

# STM4PSD – Workshop 10

## Hypothesis testing in R

It is very simple to perform a one-sample  $t$ -test in R. If the data is stored in a vector called `x`, then you can use the `t.test` function to perform the test, like so:

```
t.test(x)
```

By default, this will assume you are using a two-sided test with null hypothesis  $H_0: \mu = 0$ .

You can specify the null hypothesis by using the `mu` parameter. For example, to test the hypothesis of Question 2(b), we would use

```
x <- c(76, 88, 91, 90, 79, 96, 89, 92, 100, 87)
t.test(x, mu = 86)
```

The output will look something like this:

```
One Sample t-test

data:  x
t = 1.2418, df = 9, p-value = 0.2457
alternative hypothesis: true mean is not equal to 86
95 percent confidence interval:
 83.69913 93.90087
sample estimates:
mean of x
 88.8
```

Here we can see that it gives you the value of the test statistic, the degrees of freedom, the  $p$ -value, the 95% confidence interval, and the mean.

1. Suppose that the data from Question 1 of Workshop 9 was obtained from the following sample:

```
3.06, 2.00, 3.36, 3.60, 3.72, 4.65, 3.00, 3.74, 3.47, 4.33, 5.97, 3.11, 1.41, 3.63, 3.10, 4.76, 4.86, 7.91, 2.51
```

Use R to perform a  $t$ -test, and verify your answers from that question.

2. A company has developed a new light bulb model. The mean lifetime for their previous model was 16,000 hours. They want to know if their new model performs better than their old model.
  - (a) State appropriate hypotheses. Note that this should be a one-sided test.
  - (b) The file `lights.csv` has data for the lifetime of 300 randomly selected light bulbs (in hours). Use R to perform a `t.test` for the mean lifetime of the new model of light bulb. You can specify the type of test by using the `alternative` parameter. Check the documentation for the `t.test` function to see how it works.
  - (c) State an appropriate conclusion.
3. The term **statistical significance** refers to whether or not one rejects the null hypothesis. On the other hand, **clinical significance** refers to whether or not the deviation from the null hypothesis is actually meaningful.
  - (a) In Question 2, you should have found that increased performance of the new model was statistically significant. Does this imply that the increased performance is meaningful?
  - (b) Determine the 95% confidence interval for the mean lifetime of the new light bulb model.
  - (c) Suppose that the company intends to charge an extra 10% for the new model. Is the performance of the new model clinically significant? In other words, does your confidence interval justify this extra price?
  - (d) If the  $p$ -value was to be reported, explain why it is important to also report the confidence interval for the test.
4. Note that the `t.test` function will work only if you have the full data set. If all you have is the sample mean, sample size, and sample standard deviation, you will need to write your own code to perform the  $t$ -test.

Complete the code below to write your own  $t$ -test function. It is also found on the LMS as `blanks.R`.

```

my.t.test <- function(mean, sd, n, mu) {
  # fill in the blanks to complete the function
  df <- ... # compute the degrees of freedom
  se <- ... # compute the standard error
  t <- ... # compute the test statistic
  p.value <- ... #compute the p-value

  critical.value <- ... # calculate the coefficient of 'SE'
  lower <- ... # find the lower half of the 95% confidence interval
  upper <- ... # find the upper half of the 95% confidence interval

  # You do not need to modify the lines of code below.
  writeLines("\t\tOne Sample t-test for a two-sided hypothesis")
  writeLines("\t\t\tUsing a significance level of 0.05\n")
  writeLines(sprintf("t = %.4f, df = %.2f, p-value = %.4f", t, df, p.value))
  writeLines(sprintf("alternative hypothesis: true mean not equal to %.f4", mu))
  writeLines(sprintf("95 percent confidence interval:\n %.5f %.5f", lower, upper))
  writeLines(sprintf("sample estimates:\nmean of x\n\t%.5f", mean))
}

```

## Two-sample hypothesis testing in R

The `t.test` function in R can conduct either a paired or unpaired  $t$ -test. By default, it assumes that an unpaired  $t$ -test is desired. If the data of the first sample is in a vector  $x$ , and the data of the second sample is in a vector  $y$ , then the following commands will perform the tests:

```

t.test(x, y, paired=TRUE)      # paired t-test.
t.test(x, y)                   # unpaired t-test.

```

It is important to ensure you use the `paired` argument when needed.

5. The data from Question 5 of Workshop 9 is available on the LMS in the file `milk.csv`. Load the data into R, and then verify the output shown to you in that question, and verify your calculations from Questions 5(d) and 6.
6. Consider two different, competing battery types. For simplicity call these batteries *Battery A* and *Battery B*. Battery A is 10% more expensive, but may have a longer lifetime. To test this, 500 batteries of each type are tested and their battery life (the time until the battery is flat after continuous usage) is recorded. The data is on the LMS in a file called `BatteryLife.csv`. Load this data into R. This question will assume that the resulting data frame in R is called `BatteryLife`. We also assume that the battery life random variable for each is normal so that the  $t$ -test is appropriate to test the hypotheses

$$H_0: \mu_A = \mu_B \text{ versus } H_1: \mu_A \neq \mu_B$$

where  $\mu_A$  and  $\mu_B$  are the mean battery life for the respective batteries (measured in hours).

- (a) Execute the appropriate command in R to run the two-sample  $t$ -test.
  - (b) What is the value of the test statistic?
  - (c) What is the value of the  $p$ -value? Do you reject the null hypothesis that the batteries have an equal mean battery life based on this  $p$ -value? Explain.
  - (d) What is the sample mean battery life for Battery A? What is it for Battery B?
  - (e) What is the estimated difference between the two means?
  - (f) What is the 95% confidence interval for  $\mu_1 - \mu_2$ ?
  - (g) Do you reject the null hypothesis based on this confidence interval? Explain.
  - (h) Do you think that your confidence interval suggests that the difference in battery life is meaningful? In other words, do you think that the difference justifies the 10% extra people will pay for Battery A? Explain.
  - (i) If the  $p$ -value was to be reported, explain why it is important to also report the confidence interval for this test.
  - (j) Write a clear and simple statement that summarises your findings.
7. Two companies that produce competing household fire alarms are to be compared. A sample of 360 fire alarms are tested for faults from Company 1, of which 27 were defective. From Company 2, 215 fire alarms were sampled and 24 were defective. Let  $p_1$  denote the true proportion of defective fire alarms from Company 1 and  $p_2$  the proportion of defective units from Company 2.

- (a) By hand, calculate estimates to  $p_1$  and  $p_2$ .
- (b) Carry out a hypothesis test comparing the two proportions by using the R function `prop.test`. From the R output find and report the following:
  - i. The estimates you calculated in (a).
  - ii. The approximate 95% confidence interval for  $p_1 - p_2$ .
  - iii. The  $p$ -value for the test comparing  $p_1$  and  $p_2$ .
- (c) Using the  $p$ -value you reported above, can you reject that the proportions are equal at the  $\alpha = 0.05$  significance level? Explain.
- (d) Does your confidence interval allow you to reject that the proportions are equal? Explain.
- (e) Provide a simple statement that summarises the findings you have reported above.