

Power study on real data (B=10000) TEMPORARY

03-05-2023

Digits dataset

The aim is to compare on Digits dataset (available at <http://odds.cs.stonybrook.edu/pendigits-dataset>) the performance of three closed testing procedures, which respectively use Simes local test with and without Storey estimator for the proportion of true null hypotheses and Wilcoxon-Mann-Whitney local test.

Digits dataset consists of 6870 observations, among which $n_{inliers} = 6714$ items are inliers and the remaining $n_{outliers} = 156$ are outliers. We will denote by n, l, m respectively the train set, the calibration set and the test set size. And reproducing the same setting as in [1], we have that $n + l = n_{inliers}/2$, $l = \min\{2000, n/2\}$ and $m = \min\{2000, n/3\}$. In the case of Digits dataset we obtain

$$n + l = 6714/2, \quad l = n/3, \quad m = n/3$$

from which $n = 2517.75$, $l = 839.25$, $m = 839.25$. Arbitrarily, we choose to set

$$n = 1258, \quad l = 2099, \quad m = 420$$

in order to have exact control of type I errors at the significance level $\alpha = m/(l + 1) = 0.2$.

R functions and libraries

```
library(nout)
library(R.matlab)
library(isotree)
library(tictoc)

sim_realdata = function(B, dataset, m1, m, n, l, in_index,
                        out_index=NULL, alpha=m/(l+1), lambda = 0.5){

  m0=m-m1
  if(m1!=0 & is.null(out_index)){
    stop("Error: arg out_index must be initialized.")
  }

  # if(m!=(m1+m0)){
  #   stop("Error: equation m=m1+m0 must be verified.")
  # }

  if(m1!=0){
    tr_ind = sample(in_index, size = n)
    tr = dataset[tr_ind,]
    iso.forest = isolation.forest(tr, ndim=ncol(dataset), ntrees=10, nthreads=1,
                                scoring_metric = "depth", output_score = TRUE)
    in_index2 = setdiff(in_index, tr_ind)
```

```

crit=critWMW(m=m, n=1, alpha=alpha)

d_WMW = rep(0,B)
d_Simes = rep(0,B)
d_StoSimes = rep(0,B)
d_BH = rep(0,B)
d_StoBH = rep(0,B)

for(b in 1:B){
  cal_ind = sample(in_index2, size = 1)
  in_index3 = setdiff(in_index2, cal_ind)
  tein_ind = sample(in_index3, size = m0)
  teout_ind = sample(out_index, size = m1)

  cal = dataset[cal_ind,]
  te = dataset[c(tein_ind, teout_ind),]

  S_cal = predict.isolation_forest(iso.fo$model, cal, type = "score")
  S_te = predict.isolation_forest(iso.fo$model, te, type = "score")

  d_WMW[b] = d_mannwhitney(S_X=S_cal, S_Y=S_te, crit=crit)
  d_Simes[b] = d_Simes(S_X=S_cal, S_Y=S_te, alpha=alpha)
  d_StoSimes[b] = d_StoreySimes(S_X=S_cal, S_Y=S_te, alpha=alpha)
  d_BH[b] = d_benjhoch(S_X=S_cal, S_Y=S_te, alpha=alpha)
  d_StoBH[b] = d_StoreyBH(S_X=S_cal, S_Y=S_te, alpha=alpha)
}
}

else{
  tr_ind = sample(in_index, size = n)
  tr = dataset[tr_ind,]
  iso.fo = isolation_forest(tr, ndim=ncol(dataset), ntrees=10, nthreads=1,
                           scoring_metric = "depth", output_score = TRUE)
  in_index2 = setdiff(in_index, tr_ind)

  crit=critWMW(m=m, n=1, alpha=alpha)

  d_WMW = rep(0,B)
  d_Simes = rep(0,B)
  d_StoSimes = rep(0,B)
  d_BH = rep(0,B)
  d_StoBH = rep(0,B)

  for(b in 1:B){
    cal_ind = sample(in_index2, size = 1)
    in_index3 = setdiff(in_index2, cal_ind)
    te_ind = sample(in_index3, size = m0)

    cal = dataset[cal_ind,]
    te = dataset[te_ind,]

    S_cal = predict.isolation_forest(iso.fo$model, cal, type = "score")
    S_te = predict.isolation_forest(iso.fo$model, te, type = "score")

```

```

    d_WMW[b] = d_mannwhitney(S_X=S_cal, S_Y=S_te, crit=crit)
    d_Simes[b] = d_Simes(S_X=S_cal, S_Y=S_te, alpha=alpha)
    d_StoSimes[b] = d_StoreySimes(S_X=S_cal, S_Y=S_te, alpha=alpha)
    d_BH[b] = d_benjhoch(S_X=S_cal, S_Y=S_te, alpha=alpha)
    d_StoBH[b] = d_StoreyBH(S_X=S_cal, S_Y=S_te, alpha=alpha)
  }
}

discov = as.data.frame(cbind("d_BH"=d_BH, "d_StoBH"=d_StoBH, "d_Simes"=d_Simes,
                             "d_StoSimes"=d_StoSimes, "d_WMW"=d_WMW))
colnames(discov) = c("BH", "BHSto", "CTSim", "CTSimSto", "CTWMW")
mean.discov = apply(discov, MARGIN = 2, FUN = mean)

powerGlobalNull = as.data.frame(cbind("d_BH"=d_BH>0, "d_StoBH"=d_StoBH>0, "d_Simes"=d_Simes>0,
                                       "d_StoSimes"=d_StoSimes>0, "d_WMW"=d_WMW>0))
colnames(powerGlobalNull) = c("BH", "BHSto", "CTSim", "CTSimSto", "CTWMW")
mean.powerGlobalNull = apply(powerGlobalNull, MARGIN = 2, FUN = mean)

return(list("discoveries"=discov, "mean.discoveries" = mean.discov,
           "powerGlobalNull"=powerGlobalNull, "mean.powerGlobalNull"=mean.powerGlobalNull,
           "m1"=m1, "alpha"=alpha))
}

```

Load the data and set the parameters as described above.

```

set.seed(321)

# Initializing parameters
n = 2518
l = 2099
m = 420
alpha = m/(l+1)

data = readMat("G:\\Il mio Drive\\PHD\\Progetto di ricerca\\Conformal Inference Project\\Simulazioni\\7
dataset = cbind(data$X, data$y); colnames(dataset)[ncol(dataset)] = "y"
in_ind = which(dataset[,ncol(dataset)]==0)
out_ind = which(dataset[,ncol(dataset)]==1)

```

We fixed the train set on which we train the isolation forest algorithm and we generate $B = 1000$ calibration and test sets. For each $b = 1, \dots, B$ we compute the number of discoveries obtained by Benjamini-Hochberg procedure with and without Storey's estimator for the proportion of true null hypotheses, by closed testing using Simes local test with and without Storey estimator and by closed testing using Wilcoxon-Mann-Whitney local test.

All inliers

We now set the proportion of inliers equal to 1, so that the number of outliers $m_1 = 0$.

```

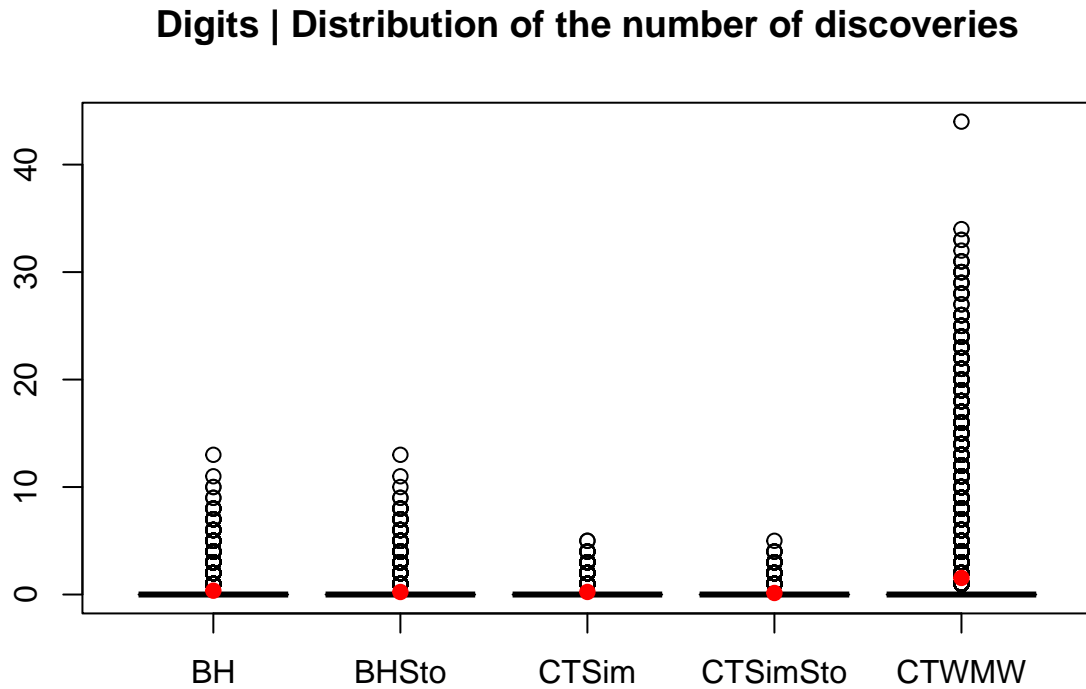
B=10^4
m1=0

tic()
res = sim_realddata(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                   alpha=alpha, l=l, n=n, m=m, m1=m1)
toc()

```

```
## 774.57 sec elapsed
```

```
boxplot(res$discoveries, main="Digits | Distribution of the number of discoveries")  
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")
```

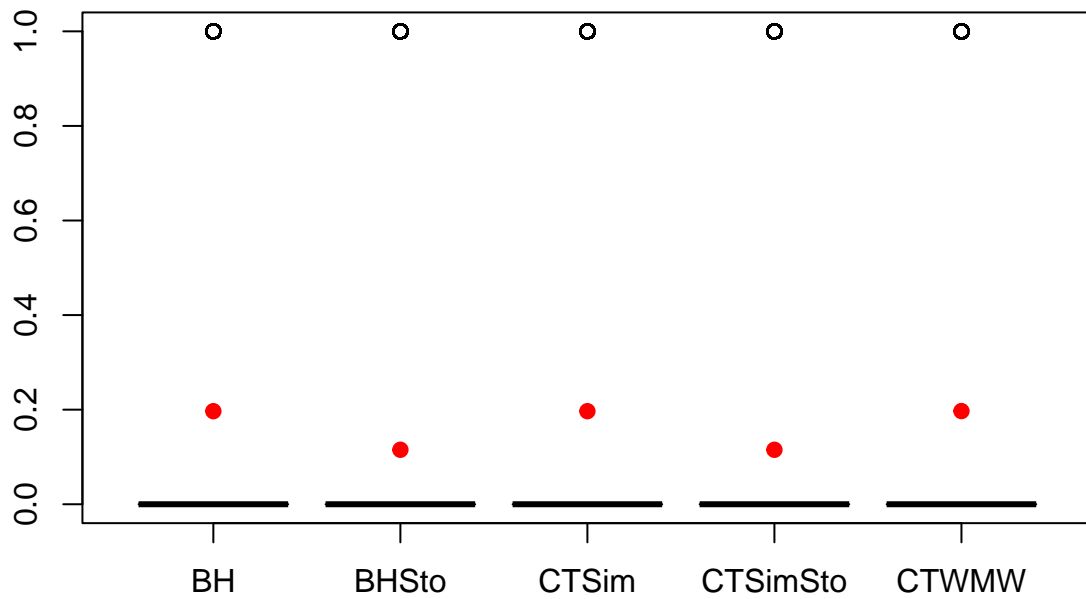


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 0.3562 0.2341 0.2396 0.1420 1.5452
```

```
boxplot(res$powerGlobalNull, main="Digits | Distribution of the power")  
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

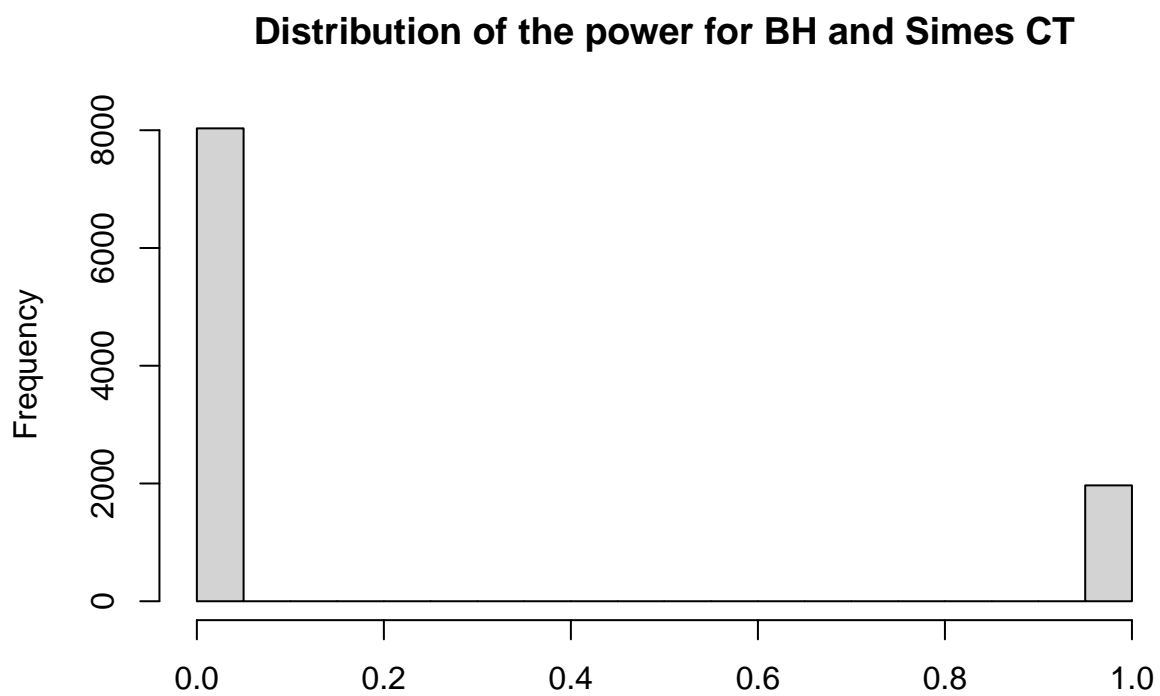
Digits | Distribution of the power



```
res$mean.powerGlobalNull
```

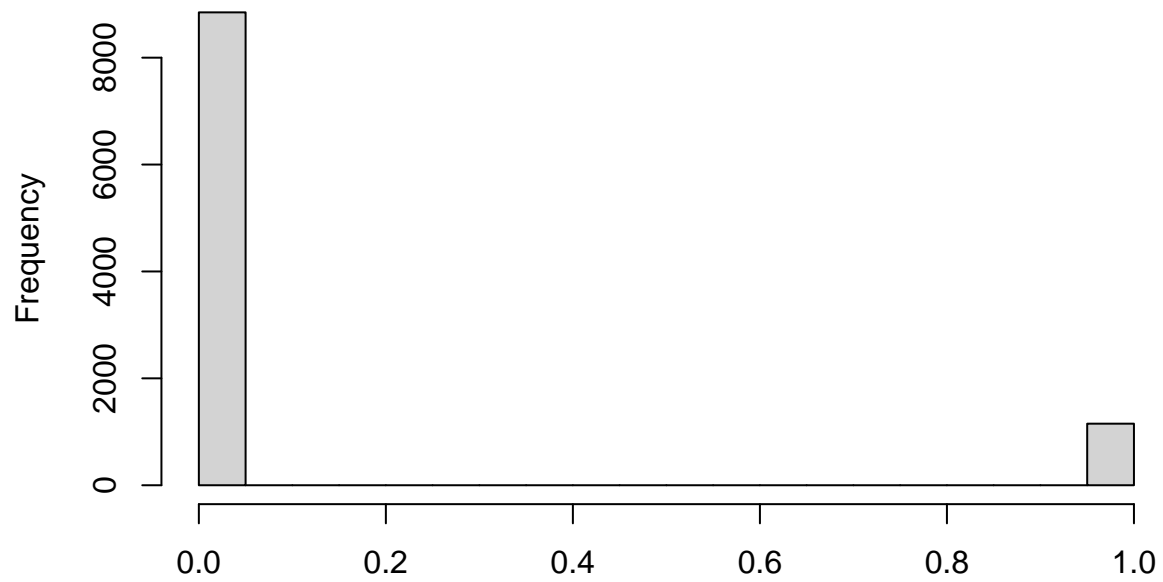
```
##      BH      BHSto      CTSim CTSimSto      CTWMW
## 0.1968 0.1152 0.1968 0.1152 0.1970
```

```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",
     main="Distribution of the power for BH and Simes CT")
```



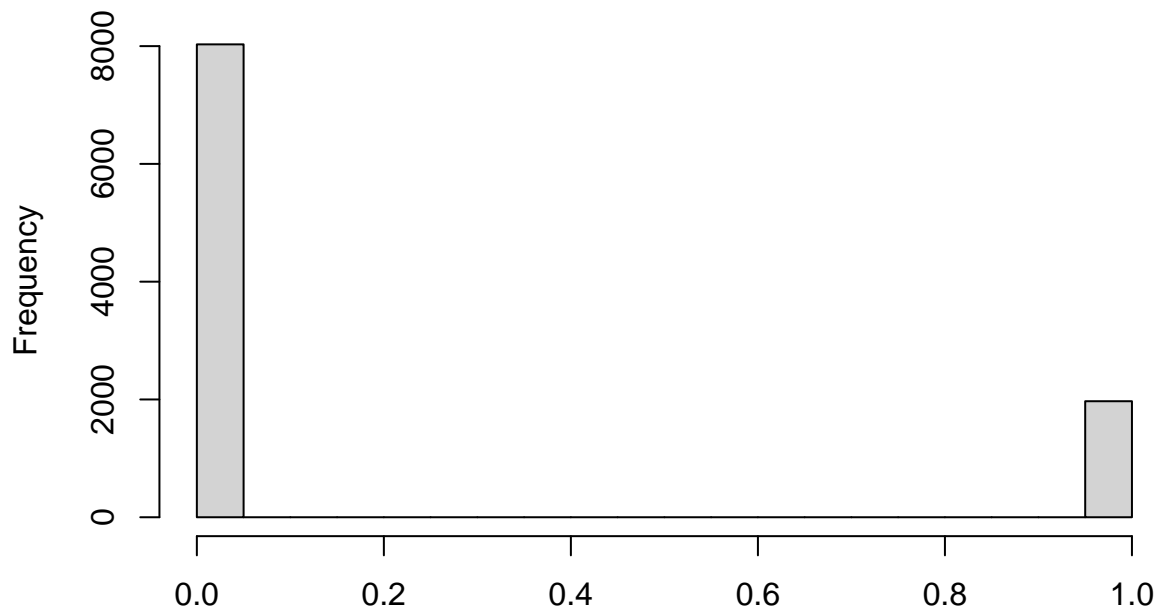
```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```

Distribution of the power for BH and Simes CT with Storey



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```

Distribution of the power for Wilcoxon–Mann–Whitney CT



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

10% outliers

We now set the proportion of inliers equal to 0.9. Referring to Digits dataset we have that the number of inliers is $m_0 = 378$ and the number of outliers is $m_1 = 42$.

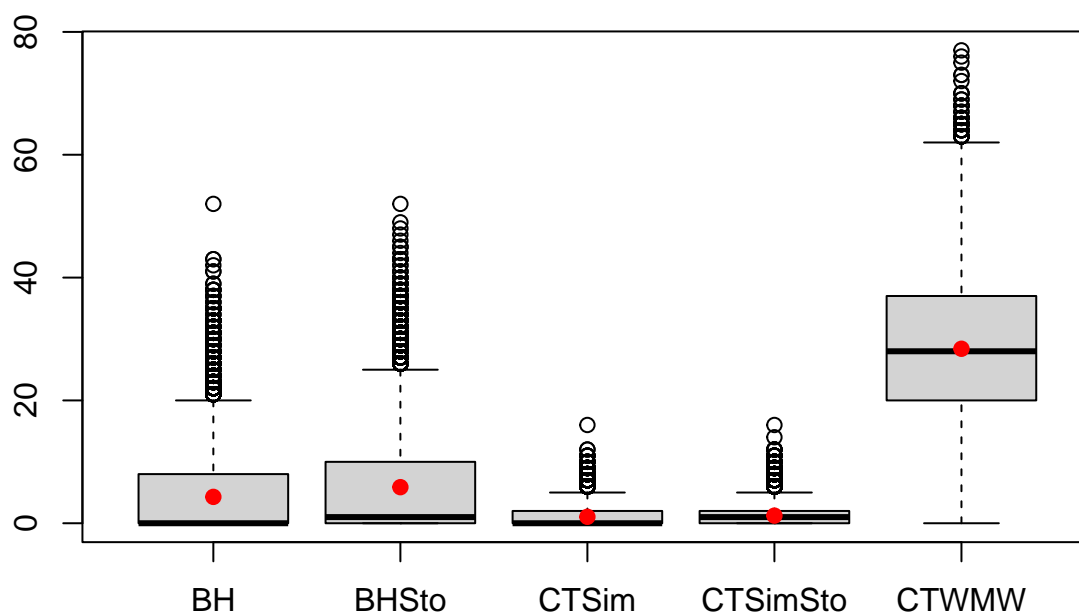
```
B=10^4
m1=round(0.1*m)

tic()
res = sim_realdData(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                    alpha=alpha,l=1, n=n, m=m, m1=m1)
toc()
```

```
## 1266.05 sec elapsed
```

```
boxplot(res$discoveries, main="Digits | Distribution of the number of discoveries")
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")
```


Digits | Distribution of the number of discoveries

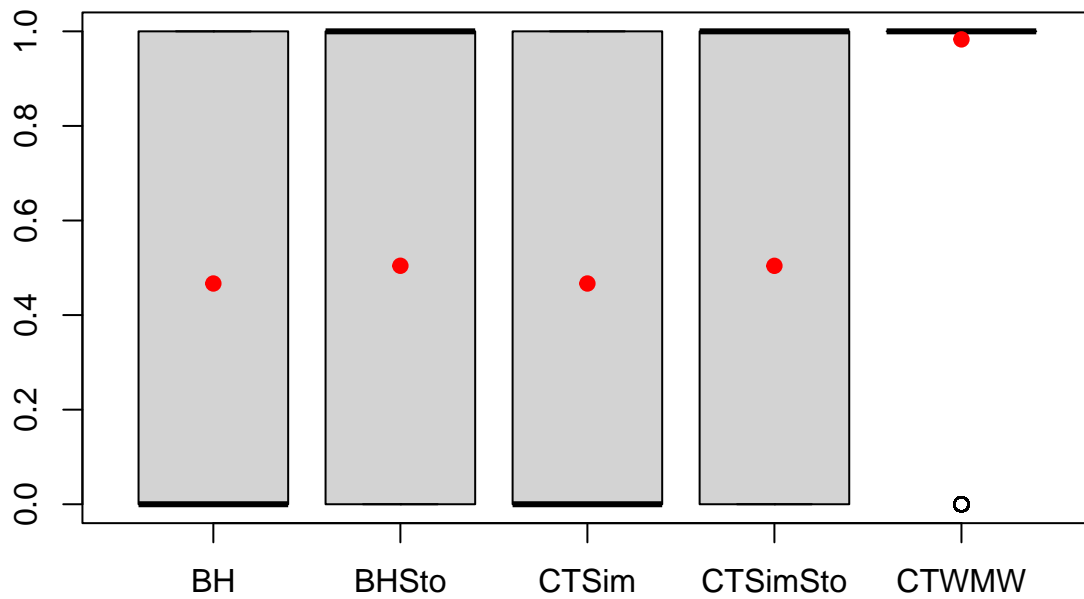


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW
##  4.2935  5.8736  1.0500  1.2383  28.4074
```

```
boxplot(res$powerGlobalNull, main="Digits | Distribution of the power")
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

Digits | Distribution of the power

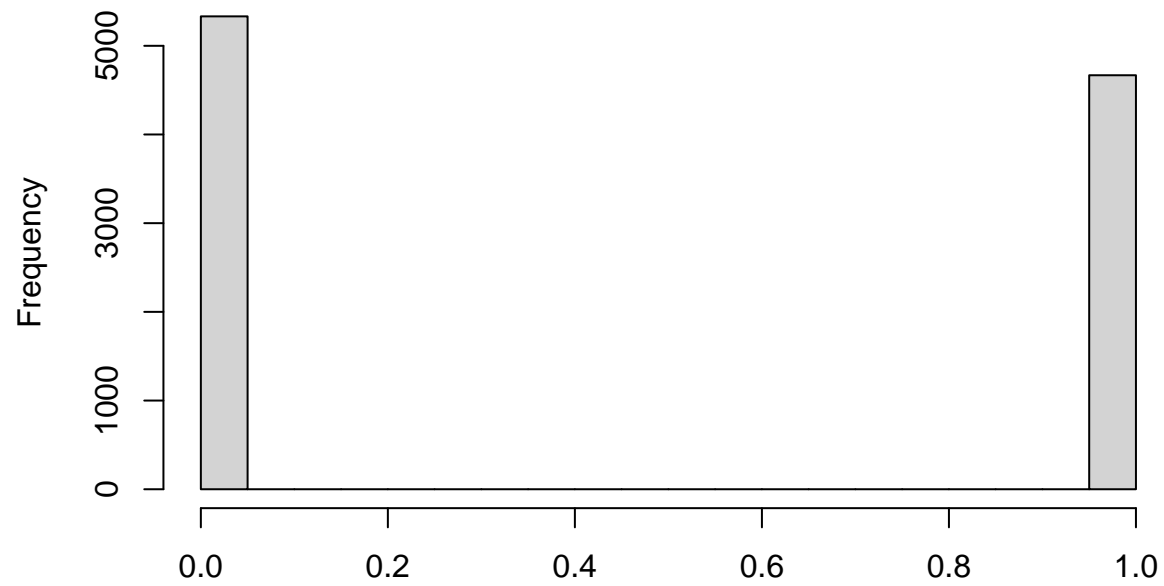


```
res$mean.powerGlobalNull
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW
## 0.4668 0.5043 0.4668 0.5042 0.9830
```

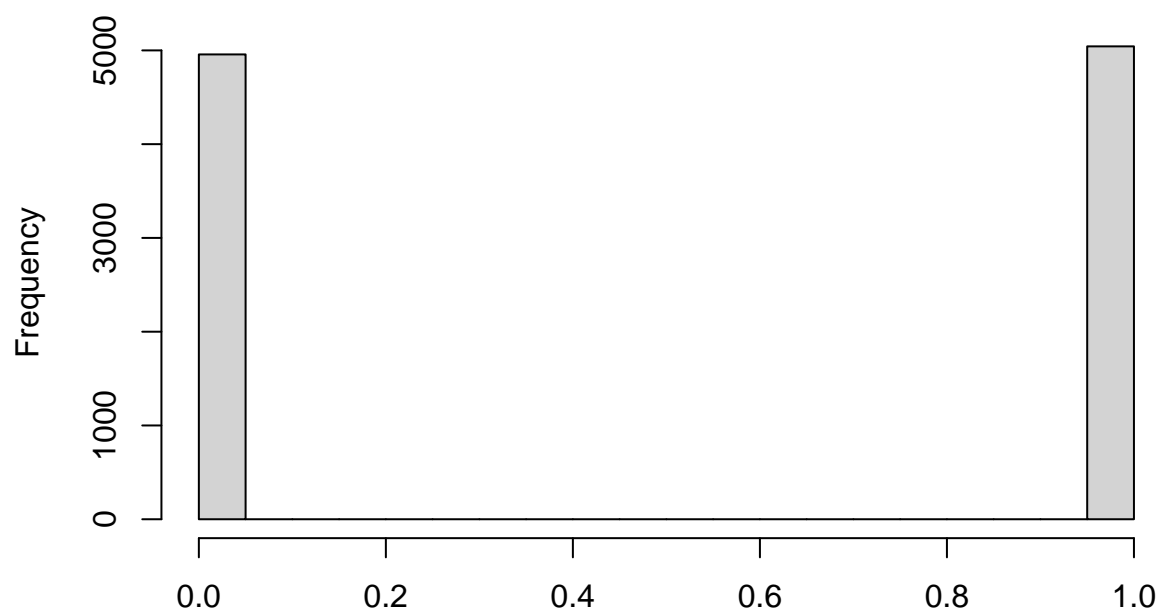
```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",
     main="Distribution of the power for BH and Simes CT")
```

Distribution of the power for BH and Simes CT

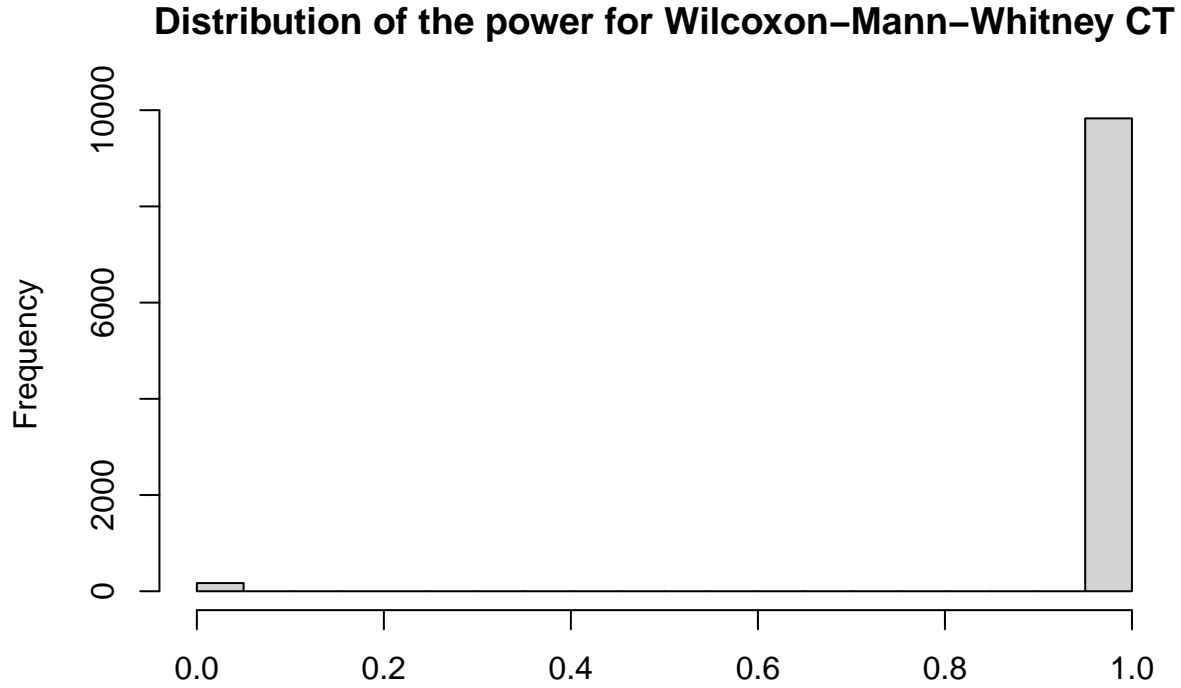


```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```

Distribution of the power for BH and Simes CT with Storey



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

Credit Card Fraud Detection dataset

The aim is to compare on Digits dataset (available at <https://www.kaggle.com/mlg-ulb/creditcardfraud>) the performance of three closed testing procedures, which respectively use Simes local test with and without Storey estimator for the proportion of true null hypotheses and Wilcoxon–Mann–Whitney local test.

Credit card dataset consists of 284807 observations, among which $n_{inliers} = 284315$ items are inliers and the remaining $n_{outliers} = 492$ are outliers. We will denote by n, l, m respectively the train set, the calibration set and the test set size. And reproducing the same setting as in [1], we have that $n + l = n_{inliers}/2$, $l = \min\{2000, n/2\}$ and $m = \min\{2000, n/3\}$. In the case of Digits dataset we obtain

$$n + l = 284315/2, \quad l = 2000, \quad m = 2000$$

from which $n = 140157.5$, $l = 2000$, $m = 2000$. Arbitrarily, we choose to set

$$n = 132158, \quad l = 9999, \quad m = 2000$$

in order to have exact control of type I errors at the significance level $\alpha = m/(l + 1) = 0.2$.

Load the data and set the parameters as described above.

```
set.seed(321)

# Initializing parameters
n = 132158
```

```

l = 9999
m = 2000
alpha = m/(l+1)

dataset = read.csv("G:\\Il mio Drive\\PHD\\Progetto di ricerca\\Conformal Inference Project\\Simulazioni\\dataset\\CreditCard.csv")
out_ind = which(dataset$Class==1)
in_ind = which(dataset$Class==0)

```

We fixed the train set on which we train the isolation forest algorithm and we generate $B = 1000$ calibration and test sets. For each $b = 1, \dots, B$ we compute the number of discoveries obtained by Benjamini-Hochberg procedure with and without Storey's estimator for the proportion of true null hypotheses, by closed testing using Simes local test with and without Storey estimator and by closed testing using Wilcoxon-Mann-Whitney local test.

All inliers

We now set the proportion of inliers equal to 1, so that the number of outliers $m_1 = 0$.

```

B=10^4
m1=0

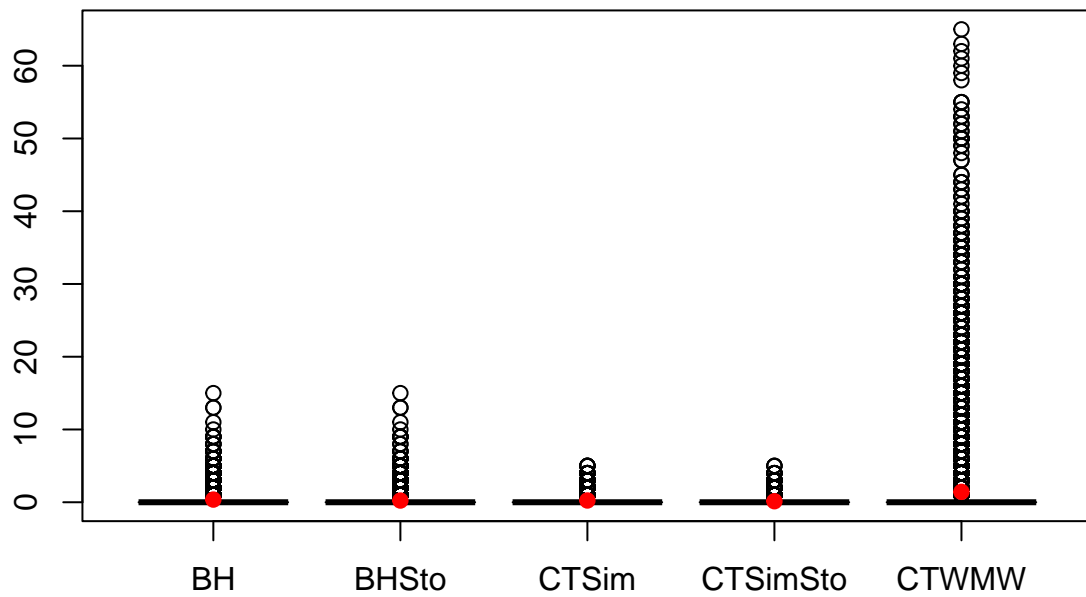
tic()
res = sim_realddata(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                    alpha=alpha, l=l, n=n, m=m, m1=m1)
toc()

## 11554.78 sec elapsed

boxplot(res$discoveries, main="CreditCard | Distribution of the number of discoveries")
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")

```

CreditCard | Distribution of the number of discoveries

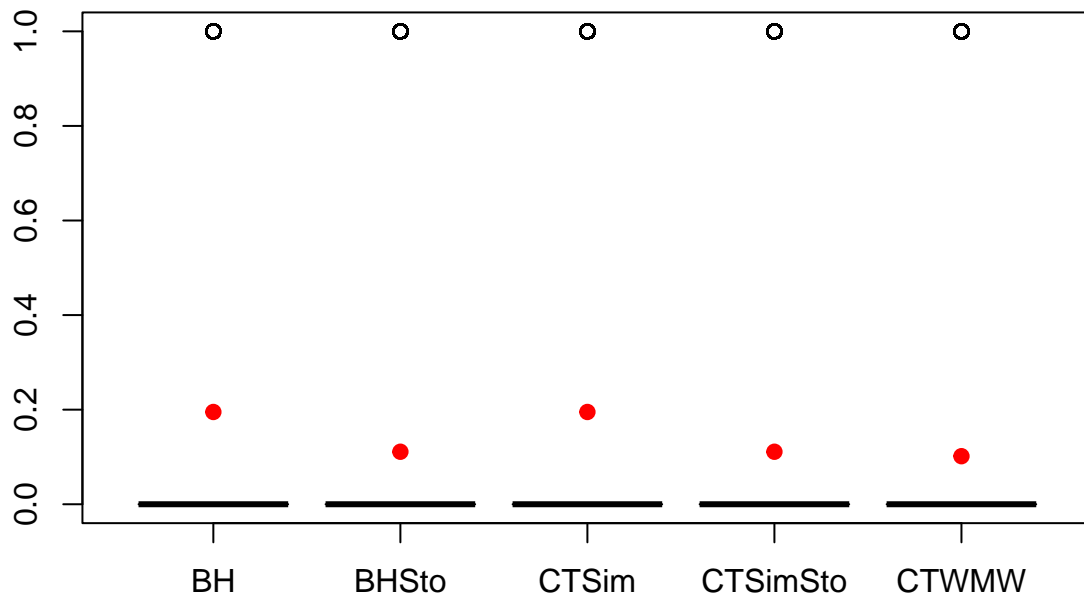


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 0.3654 0.2329 0.2414 0.1372 1.4266
```

```
boxplot(res$powerGlobalNull, main="CreditCard | Distribution of the power")  
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

CreditCard | Distribution of the power

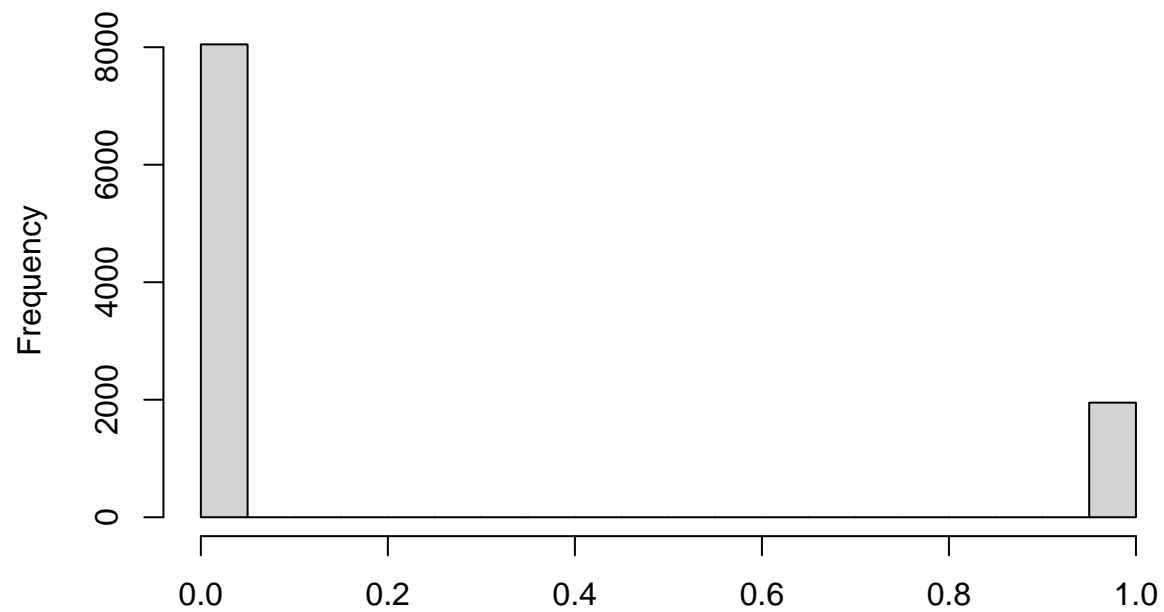


```
res$mean.powerGlobalNull
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 0.1951 0.1110 0.1951 0.1110 0.1016
```

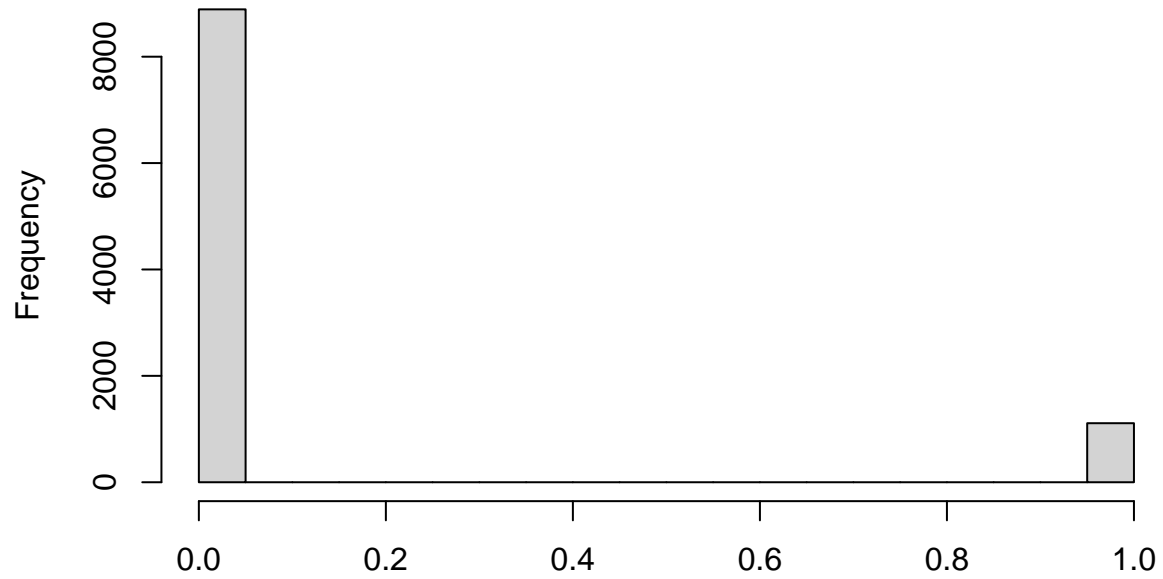
```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",  
     main="Distribution of the power for BH and Simes CT")
```


Distribution of the power for BH and Simes CT



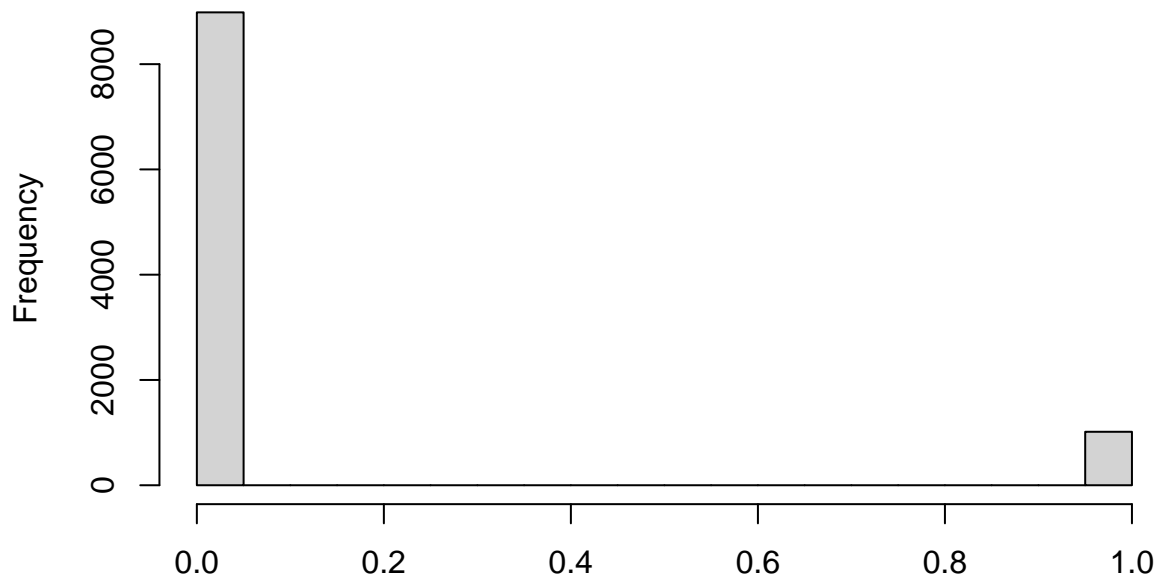
```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```

Distribution of the power for BH and Simes CT with Storey



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```

Distribution of the power for Wilcoxon–Mann–Whitney CT



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

10% outliers

We now set the proportion of inliers equal to 0.9. Referring to Digits dataset we have that the number of inliers is $m_0 = 1800$ and the number of outliers is $m_1 = 200$.

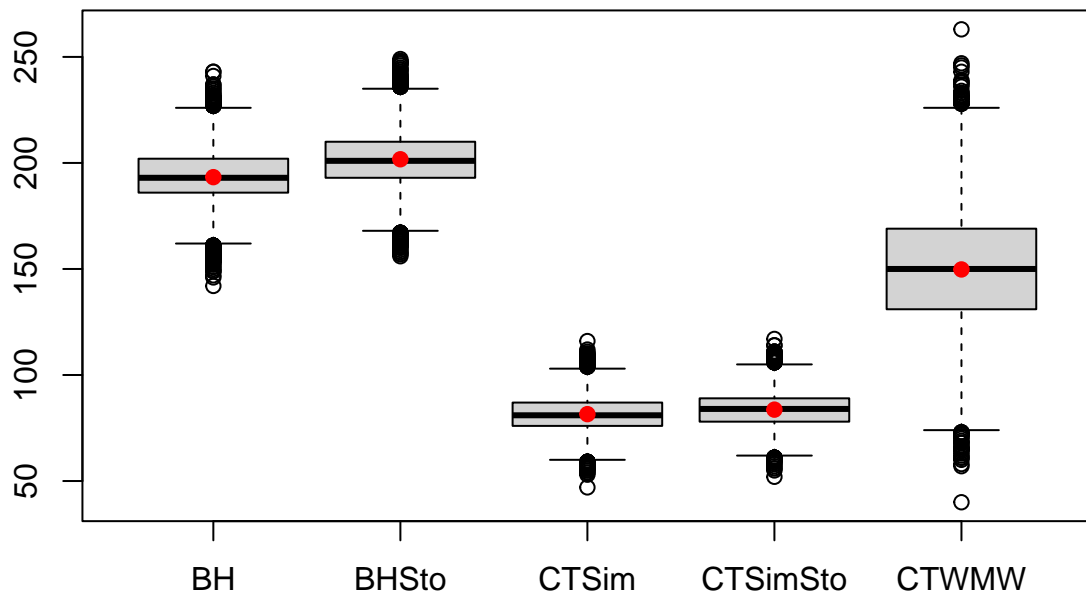
```
B=10^4
m1=round(0.1*m)

tic()
res = sim_realdData(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                   alpha=alpha,l=1, n=n, m=m, m1=m1)
toc()

## 9240.7 sec elapsed

boxplot(res$discoveries, main="CreditCard | Distribution of the number of discoveries")
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")
```

CreditCard | Distribution of the number of discoveries

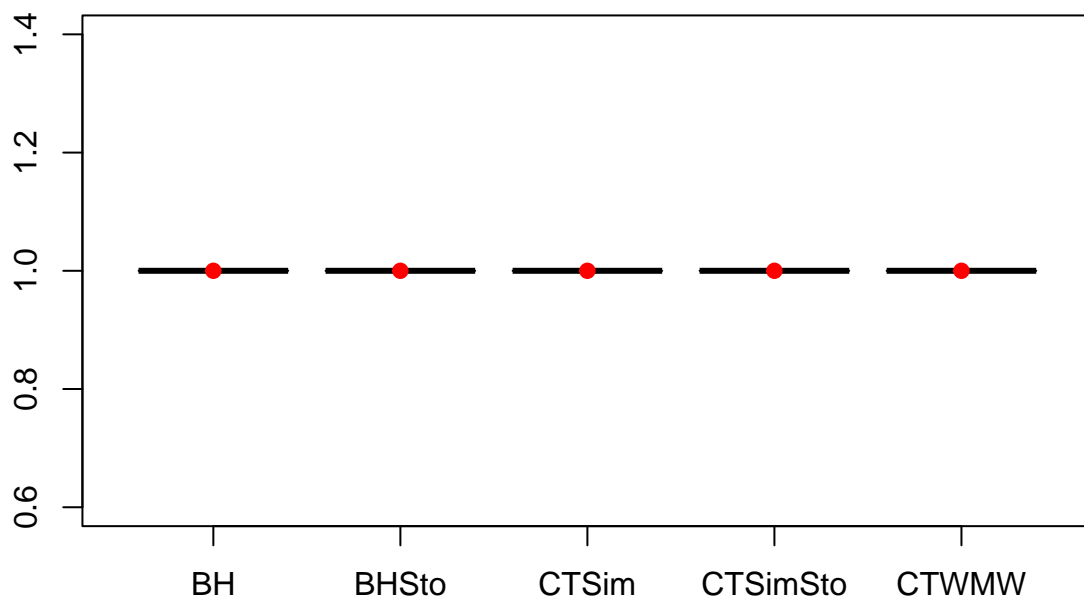


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 193.3018 201.7384  81.5388  83.6672 149.7832
```

```
boxplot(res$powerGlobalNull, main="CreditCard | Distribution of the power")  
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

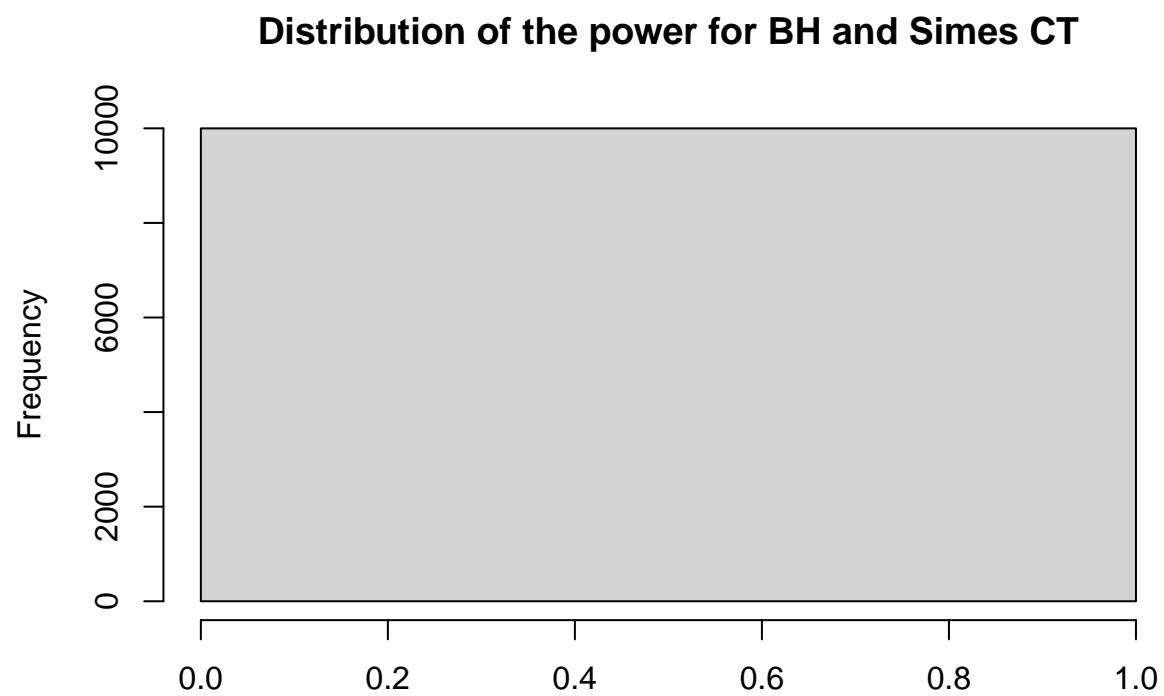
CreditCard | Distribution of the power



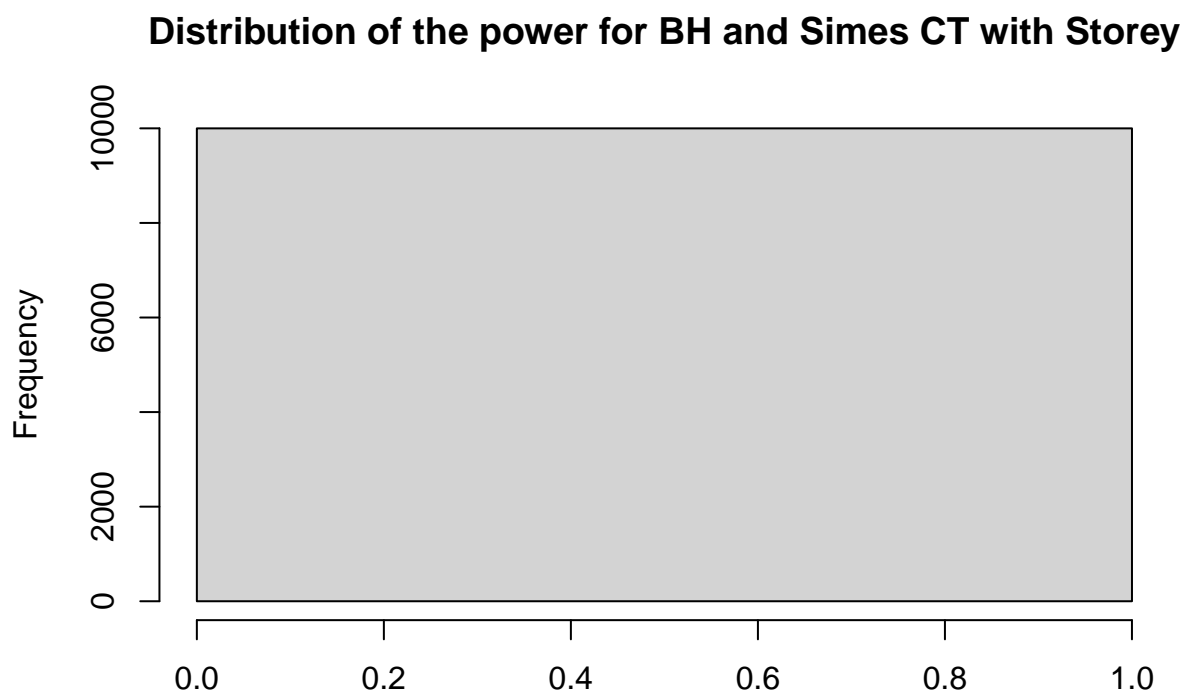
```
res$mean.powerGlobalNull
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
##      1        1        1        1        1
```

```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",  
     main="Distribution of the power for BH and Simes CT")
```

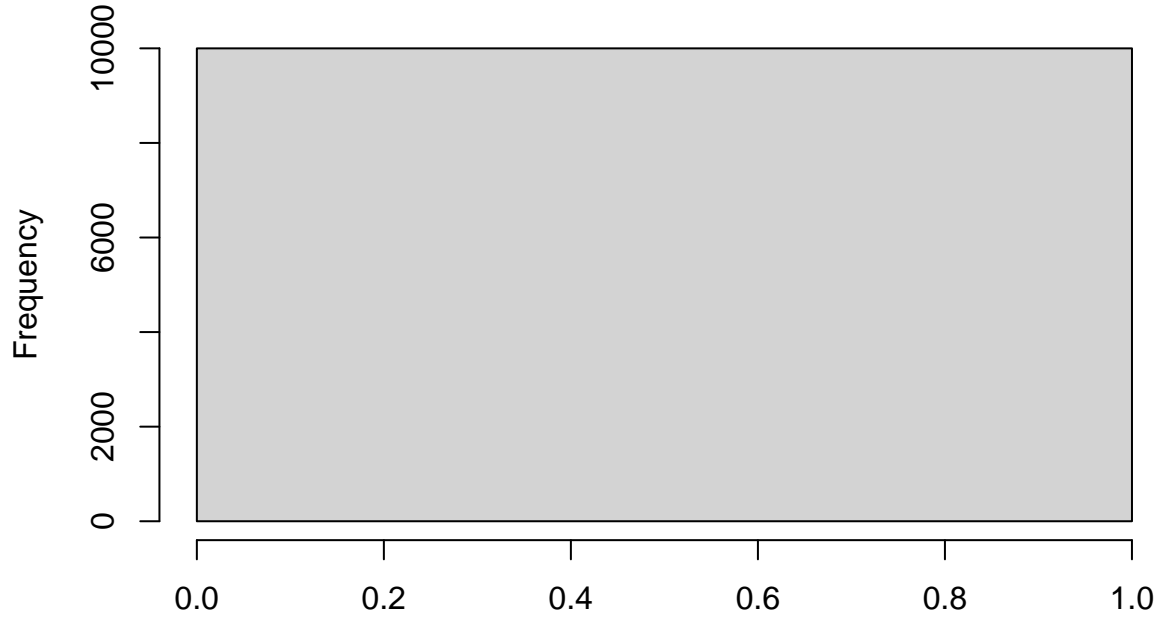


```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```

Distribution of the power for Wilcoxon–Mann–Whitney CT



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

Statlog (Shuttle) dataset

The aim is to compare on Digits dataset (available at <http://odds.cs.stonybrook.edu/shuttle-dataset>) the performance of three closed testing procedures, which respectively use Simes local test with and without Storey estimator for the proportion of true null hypotheses and Wilcoxon–Mann–Whitney local test.

Credit card dataset consists of 49097 observations, among which $n_{inliers} = 45586$ items are inliers and the remaining $n_{outliers} = 3511$ are outliers. We will denote by n, l, m respectively the train set, the calibration set and the test set size. And reproducing the same setting as in [1], we have that $n + l = n_{inliers}/2$, $l = \min\{2000, n/2\}$ and $m = \min\{2000, n/3\}$. In the case of Digits dataset we obtain

$$n + l = 45586/2, \quad l = 2000, \quad m = 2000$$

from which $n = 20793$, $l = 2000$, $m = 2000$. Arbitrarily, we choose to set

$$n = 12794, \quad l = 9999, \quad m = 2000$$

in order to have exact control of type I errors at the significance level $\alpha = m/(l + 1) = 0.2$.

Load the data and set the parameters as described above.

```
set.seed(321)

# Initializing parameters
n = 12794
```



```

l = 9999
m = 2000
alpha = m/(l+1)

data = readMat("G:\\Il mio Drive\\PHD\\Progetto di ricerca\\Conformal Inference Project\\Simulazioni\\7
dataset = cbind(data$X, data$y); colnames(dataset)[ncol(dataset)] = "y"
out_ind = which(dataset[,ncol(dataset)]==1)
in_ind = which(dataset[,ncol(dataset)]==0)

```

We fixed the train set on which we train the isolation forest algorithm and we generate $B = 1000$ calibration and test sets. For each $b = 1, \dots, B$ we compute the number of discoveries obtained by Benjamini-Hochberg procedure with and without Storey's estimator for the proportion of true null hypotheses, by closed testing using Simes local test with and without Storey estimator and by closed testing using Wilcoxon-Mann-Whitney local test.

All inliers

We now set the proportion of inliers equal to 1, so that the number of outliers $m_1 = 0$.

```

B=10^4
m1=0

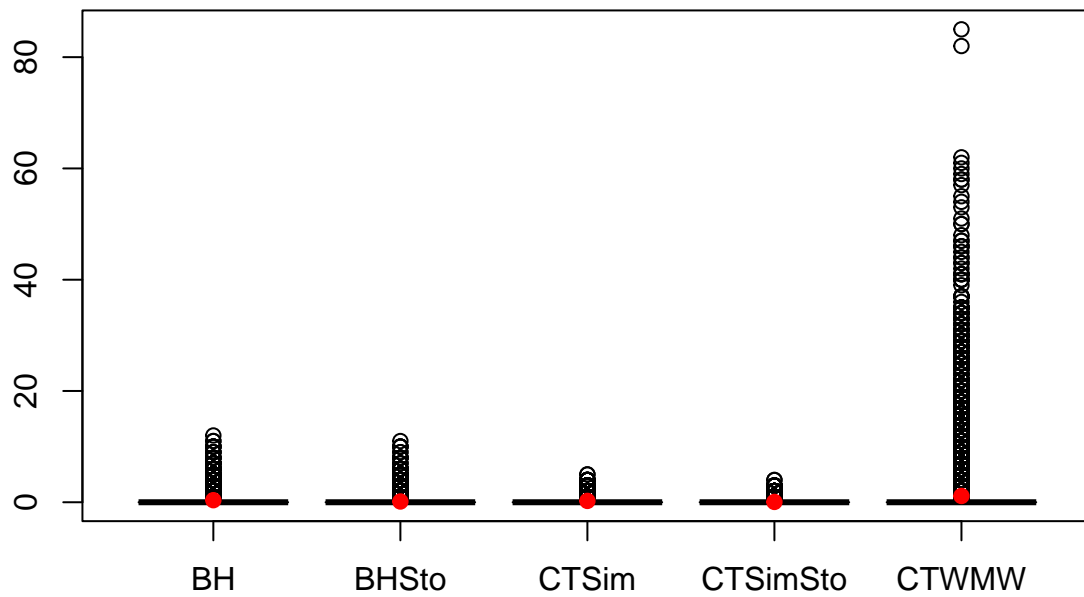
tic()
res = sim_realddata(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                    alpha=alpha, l=l, n=n, m=m, m1=m1)
toc()

## 6105.39 sec elapsed

boxplot(res$discoveries, main="Shuttle | Distribution of the number of discoveries")
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")

```

Shuttle | Distribution of the number of discoveries

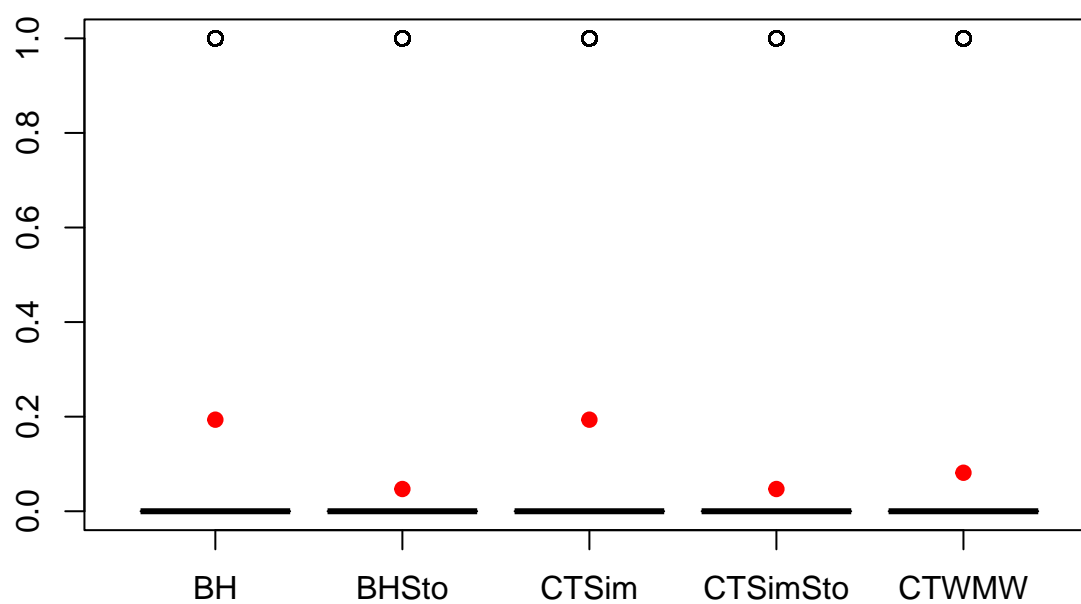


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 0.3656 0.1325 0.2393 0.0570 1.1259
```

```
boxplot(res$powerGlobalNull, main="Shuttle | Distribution of the power")  
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

Shuttle | Distribution of the power

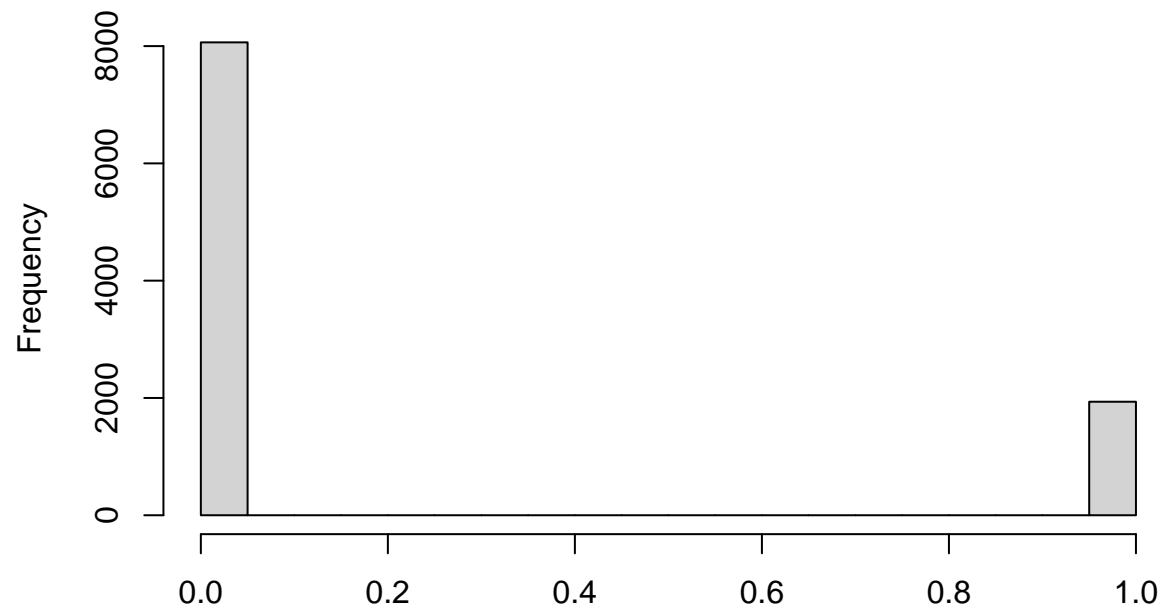


```
res$mean.powerGlobalNull
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW
## 0.1936 0.0470 0.1936 0.0470 0.0814
```

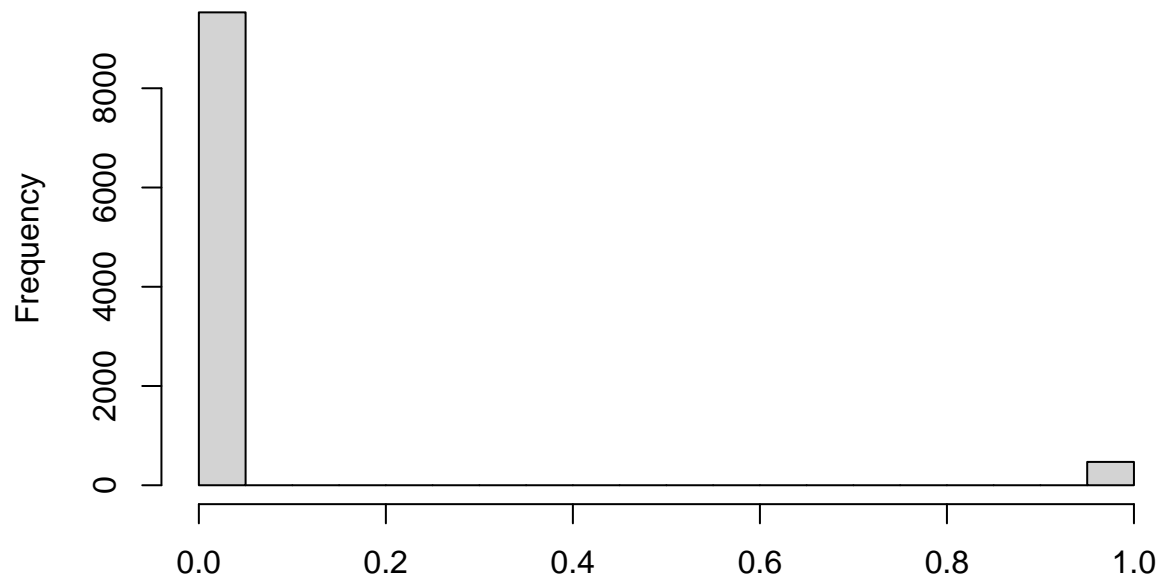
```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",
     main="Distribution of the power for BH and Simes CT")
```

Distribution of the power for BH and Simes CT



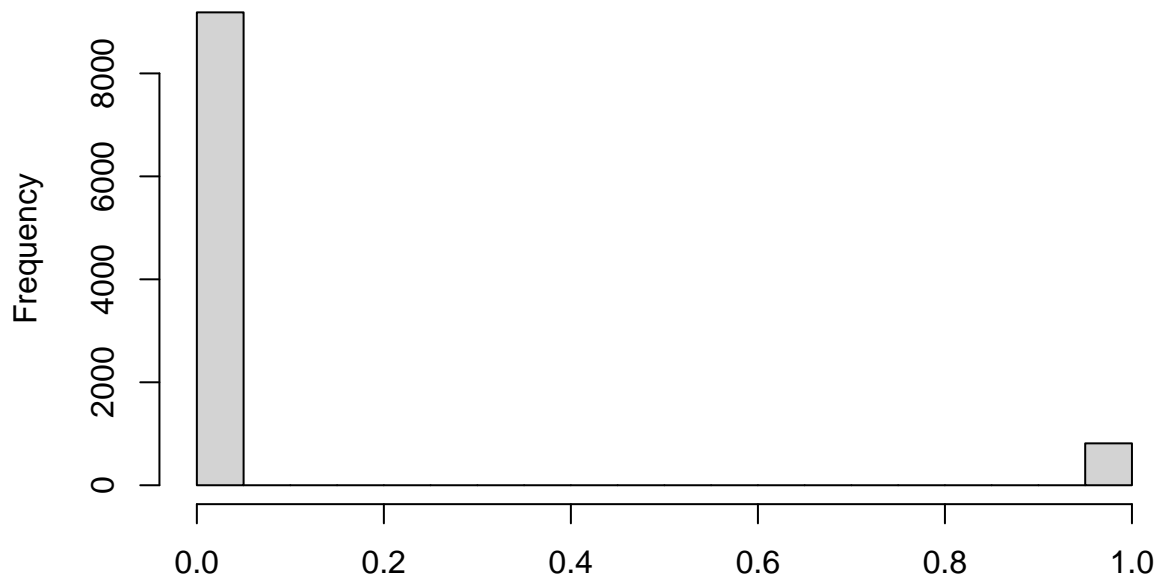
```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```

Distribution of the power for BH and Simes CT with Storey



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```

Distribution of the power for Wilcoxon–Mann–Whitney CT



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

10% outliers

We now set the proportion of inliers equal to 0.9. Referring to Digits dataset we have that the number of inliers is $m_0 = 1800$ and the number of outliers is $m_1 = 200$.

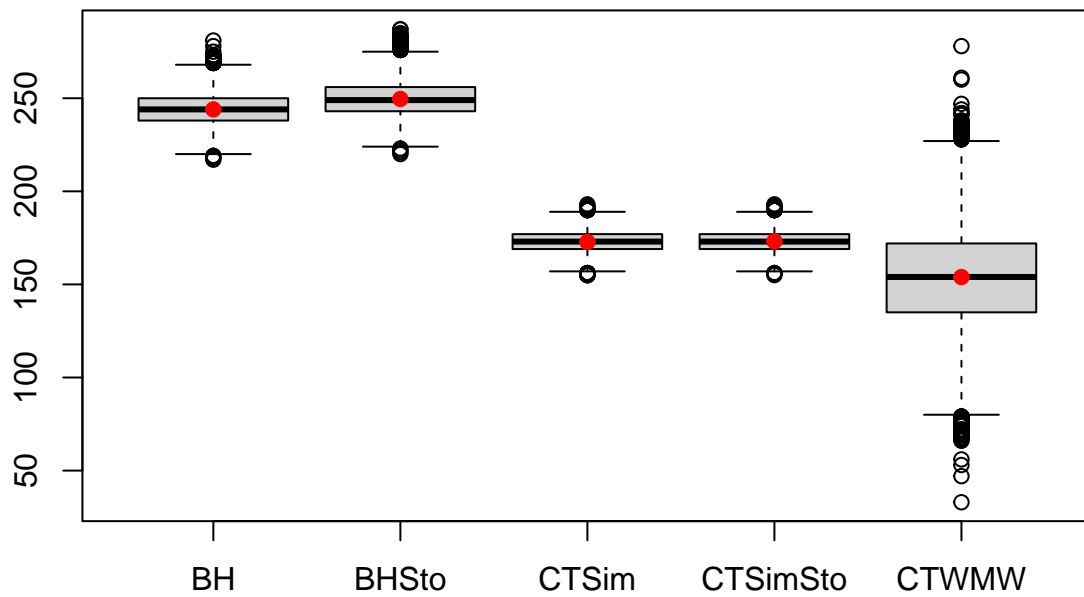
```
B=10^4
m1=round(0.1*m)

tic()
res = sim_realdData(B=B, in_index=in_ind, out_index=out_ind, dataset=dataset,
                    alpha=alpha,l=1, n=n, m=m, m1=m1)
toc()

## 6074.17 sec elapsed

boxplot(res$discoveries, main="Shuttle | Distribution of the number of discoveries")
points(x=1:5, y=res$mean.discoveries, pch=19, col="red")
```

Shuttle | Distribution of the number of discoveries

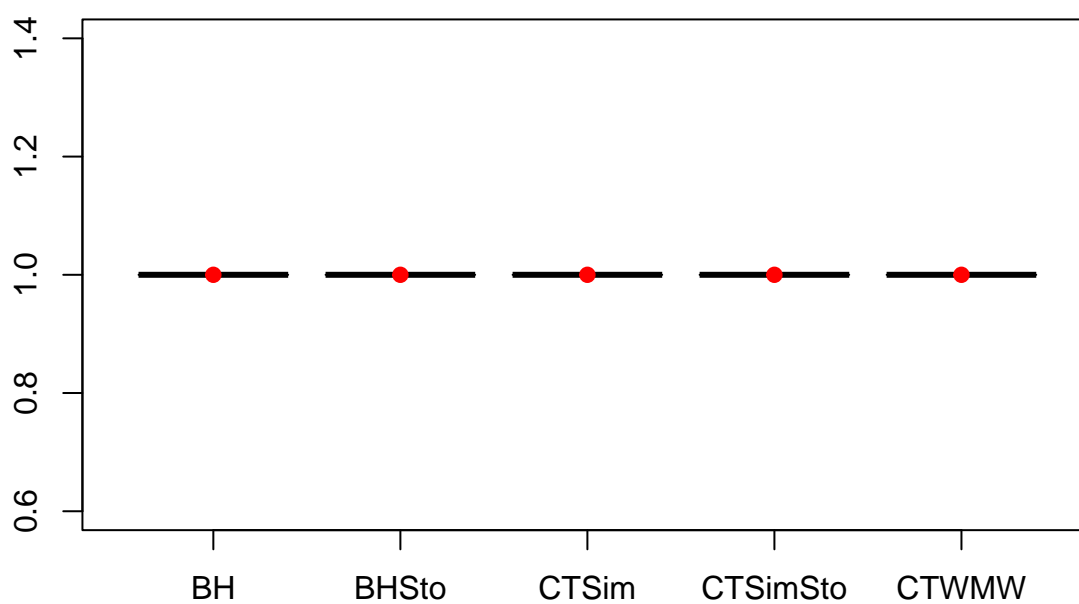


```
res$mean.discoveries
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
## 244.0126 249.6119 172.7924 173.1773 153.9530
```

```
boxplot(res$powerGlobalNull, main="Shuttle | Distribution of the power")  
points(x=1:5, y=res$mean.powerGlobalNull, pch=19, col="red")
```

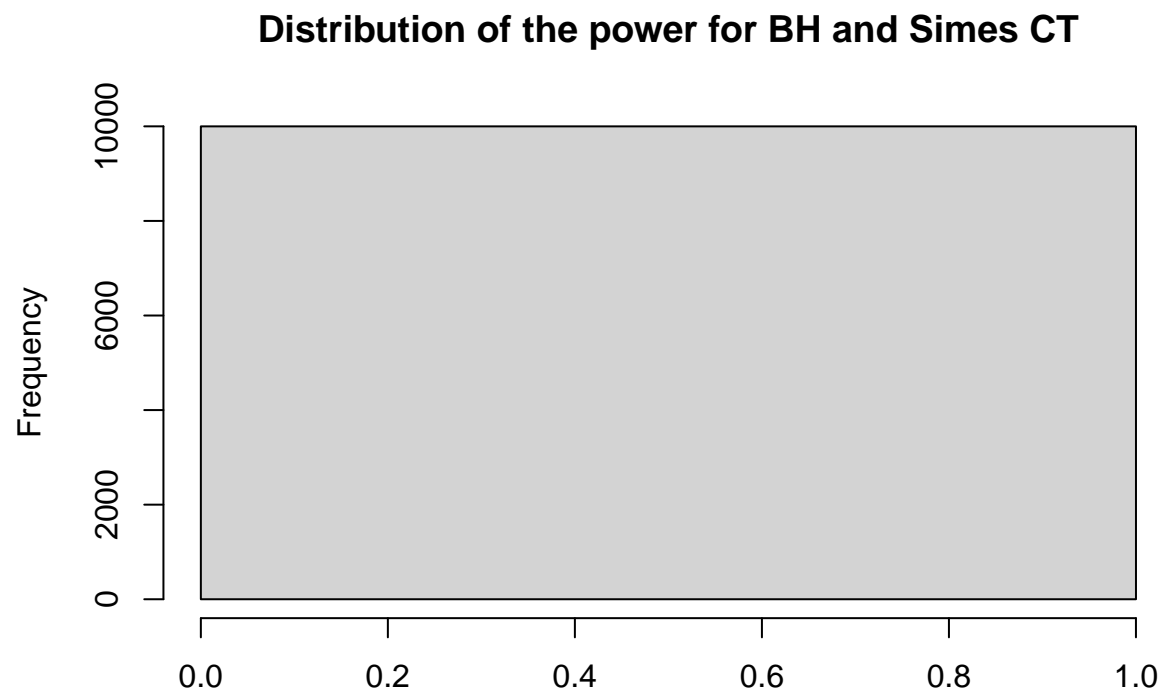
Shuttle | Distribution of the power



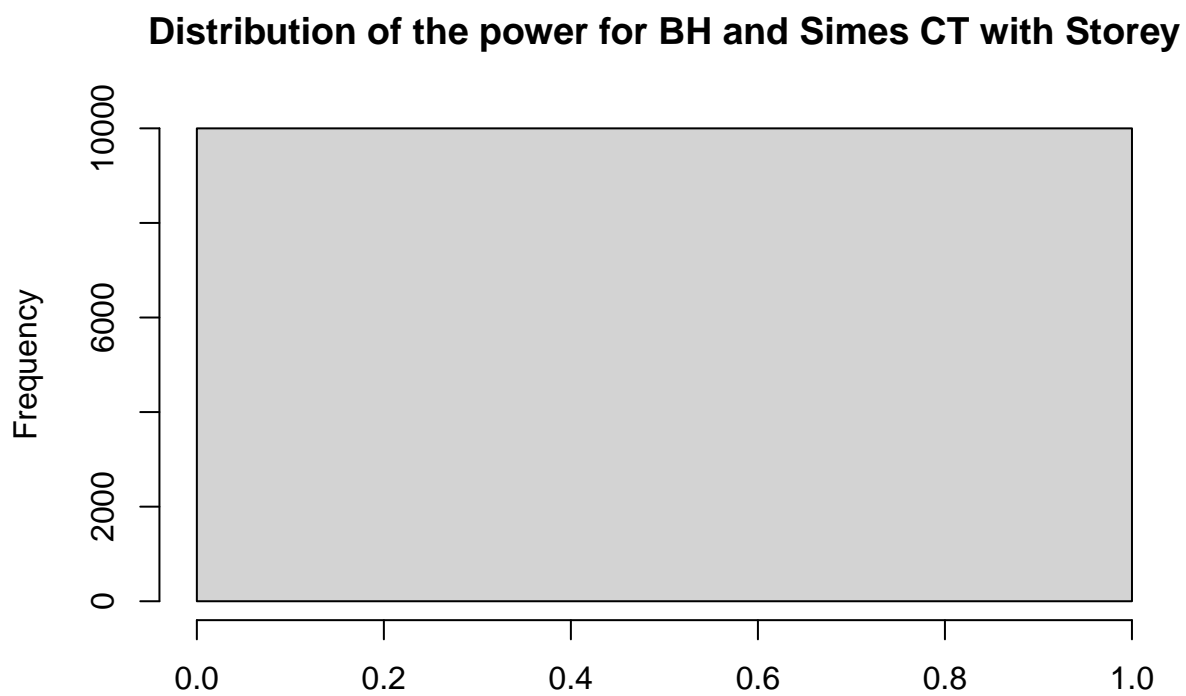
```
res$mean.powerGlobalNull
```

```
##      BH      BHSto      CTSim CTSimSto      CTWMW  
##      1        1        1        1        1
```

```
hist(as.integer(res$powerGlobalNull[[1]]), xlab="",  
     main="Distribution of the power for BH and Simes CT")
```

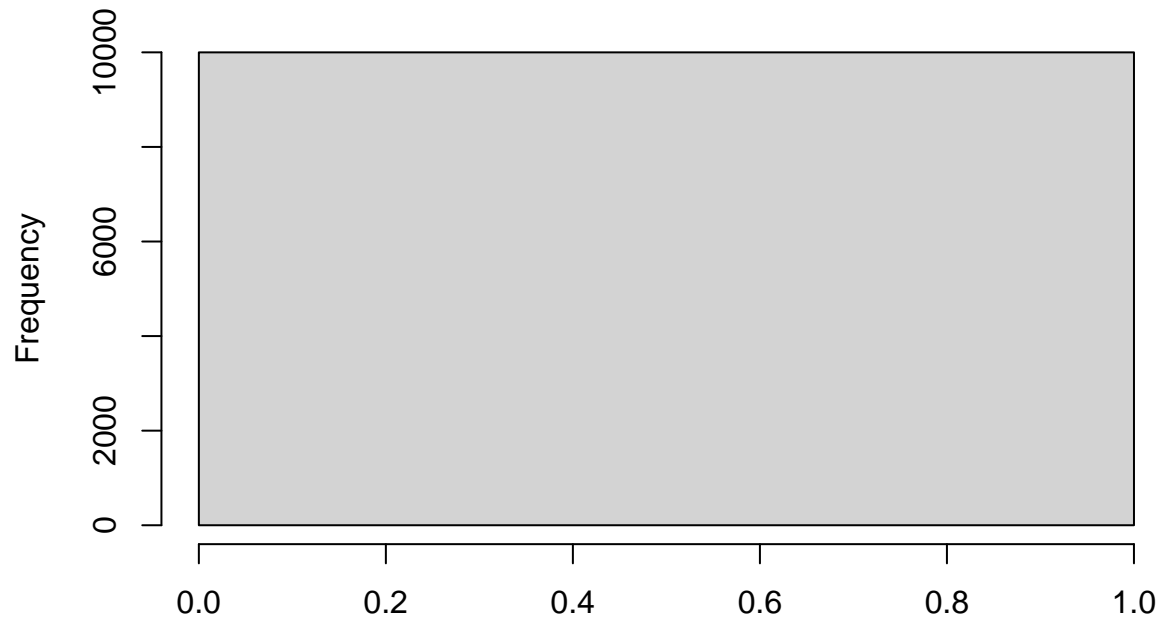



```
hist(as.integer(res$powerGlobalNull[[2]]), xlab="",  
     main="Distribution of the power for BH and Simes CT with Storey")
```



```
hist(as.integer(res$powerGlobalNull[[5]]), xlab="",  
     main="Distribution of the power for Wilcoxon-Mann-Whitney CT")
```

Distribution of the power for Wilcoxon–Mann–Whitney CT



```
# resDigits_v2 = res
# save(resDigits_v2,
#       file="C:/Users/c.magnani9/Documents/nout/trials/RealData/PowerStudy/resDigits_v2")
```

References

[1] Bates, S., E. Candes, L. Lei, Y. Romano, and M. Sesia (2023). Testing for outliers with conformal p-values. *Annals of Statistics*, {51}, 149–178.