Multiple testing: Intro & FWER¹

Mark van de Wiel mark.vdwiel@vumc.nl

Department of Epidemiology & Biostatistics, VUmc, Amsterdam

¹Some slides courtesy of Jelle Goeman

VU medisch centr

Practical notes

Study material

- 1. These slides (3 sets for multiple testing)
- "Tutorial in biostatistics: multiple hypothesis testing in genomics" by Goeman & Solari. Sections 4.3 & 4.4 are not compulsory study material.
- 3. PDF handouts: technical details.

References

- 1. [hand]: refers to PDF handout
- 2. [tut]: refers to tutorial
- 3. [exer]: refers to exercise

 Introduction Multiple Testing & Family-wise error rate (FWER)

2. False Discovery Rate (FDR) + Bayesian perspective on multiple testing

- 1. Introduction Multiple testing
- 2. FWER
- 3. Bonferroni & Holm
- 4. Permutations
- 5. R-code
- 6. Summary

Introduction Multiple Testing



Source: http://flowingdata.com/2011/04/08/statistical-significance-on-xkcd/







Differential expression

Data types Samples Continuous 3.4 4.2 8.9 7.6 2.1 8.9 5.3 6.7 (log-scale) Counts 10 4 23 42 17 $\left(\right)$ \mathbf{O} 0 AB AB AA AA AA AA AB BB Nominal 2 0 1 0 \mathbf{O} ()Group 1 Group 2

Which features are differential between groups?



Differential expression

Comparative microarray experiment

Given: matrix of (normalized) expression values



Differential expression

Top 5 genes out of 20,000

Gene	p-value
OCIAD2	5.5e-6
NEK3	6.7e-6
TAF5	7.1e-6
FOXD4L6	7.5e-6
ADIG	8.8e-6
:	:

Small p-value?

- Getting a p-value as small as 5.5e-6 is unlikely
- But is it also small if we admit that we tried 20,000 times?
- Can we reliably state that OCIAD2 is differentially expressed?
- What about NEK3?

How does multiple testing differ from single testing?

1. p-values not sufficient to control false positive rate: $p \le 0.05$ implies too many false positives

 Different error control desirable when # tests large: control proportion of false positives instead of P(at least 1 false positive)

Opportunity to learn from other features:
 When using a t-test: estimate the s.d. from all features (other lecture: limma)

- Test m null-hyposthesis: H_{01} , H_{02} , ..., H_{0m}
- Renders m p-values: p₁ ,..., p_m
- Property of p-values: $P_0(p_i \le \alpha) = \alpha$ (or $\le \alpha$) [proof: see hand-out]
- V: Number of false positives
- Under independence [tut]:

$$P(V>0) = 1 - (1 - \alpha)^{m_0}$$

Too many false positives



number of hypothesis tests at level 0.05

	True Null	False Null	
Rejected	V	U	R
Non-rejected	m ₀ - V	m ₁ - U	m-R
	m ₀	m ₁	m

Family-wise error rate (FWER)

• FWER: probability of making at least one false rejection, so

```
\mathsf{FWER} = \mathsf{P}(\mathsf{V} > 0)
```

• If m = 1 (one test), FWER reduces to type I error rate

Bonferroni's solution

- 1. Multiply p-values with m, the number of tests
- 2. Reject null-hypothesis for feature i when

$$p_i^{\text{bonf}} = m \ge p_i \le \alpha \text{ (e.g. } \alpha = 0.05)$$

Under arbitrary dependency structure:

Reject H_{0i} when $p_i^{bonf} \le \alpha$ implies FWER $\le \alpha$ [exer]

Bonferroni

- m = 5
- $\alpha = 0.05$.

Gene	p-values	Bonferroni	Reject
	p _i	Pi ^{Bonf}	H ₀
1	0.006	$p_i x m = 0.03$	Yes ◀
2	0.372	$p_i x m > 1$	No
3	0.012	$p_i x m = 0.06$	No
4	0.000	$p_i x m = 0.00$	Yes ◀
5	0.811	$p_i x m > 1$	No

Gene 1 and 4 are declared differentially expressed.

Bonferroni

- Now suppose m=1000
- Consider the same 5 genes
- α = 0.05.



Only gene 4 is declared differentially expressed.

Holm's sequential procedure

1. Reject all null-hypotheses with p-value $p_i \le \alpha/m$

2. If r null-hypotheses are rejected, reject all null-hypotheses with p-value $p_i \le \alpha/(m-r)$

3. Continu until no rejections are done

Holm's sequential procedure, adjusted p-values

1.
$$p_{(1)}^{\text{Holm}} = m p_{(1)}$$

2. $p_{(i)}^{\text{Holm}} = \max(p_{(i-1)}, (m-(i-1)) p_{(i-1)})$

3. Reject when
$$p_{(i)}^{Holm} \leq \alpha$$

Theorem

Under *arbitrary* dependency structure: Reject H_{0i} when $p_i^{Holm} \le \alpha$ implies FWER $\le \alpha$

Proof

See hand-out

Holm

- m = 5
- α = 0.05.



Gene 1,3 and 4 are declared differentially expressed.

Holm

- Now suppose m=1000
- Consider the same 5 genes
- α = 0.05.

Gene	p-values	Rank	Holm Pi ^{Holm}	Reject H ₀
1	0.006	3	p _i x 998 > 1	No
2	0.372	455	p _i x 546 > 1	No
3	0.012	44	p _i x 957 > 1	No
4	0.000	1	p _i x m = 0.00	Yes ←
5	0.811	878	p _i x 123 > 1	No

Only gene 4 is declared differentially expressed.

Bonferroni vs Holm

- Bonferroni: very simple
- Holm is more powerful and valid under same assumptions,
- •... but for high-dimensional settings (large m): little gain by using Holm.

Bonferroni and Holm in R

load("C:\\Synchr\\Onderwijs\\HighDimensional\\Slides\\FDR\\Exercises\\pvals.Rdata")

```
pvaladjHolm <- p.adjust(pvals,method="holm")
pvaladjBonf <- p.adjust(pvals,method="bonferroni")</pre>
```

```
cbind(sort(pvaladjHolm)[1:50],
sort(pvaladjBonf)[1:50])
```

Dependence



Genes are not independent, but collaborate in networks...

FWER under dependence

- Bonferroni and Holm are valid, but... likely to be conservative when many positive correlations are present.
- Very relevant in imaging data (neighboring pixels are extremely highly correlated), or DNA genomics data:



chromosomes

FWER under dependence

- For example, imagine m=100.000 variables, which are perfectly correlated withing blocks of 100.
- Bonferroni: $p_i^{bonf} = m \times p_i$,
- However, $p_i^{adj} = m/100 \times p_i$ would suffice for FWER control and is much less conservative
- m/100: effective dimension. Often hard to determine.
- Solution: permutation. Retains the correlation structure between hypotheses (genes).

X: gene matrix, Y: response (disease) labels

Original data		(X,Y)		Permuted data		(X,π(Y))	
subject	gene1	gene2	disease	subject	gene1	gene2	disease
1	3.4	3.8	1	1	3.4	3.8	1
2	5.7	1.9	1	2	5.7	1.9	0
3	2.9	3.7	1	3	2.9	3.7	0
4	3.6	1.3	1	4	3.6	1.3	1
5	1.4	4.1	0	5	1.4	4.1	1
6	1.8	3.8	0	6	1.8	3.8	0
7	2.6	4.7	0	7	2.6	4.7	1
8	3.5	2.9	0	8	3.5	2.9	0



 $FWER = P(V>0) = P_0(min_{i=1,...,m} p_i \le \alpha')$

Find α ' such that FWER $\leq \alpha$

Use permutations to find the distribution of

 $minP = min_{i=1,...,m} \ p_i$

Exchangeability condition

The null-distribution of a test statistic T can be obtained by permutation *if*, under H₀, the distribution of (X, Y) is the same as for (X, π (Y)), where π is a random permutation of response labels Y.



Single-step algorithm, permutation-equivalent of Bonferroni 1. Initiate k=1. 2. Iteration k: Randomly permute the labels Y 3. Calculate new p-values for all tests based on permuted data 4. Calculate the smallest p-value for permuted data: minP_k 5. Repeat 2 - 4 many times (say, 10,000) 6. Calculate critical value α ': the largest value of threshold t for which $P_0(\min P_k \le t) \le \alpha$ 7. Reject all H_{0i} for which $p_i \leq \alpha'$.



Westfall & Young algorithm: permutation-equivalent of Holm

- 1. Start with all hypotheses
- 2. Repeat
 - 3. Single step minP to calculate α '
 - 4. Reject hypotheses with p-value $p_i \le \alpha'$
 - 5. Remove rejected hypotheses
- 6. Stop when no more rejections occur

Why is W&Y more powerful than single-step?



31





32



Example



Summary

Familywise error rate (FWER)

- Generalizes Type I error to multiple hypotheses
- Limits the probability of any error among all inferences

FWER control methods

- Basic: Bonferroni
- Extension 1: Holm
- Extension 2: permutations
- Extension 1 & 2: Westfall & Young