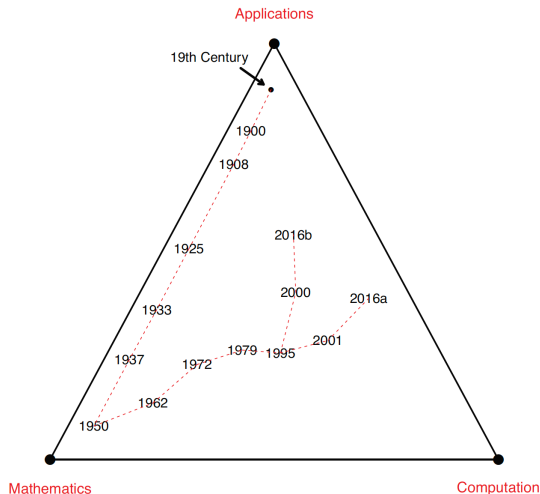


# Multiple testing

Aldo Solari  
Statistical Inference II  
PhD in Economics and Statistics  
University of Milano-Bicocca



# Where are we going?



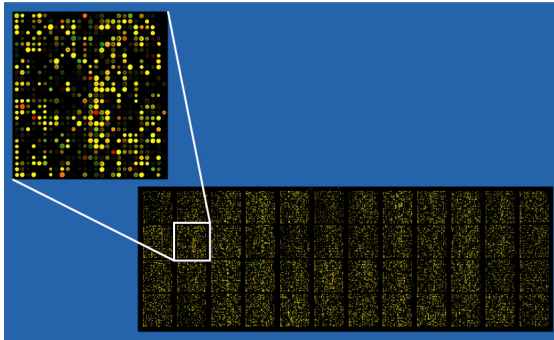
“To predict the future of statistics, we can at least examine the past to see how we’ve gotten where we are” (Efron & Hastie, 2016)

# The three eras of statistics

Efron (2012) Large-Scale Inference. Cambridge University Press

1. **The age of huge census-level data sets** were brought to bear on simple but important questions:
  - Are there more male than female births?
  - Is the rate of insanity rising?
2. **The classical period** of Pearson, Fisher, Neyman, Hotelling, and their successors, intellectual giants who developed a theory of optimal inference capable of wringing every drop of information out of a scientific experiment. The questions dealt with still tended to be simple
  - Is treatment A better than treatment B?
3. **The era of scientific mass production**, in which new technologies typified by the *microarray* allow a single team of scientists to produce *high-dimensional data*. But now the flood of data is accompanied by a deluge of questions, perhaps thousands of estimates or hypothesis tests that the statistician is charged with answering together

# Microarrays



- Biomedical devices that enabled the assessment of individual activity for thousands of genes at once
- Need to carry out thousands of simultaneous hypothesis tests, done with the prospect of finding only a few interesting genes among a haystack of null cases

# Prostate cancer data

<https://web.stanford.edu/~hastie/CASI/data.html>

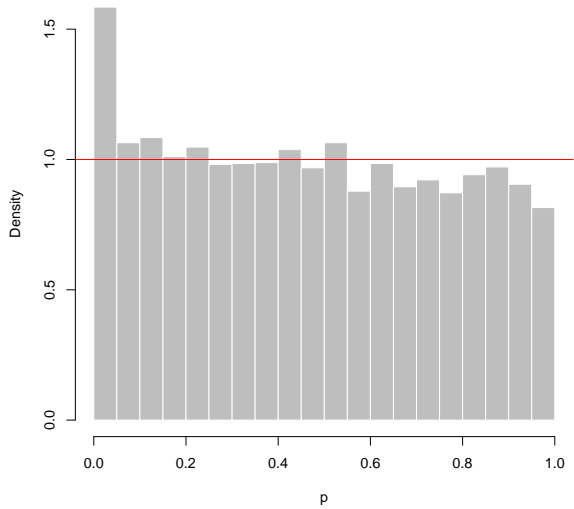
- The prostate cancer data came from a microarray study of  $n = 102$  men, 52 prostate cancer patients and 50 controls
- Each man's gene expression levels were measured on  $m = 6033$  genes, yielding a  $102 \times 6033$  matrix

$x_{ji}$  = activity of  $i$ th gene for  $j$ th subject

- The  $i$ th null hypothesis, denoted  $H_i$ , would state that the mean expression level of the  $i$ th gene is the same in both groups of patients

$$H_i : \mathbb{E}(X_i^{cancer}) = \mathbb{E}(X_i^{control})$$

- For each gene, a  $p$ -value  $p_i$  is computed



# The four eras of data

Leek (2016)

- **The era of not much data** Prior to about 1995, usually we could collect a few measurements at a time. The whole point of statistics was to try to optimally squeeze information out of a small number of samples - so you see methods like maximum likelihood and minimum variance unbiased estimators being developed
- **The era of lots of measurements on a few samples** This one hit hard in biology with the development of the microarray and the ability to measure thousands of genes simultaneously. This is the same statistical problem as in the previous era but with a lot more noise added. Here you see the development of methods for multiple testing and regularized regression to separate signals from piles of noise
- **The era of a few measurements on lots of samples** This era is overlapping to some extent with the previous one. Large scale collections of data from EMRs and Medicare are examples where you have a huge number of people (samples) but a relatively modest number of variables measured. Here there is a big focus on statistical methods for knowing how to model different parts of the data with hierarchical models and separating signals of varying strength with model calibration
- **The era of all the data on everything** This is an era that currently we as civilians don't get to participate in. But Facebook, Google, Amazon, the NSA and other organizations have thousands or millions of measurements on hundreds of millions of people. Other than just sheer computing I'm speculating that a lot of the problem is in segmentation (like in era 3) coupled with avoiding crazy overfitting (like in era 2)

# Many tests

- In a single test, the probability of making a type I error is bounded by  $\alpha$ , conventionally set at 0.05
- Problems arise, however, when researchers do not perform a single hypothesis test but many of them
- There are many ways of dealing with type I errors. We will focus on three types of multiple testing methods:
  1. those that control the *FamilyWise Error Rate* (FWER)
  2. those that control the *False Discovery Rate* (FDR)
  3. those that estimate the *False Discovery Proportion* (FDP) or make confidence intervals for it

## Reference

Goeman, Solari (2014)

Multiple Hypothesis Testing in Genomics.

Statistics in Medicine, 33, 1946–78



# Rejections

Suppose we have a collection  $\mathcal{H} = \{H_1, \dots, H_m\}$  of  $m$  null hypotheses:

- an unknown number  $m_0$  of these hypotheses is true, whereas the other  $m_1 = m - m_0$  is false. The proportion of true hypotheses is  $\pi_0 = m_0/m$
- The collection of true hypotheses is  $\mathcal{T} \subseteq \mathcal{H}$  and of false hypotheses is  $\mathcal{F} = \mathcal{H} \setminus \mathcal{T}$
- The goal of a multiple testing procedure is to choose a collection  $\mathcal{R} \subseteq \{H_1, \dots, H_m\}$  of hypotheses to reject. If we have  $p$ -values  $p_1, \dots, p_m$  for  $H_1, \dots, H_m$ , a natural choice is

$$\mathcal{R} = \{H_i : p_i \leq c\}$$

rejecting all hypotheses with a  $p$ -value below a critical value  $c$

# Errors

Ideally, the set of rejected hypotheses  $\mathcal{R}$  should coincide with the set  $\mathcal{F}$  of false hypotheses as much as possible. However, two types of error can be made:

- type I errors: the rejected hypotheses that are true hypotheses, i.e.  $\mathcal{R} \cap \mathcal{T}$
- type II errors: the false hypotheses that we failed to reject, i.e.  $\mathcal{F} \setminus \mathcal{R}$

Rejected hypotheses are sometimes called *discoveries*, hence the terms *true discovery* and *false discovery* are sometimes used for correct and incorrect rejections

# Type I errors

- Type I errors are traditionally considered more problematic than type II errors
- If a rejected hypothesis allows publication of a scientific finding, a type I error brings a false discovery, and the risk of publication of a potentially misleading scientific result
- Type II errors, on the other hand, mean missing out on a scientific result. Although unfortunate for the individual researcher, the latter is, in comparison, less harmful to scientific research as a whole

## $2 \times 2$ table

We can summarize the numbers of errors in a contingency table:

	true	false	total
rejected	$V$	$U$	$R$
not rejected	$m_0 - V$	$m_1 - U$	$m - R$
total	$m_0$	$m_1$	$m$

We can observe  $m$  and  $R = |\mathcal{R}|$ , but all quantities in the first two columns of the table are unobservable

# FDP

The false discovery proportion (FDP)  $Q$  is defined as

$$Q = \frac{V}{\max(R, 1)} = \begin{cases} V/R & \text{if } R > 0 \\ 0 & \text{otherwise,} \end{cases}$$

the proportion of false rejections among all rejections, defined as 0 if no rejections are made

# FWER and FDR control

1. Familywise error rate (FWER):

$$\text{FWER} = P(V > 0) = P(Q > 0)$$

the probability that the rejections contains any type I error

2. False discovery rate (FDR):

$$\text{FDR} = \mathbb{E}(Q)$$

the expected proportion of type I errors among the rejections

We say that FWER or FDR is *controlled* at level  $\alpha$  when the set  $\mathcal{R}$  is chosen in such a way that the corresponding aspect of the distribution of  $Q$  is guaranteed to be at most  $\alpha$ , i.e.

$$\text{FWER} \leq \alpha \quad \text{or} \quad \text{FDR} \leq \alpha$$

$$\text{FWER} \geq \text{FDR}$$

- The two error rates FDR and FWER are related. Because  $0 \leq Q \leq 1$ , we have  $Q \leq \mathbb{1}\{Q > 0\}$  and

$$E(Q) \leq P(Q > 0)$$

which means that FWER control implies FDR control

- If all hypotheses are true, FDR and FWER are identical; because  $R = V$  in this case,  $Q$  is a Bernoulli variable, and

$$E(Q) = P(Q > 0)$$

- Both FDR and FWER are proper generalizations of the concept of type I error to multiple hypotheses: if  $m = 1$ , the two error rates are identical and equal to the type I error rate

$m = 100, m_0 = 80, T_i \sim N(\mu_i, 1), \mu_i = 0$  if  $H_i$  true,  $\mu_i = 2$  otherwise

	1	2	3	4	5	6	7	8	9	10
$R$	20	17	23	16	20	16	15	17	20	17
$V$	4	5	6	5	5	3	3	5	7	4
$\mathbb{1}\{V > 0\}$	1	1	1	1	1	1	1	1	1	1
$V/R$	0.20	0.29	0.26	0.31	0.25	0.19	0.20	0.29	0.35	0.24

Reject  $H_i$  if  $p_i \leq 0.05$  gives FWER = 0.983 and FDR = 0.232



## Null $p$ -values

- All methods we will consider start from a collection of  $p$ -values  $p_1, \dots, p_m$ , one for each hypothesis tested. We call these  $p$ -values *raw* as they have not been corrected for multiple testing yet
- Assumptions on the  $p$ -values often involve only the  $p$ -values of true hypotheses. We denote these *null*  $p$ -values by

$$q_1, \dots, q_{m_0}$$

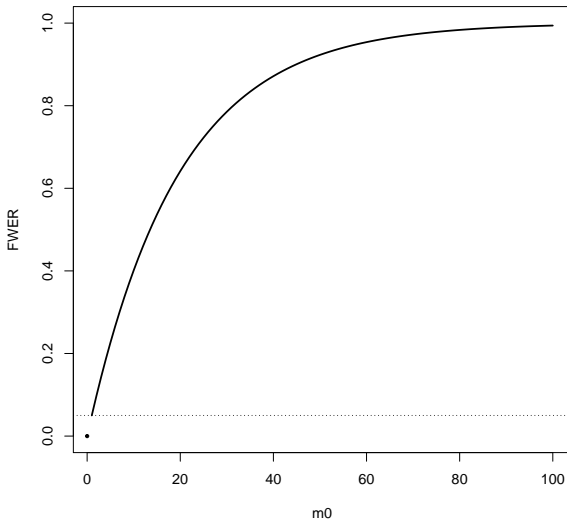
- Null  $p$ -values are assumed to be *valid* in the sense

$$P(q_i \leq u) \leq u$$

with equality when  $q_i \sim U(0, 1)$

If  $q_1, \dots, q_{m_0} \stackrel{i.i.d.}{\sim} U(0, 1)$ , then  $\mathcal{R} = \{H_i : p_i \leq 0.05\}$  has

$$\text{FWER} = 1 - (1 - 0.05)^{m_0}$$



## Expected number of type I errors

- The Per Family Error Rate (PFER) is the expected number of type I errors

$$\text{PFER} = \mathbb{E}(V)$$

- By Markov's inequality

$$P(V > 0) \leq \frac{\mathbb{E}(V)}{1}$$

we obtain

$$\text{FDR} \leq \text{FWER} \leq \text{PFER}$$

- If we consider

$$\mathcal{R} = \{H_i : p_i \leq c\}$$

then

$$V = \sum_{i=1}^{m_0} \mathbb{1}\{q_i \leq c\}$$

## Expected number of type I errors

- Suppose that  $q_i \sim U(0, 1)$

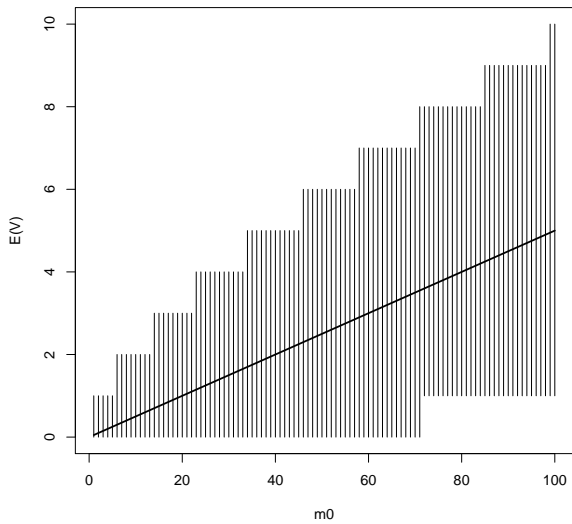
-

$$\mathbb{E}(V) = \sum_{i=1}^{m_0} \mathbb{E}(\mathbb{1}\{q_i \leq c\}) = m_0 c$$

-

$$\begin{aligned}\mathbb{V}\text{ar}(V) &= \sum_{i=1}^{m_0} \sum_{j=1}^{m_0} \mathbb{C}\text{ov}(\mathbb{1}\{q_i \leq c\} \mathbb{1}\{q_j \leq c\}) = m_0 c(1 - c) + \\ &+ 2 \sum_{i < j} \left[ \mathbb{P}(\mathbb{1}\{q_i \leq c, q_j \leq c\}) - \mathbb{P}(\mathbb{1}\{q_i \leq c\})\mathbb{P}(\mathbb{1}\{q_j \leq c\}) \right] \\ &= m_0 c(1 - c) + 2 \sum_{i < j} \left[ \mathbb{P}(\mathbb{1}\{q_i \leq c, q_j \leq c\}) - c^2 \right]\end{aligned}$$

where the first term represents the independence structure and last term the *overdispersion*



Methods for familywise error rate control

# Bonferroni method

## Theorem

*Bonferroni method rejects the hypotheses with  $p$ -value less than  $\alpha/m$*

$$\mathcal{R} = \{H_i : p_i \leq \frac{\alpha}{m}\}$$

*It controls the FWER at level  $\alpha$ .*

Proof.

$$\mathbb{P}\left(\bigcup_{i=1}^{m_0} \left\{q_i \leq \frac{\alpha}{m}\right\}\right) \leq \sum_{i=1}^{m_0} \mathbb{P}\left(q_i \leq \frac{\alpha}{m}\right) \leq m_0 \frac{\alpha}{m} \leq \alpha.$$



# Bonferroni conservativeness

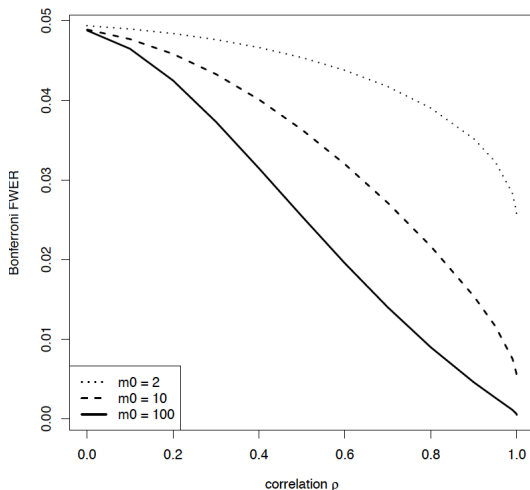
The two inequalities indicate in which cases the Bonferroni method can be *conservative*, i.e.  $\text{FWER} < \alpha$

- The right-hand one shows that Bonferroni controls the FWER at level  $\pi_0\alpha$ , where  $\pi_0 = m_0/m$ . If there are many false null hypotheses, Bonferroni will be conservative
- The left-hand inequality is due to Boole's inequality, i.e. for any collection of events  $E_1, \dots, E_k$ , we have  $P(\bigcup_{i=1}^k E_i) \leq \sum_{i=1}^k P(E_i)$ . This inequality is a strict one in all situations except the one in which all events  $\{q_i \leq \alpha/m\}$  are disjoint. With independent  $p$ -values, the conservativeness is present but very minor



# Positively correlated $p$ -values

Much more serious conservativeness can occur if  $p$ -values are positively correlated. Suppose that the correlation matrix is such that  $\{\Sigma\}_{ij} = \rho$  for  $i \neq j$



## Adjusted $p$ -values

- When testing a single hypothesis, we often do not only report whether a hypothesis was rejected, but also the corresponding  $p$ -value
- By definition, the  $p$ -value is the smallest chosen  $\alpha$ -level of the test at which the hypothesis would have been rejected
- The direct analogue of this in the context of multiple testing is the *adjusted*  $p$ -value, defined as the smallest  $\alpha$  level at which the multiple testing method would reject the hypothesis.
- For the Bonferroni procedure, this adjusted  $p$ -value is given by

$$\tilde{p}_i = \min(mp_i, 1)$$

where  $p_i$  is the raw  $p$ -value

# Sidak method

## Theorem

*Sidak method rejects*

$$\mathcal{R} = \{H_i \in \mathcal{H} : p_i \leq 1 - (1 - \alpha)^{1/m}\}$$

*If the null  $p$ -values  $q_1, \dots, q_{m_0} \stackrel{i.i.d.}{\sim} U(0, 1)$ , it controls the FWER at level  $\alpha$ .*

Proof.

$$P\left(\bigcup_{i=1}^{m_0} \{q_i \leq c\}\right) = 1 - \prod_{i=1}^{m_0} P(q_i > c) = 1 - (1 - c)^{m_0}$$
 which equals  $\alpha$  for  $c = 1 - (1 - \alpha)^{1/m_0}$ . Since we don't know  $m_0$ , we can use

$$1 - (1 - \alpha)^{1/m} \leq 1 - (1 - \alpha)^{1/m_0}$$



# Holm method

- Holm's method is a sequential variant of the Bonferroni method that always rejects at least as much as Bonferroni's method, and often a bit more, but still has valid FWER control under the same assumptions
- In the first step, all hypotheses with  $p$ -values at most  $\alpha/h_0$  are rejected, with  $h_0 = m$  just like in the Bonferroni method. Suppose this leaves  $h_1$  hypotheses unrejected. Then, in the next step, all hypotheses with  $p$ -values at most  $\alpha/h_1$  are rejected, which leaves  $h_2$  hypotheses unrejected, which are subsequently tested at level  $\alpha/h_2$ . This process is repeated until either all hypotheses are rejected, or until a step fails to result in any additional rejections

# Holm algorithm

Step 0 Begin by ordering the p-values in ascending order

$$p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$$

and let  $H_{(1)}, H_{(2)}, \dots, H_{(m)}$  be the corresponding hypotheses

Step 1 : If  $p_{(1)} \leq \alpha/m$  reject  $H_{(1)}$  and go to Step 2. Stop otherwise

Step 2 : If  $p_{(2)} \leq \alpha/(m-1)$  reject  $H_{(2)}$  and go to Step 3. Stop otherwise

...

Step  $i$  : If  $p_{(i)} \leq \alpha/(m-i+1)$  reject  $H_{(i)}$  and go to Step  $i+1$ . Stop otherwise

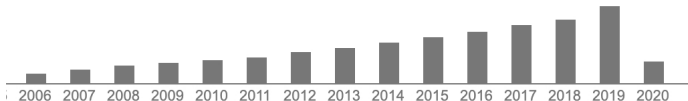
...

Step  $m$  : If  $p_{(m)} \leq \alpha$  reject  $H_{(m)}$

Methods for false discovery rate control

- If we are testing millions of hypotheses at once, and making few false discoveries is not the end of the world
- The concept of False Discovery Rate (FDR) has changed thinking about multiple testing quite radically, showing that FWER control is not only way to do of multiple testing, and stimulating the field of multiple testing enormously
- FDR was introduced by Benjamini and Hochberg in 1995, and currently has 63K citations. It is one of the most-cited research of all time

- *Controlling the false discovery rate: a practical and powerful approach to multiple testing*. Benjamini & Hochberg, JRSS-B (1995). Citations: 63K



- *Maximum likelihood from incomplete data via EM algorithm*. Dempster, Laird & Rubin, JRSS-B (1977). Citations: 60K
- *Nonparametric estimation from incomplete observations*. Kaplan & Meier, JASA (1958). Citations: 57K
- *Regression models and life-tables*. Cox, JRSS-B (1972). Citations: 51K



# Benjamini & Hochberg

1. Sort the  $p$ -values

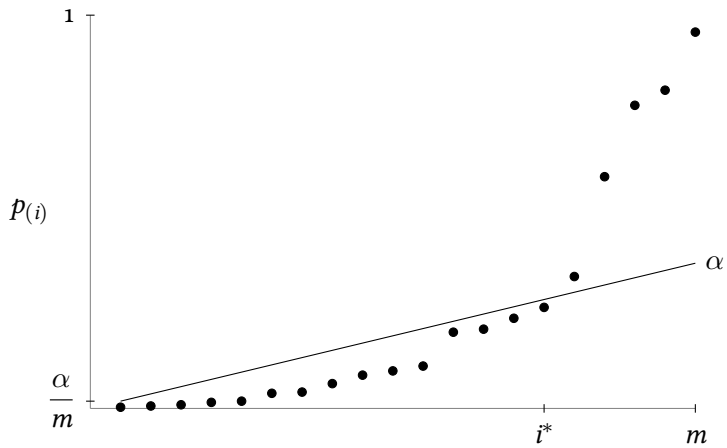
$$p_{(1)} \leq \dots \leq p_{(m)}$$

2. If  $p_{(i)} > \frac{i\alpha}{m}$  for all  $i$ , reject nothing, i.e.  $\mathcal{R} = \emptyset$
3. Otherwise, let

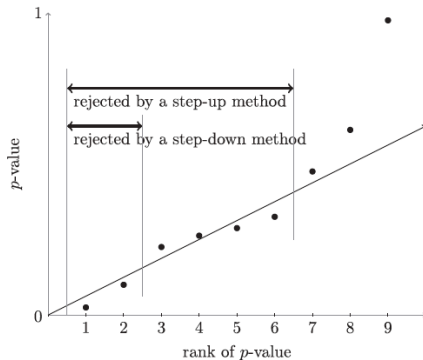
$$i^* = \max \left\{ i \in \{1, \dots, m\} : p_{(i)} \leq \frac{i\alpha}{m} \right\}$$

be the largest  $i$  for which  $p_{(i)} \leq \frac{i\alpha}{m}$

4. Reject all  $H_{(i)}$  with  $i \leq i^*$ , i.e.  $\mathcal{R} = \left\{ H_i : p_i \leq \frac{i^*\alpha}{m} \right\}$



# Step-up and step-down



- BH is a step-up method with threshold  $\frac{i\alpha}{m}$
- Holm is a step-down method with threshold  $\frac{\alpha}{m - i + 1}$

## Theorem

*For independent  $p$ -values  $p_1, \dots, p_m$  and null  $p$ -values*

*$q_1, \dots, q_{m_0} \stackrel{i.i.d.}{\sim} U(0, 1)$ , the FDR of the Benjamini-Hochberg method is exactly  $\pi_0 \alpha$ .*

## Proof (Candes and Barber version)

- The conclusion is obvious when  $m_0 = 0$ : assume  $m_0 \geq 1$
- Define  $V_i = \mathbb{1}\{H_i \text{ rejected}\}$  for each  $i \in T$  where  $T = \{i : H_i \in \mathcal{T}\}$ . We can express the FDP as

$$Q = \sum_{i \in T} \frac{V_i}{R \vee 1}$$

- We claim that

$$\mathbb{E}\left(\frac{V_i}{R \vee 1}\right) = \frac{\alpha}{m}, \quad i \in T$$

based on which we have

$$\text{FDR} = \mathbb{E}(Q) = \sum_{i \in T} \mathbb{E}\left(\frac{V_i}{R \vee 1}\right) = \sum_{i \in T} \frac{\alpha}{m} = \pi_0 \alpha$$

What remains for the proof is to show that the claim is true

## Proof - I

- When there are  $R = k$  rejections, then  $H_i$  is rejected if and only if  $p_i \leq (\alpha k)/m$ , and therefore, we have

$$V_i = \mathbb{1}\{p_i \leq (\alpha k)/m\}$$

- Suppose  $p_i \leq (\alpha k)/m$  (i.e.  $H_i$  is rejected). Let us take  $p_i$  and set its value to 0, and denote the new number of rejections by  $R(p_i \downarrow 0)$ . This new number of rejections is exactly  $R$ , because we have only reordering the first  $k$   $p$ -values, all of which remain below the threshold  $(\alpha k)/m$ . On the other hand, if  $p_i > (\alpha k)/m$ , then we do not reject  $H_i$ , and so  $V_i = 0$ . Therefore we have

$$V_i \mathbb{1}\{R = k\} = V_i \mathbb{1}\{R(p_i \downarrow 0) = k\}$$

## Proof - II

Combining the observations above and taking the expectation conditional on all  $p$ -values except for  $p_i$ , i.e.

$\mathcal{F}_i = \{p_1, \dots, p_{i-1}, p_{i+1}, \dots, p_m\}$ , we have

$$\begin{aligned}\mathbb{E}\left(\frac{V_i}{R \vee 1} | \mathcal{F}_i\right) &= \sum_{k=1}^m \frac{\mathbb{E}(\mathbb{1}\{p_i \leq (\alpha k)/m\} \mathbb{1}\{R(p_i \downarrow 0) = k\} | \mathcal{F}_i)}{k} \\ &= \sum_{k=1}^m \frac{\mathbb{1}\{R(p_i \downarrow 0) = k\} (\alpha k)/m}{k}\end{aligned}$$

where the second equality holds because knowing  $\mathcal{F}_i$  and  $p_i = 0$  makes  $\mathbb{1}\{R(p_i \downarrow 0)\}$  deterministic, and the fact that  $p_i \sim U(0, 1)$  and the  $p$ -values  $p_1, \dots, p_m$  are independent

## Proof - III

- Next, we have

$$\mathbb{E}\left(\frac{V_i}{R \vee 1} \mid \mathcal{F}_i\right) = \frac{\alpha}{m} \sum_{k=1}^m \mathbb{1}\{R(p_i \downarrow 0) = k\} = \frac{\alpha}{m}$$

after noticing that  $\sum_{k=1}^m \mathbb{1}\{R(p_i \downarrow 0) = k\} = 1$

- Since we have set  $p_i$  to 0, we must make at least one rejection - we will always reject  $H_i$ . Therefore  $R(p_i \downarrow 0) \geq 1$ , and  $R(p_i \downarrow 0)$  must take a value between 1 and  $m$
- The tower property verifies that

$$\text{FDR} = \sum_{i \in T} \mathbb{E}\left(\frac{V_i}{R \vee 1}\right) = \sum_{i \in T} \mathbb{E}\left[\mathbb{E}\left(\frac{V_i}{R \vee 1} \mid \mathcal{F}_i\right)\right] = \sum_{i \in T} \frac{\alpha}{m} = \pi_0 \alpha$$

□



BH is valid under the more general assumption of *positive regression dependence on a subset* (PDS). One case under which the PDS condition holds is one-sided test statistics that are jointly normally distributed, if all correlations between test statistics are positive.

### Theorem

*For  $p$ -values satisfying the PDS assumption, the Benjamini-Hochberg procedure controls the FDR at level  $\pi_0\alpha$ .*

# Adaptive Benjamini-Hochberg

- The Benjamini & Hochberg method, like Bonferroni, controls its error rate at level  $\pi_0\alpha$ , rather than at  $\alpha$ . This suggests the possibility of an alternative, more powerful Benjamini & Hochberg procedure that uses critical values

$$\frac{i\alpha}{\hat{\pi}_0 m}$$

rather than  $(i\alpha)/m$  if a good estimate  $\hat{\pi}_0$  of the proportion of true hypotheses  $\pi_0$  would be available

- Such procedures are called *adaptive* procedures, and many have been proposed on the basis of various estimates of  $\pi_0$
- A problem with the adaptive approach, however, is that estimates of  $\pi_0$  can have high variance, especially if  $p$ -values are strongly correlated. Naive plug-in procedures, in which this variance is not taken into account, will therefore generally not have FDR control

## $\pi_0$ estimator

$$\hat{m}_0(\lambda) = \frac{\sum_{i=1}^m \mathbb{1}\{p_i > \lambda\}}{1 - \lambda}$$

If null  $p$ -values have marginal  $U(0, 1)$  distribution, a proportion  $1 - \lambda$  is expected to be above  $\lambda$ :

$$\mathbb{E}\left(\sum_{i=1}^m \mathbb{1}\{p_i > \lambda\}\right) \geq \mathbb{E}\left(\sum_{i \in T} \mathbb{1}\{q_i > \lambda\}\right) = m_0(1 - \lambda)$$

thus  $\mathbb{E}(\hat{m}_0) \geq m_0$ .

$\hat{\pi}_0 = \frac{\hat{m}_0}{m}$  is a conservative estimator of  $\pi_0$ , i.e.  $\mathbb{E}(\hat{\pi}_0) \geq \pi_0$

# Storey method

1. Choose  $\lambda \in (0, 1)$ . Estimate  $\pi_0$  by

$$\hat{\pi}_0(\lambda) = \frac{\sum_{i=1}^m \mathbb{1}\{p_i > \lambda\} + 1}{(1 - \lambda)m}$$

2. Perform Benjamini-Hochberg procedure at level

$$\frac{\alpha}{\hat{\pi}_0(\lambda)}$$

- The addition of 1 to the numerator makes sure that  $1/\hat{\pi}_0$  is always well-defined (but it may happen  $\hat{\pi}_0 > 1$ )
- The value of  $\lambda$  is typically  $1/2$ , although  $\lambda = \alpha$  has also been advocated
- Storey method controls FDR under independence of  $p$ -values but generally not under positive dependence

# General dependence

1. It has been proved that the Benjamini-Hochberg procedure controls FDR at level  $\alpha$  also under the PDS assumption
2. If the PDS assumption is not valid, an alternative is the procedure of Benjamini & Yekutieli, which is valid under general dependence

## Example

Consider  $m = m_0 = 2$ . The two null  $p$ -values  $q_1$  and  $q_2$  are marginally  $U(0, 1)$ , but the joint distribution of  $(q_1, q_2)$  is piecewise constant with density

- $1/(1 - \alpha)$  in areas b
- $2/\alpha$  in area c
- $b(1 - b\alpha/2)$  in area a
- 0 in gray areas

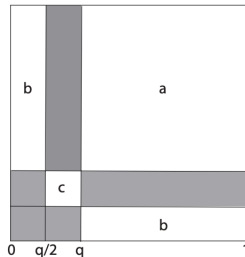


Figure 2: A piecewise constant joint distribution.

# Benjamini & Yekutieli

$$\text{FDR}[\text{BH}(\alpha)] = P(I) + P(II) + P(III)$$

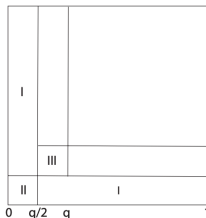


Figure 1: The BHq rejection region.

- In our example,  $\text{FDR}[\text{BH}(\alpha)] = 3\alpha/2$
- In general,  $\exists$  a worst-case joint distribution of  $p$ -values such that  $\text{FDR}[\text{BH}(\alpha)] = \alpha H_m$  with  $H_m = \sum_{j=1}^m \frac{1}{j}$  harmonic number