

Introduction to Statistical Analysis

Mark Dunning and Sarah Vowler

Last modified: 04 Nov 2015

Contents

1	Introduction	1
2	T-tests practical	1
2.1	The effect of disease on height	1
2.2	Biological processes duration	3
2.3	Blood vessel formation	5
3	Tests for categorical variables	6
3.1	Nucleotide frequency	6
3.2	Disease association	7

1 Introduction

In this practical, we will use several 'read-life' datasets to demonstrate some of the concepts you have seen in the lectures. We will guide you through how to analyse these datasets in Shiny and the kinds of questions you should be asking yourself when faced with similar data.

To answer the questions in this practical we will be using apps that we have developed using the *Shiny* add-on for the *R* statistical package. *R* is a freely-available open-source software that is popular within academic and commercial communities. The functionality within the software compares favourably with other statistical packages (SAS, SPSS and Stata). The downside is that *R* has a steep learning-curve and requires a basic familiarity with command-line software. To ease the transition we have chosen to present this course using a series of online tools that will allow you to perform statistical analysis without having to worry about learning *R*. At the same time, the *R* code required for the analysis will be recorded in the background. You will therefore be able to repeat the analysis at a later date, or pass-on to others. As you gain familiarity with *R* through other courses, you will see how the code generated by Shiny can be adapted to your own needs.

2 T-tests practical

2.1 The effect of disease on height

A scientist knows that the mean height of females in England is 165cm and wants to know whether her patients with disease X have heights that differ significantly from the population mean - we will use a one-sample t-test to test this. The data are contained in the file `diseaseX.csv` and can be analysed online at:-

<http://bioinf-rstud001:3838/OneSampleTest/>

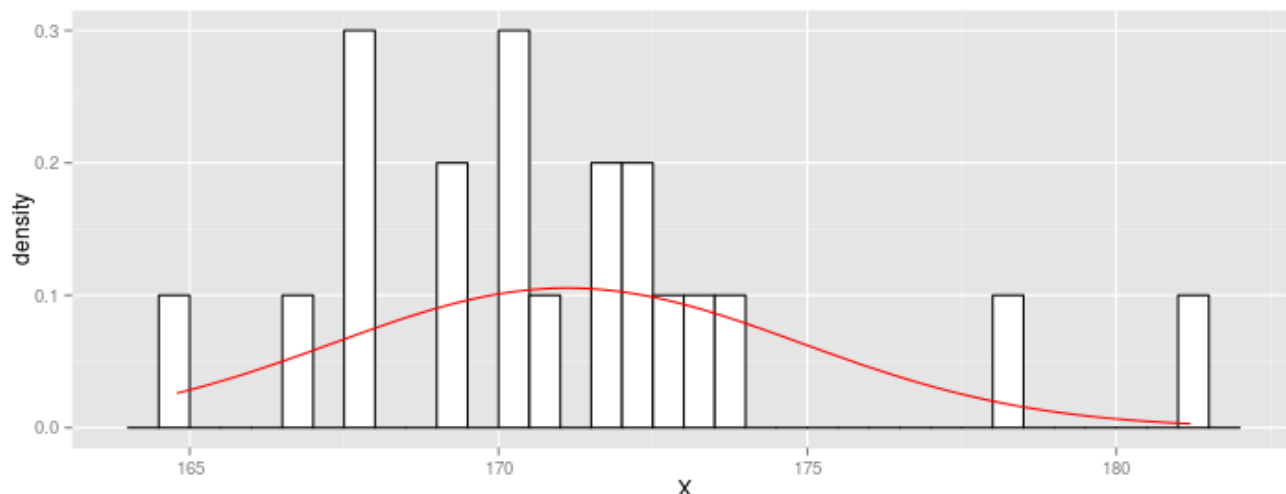
a) What are your null and alternative hypotheses?

```
## Null hypothesis: The mean height of female patients with disease X = 165cm
##      (the population mean for females)
```

```
## Alternative hypothesis: The mean height of female patients with disease X != 165cm
##      (the population mean for females)
```

To import the file `diseaseX.csv` into **Shiny** you will need to select the Choose File option and navigate to where the course data are located on your laptop. You can use the **The data** tab to check that the data has been imported correctly.

b) A histogram of the Height variable will be automatically generated for you. To view it, click on the **Data Distribution**. Do the data look normally distributed? Based on the histogram, is the one-sample t-test appropriate?



In this case the data look normally distributed. Therefore, the one-sample t-test is appropriate.

c) We are interested in knowing whether the mean height in our sample of patients with disease X is different from that of the general population. Perform a **one-sample t-test**

Remember to change the value of **True mean**.

What is the mean height in your sample? What is your value of t? What is the p-value? How do you interpret the p-value?

```
##
## One Sample t-test
##
## data: X
## t = 202.38, df = 19, p-value < 2.2e-16
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 169.3428 172.8822
## sample estimates:
## mean of x
## 171.1125

## Mean height in sample = 171.1cm (95% CI: 169.3-172.9)
## t = 7.23,
## df = 19,
## p <= 0.0001.
##
## Under the null hypothesis, the probability of observing a t-statistic as extreme as 7.23,
## is very small  $P(t \leq 7.23 \mid t \geq 7.23) < 0.0001$ . Therefore, there is strong evidence
## to reject the null hypothesis in favour of the alternative hypothesis. There is strong evidence
## to suggest that the mean height in female patients with disease X is different to the
## population mean height of females of 165cm.
```

2.2 Biological processes duration

In the file `bp_times.csv`, we have the durations of a biological process for two samples of wild-type and knock-out cells (times in seconds). We are interested in seeing whether there is a difference in the durations for the two types of cells – we shall use an **independent t-test** to compare the two cell-types.

These data can be analysed online at <http://bioinf-rstud001:3838/TwoSampleTest/>

a) What are your null and alternative hypotheses?

```
## Null hypothesis: there is no difference in the duration of the
```

```
## biological process for the two cell types.
```

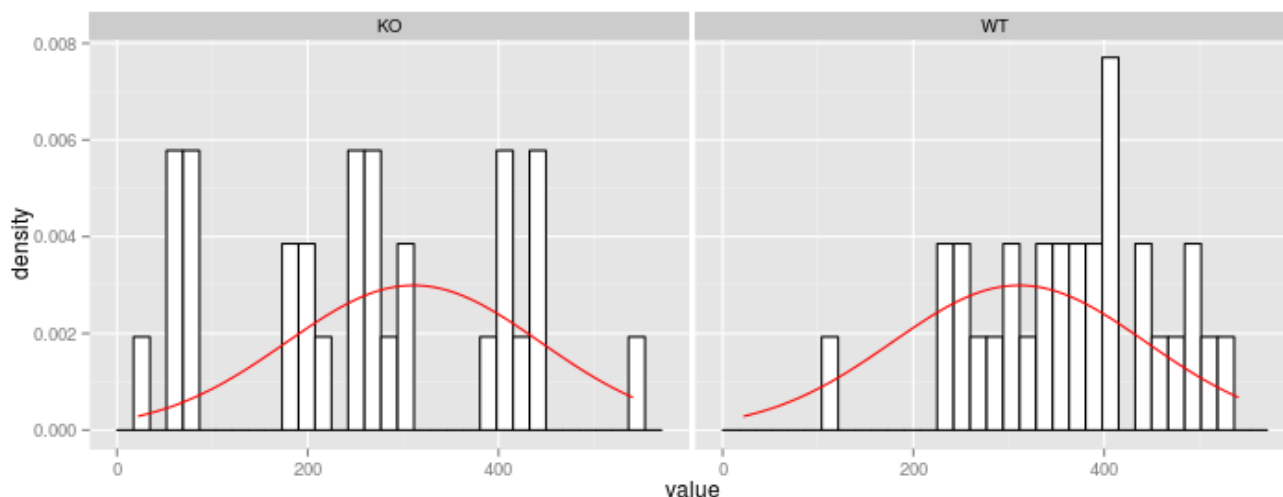
```
##
```

```
## Alternative hypothesis: there is a difference in the duration of the
```

```
## biological process for the two cell types.
```

Import the data using **Choose File** as before. Make sure that the **1st column is a factor?** checkbox is ticked.

b) Histograms to compare the two groups will be created for you automatically. Do the data look normally distributed for each cell-type? Is the independent t-test appropriate.



```
## The data do appear to be approximately normally distributed as we could
```

```
## easily draw a bell shape over each of the two histograms. The independent t-test is appropriate
```

c) Use the *Numerical statistics* analysis to compute descriptive statistics for each group.

```
## $KO
```

```
##      mean      sd   IQR 0% 25% 50% 75% 100%  n
```

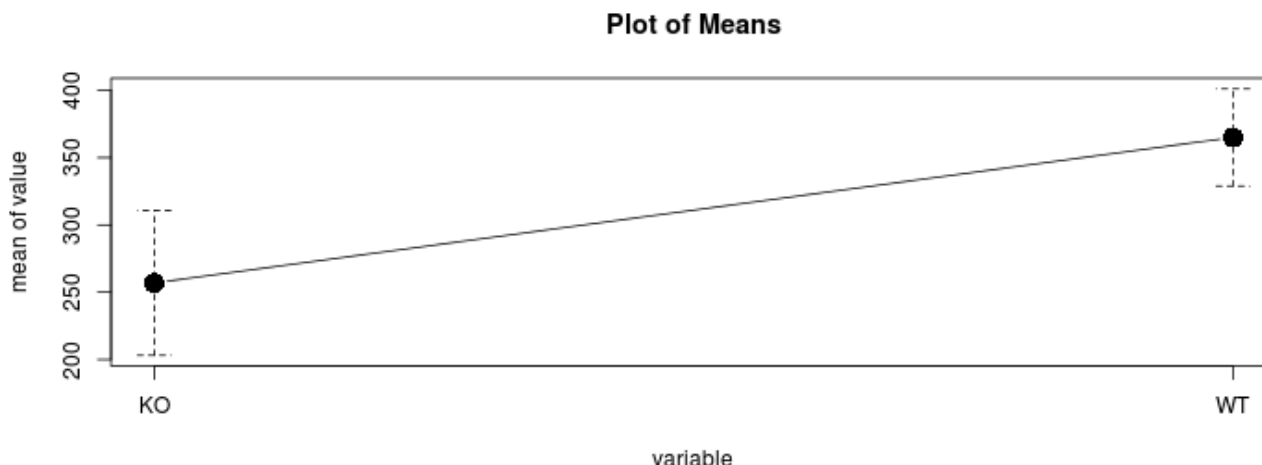
```
## 256.8333 143.9413 219.5 22 177 256 396.5 541 30
```

```
##
```

```
## $WT
```

```
##      mean      sd   IQR 0% 25% 50% 75% 100%  n
```

```
## 365.0333 96.6335 118.75 118 310.25 374 429 530 30
```



Given the distribution of your data, which statistics might you report to summarise your data? Look at and compare the 95% confidence intervals of the mean durations of the two cell-types. Do they overlap?

```
## The normality assumption seems reasonable so the mean and standard deviation
## (or standard error or 95% CI) provide a good summary of the data.
## If the data were skewed, the median and interquartile range would be more meaningful.
## The confidence intervals do not overlap.
```

- d) In order to apply the correct statistical test, we need to test to see if the variances of the two groups are comparable. This is tested for us automatically in the Shiny app. Click the **Histogram** tab to see the result.

What do you conclude from the p-value of this test. How does it influence what test to use?

```
## The test yields a p-value of 0.03568, which is sufficient evidence
## to reject the null hypothesis that the variances of the two groups are the same.
## Therefore we should apply Welch's correction.
```

- e) Use the appropriate test to compare the durations of the two groups.

Is a Welch's correction needed? What is your value of t ? What is the p-value? How do you interpret the p-value? Is this in agreement with the 95% confidence intervals?

```
##
## Welch Two Sample t-test
##
## data: value by variable
## t = -3.4183, df = 50.727, p-value = 0.00125
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -171.75411 -44.64589
## sample estimates:
## mean in group KO mean in group WT
## 256.8333 365.0333

## When we run the independent (unpaired) t test, a formal comparison of the variances
## between the two groups is automatically run. The corresponding p-value is 0.00125 which indicates
## evidence to reject the null hypothesis of equal variances between the two groups.
## Therefore, Welch's correction is needed. We can also look at the histograms and summary statistics
## to see whether we might need to use the Welch's correction. We can see that the standard deviation
## (i.e. spread of data) in each of our groups is quite different (96.63 for WT and 143.94 for KO).
## The widths of the base of our histograms are also different, though not by a very large amount.
```

Both of these suggest a Welch's correction may be needed.

2.3 Blood vessel formation

In blood plasma cancer, there is an increase in blood vessel formation in the bone marrow. A stem cell transplant can be used as a treatment for blood plasma cancer. The bone marrow micro vessel density was measured before and after treatment for 7 patients with blood plasma cancer.

We are interested in seeing whether there is a decrease in the bone marrow micro vessel density after treatment with a stem cell transplant. We will use a paired two-sample t-test to compare the before and after bone marrow micro vessel densities.

The data are contained in the file `bloodplasmacancer2.csv`. These data can be analysed online at <http://bioinf-rstud001:3838/TwoSampleTest/>

a) What are your null and alternative hypotheses?

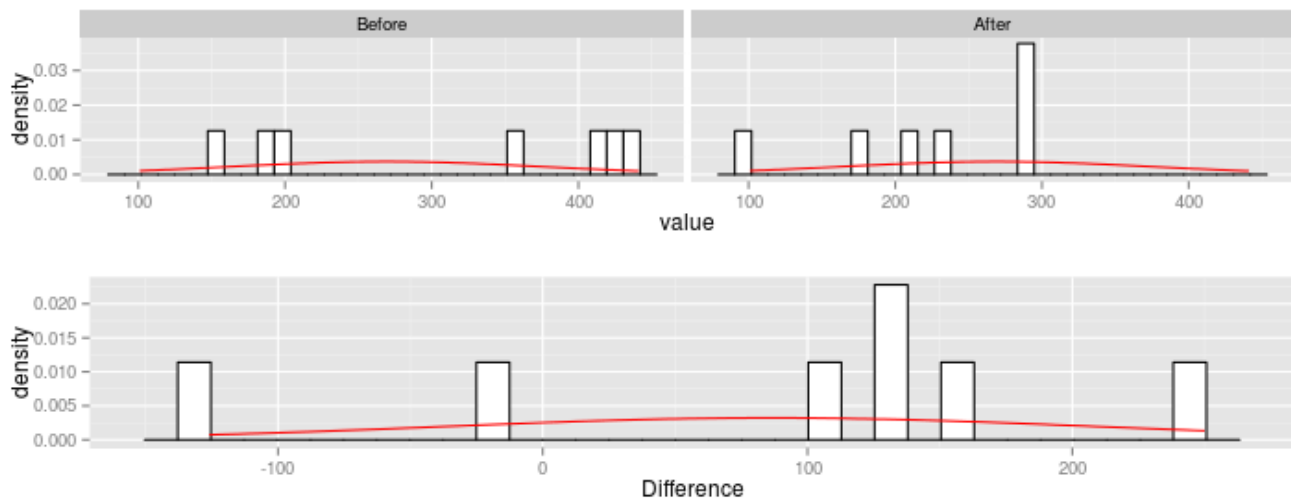
Null hypothesis: the bone marrow micro vessel density after treatment is greater than or equal to the bone marrow micro vessel density before treatment.

##

Alternative hypothesis: the bone marrow micro vessel density after treatment is less than the bone marrow micro vessel density before treatment.

Import the data and create a column of differences (after-before).

b) Plot a histogram of the differences. Do the data look normally distributed? Is the paired t-test appropriate?



From this histogram it is difficult to tell whether the differences between the densities before and after treatment are normally distributed. In situations like this, we may need to draw on the experience of similar sets of measurements.

c) We are interested in seeing whether there is a decrease in the bone marrow micro vessel density after treatment with a stem cell transplant. Is this a one-tailed or two-tailed test?

One-tailed as we are only interested in a decrease.

Usually a two-sided test is preferred unless there is a strong argument for a one-sided test. In this case our treatment is only considered to be effective if we see a reduction in the bone marrow micro vessel density after treatment. Observing an increase in bone marrow density after treatment would lead to the same action/conclusion as if no difference had been observed - the treatment might be

dropped from the research programme (but bear in mind here,
the sample size is small and so only large differences may be detected).

- d) Compare the durations before and after values. Ensure you select the one- or two-tailed test as appropriate. What is the mean difference? What is your value of t ? What is the p -value? How do you interpret the p -value?

```
##
## Paired t-test
##
## data: value by variable
## t = 1.8425, df = 6, p-value = 0.05749
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
## -4.709347      Inf
## sample estimates:
## mean of the differences
##                86.14286

## Under the null hypothesis, the probability of observing a t-statistic as extreme as 1.84, is 0.06,
## slightly greater than 0.05, our nominal significance level. Our result is borderline.
## There is insufficient evidence to reject the null hypothesis.
## Therefore, we might conclude that there is an association of a decrease in bone marrow micro vessel
## density after treatment with bone marrow transplant. It is important to note the small sample size here
## Studying just 7 patients means we will only be able to detect large differences.
```

3 Tests for categorical variables

3.1 Nucleotide frequency

In **Table 1**, we have the frequencies of the four nucleotides in two sequences. We are interested in comparing the nucleotide proportions of the two sequences.

	A	C	G	T	Total
Sequence 1	273.00	233.00	236.00	258.00	250.00
Sequence 2	281.00	246.00	244.00	229.00	250.00
Total	277.00	239.50	240.00	243.50	250.00

Table 1: Nucleotide frequencies for two sequences

- a) What are your null and alternative hypotheses?

```
## Null hypothesis. there is no association between sequence number and nucleotide.
##
```

```
## Alternative hypothesis. there is an association between sequence number and nucleotide.
```

We can analyse these data online in the [Shiny app](#), by modifying the contents of the **Enter your data as a table** box. Columns need to be separated by a '-' character, and rows by a '|'. You can check that you have entered the data correctly by looking at **The data** tab. You do not need to calculate row or column totals.

Note that you do not need to enter the totals.

What is your value of your Chi-squared statistic and its corresponding p -value? How do you interpret the result?

```
##      1      2      3      4
## 1 277 239.5 240 243.5
## 2 277 239.5 240 243.5
```

```
##
## Pearson's Chi-squared test
##
## data: .Table
## X-squared = 2.3286, df = 3, p-value = 0.5071

## Under the null hypothesis, the probability of observing a Chi-squared statistic
## as extreme as 2.3286, is 0.5071. There is no evidence to reject the null hypothesis.
## Therefore, there is no evidence of an association between sequence number and nucleotide.
```

3.2 Disease association

Table 2 gives the frequencies of wild-type and knock-out mice developing a disease thought to be associated to the absence of the knock-out gene.

	WT	KO	Total
Sequence 1	1.00	7.00	4.00
Sequence 2	9.00	3.00	6.00
Total	5.00	5.00	5.00

Table 2: Frequencies of wild-type and knock-out mice developing disease

a) What are your null and alternative hypotheses?

```
## Null hypothesis: there is no association between mouse type and disease X
##
## Alternative hypothesis: there is an association between mouse type and disease X
```

b) What are your expected frequencies?

```
##           WT KO
## Disease    4  4
## No Disease  6  6
```

Enter the data as before;\

c) Select the **Fisher's exact test** option to compare the proportion of mice in each group that developed the disease. What is your p-value? How do you interpret the result?

```
##           WT KO
## Disease    4  4
## No Disease  6  6
```

```
##
## Fisher's Exact Test for Count Data
##
## data: .Table
## p-value = 0.01977
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.0009621944 0.7209145117
## sample estimates:
## odds ratio
## 0.05788421

## p = 0.02. Under the null hypothesis, there is a small probability (p=0.02<0.05)
## of observing such an extreme distribution of the mice given the observed row and column totals.
## There is evidence to reject the null hypothesis in favour of the alternative hypothesis.
```

There is evidence of an association between mouse type and disease X.