Multiple testing: FDR + Bayesian perspective

Mark van de Wiel mark.vdwiel@vumc.nl

Department of Epidemiology & Biostatistics Department of Mathematics VUmc/VU, Amsterdam







VU medisch centrum

1.	Recap FWER
2.	False Discovery Rate (FDR): Benjamini-Hochberg
3.	FDP and extensions of FDR
4.	Bayesian multiple testing
5.	FDR: pitfalls
6.	FDR vs FWER discussion

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 FWER= 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)

Bonferroni's solution

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

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

Probability on one or more false rejections, FWER, is smaller or equal to α .

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.

Bonferroni's rule:

- Does not scale with the number of tests m
- Is too conservative for very high-dimensional data (e.g. m=100,000)

Alternative: False Discovery Rate (FDR)

FDR: expected proportion of false discoveries among all discoveries given a p-value threshold.

 $FDR = E(V/R) \approx E(V) / R = E(\#FD) / R.$

FDR and False Discovery Proportion (FDP)

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

$FDR = E(FDP) = E(V/R) \approx E(V) / R$

FDR and FDP

FDR control vs FDP estimation

- FDR control: guarantee that $E(V/R) \le \alpha$
- FDP estimation: $\hat{FDP} \leq \alpha$
- FDR control is more strict than FDP estimation

FDR, Benjamini-Hochberg

Control [hand] of FDR, Benjamini-Hochberg's stepup [tut] procedure

- Select a desired limit α of the FDR.
- Order the p-values: $p_{(1)} \le p_{(2)} \le \dots \le p_{(m)}$.
- Find largest rank r for which: m ×p_(r)/r ≤ α.
- All genes with a lower or equal rank are declared significant.



The cut-off line is given by $y = \alpha \times r / m$.

Benjamini-Hochberg: adjusted p-values

• Order the *p*-values: $p_{(1)} \le p_{(2)} \le \dots \le p_{(m)}$.

• $p'_{(r)} = m \times p_{(r)}/r$

• $p_{(r)}^{BH} = min(p'_{(r)}, p'_{(r+1)})$

m=100, five genes with lowest p-values

p _(r)	p' _(r)	$p_{(r)}^{BH}$
0.00010	0.010	0.005
0.00011	0.005	0.005
0.0052	0.173	0.173
0.024	0.600	0.600
0.066	>1	>1

Last step: guarantees that adjusted p-values are in the same order as the raw ones



FDR

Conditions for Benjamini-Hochberg (BH)

- Proven FDR-control for independent p-values [hand] or specific type of correlation (PDSS [tut])
- Simulations, however, show that the BH procedure is very robust against many types of correlations
- Alternative when one suspects very strong correlations:

Replace m/r by:
$$(m/r)\sum_{j=1}^{m} 1/j \approx (m/r) \log(m)$$

• Known as Benjamini-Yekutieli rule

Benjamini-Hochberg and Benjamini-Yekutieli in R

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

```
pvaladjBH <- p.adjust(pvals,method="BH")
pBH50 <- sort(pvaladjBH)[1:50]</pre>
```

```
pvaladjBY <- p.adjust(pvals,method="BY")
pBY50 <- sort(pvaladjBY)[1:50]</pre>
```

```
pvaladjBonf <- p.adjust(pvals,method="bonferroni")
pBonf50 <- sort(pvaladjBonf)[1:50]</pre>
```

cbind(pBH50, pBY50, pBonf50)

FDP/FDR estimation

- Estimation of FDP = V/R. Literature not consistent about FDP/FDR terminology.
- Given p-value threshold t:

$$\begin{split} & F\hat{D}P(t) = \hat{V}(t) / \hat{R}(t) = \sum_{i=1}^{m} P(p_i \leq t \,|\, H_{0i} = true) \, \hat{P}(H_{0i} = true) \, \middle/ \#(p_i \leq t) \\ & = \hat{\pi}_0 mt / \#(p_i \leq t) \end{split}$$

• Methods differ in estimation of the proportion of non-rejections: π_0 [hand,tut]

Equivalence Theorem [hand]

When using $\pi_0 = 1$, the BH rule is equivalent to rejecting all H_{0i} for which $p_i \le t$, with

$$t = max_u(F\hat{D}P(u) \le \alpha)$$

• Equivalence Theorem is nice, but crux of FDP estimators is estimation of π_0

 No guarantee for control of FDR, but generally more powerful than BH



Bayesian interpretation

- Model distribution of p-values by $F(t) = \pi_0 F_0(t) + (1 \pi_0)F_1(t)$
- bFDR = Posterior probability <u>given p-value</u> p_i:

$$\begin{split} b\mathsf{FDR} &= \mathsf{P}(\mathsf{H}_{0i} = \mathsf{true} \mid \mathsf{p}_i \leq \mathsf{t}) \\ &= \mathsf{P}(\mathsf{p}_i \leq \mathsf{t} \mid \mathsf{H}_{0i} = \mathsf{true})\mathsf{P}(\mathsf{H}_{0i} = \mathsf{true}) \big/ \mathsf{P}(\mathsf{p}_i \leq \mathsf{t}) \\ &= \pi_0 \mathsf{F}_0(\mathsf{t}) / \mathsf{F}(\mathsf{t}) \\ &\approx \hat{\pi}_0 \mathsf{mt} / \#(\mathsf{p}_i \leq \mathsf{t}) \end{split}$$

• Local FDR (lfdr) concept [tut,hand]:

$$fdr = P(H_{0i} = true | p_i = t) = \pi_0 f_0(t) / f(t)$$

Bayesian interpretation

Relationship between Ifdr and bFDR [hand]

$E_{f}[Ifdr(U) | U \leq v] = bFDR(v)$

So, at threshold v, bFDR is the mean lfdr over all thresholds $U \le v$

Why bFDR and lfdr are not really "Bayesian"?

$$bFDR = P(H_{0i} = true | p_i \le t)$$

If dr = P(H_{0i} = true | p_i = t)



These quantities still rely on p-values!

True Bayesian analogues of Ifdr and bFDR: BFDR

Recap: posterior probabilities

• Simple regression (or GLM-type) model: $Y_j = \alpha + \beta X_j + \varepsilon_j$

• Prior on β : $\pi(\beta)$

• Likelihood: $\pi(\mathbf{Y}|\beta)$

Bayes' rule

Posterior:
$$\pi(\beta \mid \mathbf{Y}) = \frac{\pi(\beta)\pi(\mathbf{Y} \mid \beta)}{\pi(\mathbf{Y})} = \frac{\pi(\beta)\pi(\mathbf{Y} \mid \beta)}{\int \pi(\beta)\pi(\mathbf{Y} \mid \beta)d\beta}$$

Bayesian inference, single test

Parameter of interest: gene expression difference between two conditions: primary and metastasis. Model:

$$\mathbf{Y}_{j} = \boldsymbol{\alpha} + \boldsymbol{\beta}_{\mathsf{P-M}} \mathbf{I}_{\{\mathbf{j} \square \mathsf{P}\}} + \boldsymbol{\varepsilon}_{j}$$



Hypothesis testing

Null-hypothesis H_0 : $|\beta| \le \delta$

Posterior null-probability: $\pi_0(\mathbf{Y}) = P(\mathbf{H}_0 | \mathbf{Y}) = \int_{-\delta}^{\delta} \pi(\beta | \mathbf{Y}) d\beta$

Recall lfdr: $Ifdr = P(H_0 | p = t) = P(H_0 | p(Y) = t)$

 $\pi_0(\mathbf{Y})$ is a Bayesian version of lfdr

Point null-hypothesis in a Bayesian setting

```
What if one aims to test H_0: \beta = 0?
```

```
\pi_0(\mathbf{Y}) = P(H_0 | \mathbf{Y}) = \pi(0 | \mathbf{Y}) = 0????
```



Prior $\pi(\beta)$ needs to have mass on 0 to avoid this

Point null-hypothesis in a Bayesian setting

If:
$$\pi(\beta) = p_0 \delta(0) + (1 - p_0) \pi'(\beta)$$

+

Then[Exer]:

$$\pi_{0}(\mathbf{Y}) = \frac{p_{0}\pi(\mathbf{Y} \mid 0)}{p_{0}\pi(\mathbf{Y} \mid 0) + (1 - p_{0})\Pi(\mathbf{Y})} = \frac{p_{0}\pi(\mathbf{Y} \mid 0)}{p_{0}\pi(\mathbf{Y} \mid 0) + (1 - p_{0})\int \pi'(\beta)\pi(\mathbf{Y} \mid \beta)d\beta}$$

 $\pi(\mathbf{Y}|\mathbf{0})$ and $\Pi(\mathbf{Y})$ are called *marginal likelihoods*.

BFDR

$$BFDR(t) = E[\pi_0(\mathbf{Y}) \mid \pi_0(\mathbf{Y}) \le t] = \frac{\sum_{i=1}^p \pi_0(\mathbf{Y}) I_{\{\pi_0(\mathbf{Y}) \le t\}}}{\sum_{i=1}^p I_{\{\pi_0(\mathbf{Y}) \le t\}}}$$

piO(Y) = lfdr = t	sum	div	BFDR
0.001	0.001	1	0.001
0.003	0.004	2	0.002
0.008	0.012	3	0.004
0.012	0.024	4	0.006

BFDR vs bFDR

$bFDR(u) = E_f[Ifdr(t) | t \le u]$

$\mathsf{BFDR}(\mathsf{u}) = \mathsf{E}\big[\pi_{0}(\mathbf{Y}) \,|\, \pi_{0}(\mathbf{Y}) \leq \mathsf{u}\big]$

Some discussion

- Bayesian FDR is only an estimate; Bayesians general are less concerned with error control
- Prior very important for performance (see ShrinkBayes lectures)
- Bayesian FDR and (frequentistic) FDR are different entities BUT can be interpreted similary for the purpose of estimation
- Very useful in <u>practice</u>: Biologists, medical researchers etc. do not have intuition for posterior probabilities, but do for (B)FDR
- The more true null features are separated from non-null features the more similar BFDR and FDR are.



FDR: Pitfalls

Subsetting property: for each subset S of all hypotheses, error control on the entire set of hypotheses implies error control for S

FWER has the subsetting property (trivial)

FDR does generally not have the subsetting property

Very important to *a priori* fix the relevant set of hypotheses to be tested

FDR: Pitfalls

m = 1000 genes on the array, α = 0.05.



Significance of gene 1 & 3 depends also on p-values of the *other* genes (the more small p-values, the better...)



FDR: Pitfalls

<u>Setting</u>: researcher knows that apoptosis-related genes are likely to be differentially expressed given two conditions <u>Interest</u>: which *other* genes are differentially expressed?

Example: 500 apoptosis genes, $m_1 = 10.000$ genes, $m_2 = 10.000-500 = 9.500$ non-apoptosis genes

p-value	Apoptosis	₽1 ^{BH}	₽₂ ^{BH}
2.3*10-6	yes	p*m ₁ /1=0.023	
5.2*10 ⁻⁶	yes	p*m ₁ /2=0.026	
7.9*10 ⁻⁶	no	p*m ₁ /3=0.026	p*m ₂ / 1=0.075
11.8*10 ⁻⁶	yes	p*m ₁ /4=0.030	
21.7*10 ⁻⁶	no	p*m ₁ /5=0.043	p*m ₂ / <mark>2=0.103</mark>

Removal of apoptosis genes leads to less discoveries!



Multiple testing

FWER or FDR?

- FDR usually preferable when screening is the pupose and validation on *independent* samples can be performed using a robust technique
- Interpretation of FDR is restricted to the context of the set of features tested. Good definition of this set is very critical
- FWER is preferable when the number of tests is small and no further validation is available
- FWER can be generalized to more powerful and liberal procedures, e.g. kFWER: $P(V > k) \le \alpha$.