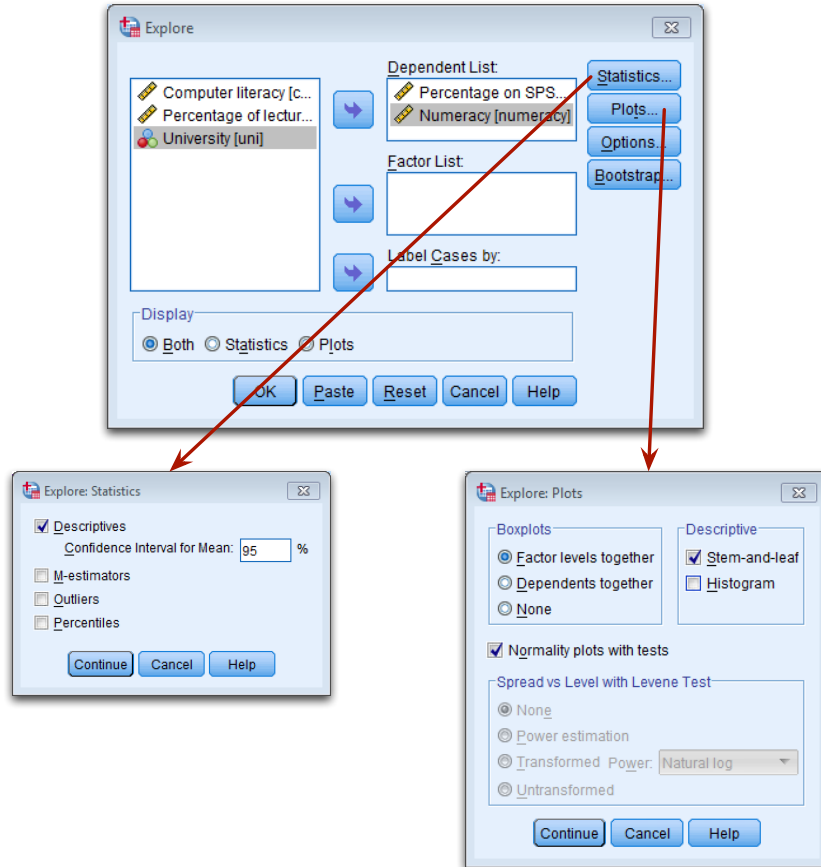b. University = Sussex University

Computer Literacy



Percentage of Lectures Attended

# Task 7

*Conduct and interpret a K-S test for **numeracy** and **exam**.*

The Kolmogorov–Smirnov (K-S) test can be accessed through the *explore* command (Analyze Descriptive Statistics ▶ Explore… ). First, enter any variables of interest in the box labelled *Dependent List* by highlighting them on the left-hand side and transferring them by clicking on . For this example, just select the exam scores and numeracy scores. It is also possible to select a factor (or grouping variable) by which to split the output (so if you select **Uni** and transfer it to the box labelled *Factor List*, SPSS will produce exploratory analysis for each group — a bit like the *split file* command). If you click on Statistics… a dialog box appears, but the default option is fine (it will produce means, standard deviations and so on). The more interesting option for our current purposes is accessed by clicking on Plots… . In this dialog box select the option ✔ Normality plots with tests , and this will produce both the K-S test and some graphs called *normal Q-Q plots*. Click on Continue to return to the main dialog box and then click on OK to run the analysis.

**Dialog boxes for the *explore* command**

The first table produced by SPSS contains descriptive statistics (mean, etc.) and should have the same values as the tables obtained using the frequencies procedure. The important table is that of the K-S test, which includes the test statistic itself, the degrees of freedom (which should equal the sample size) and the significance value of this test. Remember that a significant value (*Sig.* less than .05) indicates a deviation from normality. For both numeracy and SPSS exam scores, the K-S test is highly significant, indicating that both distributions are not normal. This result is likely to reflect the bimodal distribution found for exam scores, and the positively skewed distribution observed in the numeracy scores. However, these tests confirm that these deviations were *significant*. (But bear in mind that the sample is fairly big.)

**Tests of Normality**

| | Kolmogorov-Smirnov[a] | | | Shapiro-Wilk | | |
|---|---|---|---|---|---|---|
| | Statistic | df | Sig. | Statistic | df | Sig. |
| Percentage on SPSS exam | .102 | 100 | .012 | .961 | 100 | .005 |
| Numeracy | .153 | 100 | .000 | .924 | 100 | .000 |

a. Lilliefors Significance Correction

**SPSS Output**

We can report the results in the SPSS output in the following way:

- The percentages on the SPSS exam, $D(100) = 0.10$, $p < .05$, and the numeracy scores, $D(100) = 0.15$, $p < .001$, were both significantly non-normal.

As a final point, bear in mind that when we looked at the exam scores for separate groups, the distributions seemed quite normal; now if we'd asked for separate tests for the two universities (by placing **Uni** in the box labelled *Factor List*) the K-S test might not have been significant. In fact if you try this out, you'll find that the percentages on the SPSS exam are indeed normal within the two groups (the values in the *Sig.* column are greater than .05). This is important because if our analysis involves comparing groups, then what's important is not the overall distribution but the distribution in each group.
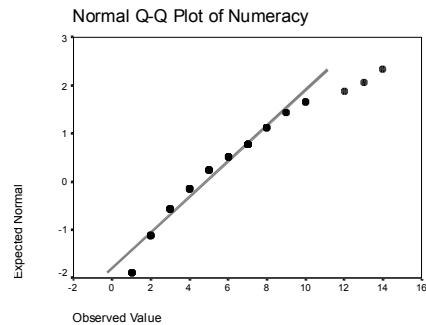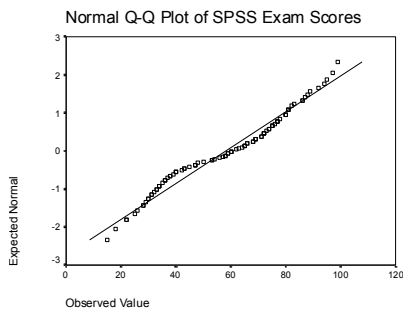
**Tests of Normality**

| | University | Kolmogorov-Smirnov[a] | | | Shapiro-Wilk | | |
|---|---|---|---|---|---|---|---|
| | | Statistic | df | Sig. | Statistic | df | Sig. |
| Percentage on SPSS exam | Duncetown University | .106 | 50 | .200* | .972 | 50 | .283 |
| | Sussex University | .073 | 50 | .200* | .984 | 50 | .715 |
| Numeracy | Duncetown University | .183 | 50 | .000 | .941 | 50 | .015 |
| | Sussex University | .155 | 50 | .004 | .932 | 50 | .007 |

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

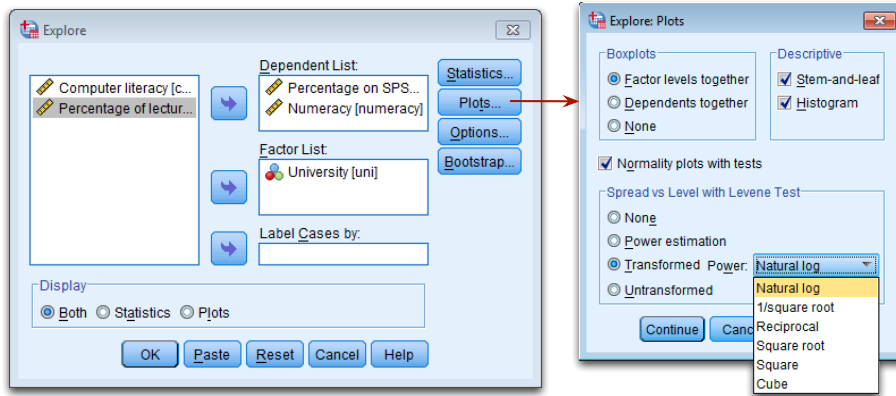a. Lilliefors Significance Correction

**SPSS Output**



**Normal Q-Q plots of numeracy and SPSS exam scores**

SPSS also produces a normal Q-Q plots. In both of the variables analysed we already know that the data are not normal, and these plots confirm this observation because the dots deviate substantially from the line. It is noteworthy that the deviation is greater for the numeracy scores, and this is consistent with the higher significance value of this variable on the K-S test.

## Task 8

*Conduct and interpret a Levene's test for **numeracy** and **exam**.*

We can get Levene's test using the *Explore* menu: Analyze Descriptive Statistics ▸ Explore... . Transfer the SPSS exam scores and the numeracy scores from the list on the left-hand side to the box labelled *Dependent List* by clicking on the ⮕ next to this box, and because we want to split the output by the grouping variable to compare the variances, select the variable **Uni** and transfer it to the box labelled *Factor List* by clicking on the appropriate ⮕. Then click on Plots... to open another dialog box. To get Levene's test we need to select one of the options where it says *Spread vs. level with Levene's test*. If you select ⦿ Untransformed Levene's test is carried out on the raw data (a good place to start). When you've finished with this dialog box click on Continue to return to the main *Explore* dialog box and then click on OK to run the analysis.

**Exploring groups of data and obtaining Levene's test**

Levene's test is non-significant for the SPSS exam scores (values in the *Sig.* column are more than .05), indicating that the variances are not significantly different (i.e. they are similar and the homogeneity of variance assumption is tenable). However, for the numeracy scores, Levene's test is significant (values in the *Sig.* column are less than .05), indicating that the variances are significantly different (i.e., they are not the same and the homogeneity of variance assumption has been violated). We can also calculate the variance ratio. For SPSS exam scores the variance ratio is 158.48/104.14 = 1.52 and for numeracy scores the value is 9.43/4.27 = 2.21. Our group sizes are 50 and we're comparing two variances so the critical value is (from the table in the additional material) approximately 1.67. These ratios concur with Levene's test: variances are significantly different for numeracy scores (2.21 is bigger than 1.67) but not for SPSS exam scores (1.52 is smaller than 1.67). We could report these findings as follows:

- For the percentage on the SPSS exam, the variances were equal for Duncetown and Sussex University students, $F(1, 98) = 2.58$, *ns*, but for numeracy scores the variances were significantly different in the two groups, $F(1, 98) = 7.37$, $p < .01$.
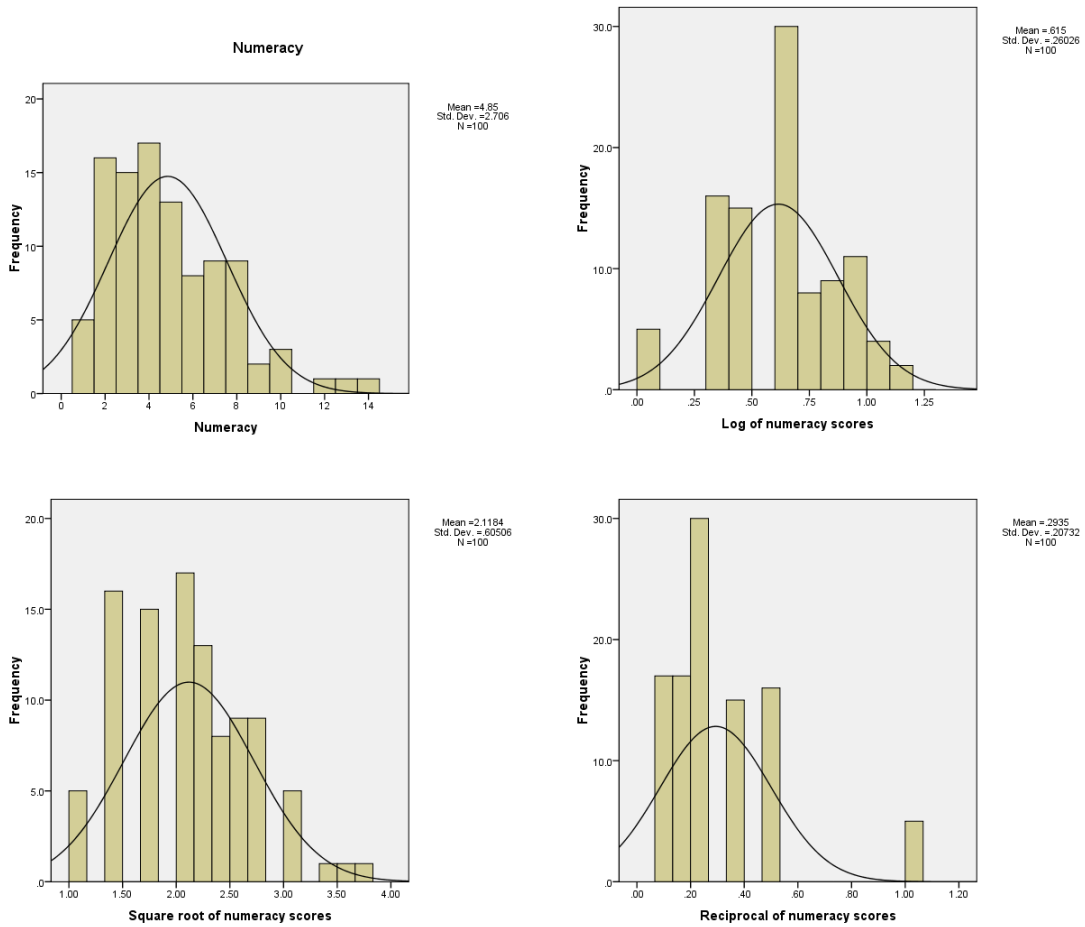
**Test of Homogeneity of Variance**

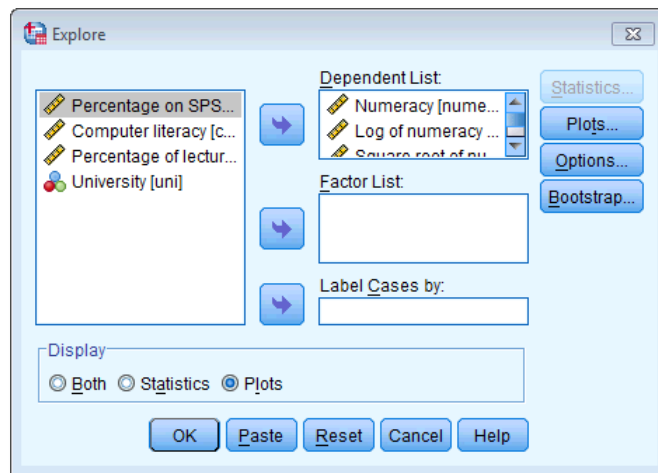| | | Levene Statistic | df1 | df2 | Sig. |
|---|---|---|---|---|---|
| Percentage on SPSS exam | Based on Mean | 2.584 | 1 | 98 | .111 |
| | Based on Median | 2.089 | 1 | 98 | .152 |
| | Based on Median and with adjusted df | 2.089 | 1 | 94.024 | .152 |
| | Based on trimmed mean | 2.523 | 1 | 98 | .115 |
| Numeracy | Based on Mean | 7.368 | 1 | 98 | .008 |
| | Based on Median | 5.366 | 1 | 98 | .023 |
| | Based on Median and with adjusted df | 5.366 | 1 | 83.920 | .023 |
| | Based on trimmed mean | 6.766 | 1 | 98 | .011 |

**SPSS output**

# Task 9

*Transform the **numeracy** scores (which are positively skewed) using one of the transformations described in this chapter. Do the data become normal?*

These are the original histogram and those of the transformed scores (I've included all three transformations discussed in the chapter):



None of these histograms appear to be normal. Below is the table of results from the K–S test, all of which are significant. The only conclusion is that although the square root transformation does the best job of normalizing the data, none of these transformations actually works!
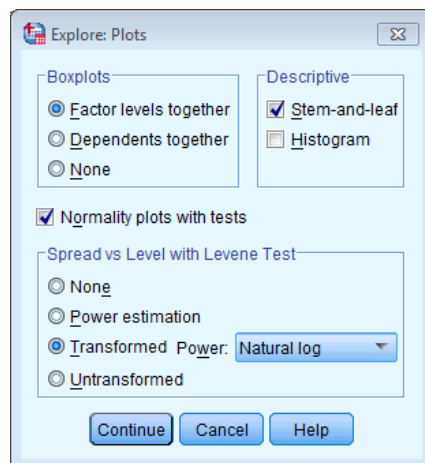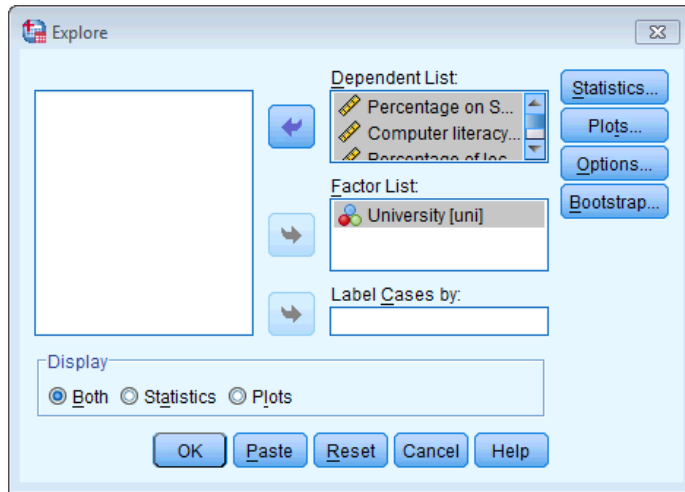
**Tests of Normality**

| | Kolmogorov-Smirnov[a] | | | Shapiro-Wilk | | |
|---|---|---|---|---|---|---|
| | Statistic | df | Sig. | Statistic | df | Sig. |
| Numeracy | .153 | 100 | .000 | .924 | 100 | .000 |
| Log of numeracy scores | .120 | 100 | .001 | .959 | 100 | .003 |
| Square root of numeracy scores | .108 | 100 | .006 | .970 | 100 | .020 |
| Reciprocal of numeracy scores | .223 | 100 | .000 | .763 | 100 | .000 |

a. Lilliefors Significance Correction

## Task 10

*Use the* explore *command to see what effect a natural log transformation would have on the four variables measured in* **SPSSExam.sav***.*

The completed dialog boxes should look like this:





The SPSS output below shows Levene's test on the log-transformed scores. Compare this table to the one in the book (which was conducted on the untransformed SPSS exam scores and numeracy). To recap the book chapter, for the untransformed scores Levene's test was non-significant for the SPSS exam scores (the value in the column labelled *Sig.* was .111,  more than .05) indicating that the

variances were not significantly different (i.e., the homogeneity of variance assumption is tenable). However, for the numeracy scores, Levene's test was significant (the value in the column labelled *Sig.* was .008, less than .05) indicating that the variances were significantly different (i.e. the homogeneity of variance assumption was violated).

For the log-transformed scores (below), the problem has been reversed: Levene's test is now significant for the SPSS exam scores (values in the column labelled *Sig.* are less than .05) but is no longer significant for the numeracy scores (values in the column labelled *Sig.* are more than .05). This reiterates my point from the book chapter that transformations are often not a magic solution to problems in the data!

**Test of Homogeneity of Variance**

| | | Levene Statistic | df1 | df2 | Sig. |
|---|---|---|---|---|---|
| Percentage on SPSS exam | Based on Mean | 25.055 | 1 | 98 | .000 |
| | Based on Median | 24.960 | 1 | 98 | .000 |
| | Based on Median and with adjusted df | 24.960 | 1 | 64.454 | .000 |
| | Based on trimmed mean | 25.284 | 1 | 98 | .000 |
| Computer literacy | Based on Mean | .003 | 1 | 98 | .959 |
| | Based on Median | .037 | 1 | 98 | .847 |
| | Based on Median and with adjusted df | .037 | 1 | 83.139 | .847 |
| | Based on trimmed mean | .004 | 1 | 98 | .953 |
| Percentage of lectures attended | Based on Mean | 5.830 | 1 | 98 | .018 |
| | Based on Median | 3.105 | 1 | 98 | .081 |
| | Based on Median and with adjusted df | 3.105 | 1 | 76.343 | .082 |
| | Based on trimmed mean | 4.760 | 1 | 98 | .032 |
| Numeracy | Based on Mean | .211 | 1 | 98 | .647 |
| | Based on Median | .279 | 1 | 98 | .598 |
| | Based on Median and with adjusted df | .279 | 1 | 97.664 | .598 |
| | Based on trimmed mean | .238 | 1 | 98 | .627 |