

MHT Simulation

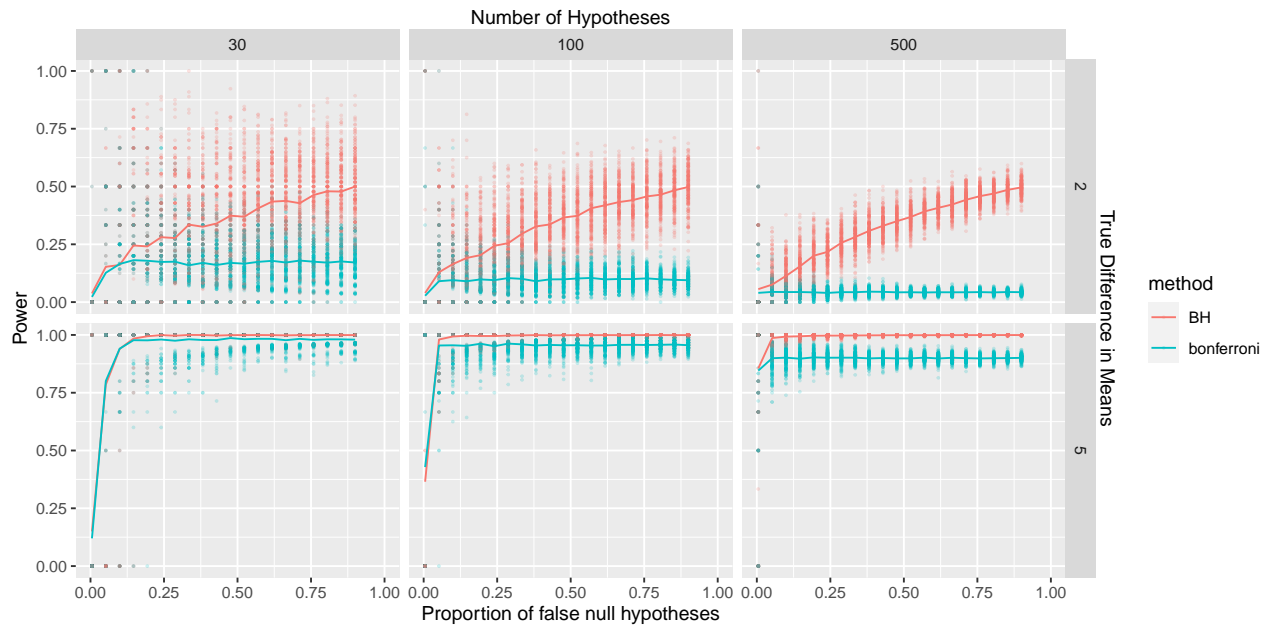
Connor Demorest

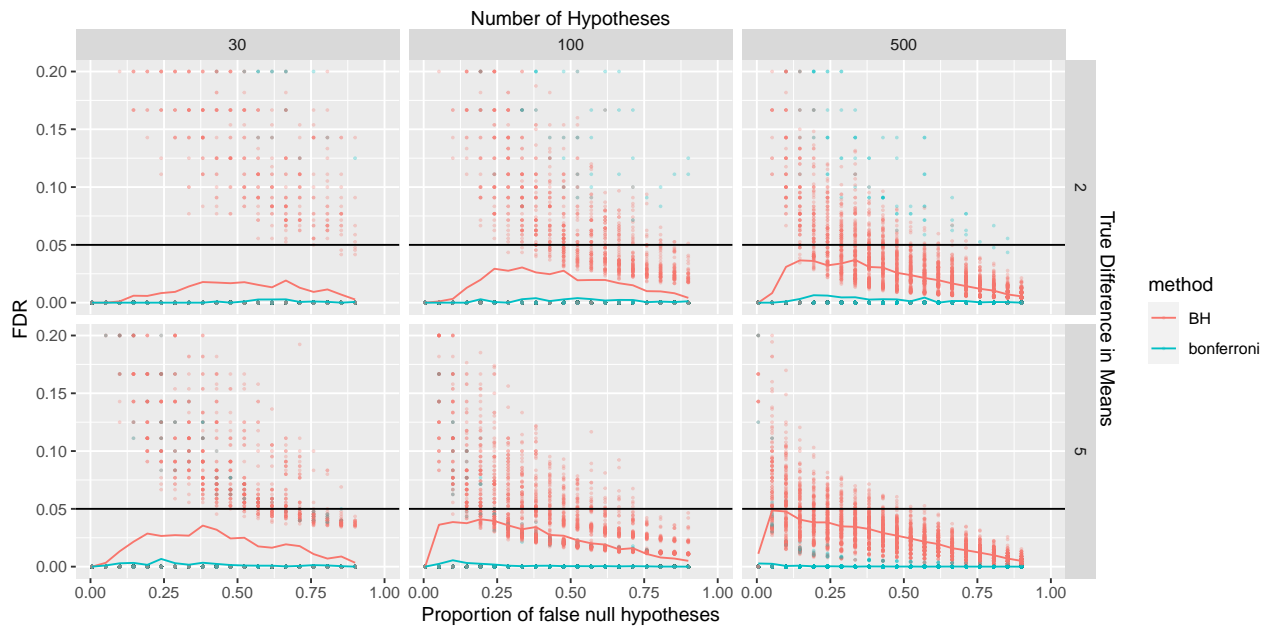
9/4/2021

Simulation study:

- 1) Take 100 samples from $N(0,1)$ and $N(5,1)$ mixed with proportion p in $(0,1)$. Take $\alpha = 0.05$. $X_i = (1 - p) * N(0, 1) + p * N(3, 1)$
- 2) Find p values of X
- 3) Use BH, etc
- 4) Get FDP and Power
- 5) Repeat many times to find FDR and avg power across values of p

As a function:





Using Lfdr and PFDR

Code:

```
knitr::opts_chunk$set(echo = F, fig.width = 10, fig.height = 5,
  cache = TRUE)

library(tidyverse)
p = 0.4
a = 0.05
num_samples = 100
# True if true mean != 0
A = rbernoulli(num_samples, p)
t_nulls = rnorm(num_samples, 0, 1)
f_nulls = rnorm(num_samples, 2, 1)
samples = ifelse(A, f_nulls, t_nulls)
p_vals = pnorm(samples, lower.tail = F)

# Test discoveries
B = (p_vals <= a)
# FP = A' & B
## False_positive = (A & B)
# FDP = P(A' & B)/P(B)
FDP = ifelse(is.na(sum(!A & B)/sum(B)), 0, sum(!A & B)/sum(B))

# Power = P(reject Null | Null false) = P(reject & Null false)/P(Null false)
Power = ifelse(is.na(sum(A & B)/sum(A)), 0, sum(A & B)/sum(A))

# FNR = ?????? :)
# m === number of total hypotheses
func = function(p = 0.4, a = 0.05, m = 1000, alt_hyp = 3, method = "none") {
  # A === True discoveries
  A = rbernoulli(m, p)
  t_nulls = rnorm(m, 0, 1)
```

```

f_nulls = rnorm(m, alt_hyp, 1)
samples = ifelse(A, f_nulls, t_nulls)
p_vals = pnorm(samples, lower.tail = F)
adj_p_vals = p.adjust(p = p_vals, method = method)

# Test discoveries
B = (adj_p_vals <= a)
# FP = A' & B
False_positive = ifelse(!A & B, 1, 0)
# FDP = P(A' & B)/P(B)
FDP = ifelse(is.na(sum(!A & B)/sum(B)), 0, sum(!A & B)/sum(B))

# Power = P(reject Null | Null false) = P(reject & Null false)/P(Null false)
Power = ifelse(is.na(sum(A & B)/sum(A)), 0, sum(A & B)/sum(A))
return(c(FDP, Power))
}

# New function: do.many, where it just does 'many' and returns averages?
# Calculate FDR and Avg Power

# many = replicate(1000, func(method = "none")) %>% apply(., 1, FUN = mean) %>% print

df = expand_grid(p = rep(seq(0.005, 0.9, length.out = 20), each = 200),
                method = c("bonferroni", "BH"),
                alt_hyp = c(2,5),
                m = c(30, 100, 500))
many = t(mapply(func, p = df$p, method = df$method, alt_hyp = df$alt_hyp, m = df$m)) %>%
  as_tibble %>%
  rename(FDR = V1, Power = V2) %>%
  mutate(df)

# do_many = function(num_reps = 1000, p = 0.4, a = 0.05, num_samples = 1000, method = "none") {
#   many = replicate(num_reps, func(p = p, a = a, num_samples = num_samples, method = method)) %>%
#   apply(., 1, FUN = mean)
#   return(many)
# }

#
# df = expand_grid(p = seq(0.1, 0.9, by = 0.01), method = c("none", "BH"))
# df2 = t(mapply(do_many, p = df$p, method = df$method, num_reps = 100)) %>%
# as_tibble %>%
#   rename(FDR = V1, Power = V2) %>%
#   mutate(df)

# Increased p gives more true discoveries
ggplot(data = many, aes(x = p, y = Power, color = method)) +
  geom_point(size = 0.3, alpha = 0.2) +
  stat_summary(fun = "mean", geom = "line") +
  facet_grid(alt_hyp ~ m) +
  labs(x = "Proportion of false null hypotheses") +
  scale_y_continuous(sec.axis = sec_axis(~ ., name = "True Difference in Means",
                                          breaks = NULL, labels = NULL), limits = c(0,1)) +
  scale_x_continuous(sec.axis = sec_axis(~ ., name = "Number of Hypotheses",
                                          breaks = NULL, labels = NULL), limits = c(0,1))

```

```

ggplot(data = many, aes(x = p, y = FDR, color = method)) +
  geom_point(size = 0.3, alpha = 0.3) +
  stat_summary(fun = "mean", geom = "line") +
  facet_grid(alt_hyp ~ m) +
  geom_abline(intercept = 0.05, slope = 0) +
  labs(x = "Proportion of false null hypotheses") +
  scale_y_continuous(sec.axis = sec_axis(~ ., name = "True Difference in Means",
                                           breaks = NULL, labels = NULL), limits = c(0,0.2)) +
  scale_x_continuous(sec.axis = sec_axis(~ ., name = "Number of Hypotheses",
                                           breaks = NULL, labels = NULL), limits = c(0,1))

# Basically re-using code from above
# p known here!
p = 0.4
alpha = 0.05
num_samples = 1000
## True if true mean != 0
# theta = rbernoulli(num_samples, p)
# f0 = rnorm(num_samples, 0, 1)
# f1 = rnorm(num_samples, 5, 1)
# X = ifelse(theta, f1, f0)
# f = (1-p)*dnorm(X, 0, 1) + p*dnorm(X, 5, 1)

# Lfdr = (1-pi)f0(X)/f(x)
PFDR_fun = function(p = 0.4, alpha = 0.05, num_samples = 1000, alt_hyp = 4) {
  df = tibble(theta = rbernoulli(num_samples, p),
              X = ifelse(theta, rnorm(num_samples, alt_hyp, 1), rnorm(num_samples, 0, 1)),
              f = (1-p)*dnorm(X, 0, 1) + p*dnorm(X, alt_hyp, 1),
              Lfdr = (1-p)*dnorm(X, 0, 1)/f) %>%
    arrange(Lfdr) %>%
    mutate(q = cumsum(Lfdr)/(1:length(Lfdr)),
           `q<a` = q < alpha, # Total discoveries
           Type1error = `q<a`&!theta, # Type 1 errors
           Type2error = theta&!`q<a` %>% # Type 2 errors
           summarise(theta = sum(theta), disc = sum(`q<a`), t1 = sum(Type1error), t2 = sum(Type2error))
    return(df)
}

# PFDR_fun()
reps = plyr::rdply(1000, PFDR_fun(), .id = NULL)

# PFDR = rep(NA, length(Lfdr))
# Q = rep(NA, length(Lfdr))
# for(t in 1:length(Lfdr)) {
#   Q[t] = (sum(sort(Lfdr)[1:t]/t))# %>% print
#   PFDR[t] = ((sum(sort(Lfdr)[1:t]/t)) < alpha)# %>% print
# }

```