

# Assignment 8

EMATM0061: Statistical Computing and Empirical Methods, TB1, 2022

Dr. Rihuan Ke

## Introduction

This is the sixth assignment for Statistical Computing and Empirical Methods (Unit EMATM0061) on the MSc in Data Science & MSc in Financial Technology with Data Science. This assignment is mainly based on Lectures 14 and 15 (see the Blackboard).

The submission deadline for this assignment is 23:59, 28 November 2022. Note that this assignment will not count towards your final grade. However, it is recommended that you try to answer the questions to gain a better understanding of the concepts.

## Create an R Markdown for the assignment

It is a good practice to use R Markdown to organize your code and results. You can start with the template called `Assignment08_Template.Rmd` which can be downloaded via Blackboard.

If you are considering submitting your solutions, please generate a PDF file. For example, you can choose the “PDF” option when creating the R Markdown file (note that this option may require Tex to be installed on your computer), or use R Markdown to output an HTML and print it as a PDF file in a browser, or use your own way of creating a PDF file that contains your solutions.

*Only a PDF file will be accepted in the submission of this assignment.* To submit the assignment, please visit the “Assignment” tab on the Blackboard page, where you downloaded the assignment.

## Wish to know more about a particular question?

You may want to ask a question during the computer lab.

Alternatively, we are collecting questions about this assignment that need to be addressed, through the following form. And this can be done either during the labs or outside the lab sessions. So, If you found a question in this assignment interesting but had difficulty in a particular step when trying to develop your answer, please put your remark in the form via the following link. A brief description of the difficulty would be very helpful. Giving your remark is optional, but we aim to know the most common questions that you might want to get some support.

<https://forms.office.com/r/xQM4yX45GC>

## Load packages

Some of the questions in this assignment require the tidyverse package. If it hasn't been installed on your computer, please use `install.packages()` to install them first.

To load the tidyverse package:

```
library(tidyverse)
```

## 1. Obstacles to valid scientific inference

(Q1)

For each of the following give

- (A) an explanation of the concept and
- (B) an example of a situation (real or hypothetical) where they create a barrier to drawing scientific conclusions based on data.

You are encouraged to discuss these concepts with your colleagues.

1. Measurement distortions
2. Selection bias
3. Confounding variables

## 2. paired t-test and effect size

The Barley data set gives the yields of two types of barley - Glabron and Velvet across twelve different fields. The data is paired as yields are given for both types of barley across each of the twelve fields.

```
library(PairedData) # you might need to install the package first
data("Barley")
detach('package:PairedData', unload=TRUE)
detach('package:MASS', unload=TRUE)
# unload package because it contains another select() function

head(Barley, 4)
```

```
##   Farm Glabron Velvet
## 1  F01      49      42
## 2  F02      47      47
## 3  F03      39      38
## 4  F04      37      32
```

(Q1)

Carry out a paired t-test to determine whether there is a difference in average yield between the two types of barley. Use a significance level of 0.01. You can use the `t.test()` function.

(Q2)

Compute the effect size using Cohen's d statistic.

(Q3)

What assumptions are required for the paired t-test? Are these assumptions justified in this case?

## 3. Implementing unpaired t-test

In this question the goal is to create a function called `t_test_function()` which implements an unpaired Student's t-test, in order to test for a difference of population means between two unpaired samples from two distributions. Your function will play a similar role to the following standard R function:

```
t.test(body_mass_g~species, data=peng_AC, var.equal = TRUE)
```

Begin by creating a data frame called “peng\_AC” which is a subset of the Palmer penguins data set consisting of all those penguins which belong to either the “Adelie” or the “Chinstrap” species of penguins.

```
library(palmerpenguins)
peng_AC<-penguins %>%
  drop_na(species,body_mass_g) %>%
  filter(species !="Gentoo")
head(peng_AC %>% select(species, flipper_length_mm, body_mass_g), 5)
```

```
## # A tibble: 5 x 3
##   species flipper_length_mm body_mass_g
##   <fct>         <int>         <int>
## 1 Adelie           181           3750
## 2 Adelie           186           3800
## 3 Adelie           195           3250
## 4 Adelie           193           3450
## 5 Adelie           190           3650
```

### (Q1)

First, try to understand what the following piece of code does.

```
val_col <- "body_mass_g"
group_col <- "species"
data <- peng_AC

data_new <- data %>%
  # rename the columns; note that you can not drop the "!!" (why?)
  rename(group=(!group_col),val=(!val_col))%>%
  group_by(group) %>%
  drop_na(val) %>%
  summarise(mn=mean(val))

data_new
```

```
## # A tibble: 2 x 2
##   group      mn
##   <fct>    <dbl>
## 1 Adelie  3701.
## 2 Chinstrap 3733.
```

```
data_new$mn[2]
```

```
## [1] 3733.088
```

Now, let’s create the function `t_test_function()`. Your function should take in the following arguments:

1. “data” - A data frame argument,
2. “val\_col”- A string argument. This argument is for the column name for a continuous variable (e.g., the body mass column `body_mass_g`),
3. “group\_col” - A string argument. This argument is for the column name for a binary variable (e.g., the species column `species`).

The function should carry out an unpaired Student’s t test based on the value of the continuous variable with column name “val\_col”:

1. The function should begin by partitioning the sample into two groups based on the value of the binary variable named “group\_col”. For example, suppose that the column “group\_col” contains entries

- “Adelie” and “Chinstrap”, then you should partition the sample into two groups corresponding to “Adelie” and “Chinstrap” respectively. You can use the `group_by()` function.
2. Your function should then compute the sample mean, sample variance and sample size for each of these two groups, based upon the variable within the column named “`var_col`”.
  3. Your function should compute a test statistic for the Student’s unpaired t-test (Lecture 20). In addition, the function should compute the corresponding p-value. Finally, your function should compute an estimate for the effect size.
  4. Your function should return a data frame containing the test statistic, p-value and effect size.

Your function should have the following output:

```
t_test_function(data=peng_AC, val_col="body_mass_g", group_col="species")

##           t_stat effect_size      p_val
## 1 -0.5080869 -0.07420226 0.6119085
```

### (Q2)

As an additional challenge you can modify your function so that it takes a fourth argument called “`var_equal`” which takes a Boolean value. If the input of “`var_equal`” is the Boolean “TRUE” your function should compute the test statistic and p-value for an unpaired Student’s t-test. If, on the other hand, the input of “`var_equal`” is the Boolean “FALSE” your function should compute the test statistic and p-value for Welch’s t-test. Your function should have the following output:

```
t_test_function(data=peng_AC, val_col="body_mass_g", group_col="species", val_equal=FALSE)
```

You can compare the output of your function with R’s inbuilt `t.test()` function.

## 4. Useful concepts in statistical hypothesis testing

This question is about the basic concepts in statistical hypothesis testing.

### (Q1)

Explain the following concepts:

1. Null hypothesis
2. Alternative hypothesis
3. Test statistic
4. Type 1 error
5. Type 2 error
6. The size of a test
7. The power of a test
8. The significance level
9. The p-value
10. Effect size

### (Q2)

- (1). Is the p-value the probability that the null hypothesis is true?
- (2). If I conduct a statistical hypothesis test, and my p-value exceeds the significance level, do I have good evidence that the null hypothesis is true?

## 5. Investigating test size for an unpaired Student's t-test

In this question, we shall investigate the performance of an unpaired Student's t-test from the perspective of test size. Recall a Type 1 error occurs when we reject the null hypothesis when the null hypothesis is true. The size of a test is the probability of a Type 1 error. A key property of valid statistical hypothesis tests with a given significance level is that the size of the test does not exceed the significance level.

Note that we can apply unpaired Student's t-test with significance level  $\alpha$  on a samples `sample_0`, `sample_1`:

```
t.test(sample_0,sample_1,var.equal = TRUE, conf.level = 1-alpha)
```

We can apply an unpaired Student's t-test and extract the p-value as follows:

```
t.test(sample_0,sample_1,var.equal = TRUE)$p.value
```

Notice that the significance level wasn't supplied as an argument. Is this a problem?

The following code checks the size of an unpaired Student's t-test with a significance level of  $\alpha = 0.05$ .

```
num_trials<-10000
sample_size<-30
mu_0<-1
mu_1<-1
sigma_0<-3
sigma_1<-3
alpha<-0.05
set.seed(0) # set random seed for reproducibility

single_alpha_test_size_simulation_df <- data.frame(trial=seq(num_trials)) %>%
  # generate random Gaussian samples
  mutate(sample_0=map(.x=trial,.f=~rnorm(n=sample_size,mean=mu_0,sd=sigma_0)),
          sample_1=map(.x=trial,.f=~rnorm(n=sample_size,mean=mu_1,sd=sigma_1))) %>%
  # generate p values
  mutate(p_value=pmmap(.l=list(trial,sample_0,sample_1),
                           .f=~t.test(..2,..3,var.equal = TRUE)$p.value))%>%
  # type I error
  mutate(type_1_error=p_value<alpha)

single_alpha_test_size_simulation_df %>%
  pull(type_1_error) %>%
  mean() # estimate of coverage probability
```

```
## [1] 0.0502
```

Check that you understand the above code.

(Q1)

Modify the above code to explore how the size of the test varies as a function of the significance level  $\alpha$ . You might want to use visualization.

## 6. The statistical power of an unpaired t-test

In this question, we shall investigate the performance of an unpaired Student's t-test from the perspective of statistical power. Recall that the statistical power of a test is the probability of correctly rejecting the null hypothesis when an alternative hypothesis holds.

Consider a setting in which we have two samples i.i.d with Gaussian distribution. The first sample consists of  $n_0$  observations with a population mean  $\mu_0$  and population variance  $\sigma_0^2$ . The second sample consists of  $n_1$  observations with a population mean  $\mu_1$  and population variance  $\sigma_1^2$ .

The following code checks the statistical power of an unpaired Student's t-test in sample sizes  $n_0 = n_1 = 30$ ,  $\mu_0 = 3$ ,  $\mu_1 = 4$ ,  $\sigma_0 = \sigma_1 = 1$  and with a significance level of  $\alpha = 0.05$ .

```
num_trials<-10000

n_0<-30
n_1<-30
mu_0<-3
mu_1<-4
sigma_0<-2
sigma_1<-2

alpha<-0.05
set.seed(0) # set random seed for reproducibility

data.frame(trial=seq(num_trials)) %>%
  # generate random Gaussian samples
  mutate(sample_0 = map(.x=trial,.f=~ rnorm(n=n_0,mean=mu_0,sd=sigma_0)),
          sample_1 = map(.x=trial,.f=~ rnorm(n=n_1,mean=mu_1,sd=sigma_1))) %>%
  # for each sample, generate p value; check examples of pmap() with ?map
  mutate(p_value=pmap(.l = list(trial,sample_0,sample_1),
                          .f =~ t.test(..2, ..3, var.equal = TRUE)$p.value)) %>%
  # estimate of coverage probability
  mutate(reject_null = p_value<alpha ) %>%
  # extract a column
  pull(reject_null) %>%
  # compute probability
  mean()

## [1] 0.4862
```

(Q1)

Conduct a simulation study to explore how the statistical power varies as a function of the significance level.

(Q2)

Conduct a simulation study to explore how the statistical power varies as a function of the difference in means  $\mu_1 - \mu_0$ .

(Q3)

Conduct a simulation study to explore how the statistical power varies as a function of the population standard deviation  $\sigma = \sigma_0 = \sigma_1$ .

(Q4)

Conduct a simulation study to explore how the statistical power varies as a function of the sample size  $n = n_0 = n_1$

## 7. (\*Optional) Comparing the paired and unpaired t-tests on paired data

The aim of this question is to explore the benefits of using a paired test when a natural pairing is available. Consider a situation in which we have two i.i.d. samples  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$ .

Suppose that  $X_1, \dots, X_n \sim \mathcal{N}(\mu_X, \sigma_X^2)$  and for each  $i = 1, \dots, n$ , we have  $Y_i = X_i + Z_i$  where  $Z_1, \dots, Z_n \sim \mathcal{N}(\mu_Z, \sigma_Z^2)$  are independent and identically distributed random variables. It follows that  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$  are independent and identically distributed with  $\mu_Y = \mu_X + \mu_Z$  and  $\sigma_Y^2 = \sigma_X^2 + \sigma_Z^2$ .

In this situation we only observe the two samples  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$ . We are interested in performing a statistical hypothesis test to see if  $\mu_X \neq \mu_Y$ . We have two options here. We could either (1) use the pairing and apply a paired test or (2) ignore the pairing and use an unpaired test. In the console run `?t.test()` to see how to carry out an unpaired and a paired test within R.

### (Q1)

Conduct a simulation study to explore the statistical power of these two approaches. You may want to consider a setting in which  $n = 30, \mu_X = 10, \sigma_X = 5, \mu_Z = 1$  and  $\sigma_Z = 1$ . Consider a range of different significance levels  $\alpha$