

Assignment 1, EDDA 2017

Fabio Curi Paixao (2592802) Arash Parnia (2591051) - Group 22

12 April 2017

Introduction

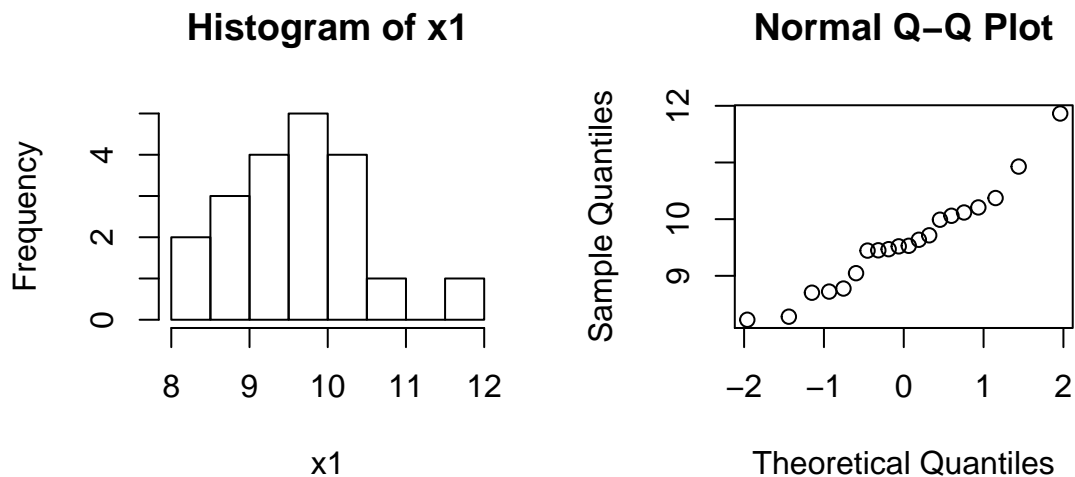
In the present document, the results for the first assignment of the EDDA course are presented.

Exercise 1

```
load(file="assign1.RData")
```

The histogram and QQ-plot for x1, x2, x3, x4 and x5 are shown here-after.

```
par(mfrow=c(1,2));  
hist(x1); qqnorm(x1)
```

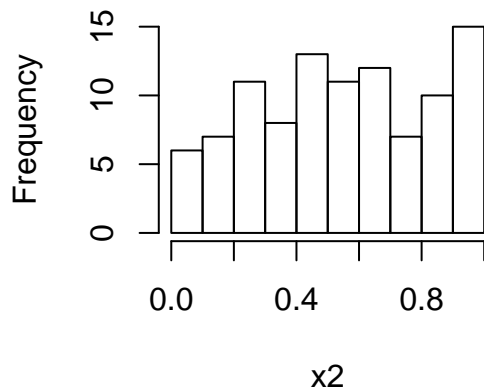


```
shapiro.test(x1);
```

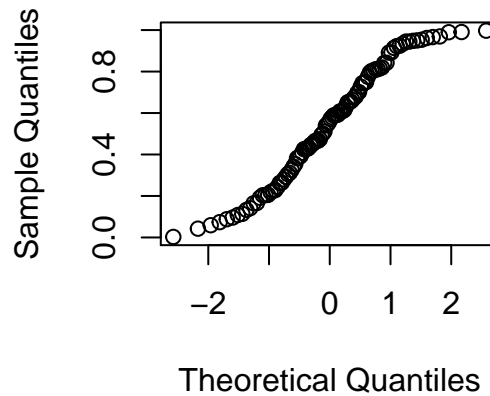
```
##  
## Shapiro-Wilk normality test  
##  
## data: x1  
## W = 0.95401, p-value = 0.432
```

```
par(mfrow=c(1,2));  
hist(x2); qqnorm(x2);
```

Histogram of x2



Normal Q-Q Plot

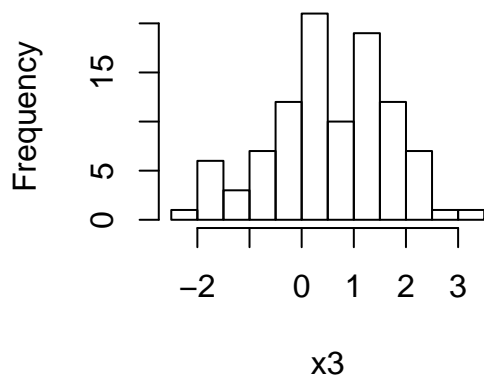


```
shapiro.test(x2);
```

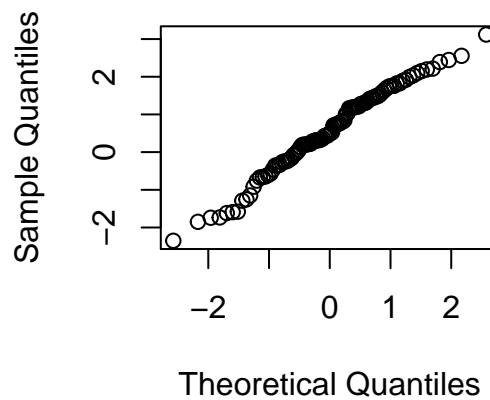
```
##  
## Shapiro-Wilk normality test  
##  
## data: x2  
## W = 0.95825, p-value = 0.003022
```

```
par(mfrow=c(1,2));  
hist(x3); qqnorm(x3);
```

Histogram of x3



Normal Q-Q Plot

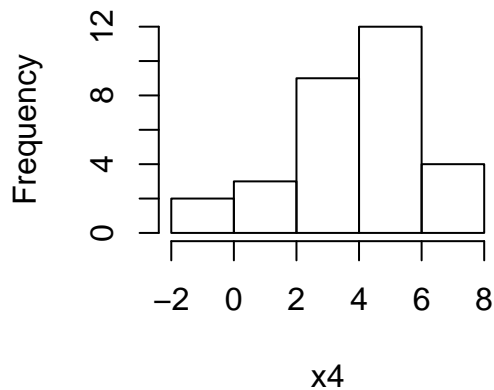


```
shapiro.test(x3);
```

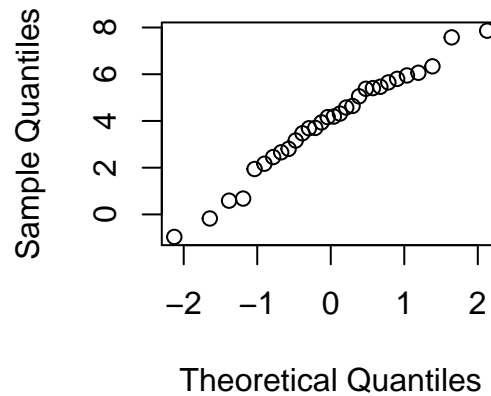
```
##  
## Shapiro-Wilk normality test  
##  
## data: x3  
## W = 0.98461, p-value = 0.2975
```

```
par(mfrow=c(1,2));  
hist(x4); qqnorm(x4);
```

Histogram of x4



Normal Q-Q Plot

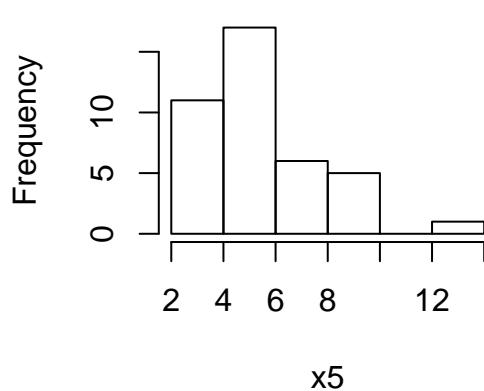


```
shapiro.test(x4);
```

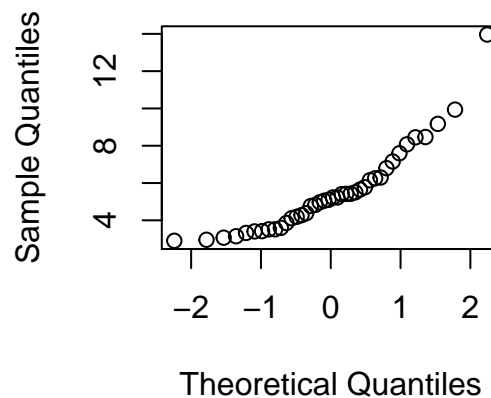
```
##  
##  Shapiro-Wilk normality test  
##  
## data:  x4  
## W = 0.97395, p-value = 0.6517
```

```
par(mfrow=c(1,2));  
hist(x5); qqnorm(x5);
```

Histogram of x5



Normal Q-Q Plot



```
shapiro.test(x5);
```

```
##  
##  Shapiro-Wilk normality test  
##  
## data:  x5  
## W = 0.86828, p-value = 0.0002581
```

From the five histogram and QQ-plots here above, we should be able to identify whether the data follows a normal distribution. A plot using qqnorm of a sample from a normal distribution will show approximately a straight line, and a deviation from a line indicates that the sample was not

taken from a normal population. We have that x1, x2, x3, x4 and x5 have sizes of 20, 100, 100, 30 and 40, respectively.

The elements x2, x3 and x4 look very much like a QQ-plot of a normal distribution. Even though x1 has an uprising behavior, it could still have been sampled from a non-normal distribution as it does not exactly follow a straight line behaviour. Finally, the x5 curve looks more like an exponential one rather than a straight line, thus it does not invite the thought that it has been sampled from a normal distribution.

The Shapiro test suggests to reject the null hypothesis that the distribution is normal in cases x2 and x5. `### Exercise 2`

Part 1

```
mu=nu=180
m=n=30
sd=10
B=1000
p=numeric(B)
for (b in 1:B) {x=rnorm(m,mu,sd)
y=rnorm(n,nu,sd)
p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
```

The 1000 p-values are stored within “p”. The number of elements smaller than 5% are:

```
length(p[p<0.05])
```

```
## [1] 44
```

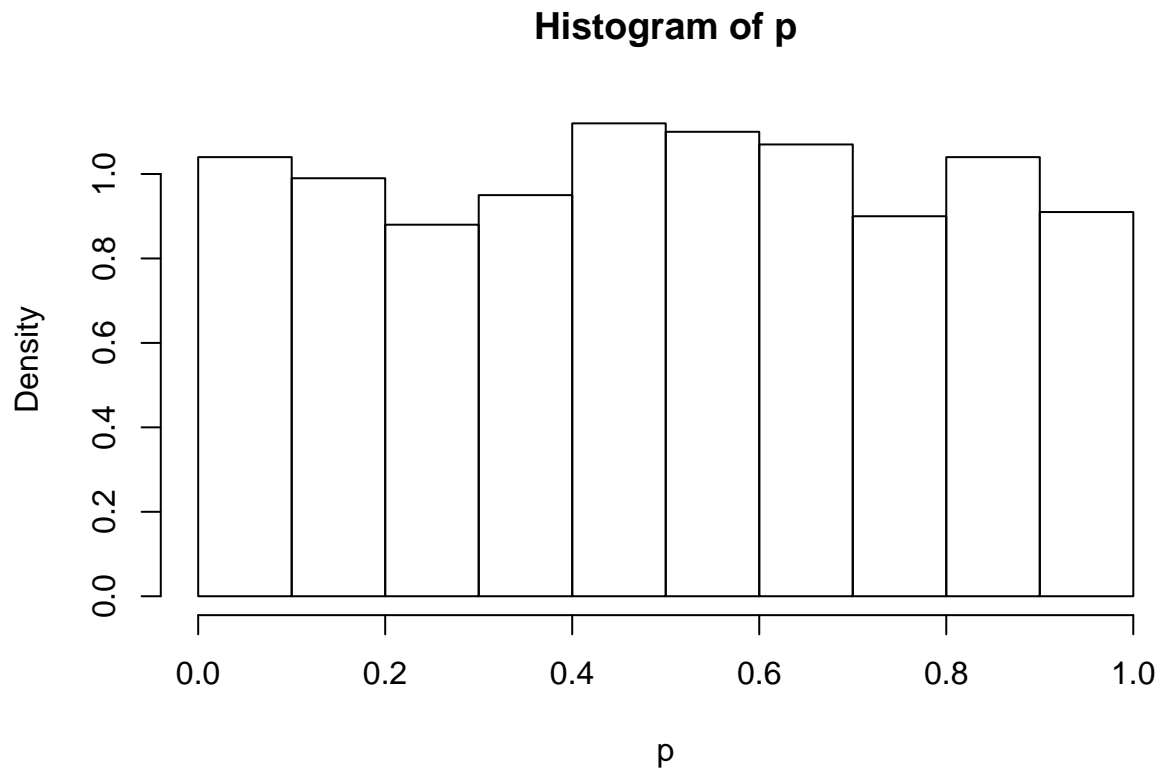
The number of elements smaller than 10% are:

```
length(p[p<0.1])
```

```
## [1] 104
```

The distribution of the p-values is:

```
hist(p,prob=TRUE)
```



Part 2

```
mu=nu=180
m=n=30
sd=1
B=1000
p=numeric(B)
for (b in 1:B) {x=rnorm(m,mu,sd)
y=rnorm(n,nu,sd)
p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
```

The 1000 p-values are stored within “p”. The number of elements smaller than 5% are:

```
length(p[p<0.05])
```

```
## [1] 58
```

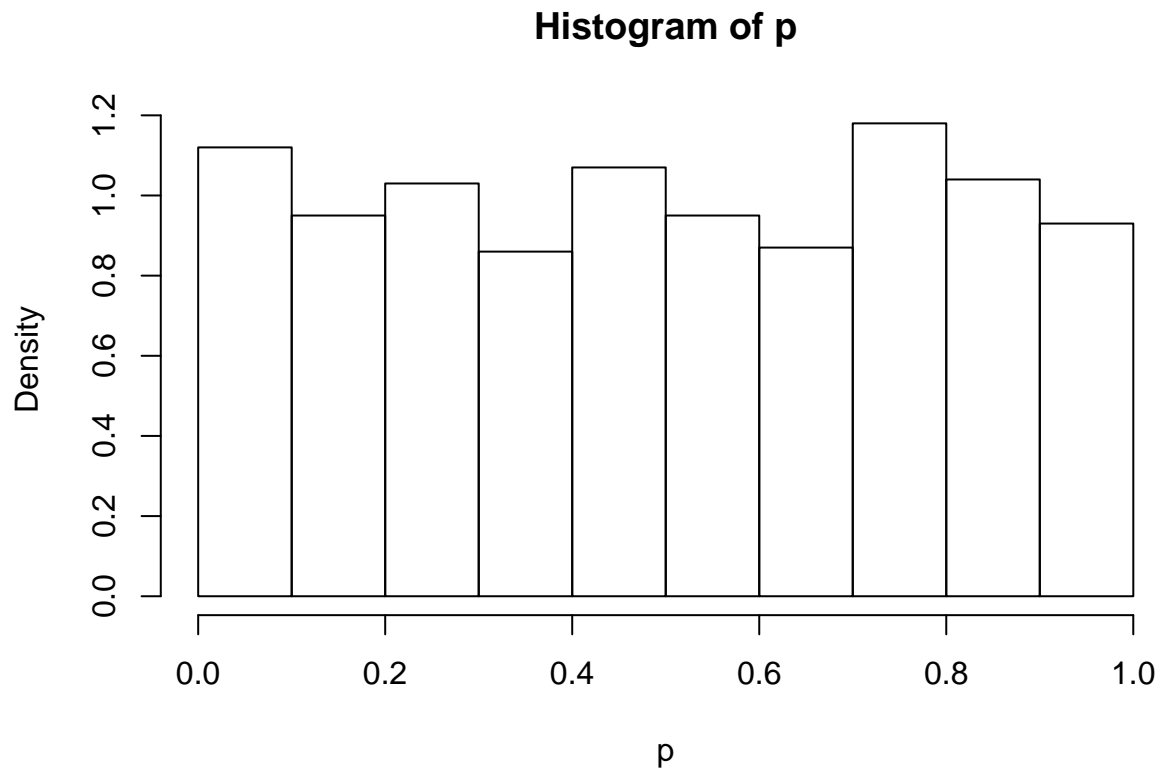
The number of elements smaller than 10% are:

```
length(p[p<0.1])
```

```
## [1] 112
```

The distribution of the p-values is:

```
hist(p,prob=TRUE)
```



Part 3

```
mu=180
nu=175
m=n=30
sd=6
B=1000
p=numeric(B)
for (b in 1:B) {x=rnorm(m,mu,sd)
y=rnorm(n,nu,sd)
p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
```

The 1000 p-values are stored within “p”. The number of elements smaller than 5% are:

```
length(p[p<0.05])
```

```
## [1] 887
```

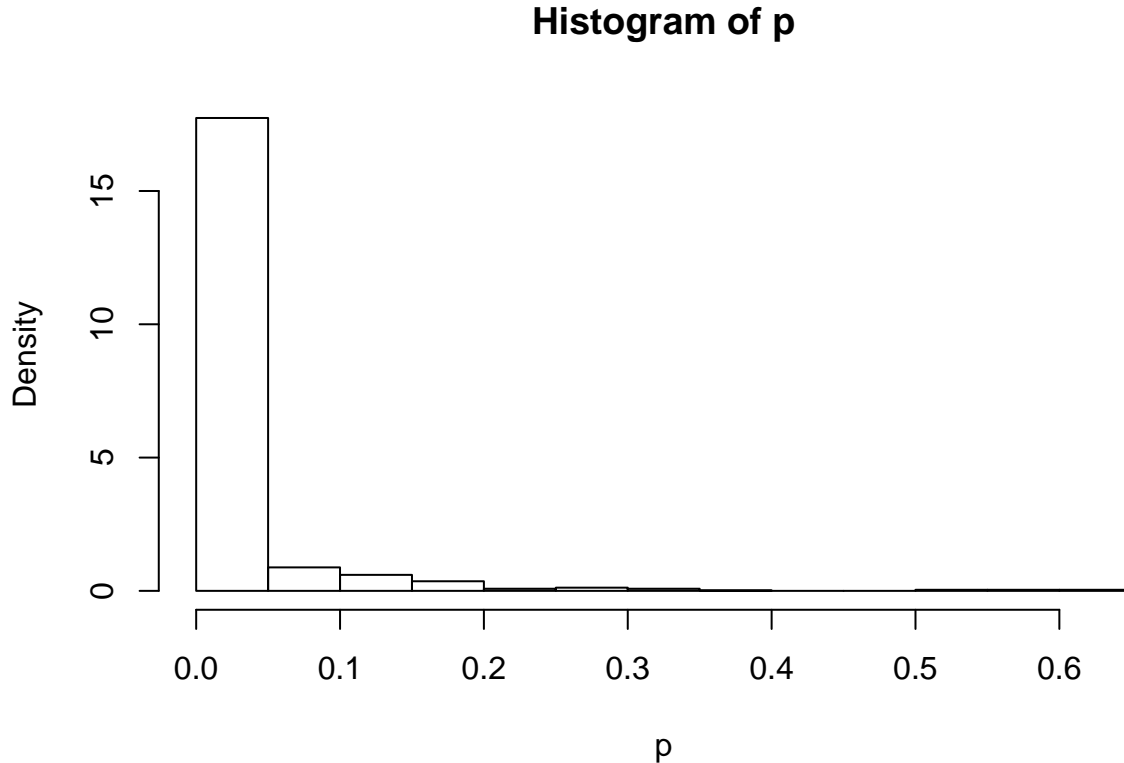
The number of elements smaller than 10% are:

```
length(p[p<0.1])
```

```
## [1] 931
```

The distribution of the p-values is:

```
hist(p,prob=TRUE)
```



Part 4

The test here is to check whether the null hypothesis that the population means of the two populations are equal. The two first histograms show the frequencies of the p-values for two randomly generated x and y arrays with center values of 180, size 30 and differing only by their standard deviations. For $sd=10$, the lowest p-values are seen more often than $sd=1$. Analogically, for $sd=1$, the highest p-values are seen more often than $sd=10$. This is explained by the fact that for smaller standard deviation values, the two arrays x and y have their values within a narrower interval, which increases the probability that the H_0 hypothesis is true. ### Exercise 3

Parts 1 & 2 & 3

Here after we have the code for computing the power of a test with respect to the different nu values. The power of a test is 1 minus the probability of an error of the second kind, which is the p-value. Thus, the power of a test is here defined as $1 - p\text{-value}$.

```
mu=180
m=n=30
sd=5
nu = seq(175,185,by=0.1)
size=length(nu)
p=numeric(size)
for (b in 1:size) {x=rnorm(m,mu,sd)
y=rnorm(n,nu[b],sd)
```

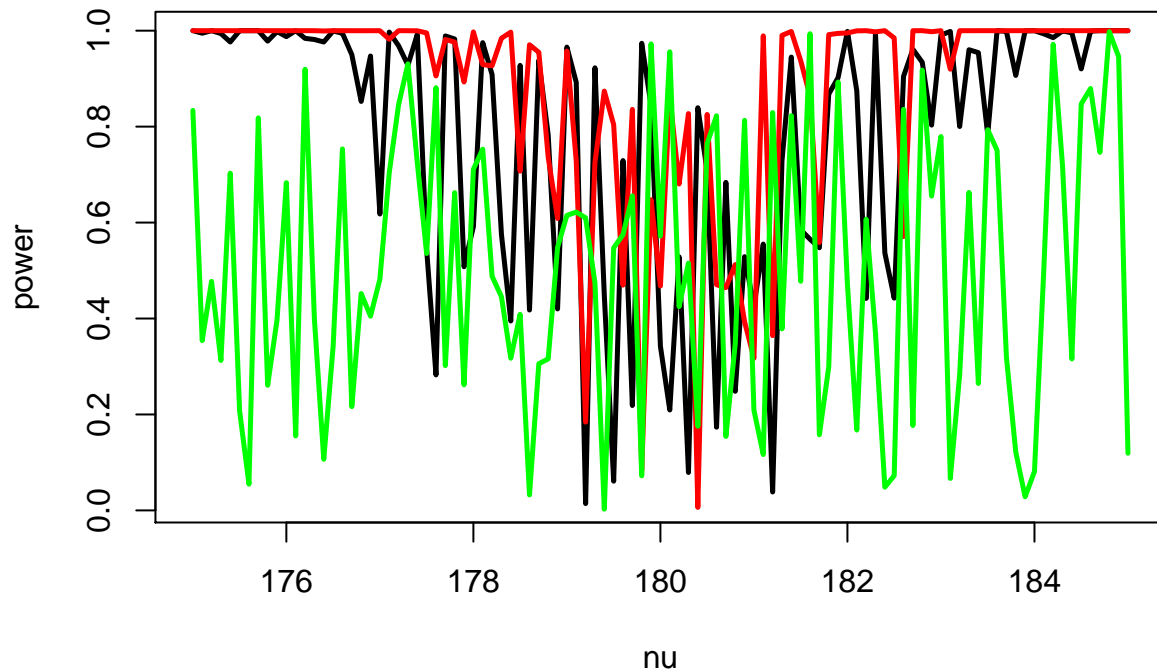
```

p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
power=1-p
plot(nu,power,type="l",lwd=2.5)

mu=180
m=n=100
sd=5
nu = seq(175,185,by=0.1)
size=length(nu)
p=numeric(size)
for (b in 1:size) {x=rnorm(m,mu,sd)
y=rnorm(n,nu[b],sd)
p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
power=1-p
lines(nu,power,type="l",col="red",lwd=2.5)

mu=180
m=n=30
sd=100
nu = seq(175,185,by=0.1)
size=length(nu)
p=numeric(size)
for (b in 1:size) {x=rnorm(m,mu,sd)
y=rnorm(n,nu[b],sd)
p[b]=t.test(x,y,var.equal=TRUE)[[3]]}
power=mean(p<0.05)
mean_thres=mean(p[p<0.05])
power=1-p
lines(nu,power,type="l",col="green",lwd=2.5)

```

The plots in the previous Figure represent the following cases:

Black: $\mu=180$; $m=n=30$; $sd=5$;

Red: $\mu=180$; $m=n=100$; $sd=5$;

Green: $\mu=180$; $m=n=30$; $sd=100$;

Comparing the black and red curves, the principal difference is the length of the sample, which is more than three times bigger in the second case. As the samples are randomly taken from an uniform distribution, the higher the amount of the sample, the better the power of the test in this case, as there will be more numbers randomly generated with central values around 180 and sd of 5. Furthermore, the power of the test is overall lower for the green case, which makes use of a higher standard deviation. Thus, the mean values of x and y are a lot more unlikely to match (which is the hypothesis H_0).

section extra

in this section we used fuctions and a different way to implement part 2 and 3

```
twoSampleT <- function(mu,nu,m,n,sd,pValue){
  B=1000
  p=numeric(B)
  c = 0
  for (b in 1:B) {
    x=rnorm(m,mu,sd)
    y=rnorm(n,nu,sd)
    p[b]=t.test(x,y,var.equal=TRUE)[[3]]
    if (p[b] < pValue){ c= c+1}
  }
  hist(p,main=c("with",pValue ))
}
```

```
twoSampleTNU <- function(mu,nuMin,nuMax,nuStep,m,n,sd,pValue){
  B=1000
  p=numeric(B)
  power = numeric()
  nu = seq(nuMin,nuMax,by=nuStep)
  for (i in 1:length(nu)){
    for (b in 1:B) {
      x=rnorm(m,mu,sd)
      y=rnorm(n,nu[i],sd)
      p[b]=t.test(x,y,var.equal=TRUE)[[3]]
    }
    power[i]=mean(p<pValue)
  }
  plot(nu,power,xlab = "nu",ylab = "power" ,type="l",lwd=2.5)
}
```

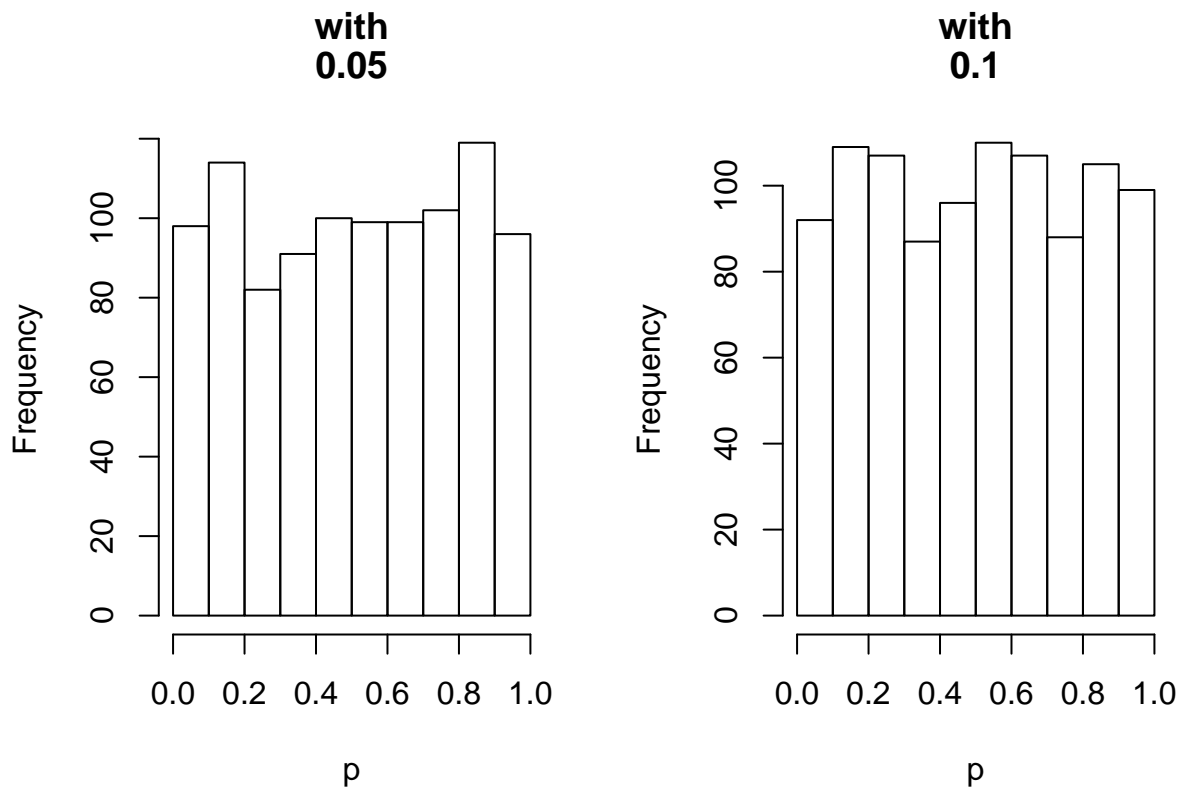
1. Set $\mu=\nu=180$, $m=n=30$ and $sd=10$. Repeat the script 1000 times, and record the 1000 p-values. How many p-values are smaller than 5%? How many are smaller than 10%? What is the distribution of the p-values (make a histogram)?

```
par(mfrow=c(1,2))
print(twoSampleT(180,180,30,30,10,0.05))
```

```
## $breaks
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
##
```

```
## $counts
## [1] 98 114 82 91 100 99 99 102 119 96
##
## $density
## [1] 0.98 1.14 0.82 0.91 1.00 0.99 0.99 1.02 1.19 0.96
##
## $mids
## [1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

```
print(twoSampleT(180,180,30,30,10,0.10))
```



```
## $breaks
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
##
## $counts
## [1] 92 109 107 87 96 110 107 88 105 99
##
```

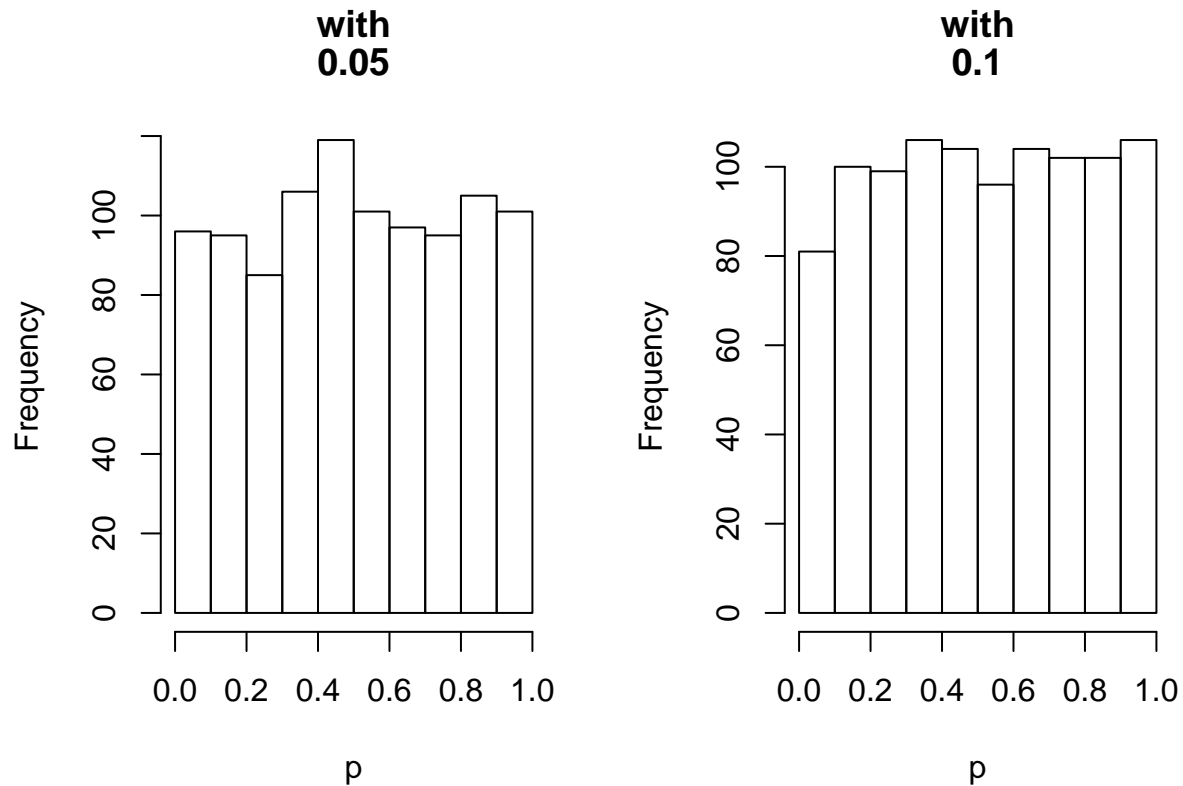
```
## $density
## [1] 0.92 1.09 1.07 0.87 0.96 1.10 1.07 0.88 1.05 0.99
##
## $mids
## [1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

2. Set $\mu=\nu=180$, $m=n=30$ and $sd=1$. Answer the same questions.

```
par(mfrow=c(1,2))
print(twoSampleT(180,180,30,30,1,0.05))
```

```
## $breaks
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
##
## $counts
## [1] 96 95 85 106 119 101 97 95 105 101
##
## $density
## [1] 0.96 0.95 0.85 1.06 1.19 1.01 0.97 0.95 1.05 1.01
##
## $mids
## [1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

```
print(twoSampleT(180,180,30,30,1,0.10))
```



```
## $breaks
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
##
## $counts
## [1] 81 100 99 106 104 96 104 102 102 106
##
## $density
## [1] 0.81 1.00 0.99 1.06 1.04 0.96 1.04 1.02 1.02 1.06
##
## $mids
## [1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

3. Set $\mu=180$, $\nu=175$, $m=n=30$ and $sd=6$. Answer the same questions.

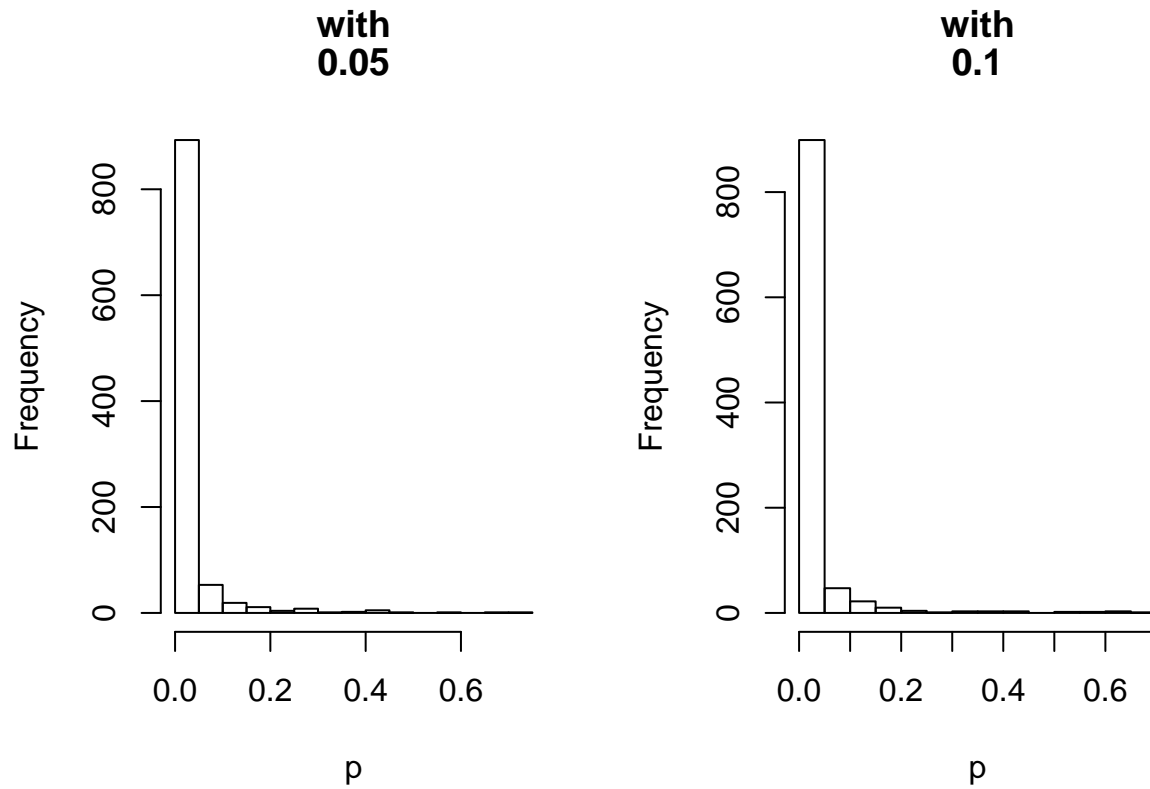
```

par(mfrow=c(1,2))
print(twoSampleT(180,175,30,30,6,0.05))

## $breaks
## [1] 0.00 0.05 0.10 0.15 0.20 0.25 0.30 0.35 0.40 0.45 0.50 0.55 0.60 0.65
## [15] 0.70 0.75
##
## $counts
## [1] 893 53 19 11 4 8 1 2 5 1 0 1 0 1 1
##
## $density
## [1] 17.86 1.06 0.38 0.22 0.08 0.16 0.02 0.04 0.10 0.02 0.00
## [12] 0.02 0.00 0.02 0.02
##
## $mids
## [1] 0.025 0.075 0.125 0.175 0.225 0.275 0.325 0.375 0.425 0.475 0.525
## [12] 0.575 0.625 0.675 0.725
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"

print(twoSampleT(180,175,30,30,6,0.10))

```



```
## $breaks
## [1] 0.00 0.05 0.10 0.15 0.20 0.25 0.30 0.35 0.40 0.45 0.50 0.55 0.60 0.65
## [15] 0.70
##
## $counts
## [1] 899 47 22 10 4 1 3 3 3 0 2 2 3 1
##
## $density
## [1] 17.98 0.94 0.44 0.20 0.08 0.02 0.06 0.06 0.06 0.00 0.04
## [12] 0.04 0.06 0.02
##
## $mids
## [1] 0.025 0.075 0.125 0.175 0.225 0.275 0.325 0.375 0.425 0.475 0.525
## [12] 0.575 0.625 0.675
##
## $xname
## [1] "p"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

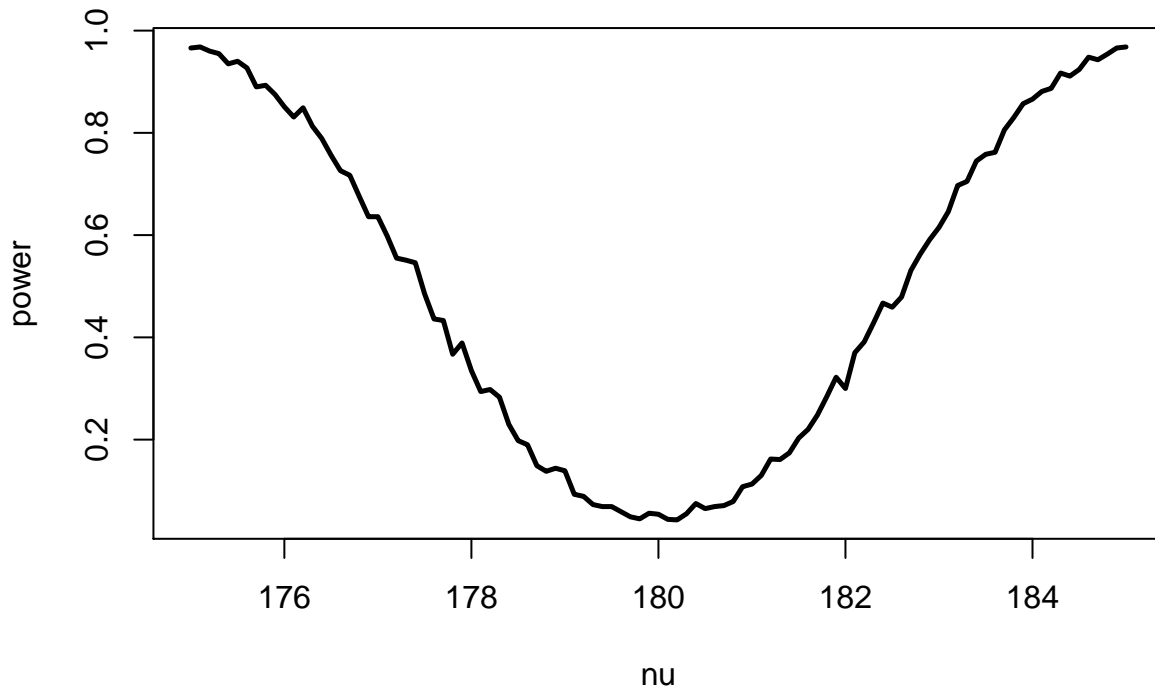
4. Explain the findings.

The test here is to check whether the null hypothesis that the population means of the two populations are equal. The two first histograms show the frequencies of the p-values for two randomly generated x and y arrays with center values of 180, size 30 and differing only by their standard deviations. For $sd=10$, the lowest p-values are seen more often than $sd=1$. Analogically, for $sd=1$, the highest p-values are seen more often than $sd=10$. This is explained by the fact that for smaller standard deviation values, the two arrays x and y have their values within a narrower interval, which increases the probability that the H_0 hypothesis is true.

EXERCISE 3

1. Set $\mu=180$, $m=n=30$ and $sd=5$. Calculate the power of the t-test for every value of ν in the grid `seq(175,185,by=0.1)`. Plot the power as a function of ν .

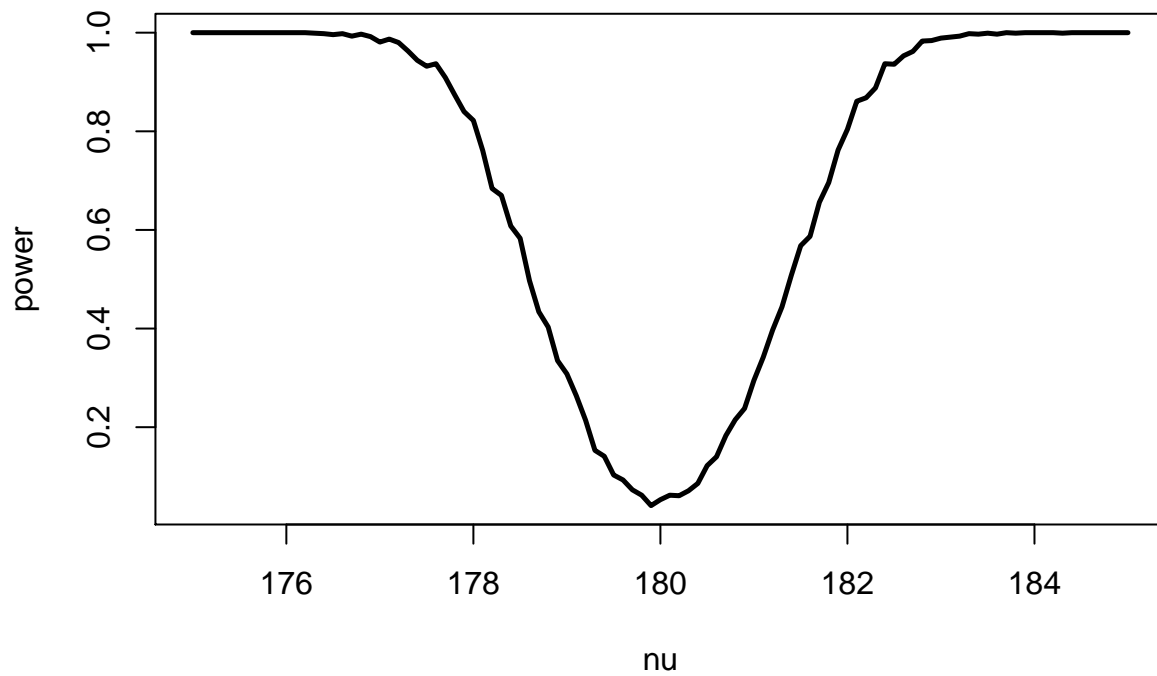
```
print(twoSampleTNU(180,175,185,0.1,30,30,5,0.05))
```



NULL

2. Set $\mu=180$, $m=n=100$ and $sd=5$. Repeat the preceding exercise. Add the plot to the preceding plot.

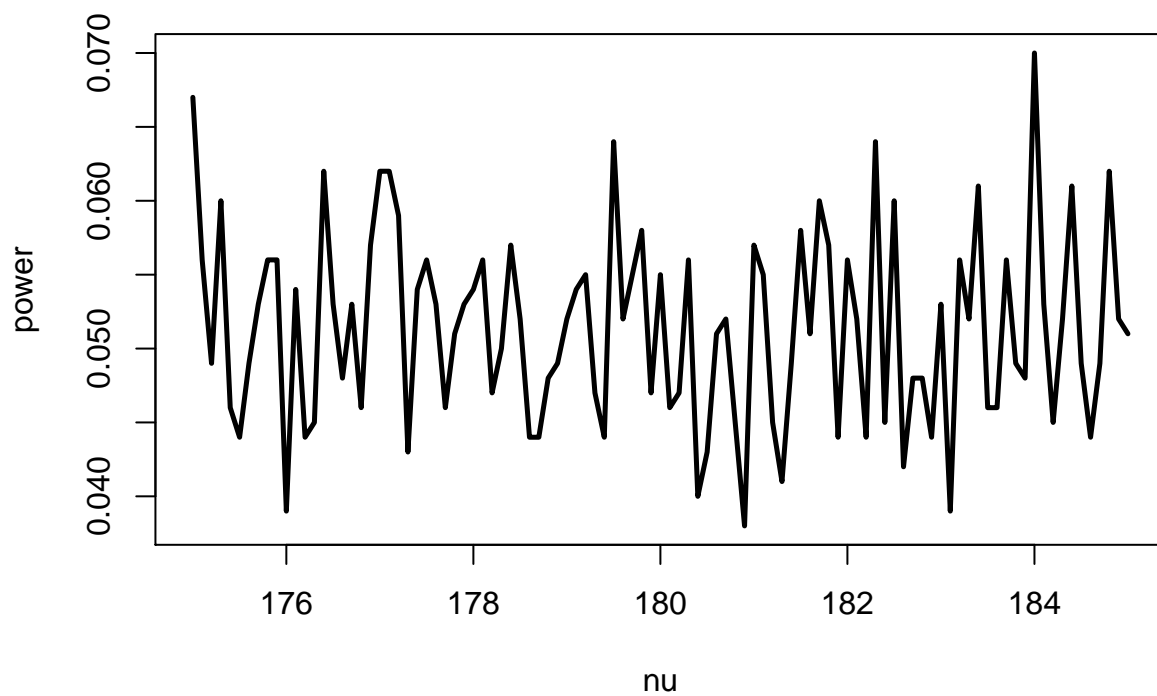
```
print(twoSampleTNU(180,175,185,0.1,100,100,5,0.05))
```

NULL

3. Set $\mu=180$, $m=n=30$ and $sd=100$. Repeat the preceding exercise.

```
print(twoSampleTNU(180,175,185,0.1,30,30,100,0.05))
```



NULL

4. Explain the findings.

The plots in the previous Figure represent the following cases:

Black: $\mu=180$; $m=n=30$; $sd=5$; Red: $\mu=180$; $m=n=100$; $sd=5$; Green: $\mu=180$; $m=n=30$; $sd=100$;

Comparing the black and red curves, the principal difference is the length of the sample, which is more than three times bigger in the second case. As the samples are randomly taken from an uniform distribution, the higher the amount of the sample, the better the power of the test in this case, as there will be more numbers randomly generated with central values around 180 and sd of 5. Furthermore, the power of the test is overall lower for the green case, which makes use of a higher standard deviation. Thus, the mean values of x and y are a lot more unlikely to match (which is the hypothesis H_0).