

Comparison between different local tests: Simes, Simes with Storey and Wilcoxon-Mann-Whitney

20-04-2023

The aim is to compare 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.

We consider a null distribution F from which inliers come and we consider outliers which come from the alternative distribution

$$F^k$$

with $k > 0$, especially if $k \in \mathbb{N}_{>0}$ we know that F^k is the distribution of the random variable defined as the maximum of k observations drawn from F . So, we consider a sample drawn from the mixture distribution

$$G = (1 - \theta)F + \theta F^k$$

where $\theta \in [0, 1]$ is the proportion of inliers.

Since we deal with conformal p -values, we are interested not really in the sample distribution, but rather in the scores sample distribution. In our simulation study we draw n observations for the train set, l for the calibration set and mk for the test set. All observations are drawn from a d -multivariate standard normal distribution with $d = 3$ and using the algorithm of *isolation forest* trained on the training samples we compute the scores for the calibration and test samples. In order to generate scores associated to outlier observations, we consider the m blocks of k observations which create the test set. For each block i with $i = 1, \dots, m$, we draw from a Bernoulli random variable with success probability equal to θ . If the value is 1 then the i -th observation of the test set is inlier and we randomly sample its score value from the k scores of the i -th block, otherwise it is an outlier and its score value will be the maximum of the scores of the i -th block.

```
library(mvtnorm)
library(nout)
library(isotree)
```

```
## Warning: il pacchetto 'isotree' è stato creato con R versione 4.1.3
```

```
d_benjhoch = function(S_Y, S_X, alpha = 0.1){
  m = length(S_Y)
  n = length(S_X)
  pval = sapply(1:m, function(i) (1+sum(S_X >= S_Y[i]))/(n+1))
  d = sum(stats::p.adjust(pval, "BH") <= alpha)
  return(d)
}
```

```
d_StoreyBH = function(S_Y, S_X, alpha = 0.1, lambda=0.5){
```

```

m = length(S_Y)
n = length(S_X)
pval = sort(sapply(1:m, function(i) (1+sum(S_X >= S_Y[i]))/(n+1)), decreasing=FALSE)
pi0Sto = (1+sum(pval>lambda))/(m*(1-lambda))
d = sum(stats::p.adjust(pval,"BH")<=alpha/pi0Sto)
return(d)
}

```

```

scores_from_mixture = function(k, raw_scores, theta){

  if(theta>1 || theta<0){
    stop("Error: argument theta should in [0,1] interval")
  }

  ll = length(raw_scores)
  if(ll<k){
    stop("Error: length of raw_scores is smaller than k.")
  }

  quotient = ll%%k # is m
  remainder = ll%%k

  if(remainder != 0){
    cat("Warning: length of raw_scores is not a multiple of k. Last ",
        remainder, "elements of raw_scores will not be used.")
  }

  usable.raw_scores = raw_scores[1:(ll-remainder)]

  m1 = ifelse((theta*m)%1!=0, round(theta*m), theta*m)

  scores = rep(0, times = quotient)
  outlier = rep(0, times = quotient)

  if(m1==0){
    for(i in 0:(m-1)){
      scores[i+1] = sample(usable.raw_scores[(i*k+1):(i*k+k)], size=1)
    }
  }

  if(m1==m){
    for(i in 0:(m-1)){
      scores[i+1] = max(usable.raw_scores[(i*k+1):(i*k+k)])
      outlier[i+1]=T
    }
  }

  if(0<m1 & m1<m){
    for(i in 0:(m1-1)){

```

```

    scores[i+1] = max(usable.raw_scores[(i*k+1):(i*k+k)])
    outlier[i+1]=T
  }

  for(i in m1:(m-1)){
    scores[i+1] = sample(usable.raw_scores[(i*k+1):(i*k+k)], size=1)
  }
}

return(list("scores"=scores, "outlier"=outlier))
}

simuLMPI = function(B=10^4, n, l, m, d = 3, k = 2, theta, alpha = m/(l+1)){

  train = mvtnorm::rmvnorm(n=n, mean=rep(0,d))
  iso.fo = isotree::isolation.forest(train, ndim=d, ntrees=10, nthreads=1,
                                     scoring_metric = "depth", output_score = TRUE)

  crit=critWMW(m=m, n=n, 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 = mvtnorm::rmvnorm(n=1, mean=rep(0,d))
    te = mvtnorm::rmvnorm(n=k*m, mean=rep(0,d))

    S_cal = isotree::predict.isolation_forest(iso.fo$model, cal, type = "score")
    rawS_te = isotree::predict.isolation_forest(iso.fo$model, te, type = "score")
    gen.te.score = scores_from_mixture(k=k, raw_scores=rawS_te, theta=theta)
    S_te = gen.te.score$scores

    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,
        "theta"=theta, "alpha"=alpha))
}

```

K=2

```

set.seed(321)

# Initializing parameters
B=10^5
n = 19
l = 19
m = 2
d = 3
k = 2
alpha = m/(l+1)
mis = seq(from=0, to=m, by=1)
thetas = mis/m

# Results
res = lapply(thetas, function(theta) simulMPI(B=B, n=n, l=l, m=m, d = d,
        k = k, theta, alpha = m/(l+1)))

# Storing results
store_res = list("mean.discov" = matrix(nrow=length(thetas), ncol = 5),
        "mean.powerGlobalNull" = matrix(nrow=length(thetas), ncol = 5))
row.names = rep(NA, times=length(thetas))
for(i in 1:length(thetas)){
    row.names[i] = paste("theta =",thetas[i])
}
rownames(store_res$mean.discov) = row.names
colnames(store_res$mean.discov) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")
rownames(store_res$mean.powerGlobalNull) = row.names
colnames(store_res$mean.powerGlobalNull) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")

for(i in 1:length(res)){
    store_res$mean.discov[i,] = res[[i]]$mean.discov
    store_res$mean.powerGlobalNull[i,] = res[[i]]$mean.powerGlobalNull
}

store_res$mean.discov

##           BH   StoBH   Simes StoSimes   WMW
## theta = 0   0.10593 0.06235 0.10593  0.05297 0.10712
## theta = 0.5 0.16146 0.10741 0.16146  0.09073 0.17947
## theta = 1   0.22360 0.17822 0.22360  0.14943 0.29363

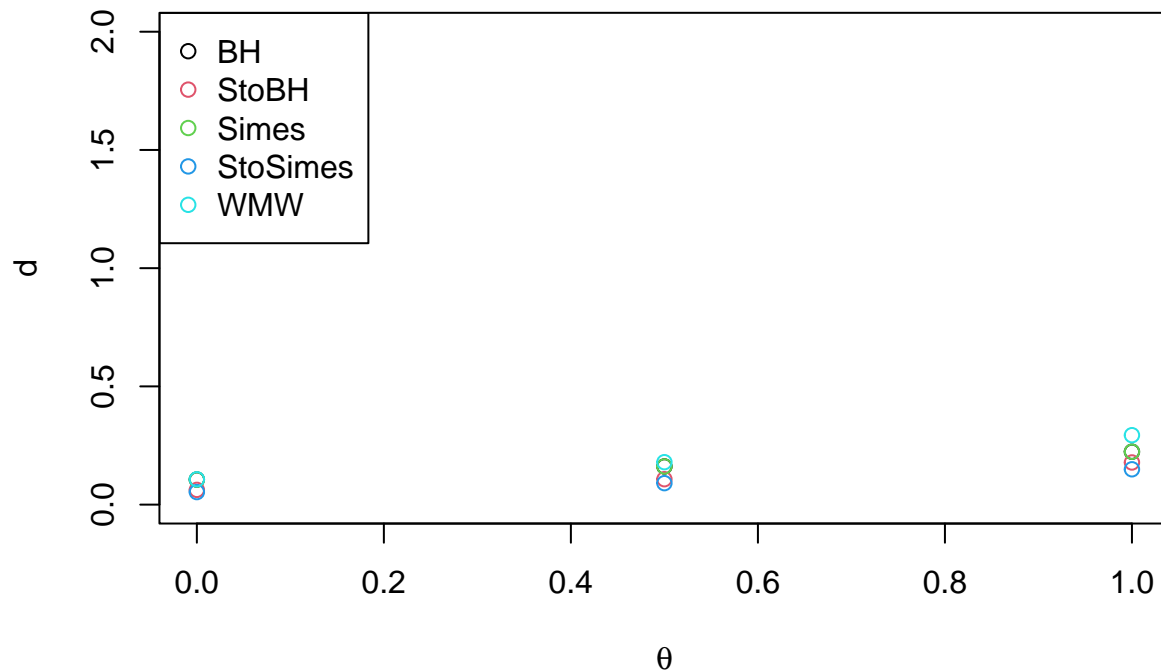
```

```
store_res$mean.powerGlobalNull
```

```
##           BH   StoBH   Simes StoSimes   WMW
## theta = 0  0.09226 0.04868 0.09226  0.04868 0.09345
## theta = 0.5 0.13632 0.08227 0.13632  0.08227 0.15433
## theta = 1   0.17998 0.13460 0.17998  0.13460 0.25001
```

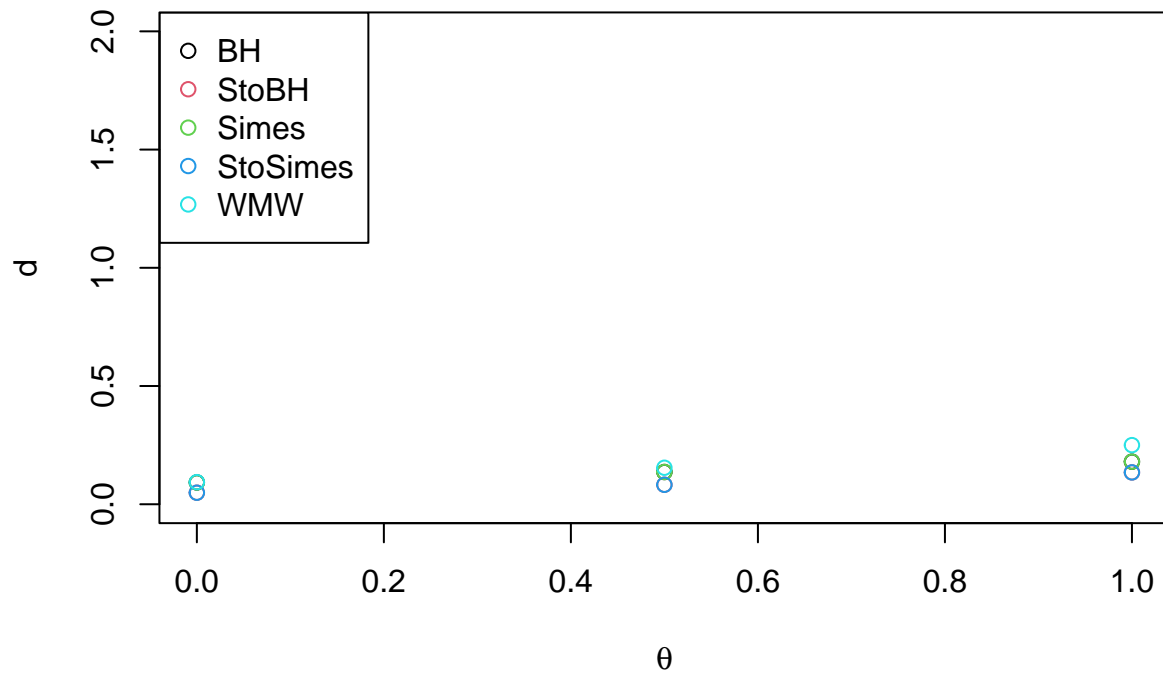
```
plot(x = thetas, y = store_res$mean.discov[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the number of discoveries on B replications")
points(x = thetas, y = store_res$mean.discov[,2], col = 2)
points(x = thetas, y = store_res$mean.discov[,3], col = 3)
points(x = thetas, y = store_res$mean.discov[,4], col = 4)
points(x = thetas, y = store_res$mean.discov[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.discov), col = c(1,2,3,4,5))
```

Mean of the number of discoveries on B replications



```
plot(x = thetas, y = store_res$mean.powerGlobalNull[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the power on B replications")
points(x = thetas, y = store_res$mean.powerGlobalNull[,2], col = 2)
points(x = thetas, y = store_res$mean.powerGlobalNull[,3], col = 3)
points(x = thetas, y = store_res$mean.powerGlobalNull[,4], col = 4)
points(x = thetas, y = store_res$mean.powerGlobalNull[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.powerGlobalNull), col = c(1,2,3,4,5))
```

Mean of the power on B replications



K=3

```
set.seed(321)

# Initializing parameters
B=10^5
n = 19
l = 19
m = 2
d = 3
k = 3
alpha = m/(n+1)
mis = seq(from=0, to=m, by=1)
thetas = mis/m

# Results
res = lapply(thetas, function(theta) simulMPI(B=B, n=n, l=l, m=m, d = d,
                                              k = k, theta, alpha = m/(l+1)))

# Storing results
store_res = list("mean.discov" = matrix(nrow=length(thetas), ncol = 5),
                 "mean.powerGlobalNull" = matrix(nrow=length(thetas), ncol = 5))
row.names = rep(NA, times=length(thetas))
for(i in 1:length(thetas)){
  row.names[i] = paste("theta =", thetas[i])
}
```

```

}
rownames(store_res$mean.discov) = row.names
colnames(store_res$mean.discov) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")
rownames(store_res$mean.powerGlobalNull) = row.names
colnames(store_res$mean.powerGlobalNull) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")

for(i in 1:length(res)){
  store_res$mean.discov[i,] = res[[i]]$mean.discov
  store_res$mean.powerGlobalNull[i,] = res[[i]]$mean.powerGlobalNull
}

store_res$mean.discov

##              BH   StoBH   Simes StoSimes   WMW
## theta = 0    0.10359 0.06030 0.10359  0.05102 0.10451
## theta = 0.5  0.21311 0.14517 0.21311  0.12188 0.23430
## theta = 1    0.33588 0.30422 0.33588  0.24832 0.49315

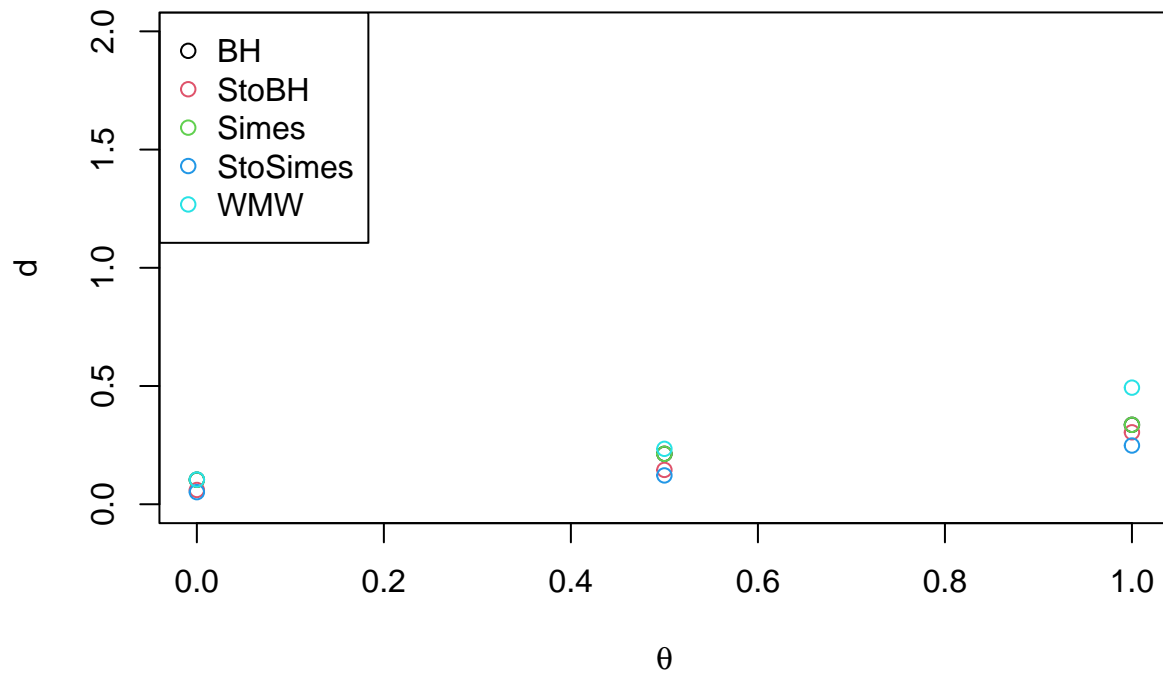
store_res$mean.powerGlobalNull

##              BH   StoBH   Simes StoSimes   WMW
## theta = 0    0.09038 0.04709 0.09038  0.04709 0.09130
## theta = 0.5  0.17789 0.10995 0.17789  0.10995 0.19908
## theta = 1    0.25003 0.21837 0.25003  0.21837 0.40730

plot(x = thetas, y = store_res$mean.discov[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the number of discoveries on B replications")
points(x = thetas, y = store_res$mean.discov[,2], col = 2)
points(x = thetas, y = store_res$mean.discov[,3], col = 3)
points(x = thetas, y = store_res$mean.discov[,4], col = 4)
points(x = thetas, y = store_res$mean.discov[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.discov), col = c(1,2,3,4,5))

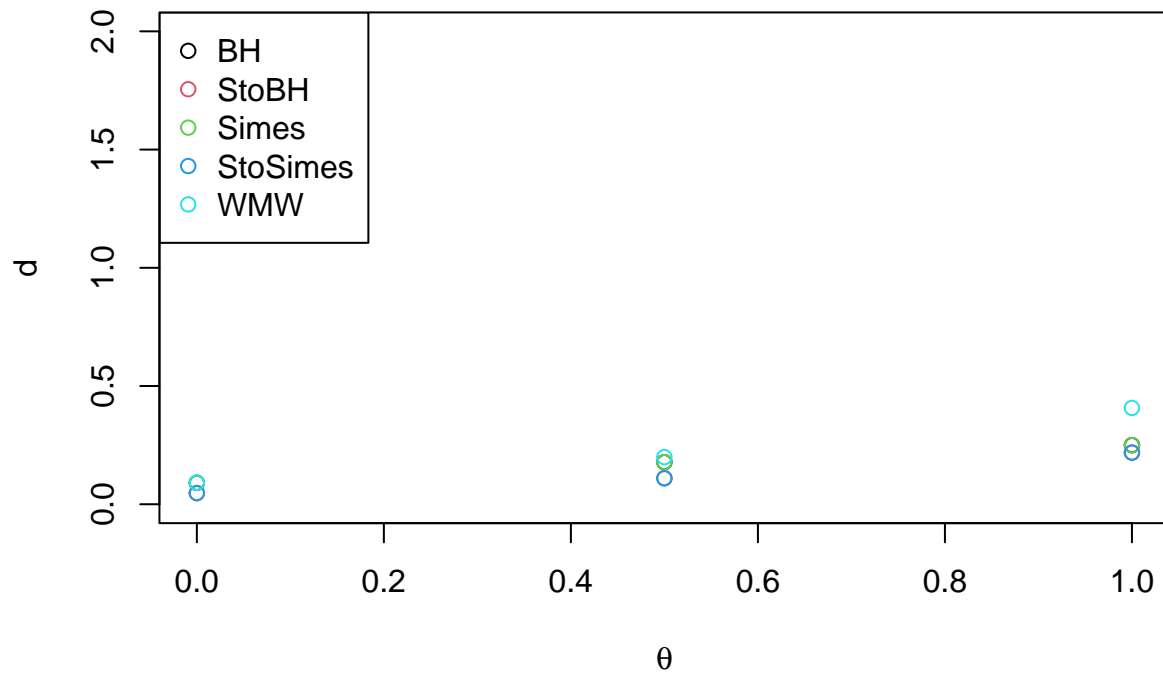
```

Mean of the number of discoveries on B replications



```
plot(x = thetas, y = store_res$mean.powerGlobalNull[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the power on B replications")
points(x = thetas, y = store_res$mean.powerGlobalNull[,2], col = 2)
points(x = thetas, y = store_res$mean.powerGlobalNull[,3], col = 3)
points(x = thetas, y = store_res$mean.powerGlobalNull[,4], col = 4)
points(x = thetas, y = store_res$mean.powerGlobalNull[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.powerGlobalNull), col = c(1,2,3,4,5))
```


Mean of the power on B replications



K=5

```
set.seed(321)

# Initializing parameters
B=10^5
n = 19
l = 19
m = 2
d = 3
k = 5
alpha = m/(l+1)
mis = seq(from=0, to=m, by=1)
thetas = mis/m
# Results
res = lapply(thetas, function(theta) simulMPI(B=B, n=n, l=l, m=m, d = d,
                                              k = k, theta, alpha = m/(l+1)))

# Storing results
store_res = list("mean.discov" = matrix(nrow=length(thetas), ncol = 5),
                 "mean.powerGlobalNull" = matrix(nrow=length(thetas), ncol = 5))
row.names = rep(NA, times=length(thetas))
for(i in 1:length(thetas)){
  row.names[i] = paste("theta =", thetas[i])
}
```

```

rownames(store_res$mean.discov) = row.names
colnames(store_res$mean.discov) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")
rownames(store_res$mean.powerGlobalNull) = row.names
colnames(store_res$mean.powerGlobalNull) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")

for(i in 1:length(res)){
  store_res$mean.discov[i,] = res[[i]]$mean.discov
  store_res$mean.powerGlobalNull[i,] = res[[i]]$mean.powerGlobalNull
}

store_res$mean.discov

##           BH   StoBH   Simes StoSimes   WMW
## theta = 0   0.10679 0.06158 0.10679  0.05258 0.10495
## theta = 0.5 0.28975 0.19300 0.28975  0.16252 0.29814
## theta = 1   0.54767 0.53166 0.54767  0.42445 0.81167

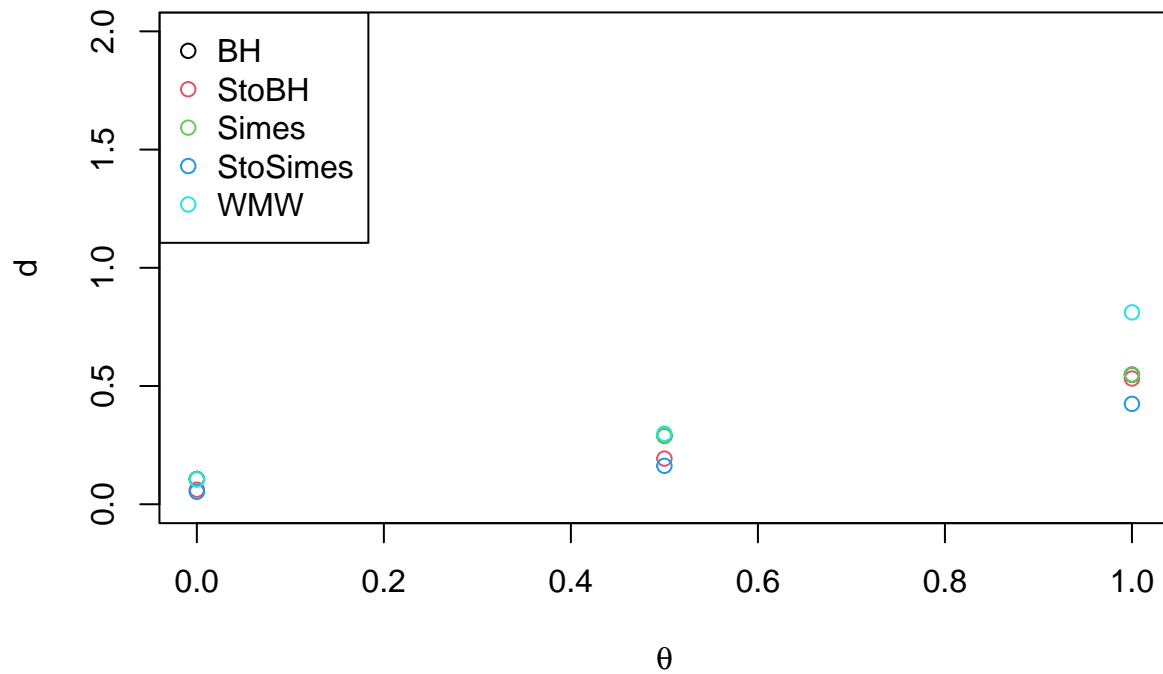
store_res$mean.powerGlobalNull

##           BH   StoBH   Simes StoSimes   WMW
## theta = 0   0.09370 0.04849 0.09370  0.04849 0.09186
## theta = 0.5 0.24255 0.14580 0.24255  0.14580 0.25094
## theta = 1   0.37586 0.35985 0.37586  0.35985 0.63986

plot(x = thetas, y = store_res$mean.discov[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the number of discoveries on B replications")
points(x = thetas, y = store_res$mean.discov[,2], col = 2)
points(x = thetas, y = store_res$mean.discov[,3], col = 3)
points(x = thetas, y = store_res$mean.discov[,4], col = 4)
points(x = thetas, y = store_res$mean.discov[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.discov), col = c(1,2,3,4,5))

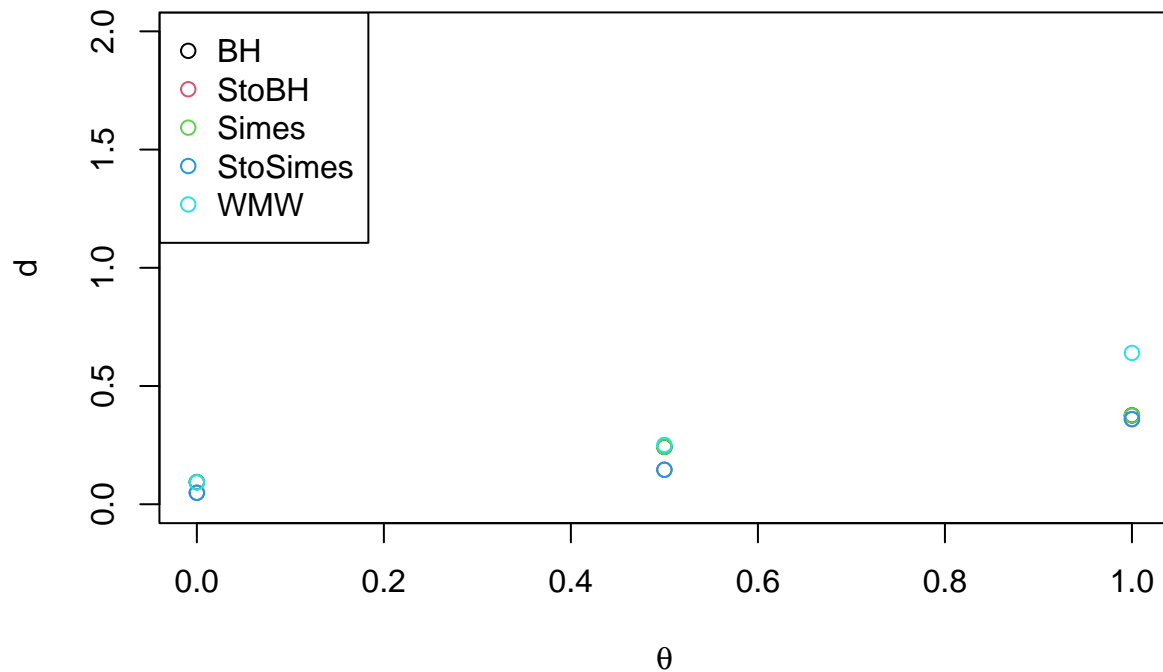
```

Mean of the number of discoveries on B replications



```
plot(x = thetas, y = store_res$mean.powerGlobalNull[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the power on B replications")
points(x = thetas, y = store_res$mean.powerGlobalNull[,2], col = 2)
points(x = thetas, y = store_res$mean.powerGlobalNull[,3], col = 3)
points(x = thetas, y = store_res$mean.powerGlobalNull[,4], col = 4)
points(x = thetas, y = store_res$mean.powerGlobalNull[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.powerGlobalNull), col = c(1,2,3,4,5))
```

Mean of the power on B replications



K=10

```
set.seed(321)

# Initializing parameters
B=10^5
n = 19
l = 19
m = 2
d = 3
k = 10
alpha = m/(l+1)
mis = seq(from=0, to=m, by=1)
thetas = mis/m
# Results
res = lapply(thetas, function(theta) simulMPI(B=B, n=n, l=l, m=m, d = d,
                                              k = k, theta, alpha = m/(l+1)))

# Storing results
store_res = list("mean.discov" = matrix(nrow=length(thetas), ncol = 5),
                 "mean.powerGlobalNull" = matrix(nrow=length(thetas), ncol = 5))
row.names = rep(NA, times=length(thetas))
for(i in 1:length(thetas)){
  row.names[i] = paste("theta =", thetas[i])
}
```

```

rownames(store_res$mean.discov) = row.names
colnames(store_res$mean.discov) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")
rownames(store_res$mean.powerGlobalNull) = row.names
colnames(store_res$mean.powerGlobalNull) = c("BH", "StoBH", "Simes", "StoSimes", "WMW")

for(i in 1:length(res)){
  store_res$mean.discov[i,] = res[[i]]$mean.discov
  store_res$mean.powerGlobalNull[i,] = res[[i]]$mean.powerGlobalNull
}

store_res$mean.discov

##              BH   StoBH   Simes StoSimes   WMW
## theta = 0    0.10450 0.06109 0.10450  0.05167 0.10479
## theta = 0.5  0.42185 0.26844 0.42185  0.22660 0.37136
## theta = 1    0.93889 0.93671 0.93889  0.73183 1.24884

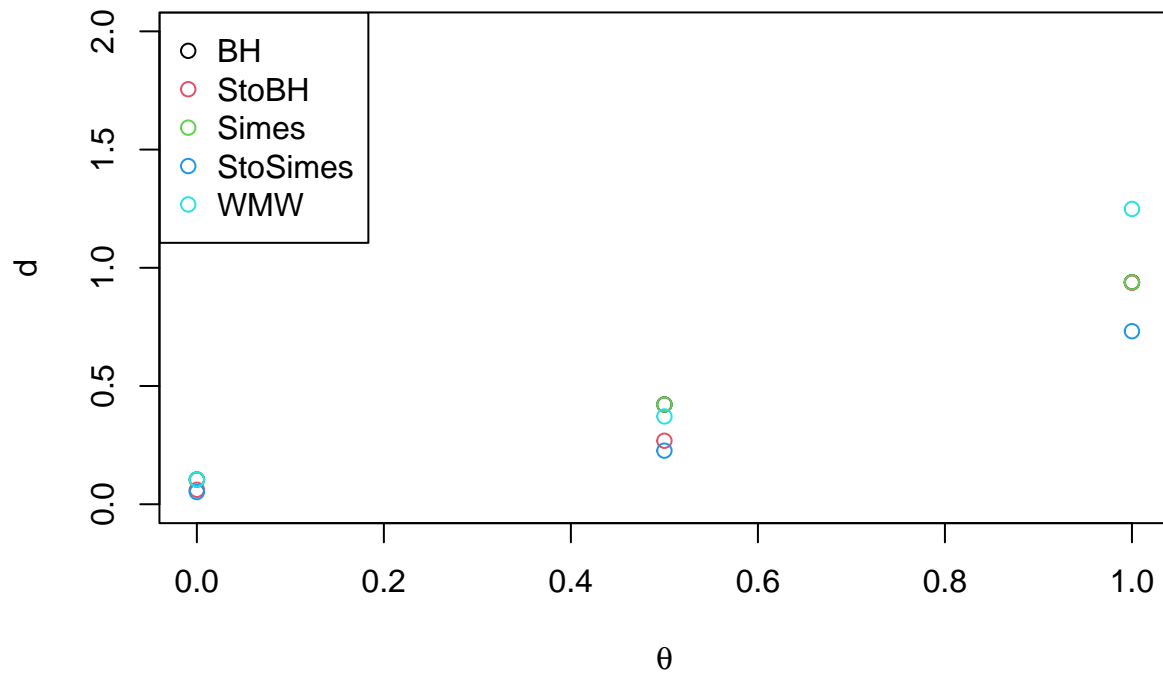
store_res$mean.powerGlobalNull

##              BH   StoBH   Simes StoSimes   WMW
## theta = 0    0.09112 0.04771 0.09112  0.04771 0.09141
## theta = 0.5  0.35480 0.20139 0.35480  0.20139 0.30431
## theta = 1    0.57199 0.56981 0.57199  0.56981 0.88194

plot(x = thetas, y = store_res$mean.discov[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the number of discoveries on B replications")
points(x = thetas, y = store_res$mean.discov[,2], col = 2)
points(x = thetas, y = store_res$mean.discov[,3], col = 3)
points(x = thetas, y = store_res$mean.discov[,4], col = 4)
points(x = thetas, y = store_res$mean.discov[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.discov), col = c(1,2,3,4,5))

```

Mean of the number of discoveries on B replications



```
plot(x = thetas, y = store_res$mean.powerGlobalNull[,1], col = 1, ylab = "d",
     xlab = expression(theta), ylim=c(0,m),
     main = "Mean of the power on B replications")
points(x = thetas, y = store_res$mean.powerGlobalNull[,2], col = 2)
points(x = thetas, y = store_res$mean.powerGlobalNull[,3], col = 3)
points(x = thetas, y = store_res$mean.powerGlobalNull[,4], col = 4)
points(x = thetas, y = store_res$mean.powerGlobalNull[,5], col = 5)
legend("topleft", pch = 1, legend = colnames(store_res$mean.powerGlobalNull), col = c(1,2,3,4,5))
```

Mean of the power on B replications

