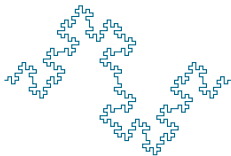


High-Throughput Sequencing Course

Multiple Testing

Biostatistics and Bioinformatics



Summer 2018

INTRODUCTION

- ▶ You have previously considered the significance of a single gene

INTRODUCTION

- ▶ You have previously considered the significance of a single gene
- ▶ The analysis of high-dimensional data is concerned with assessing the significance of multiple loci/genes
 - ▶ Microarray : 20,000-50,000 probe sets
 - ▶ GWAS: 500,000-5,000,000 typed SNPs
 - ▶ RNA-Seq: 22,000 genes (humans)

INTRODUCTION

- ▶ You have previously considered the significance of a single gene
- ▶ The analysis of high-dimensional data is concerned with assessing the significance of multiple loci/genes
 - ▶ Microarray : 20,000-50,000 probe sets
 - ▶ GWAS: 500,000-5,000,000 typed SNPs
 - ▶ RNA-Seq: 22,000 genes (humans)
- ▶ This leads to the *Multiple Testing* problem

INTRODUCTION

- ▶ You have previously considered the significance of a single gene
- ▶ The analysis of high-dimensional data is concerned with assessing the significance of multiple loci/genes
 - ▶ Microarray : 20,000-50,000 probe sets
 - ▶ GWAS: 500,000-5,000,000 typed SNPs
 - ▶ RNA-Seq: 22,000 genes (humans)
- ▶ This leads to the *Multiple Testing* problem
- ▶ Before we address this problem, let's quickly review the single hypothesis case

INTRODUCTION: TYPE I AND II ERRORS

Recall that in hypothesis testing with a single hypothesis (gene), errors can be classified as:

- ▶ Type I error - rejecting H_0 when it is true
- ▶ Type II error - not rejecting H_0 when it is false

INTRODUCTION: TYPE I AND II ERRORS

Recall that in hypothesis testing with a single hypothesis (gene), errors can be classified as:

- ▶ Type I error - rejecting H_0 when it is true
- ▶ Type II error - not rejecting H_0 when it is false

We can define probabilities associated with each of these errors:

$$\alpha = Pr(\text{reject } H_0 | H_0 \text{ true})$$

$$\beta = Pr(\text{do not reject } H_0 | H_0 \text{ false})$$

INTRODUCTION: TYPE I AND II ERRORS

Decision	H_0 True	H_0 False
Do not reject H_0	$1 - \alpha$	β
Reject H_0	α	$1 - \beta$

INTRODUCTION: TYPE I AND II ERRORS

- ▶ Always a compromise between type I and type II error

INTRODUCTION: TYPE I AND II ERRORS

- ▶ Always a compromise between type I and type II error
- ▶ An α -level test has probability α of rejecting H_0 when H_0 is true

INTRODUCTION: TYPE I AND II ERRORS

- ▶ Always a compromise between type I and type II error
- ▶ An α -level test has probability α of rejecting H_0 when H_0 is true
- ▶ Usually try to use the most powerful (smallest β) test for a given α -level

INTRODUCTION: TYPE I AND II ERRORS

- ▶ Always a compromise between type I and type II error
- ▶ An α -level test has probability α of rejecting H_0 when H_0 is true
- ▶ Usually try to use the most powerful (smallest β) test for a given α -level \rightarrow *control* type I error

INTRODUCTION: MULTIPLE TESTS

- ▶ Single-hypothesis type-I-error-control proves inadequate when experiment is characterized by a collection of hypotheses

INTRODUCTION: MULTIPLE TESTS

- ▶ Single-hypothesis type-I-error-control proves inadequate when experiment is characterized by a collection of hypotheses
 - ← Genomics

INTRODUCTION: MULTIPLE TESTS

- ▶ Single-hypothesis type-I-error-control proves inadequate when experiment is characterized by a collection of hypotheses
 ← Genomics
- ▶ Why?

INTRODUCTION: MULTIPLE TESTS

- ▶ Assume we have m hypotheses we wish to test as part of a given experiment

INTRODUCTION: MULTIPLE TESTS

- ▶ Assume we have m hypotheses we wish to test as part of a given experiment
 - ▶ These hypotheses could correspond to m genes that we are investigating for differential expression between two groups

INTRODUCTION: MULTIPLE TESTS

- ▶ Assume we have m hypotheses we wish to test as part of a given experiment
 - ▶ These hypotheses could correspond to m genes that we are investigating for differential expression between two groups
- ▶ Assume that each hypothesis is tested using a α -level test

INTRODUCTION: MULTIPLE TESTS

- ▶ Assume we have m hypotheses we wish to test as part of a given experiment
 - ▶ These hypotheses could correspond to m genes that we are investigating for differential expression between two groups
- ▶ Assume that each hypothesis is tested using a α -level test
- ▶ Assume that the tests are *INDEPENDENT* and that the null is true for each of the hypothesis

INTRODUCTION: MULTIPLE TESTS

What is the probability of rejecting a null hypothesis even though all hypotheses are actually null?

INTRODUCTION: MULTIPLE TESTS

What is the probability of rejecting a null hypothesis even though all hypotheses are actually null?

- For a given test, we have an α chance of rejecting its null

INTRODUCTION: MULTIPLE TESTS

What is the probability of rejecting a null hypothesis even though all hypotheses are actually null?

- ▶ For a given test, we have an α chance of rejecting its null
- ▶ Therefore we have a $1 - \alpha$ chance of not rejecting

INTRODUCTION: MULTIPLE TESTS

What is the probability of rejecting a null hypothesis even though all hypotheses are actually null?

- ▶ For a given test, we have an α chance of rejecting its null
- ▶ Therefore we have a $1 - \alpha$ chance of not rejecting
- ▶ We have m independent α -level tests, so the chance of not rejecting any of them is:

$$(1 - \alpha)^m$$

INTRODUCTION: MULTIPLE TESTS

What is the probability of rejecting a null hypothesis even though all hypotheses are actually null?

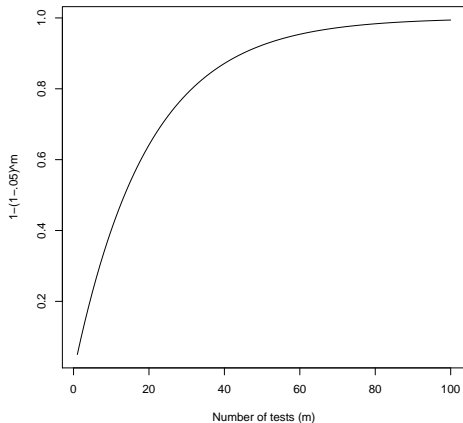
- ▶ For a given test, we have an α chance of rejecting its null
- ▶ Therefore we have a $1 - \alpha$ chance of not rejecting
- ▶ We have m independent α -level tests, so the chance of not rejecting any of them is:

$$(1 - \alpha)^m$$

- ▶ Therefore, the chance of rejecting at least one hypothesis is:

$$1 - (1 - \alpha)^m$$

INTRODUCTION: MULTIPLE TESTS



INTRODUCTION: NOTATION

- ▶ Gene j (among the m genes) is either associated with the outcome or not
- ▶ The truth is unknown to us
- ▶ The null hypothesis for gene j is denoted by H_j
- ▶ H_j : gene j is not associated with the outcome of interest
- ▶ The alternative hypothesis is denoted by \bar{H}_j
- ▶ \bar{H}_j : gene j is associated with the outcome of interest
- ▶ H_j and \bar{H}_j are called *marginal* or *local* hypotheses

INTRODUCTION: MARGINAL AND GLOBAL HYPOTHESES

- ▶ H_j and \bar{H}_j are called *marginal* or *local* hypotheses
- ▶ A global null hypothesis: None of the m genes is associated with the outcome
- ▶ A global alternative: At least one of the m genes is associated with the outcome
- ▶ Using notation
 - ▶ Global Null: $\mathbb{H}_0 : H_1 \text{ and } H_2 \text{ and } \dots H_m$
 - ▶ Global Alternative: $\mathbb{H}_1 : \bar{H}_1 \text{ or } \bar{H}_2 \text{ or } \dots \bar{H}_m$

INTRODUCTION: UNADJUSTED VS ADJUSTED P -VALUES

- ▶ Suppose that we only test a single gene, say gene j , among the m genes
- ▶ Let p_j (lower case p) denote P -value corresponding to H_j
- ▶ p_j is called the *marginal* or *unadjusted* P -value
- ▶ If m hypotheses are tested, inference on H_j on the basis of p_j is inappropriate
- ▶ The P -value for H_j has to account for testing the other $m - 1$ hypotheses
- ▶ We will denote the *adjusted* P -value by P_j (upper case P)
- ▶ When testing multiple genes, using the marginal P -value is inappropriate
- ▶ Why?

INTRODUCTION: MORE NOTATION

- ▶ Suppose that gene j is not associated with the outcome of interest (H_j is true)
 - ▶ Then
 - ▶ Decision rule rejects \rightarrow False-Positive (FP)
 - ▶ Decision rule fails to reject \rightarrow True-Negative (TN)
- ▶ Suppose that gene j is associated with the outcome of interest (H_j is false)
 - ▶ Decision rule rejects \rightarrow True-Positive (TP)
 - ▶ Decision rule fails to reject \rightarrow False-Negative (FN)

INTRODUCTION: SUMMARIZING A MULTIPLE TESTING PROCEDURE

- The results from any multiple testing procedure can be summarized as the following table

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- Notation:
 - m : Number of tests, m_0, m_1 number of null/true genes
 - R : Number of genes rejected according to the decision rule
 - A : Number of genes accepted according to the decision rule
 - R_0/R_1 number of TN/FP
 - A_0/A_1 number of FN/TP

INTRODUCTION: EXAMPLE

- Results from an analysis based on $m = 10$ genes:

##	gene	truth	pvalue
## 1	gene1	0	0.29070
## 2	gene2	1	0.61630
## 3	gene3	1	0.00320
## 4	gene4	0	0.01641
## 5	gene5	0	0.25150
## 6	gene6	0	0.58450
## 7	gene7	0	0.22890
## 8	gene8	1	0.12630
## 9	gene9	0	0.26080
## 10	gene10	0	0.04980

- Investigator decides to use following decision rule: Any gene with a corresponding unadjusted P -value of less than 0.05 will be rejected.
- Reject H_j if $p_j < 0.05$ or accept H_j otherwise

EXERCISE: FILL IN THE 2X2 TABLE

	Accept	Reject	Total
Truth Null	$A_0 = ?$	$R_0 = ?$	$m_0 = ?$
Alt.	$A_1 = ?$	$R_1 = ?$	$m_1 = ?$
	$A = ?$	$R = ?$	$m = ?$

EXAMPLE: FILL IN THE 2X2 TABLE

	Accept	Reject	Total
Truth Null	$A_0 = 5$	$R_0 = 2$	$m_0 = 7$
Alt.	$A_1 = 2$	$R_1 = 1$	$m_1 = 3$
	$A = 7$	$R = 3$	$m = 10$

- ▶ $m_0 = 7$ and $m_1 = 3$
- ▶ $R = 3$ will be rejected based on the decision rule
- ▶ Consequently $A = m - R = 7$ will be accepted
- ▶ $R_0 = 2, R_1 = 1, A_0 = 5$ and $A_1 = 2$

THE TRUTH

- What know or observe is this

```
##      gene  pvalue
## 1  gene1 0.29070
## 2  gene2 0.61630
## 3  gene3 0.00320
## 4  gene4 0.01641
## 5  gene5 0.25150
## 6  gene6 0.58450
## 7  gene7 0.22890
## 8  gene8 0.12630
## 9  gene9 0.26080
## 10 gene10 0.04980
```

- and not (truth column is not known to us):

```
dat
##      gene truth  pvalue
## 1  gene1     0 0.29070
## 2  gene2     1 0.61630
## 3  gene3     1 0.00320
## 4  gene4     0 0.01641
## 5  gene5     0 0.25150
## 6  gene6     0 0.58450
## 7  gene7     0 0.22890
## 8  gene8     1 0.12630
## 9  gene9     0 0.26080
## 10 gene10    0 0.04980
```

EXAMPLE: FILL IN THE 2X2 TABLE (BASED ON WHAT WE OBSERVE)

- We can only fill in the bottom row of the table

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	$A = 7$	$R = 3$	$m = 10$

- The remaining quantities are fixed unknown quantities or unobservable random variables.

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- m is a known constant

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- ▶ m is a known constant
- ▶ m_0 and m_1 are unknown constants

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- ▶ m is a known constant
- ▶ m_0 and m_1 are unknown constants
- ▶ R and A are determined on the basis of applying the decision rule to the data

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- ▶ m is a known constant
- ▶ m_0 and m_1 are unknown constants
- ▶ R and A are determined on the basis of applying the decision rule to the data
- ▶ They are *observable* random quantities

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- ▶ m is a known constant
- ▶ m_0 and m_1 are unknown constants
- ▶ R and A are determined on the basis of applying the decision rule to the data
- ▶ They are *observable* random quantities
- ▶ The true states of the genes of the genes are unknown

COMMENTS

	Accept	Reject	Total
Truth Null	A_0	R_0	m_0
Alt.	A_1	R_1	m_1
	A	R	m

- ▶ m is a known constant
- ▶ m_0 and m_1 are unknown constants
- ▶ R and A are determined on the basis of applying the decision rule to the data
- ▶ They are *observable* random quantities
- ▶ The true states of the genes of the genes are unknown
- ▶ A_0, A_1, R_0 and R_1 are *unobservable* random quantities

INTRODUCTION: MULTIPLE TESTING PROBLEM

- ▶ Control error rate(s) in multiple testing context

INTRODUCTION: MULTIPLE TESTING PROBLEM

- ▶ Control error rate(s) in multiple testing context
- ▶ Multiple testing methods are designed to control a particular error rate

INTRODUCTION: MULTIPLE TESTING PROBLEM

- ▶ Control error rate(s) in multiple testing context
- ▶ Multiple testing methods are designed to control a particular error rate
- ▶ Multiple error rates exist

INTRODUCTION: MULTIPLE TESTING PROBLEM

- ▶ Control error rate(s) in multiple testing context
- ▶ Multiple testing methods are designed to control a particular error rate
- ▶ Multiple error rates exist \rightarrow need to choose error rate to control and then method to control it

INTRODUCTION: ERROR RATES

- ▶ **Family-wise error rate (FWER)**: the probability of at least one type I error

INTRODUCTION: ERROR RATES

- ▶ **Family-wise error rate (FWER)**: the probability of at least one type I error
- ▶ **False discovery rate (FDR)**: the expected proportion of type I errors among the rejected hypotheses.

FAMILY-WISE ERROR RATE (FWER)

- Probability of committing at least one false-rejection (among m) given that *all* genes are null

FAMILY-WISE ERROR RATE (FWER)

- ▶ Probability of committing at least one false-rejection (among m) given that *all* genes are null
- ▶ $\text{FWER} = P(R \geq 1 | m = m_0)$

FAMILY-WISE ERROR RATE (FWER)

- ▶ Probability of committing at least one false-rejection (among m) given that *all* genes are null
- ▶ $\text{FWER} = P(R \geq 1 | m = m_0)$
- ▶ Note that when $m = 1$ (single gene), this definition is identical to the type I error we have previously considered

CONTROLLING FWER: SIDAK'S METHOD

Recall that we showed that with m independent α -level tests:

$$\text{FWER} = 1 - (1 - \alpha)^m$$

CONTROLLING FWER: SIDAK'S METHOD

Recall that we showed that with m independent α -level tests:

$$\text{FWER} = 1 - (1 - \alpha)^m$$

Therefore,

$$\begin{aligned} 1 - \text{FWER} &= (1 - \alpha)^m \\ (1 - \text{FWER})^{1/m} &= 1 - \alpha \\ 1 - (1 - \text{FWER})^{1/m} &= \alpha \end{aligned}$$

CONTROLLING FWER: SIDAK'S METHOD

Recall that we showed that with m independent α -level tests:

$$\text{FWER} = 1 - (1 - \alpha)^m$$

Therefore,

$$\begin{aligned} 1 - \text{FWER} &= (1 - \alpha)^m \\ (1 - \text{FWER})^{1/m} &= 1 - \alpha \\ 1 - (1 - \text{FWER})^{1/m} &= \alpha \end{aligned}$$

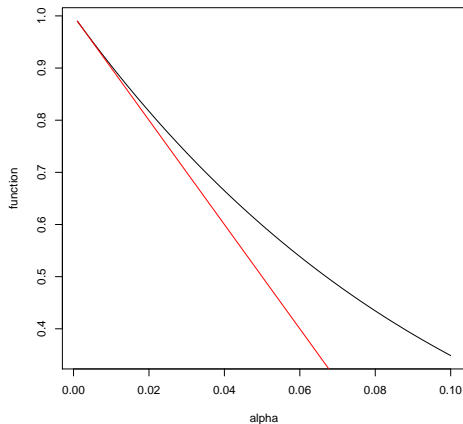
This suggests that we can control FWER by choosing α for each individual test to be $1 - (1 - \text{FWER})^{1/m}$

CONTROLLING FWER: BONFERRONI

Note that $(1 - x)^m \approx 1 - mx$ (for x close to zero)

CONTROLLING FWER: BONFERRONI

Note that $(1 - x)^m \approx 1 - mx$ (for x close to zero)



CONTROLLING FWER: BONFERRONI

Note that $(1 - x)^m \approx 1 - mx$ (for x close to zero)

Therefore

$$1 - \text{FWER} = (1 - \alpha)^m$$

$$1 - \text{FWER} \approx 1 - m\alpha$$

$$\text{FWER} \approx m\alpha$$

CONTROLLING FWER: BONFERRONI

Note that $(1 - x)^m \approx 1 - mx$ (for x close to zero)

Therefore

$$1 - \text{FWER} = (1 - \alpha)^m$$

$$1 - \text{FWER} \approx 1 - m\alpha$$

$$\text{FWER} \approx m\alpha$$

This suggests that we can control FWER by choosing α for each individual test to be $\frac{\text{FWER}}{m}$

CONTROLLING FWER: BONFERRONI

- ▶ The Bonferroni adjusted *P-value* is defined as

$$P_j = m \times p_j$$

- ▶ Technical note: P_j as defined above could be larger than 1 so a more technically rigorous definition is

$$P_j = \min\{m \times p_j, 1\}$$

- ▶ In other words, if $m \times p_j$ is larger than 1, then truncate P_j at 1.

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches
- ▶ Some power can be gained by sequential methods

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches
- ▶ Some power can be gained by sequential methods
- ▶ The simplest sequential method is Holm's method:

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches
- ▶ Some power can be gained by sequential methods
- ▶ The simplest sequential method is Holm's method:
 - ▶ Order the unadjusted P -values such that $p_1 \leq p_2 \leq \dots \leq p_m$

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches
- ▶ Some power can be gained by sequential methods
- ▶ The simplest sequential method is Holm's method:
 - ▶ Order the unadjusted P -values such that $p_1 \leq p_2 \leq \dots \leq p_m$
 - ▶ To control FWER, the step-down Holm adjusted P -values are

$$P_j = \min\{(m - j + 1) \times p_j, 1\}$$

CONTROLLING FWER: HOLM'S

- ▶ Sidak and Bonferroni are one-step approaches
- ▶ Some power can be gained by sequential methods
- ▶ The simplest sequential method is Holm's method:
 - ▶ Order the unadjusted P -values such that $p_1 \leq p_2 \leq \dots \leq p_m$
 - ▶ To control FWER, the step-down Holm adjusted P -values are

$$P_j = \min\{(m - j + 1) \times p_j, 1\}$$

- ▶ Note that every unadjusted P -value is not multiplied by same factor

CONTROLLING FWER: PERMUTATION

- ▶ These methods work well when tests are independent

CONTROLLING FWER: PERMUTATION

- ▶ These methods work well when tests are independent (not generally the case in genomics experiments)

CONTROLLING FWER: PERMUTATION

- ▶ These methods work well when tests are independent (not generally the case in genomics experiments)
- ▶ When tests are correlated, these methods are conservative

CONTROLLING FWER: PERMUTATION

- ▶ These methods work well when tests are independent (not generally the case in genomics experiments)
- ▶ When tests are correlated, these methods are conservative
- ▶ Permutation approaches are useful in this context

CONTROLLING FWER: PERMUTATION

Assume we are interested in assessing differential expression between 2 groups :

1. Compute minimum unadjusted P -value for all genes from the observed data (call it p_1)
2. Randomly permute the group labels
 - ▶ *Breaks* relationship between group and expression
 - ▶ Reflects sample from global null hypothesis
3. Compute minimum P -value from data set generated in 2 (call it p_1^1)
4. Repeat 2 and 3 B times to get $p_1^1, p_1^2, \dots, p_1^B$
5. Compute the proportion of $p_1^1, p_1^2, \dots, p_1^B$ that are $\leq p_1$
6. This proportion is the permutation adjusted P_1

CONTROLLING FWER: PERMUTATION

- ▶ Can get P_2, \dots, P_m in a similar way (a Holm's-like permutation procedure)

CONTROLLING FWER: PERMUTATION

- ▶ Can get P_2, \dots, P_m in a similar way (a Holm's-like permutation procedure)
- ▶ Correlation among the genes is accounted for

FALSE DISCOVERY RATE (FDR)

- ▶ Consider the quantity $\frac{R_0}{R}$
- ▶ This is the proportion of false discoveries among the genes rejected
- ▶ This is an *unobservable* random quantity (R_0 is not observable)
- ▶ In the FDR framework is based on controlling the *expected* value of this ratio
- ▶ $\text{FDR} \equiv E[\frac{R_0}{R}]$
 - ▶ Expectation is set to zero if $R = 0$,
therefore $\text{FDR} = E[\frac{R_0}{R} | R > 0] Pr(R > 0)$
- ▶ Note that when $m_0 = m$ (none of the genes are true),
FWER=FDR

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$
3. Reject null hypotheses affiliated with p_1, p_2, \dots, p_j

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$
3. Reject null hypotheses affiliated with p_1, p_2, \dots, p_j

This procedure will control FDR at $\alpha m_0/m$ when the tests are independent and continuous

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$
3. Reject null hypotheses affiliated with p_1, p_2, \dots, p_j

This procedure will control FDR at $\alpha m_0/m$ when the tests are independent and continuous

Benjamini and Yekutieli (2001) proposed a modification of the BH procedure that always controls FDR (no larger than $\alpha m_0/m$)

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$
3. Reject null hypotheses affiliated with p_1, p_2, \dots, p_j

This procedure will control FDR at $\alpha m_0/m$ when the tests are independent and continuous

Benjamini and Yekutieli (2001) proposed a modification of the BH procedure that always controls FDR (no larger than $\alpha m_0/m$)

← can be quite conservative

CONTROLLING FDR: BENJAMINI AND HOCHBERG

1. Order the unadjusted P -values $p_1 \leq p_2 \leq \dots \leq p_m$
2. Find the largest j such that $p_j \leq j\alpha/m$
3. Reject null hypotheses affiliated with p_1, p_2, \dots, p_j

This procedure will control FDR at $\alpha m_0/m$ when the tests are independent and continuous

Benjamini and Yekutieli (2001) proposed a modification of the BH procedure that always controls FDR (no larger than $\alpha m_0/m$)

← can be quite conservative

Note that when $m_0 = m$ (i.e., all hypotheses are null), these procedures maintain FWER at α

CONTROLLING pFDR: Q-VALUES

$$\text{FDR} = E\left[\frac{R_0}{R} | R > 0\right] Pr(R > 0)$$

CONTROLLING pFDR: Q-VALUES

$$\text{FDR} = E\left[\frac{R_0}{R} \mid R > 0\right] \Pr(R > 0)$$

$$\text{pFDR} = E\left[\frac{R_0}{R} \mid R > 0\right] \leftarrow \textit{positive FDR}$$

CONTROLLING pFDR: Q-VALUES

$$\text{FDR} = E\left[\frac{R_0}{R} | R > 0\right] Pr(R > 0)$$

$$\text{pFDR} = E\left[\frac{R_0}{R} | R > 0\right] \leftarrow \text{positive FDR}$$

- Since $Pr(R > 0)$ is often ~ 1 in most genomics experiments, FDR and pFDR are ver similar

CONTROLLING pFDR: Q-VALUES

- ▶ q-value is the minimum pFDR for which the affiliated hypothesis is rejected

CONTROLLING pFDR: Q-VALUES

- ▶ q-value is the minimum pFDR for which the affiliated hypothesis is rejected
- ▶ q-value can be interpreted as the expected proportion of false positives incurred when calling that test significant

GENOME-WIDE SIGNIFICANCE

- ▶ In GWAS papers, $\alpha = 5 \times 10^{-8}$ is typically considered the threshold for genome-wide significance
- ▶ It is based on a Bonferroni correction: If you consider testing $m = 1,000,000$ SNPs at the FWER level of 0.05, then each SNP should be tested at the

$$\alpha = \frac{0.05}{1,000,000} = 5 \times 10^{-8},$$

level

- ▶ Suppose that the unadjusted P -value for a SNP is 5×10^{-7}
- ▶ Is this "reaching" genome-wide significance?
- ▶ The term "suggestive" is also used

”REACHING” GENOME-WIDE SIGNIFICANCE

- ▶ Suppose that your $m = 1,000,000$ SNPs are independent
- ▶ The adjusted P -value is

$$P = 5 \times 10^{-7} \times m = 5 \times 10^{-7} \times 10^6 = 0.5,$$

- ▶ This is off by an order of magnitude ($0.5 = 0.05 \times 10$)
- ▶ It is not ”reaching”
- ▶ Note: Due to linkage disequilibrium among SNPs the adjusted P -value is likely to be smaller than 0.5
- ▶ The point is that while 5×10^{-7} is small number, it may not be small enough when testing a large number of hypotheses

CONCLUSIONS

- ▶ Multiple testing *must* be accounted for when testing for associations in the context of high-dimensional data
- ▶ FWER and FDR are the two common frameworks for quantifying error
- ▶ Error rate estimates can be used to compute 'adjusted' p-values
- ▶ Resampling-based methods can increase power in controlling error when sample sizes are sufficient for their use.
- ▶ When large-scale patterns of differential expression are observed, it is important to consider if such effects are biologically reasonable, and if technical factors can be attributed to the variation.