

Lab 10

Sharleen Price spp2122

April 5, 2018

Instructions

Make sure that you upload an RMarkdown file to the canvas page (this should have a .Rmd extension) as well as the PDF or HTML output after you have knitted the file. The files you upload to the Canvas page should be updated with commands you provide to answer each of the questions below. You can edit this file directly to produce your final solutions.

Goal

The goal of this lab is to investigate the empirical behavior of a common hypothesis testing procedure through simulation using R. We consider the traditional two-sample t-test.

Two-Sample T-Test

Consider an experiment testing if a 35 year old male's heart rate statistically differs between a control group and a dosage group. Let X denote the control group and let Y denote the drug group. One common method used to solve this problem is the two-sample t-test. The null hypothesis for this study is:

$$H_0 : \mu_1 - \mu_2 = \Delta_0,$$

where Δ_0 is the hypothesized value. The assumptions of the two sample pooled t-test follow below:

Assumptions

1. X_1, X_2, \dots, X_m is a random sample from a normal distribution with mean μ_1 and variance σ_1^2 .
2. Y_1, Y_2, \dots, Y_n is a random sample from a normal distribution with mean μ_2 and variance σ_2^2 .
3. The X and Y samples are independent of one another.

Procedure

The test statistic is

$$t_{calc} = \frac{\bar{x} - \bar{y} - \Delta_0}{\sqrt{\frac{s_1^2}{m} + \frac{s_2^2}{n}}},$$

where \bar{x}, \bar{y} are the respective sample means and s_1^2, s_2^2 are the respective sample standard deviations.

The approximate degrees of freedom is

$$df = \frac{\left(\frac{s_1^2}{m} + \frac{s_2^2}{n}\right)^2}{\frac{(s_1^2/m)^2}{m-1} + \frac{(s_2^2/n)^2}{n-1}}$$

Under the null hypothesis, t_{calc} has a student's t-distribution with df degrees of freedom.

Rejection rules

Alternative Hypothesis	P-value calculation
$H_A : \mu_1 - \mu_2 > \Delta_0$ (upper-tailed)	$P(t_{calc} > T)$
$H_A : \mu_1 - \mu_2 < \Delta_0$ (lower-tailed)	$P(t_{calc} < T)$
$H_A : \mu_1 - \mu_2 \neq \Delta_0$ (two-tailed)	$2 * P(t_{calc} > T)$

Reject H_0 when:

$$Pvalue \leq \alpha$$

Tasks

- 1) Using the **R** function **t.test**, run the two sample t-test on the following simulated dataset. Note that the **t.test** function defaults a two-tailed alternative. Also briefly interpret the output.

```
set.seed(5)
sigma=5
Control <- rnorm(30,mean=10,sd=sigma)
Dosage <- rnorm(35,mean=12,sd=sigma)
?t.test()
t.test(Control, Dosage)

##
##  Welch Two Sample t-test
##
## data:  Control and Dosage
## t = -1.9684, df = 62.014, p-value = 0.05349
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -4.96460632  0.03821408
## sample estimates:
## mean of x mean of y
##  10.05649  12.51969
```

Since the p-value is 0.05349 and since the difference in means are not equal to zero, this indicates that we would reject the null hypothesis and take the alternative hypothesis.

- 2) Write a function called **empirical.size** that simulates **R** different samples of X for control and **R** different samples of Y for the drug group and computes the proportion of test statistics that fall in the rejection region. The function should include the following:

- Inputs:
 - **R** is the number of simulated data sets (simulated test statistics). Let **R** have default 10,000.
 - Parameters **mu1**, **mu2**, **sigma1** and **sigma2** which are the respective true means and true standard deviations of X & Y . Let the parameters have respective defaults **mu1=0**, **mu2=0**, **sigma1=1** and **sigma2=1**.
 - Sample sizes **n** and **m** defaulted at **m=n=30**.
 - **level** is the significance level as a decimal with default at $\alpha = .05$.

- **value** is the hypothesized value defaulted at 0.
- The output should be a **list** with the following labeled elements:
 - **statistic.list** vector of simulated t-statistics (this should have length **R**).
 - **pvalue.list** vector of empirical p-values (this should have length **R**).
 - **empirical.size** is a single number that represents the proportion of simulated test statistics that fell in the rejection region.

I started the function below:

```
emperical.size <- function(R=10000,
                           mu1=0,mu2=0,
                           sigma1=1,sigma2=1,
                           m=30,n=30,
                           level=.05,
                           value=0,
                           direction="two.sided") {

  #Define empty lists
  statistic.list <- rep(0,R)
  pvalue.list <- rep(0,R)

  for (i in 1:R)
  {

    # Sample realized data
    Control <- rnorm(m,mean=mu1,sd=sigma1)
    Dosage <- rnorm(n,mean=mu2,sd=sigma2)

    # Testing values
    testing.procedure <- t.test(Control, Dosage, alternative = direction, mu=value, conf.level = level)
    statistic.list[i] <- testing.procedure$statistic
    pvalue.list[i] <- testing.procedure$p.value
  }

  size.list <- mean(pvalue.list<=level)
  return(list(statistic.list=statistic.list, pvalue.list=pvalue.list, emperical.size=size.list))

}

emperical.size(R=10,mu1=10,mu2=12,sigma1=5,sigma2=5)
```

```
## $statistic.list
## [1] -1.5594821 -1.6265940 -0.3916181 -3.0267377 -0.6315979  0.5023321
## [7]  0.1796087 -2.7168991 -2.1448224 -0.7961500
##
## $pvalue.list
## [1] 0.124574455 0.109418121 0.697191662 0.003707140 0.530363425
## [6] 0.617452414 0.858106351 0.008684406 0.036512133 0.429195763
##
## $emperical.size
## [1] 0.3
```

Evaluate your function with the following inputs **R=10,mu1=10,mu2=12,sigma1=5** and **sigma2=5**.

3) Assuming the null hypothesis

$$H_0 : \mu_1 - \mu_2 = 0$$

is true, compute the empirical size using 10,000 simulated data sets. Use the function **emperical.size** to accomplish this task and store the object as **sim**. Output the empirical size quantity **sim\$size**. Comment on this value. What is it close to?

Note: use **mu1=mu1=10** (i.e., the null is true). Also set **sigma1=5,sigma2=5** and **n=m=30**.

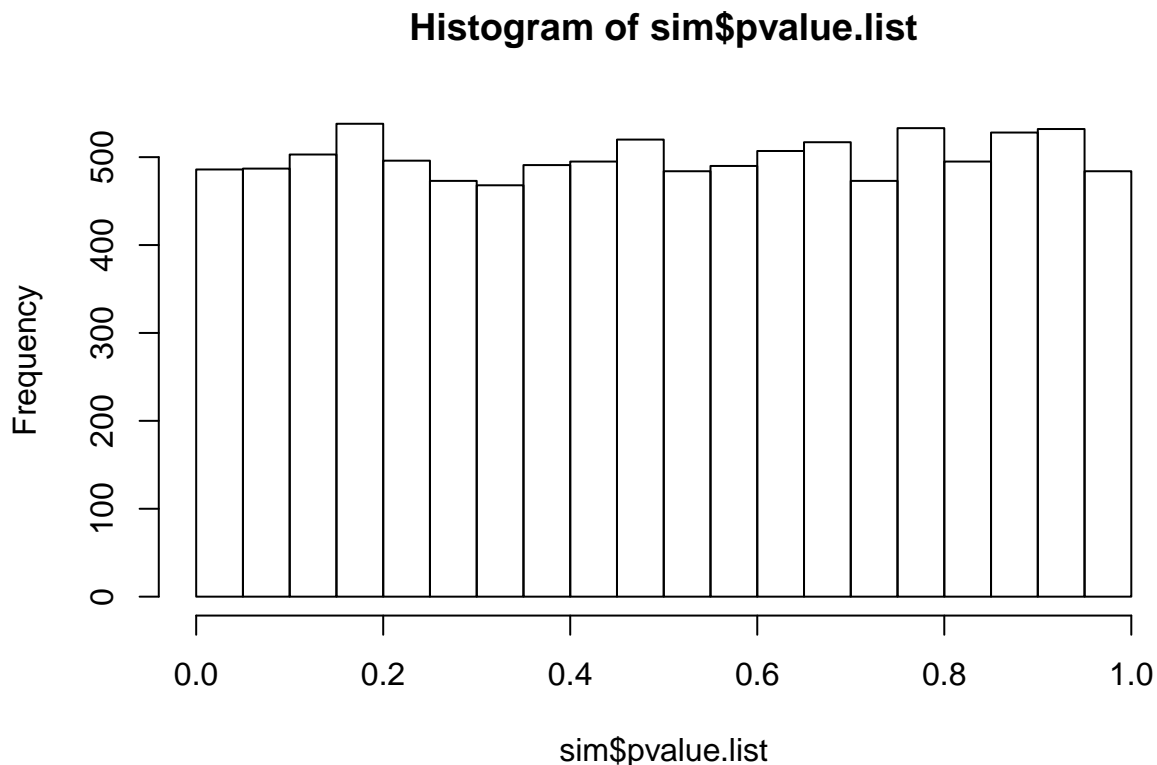
```
sim<-empercal.size(R=10000,mu1=10,mu2=10,sigma1=5,sigma2=5, m=30,n=30)
sim$empercal.size
```

```
## [1] 0.0486
```

This value of empirical size is close to 0 because the amount of time the p-value is less than level/alpha is close to zero so therefore, the null hypothesis is rarely rejected. This makes sense since we set our mu1 and mu2 to be equal and so their difference is zero.

- 4) Plot a histogram of the simulated P-values, i.e., **hist(sim\$pvalue.list)**. What is the probability distribution shown from this histogram? Does this surprise you?

```
hist(sim$pvalue.list)
```

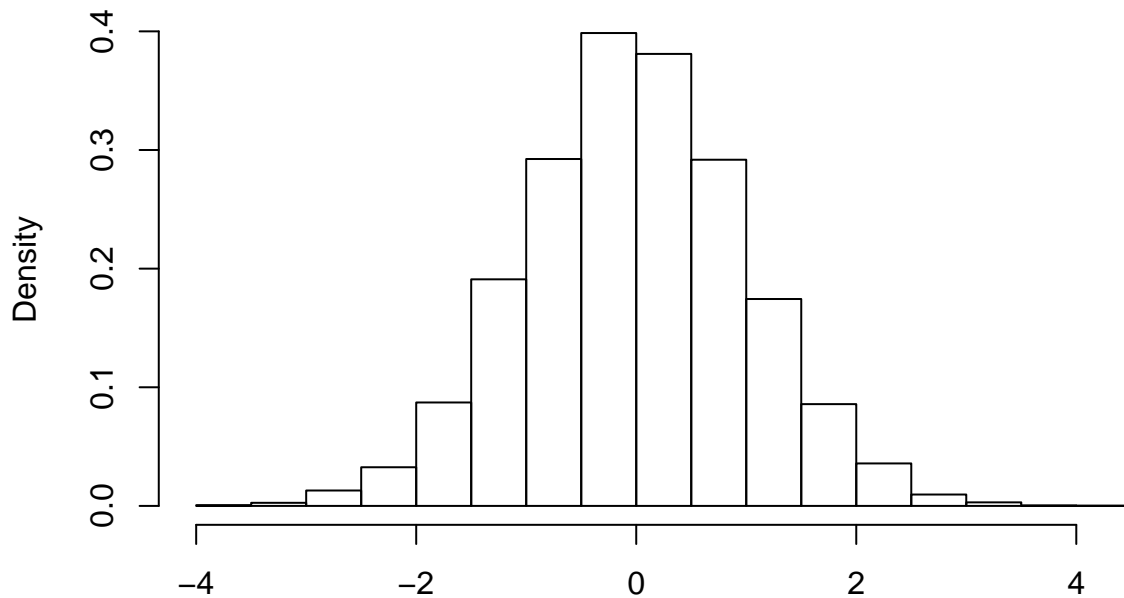


This histogram shows that there is an even distribution of p-values. This is surprising because I would think that since the two samples are the same, we would expect primarily high p-values but this shows that we can not conclude whether or not two identical samples are correlated.

- 5) Plot a histogram illustrating the empirical sampling sampling of the t-statistic, i.e., **hist(sim\$statistic.list,probability =TRUE)**. What is the probability distribution shown from this histogram?

```
hist(sim$statistic.list,probability =TRUE)
```

Histogram of sim\$statistic.list



sim\$statistic.list

This is a student's t-distribution since we are specifically looking at the case where we accept the null hypothesis.

6) Run the following four lines of code:

```
emperical.size(R=1000,mu1=10,mu2=10,sigma1=5,sigma2=5)$emperical.size
```

```
## [1] 0.055
```

```
emperical.size(R=1000,mu1=10,mu2=12,sigma1=5,sigma2=5)$emperical.size
```

```
## [1] 0.315
```

```
emperical.size(R=1000,mu1=10,mu2=14,sigma1=5,sigma2=5)$emperical.size
```

```
## [1] 0.847
```

```
emperical.size(R=1000,mu1=10,mu2=16,sigma1=5,sigma2=5)$emperical.size
```

```
## [1] 0.998
```

This makes sense since the mean increases which means that the hypothesis value is further from zero which indicates that the null hypothesis will be rejected at a higher rate.

7) Run the following four lines of code:

```
emperical.size(R=10000,mu1=10,mu2=12,sigma1=10,sigma2=10,m=10,n=10)$emperical.size
```

```
## [1] 0.0658
```

```
emperical.size(R=10000,mu1=10,mu2=12,sigma1=10,sigma2=10,m=30,n=30)$emperical.size
```

```
## [1] 0.1242
```

```
emperical.size(R=10000,mu1=10,mu2=12,sigma1=10,sigma2=10,m=50,n=50)$emperical.size
```

```
## [1] 0.1702
```

```
emperical.size(R=10000,mu1=10,mu2=12,sigma1=10,sigma2=10,m=100,n=100)$emperical.size
```

```
## [1] 0.2896
```

The empirical size increase as sample sizes increases since it results in the amount of times that the null hypothesis gets rejected to also increase. The sample sizes will naturally cause the mean and the level/alpha value to stray further apart and cause less points to be above the level/alpha levels.