Improving the Performance of the FDR Procedure Using an Estimator for the Number of True Null Hypotheses

Amit Zeisel, Or Zuk, Eytan Domany

W.I.S.

June 15, 2009

Amit Zeisel, Or Zuk, Eytan Domany (W.I.S.)Improving the Performance of the FDR Proce

June 15, 2009 1 / 17

The vast use of high throughput technologies involves testing thousands of hypotheses simultaneously. The field of multiple testing deals with developing methods to determine the level of significance in such a scenario.

Multiple testing

m null hypotheses $H_{0,i}$, $\forall i = 1, 2, \dots, m$.

Example: for *m* variables $H_{0,i}: \mu_i^A = \mu_i^B$. Calculate p-values p_i , and set a threshold for significance. The set of random variables (U, V, T, S), and parameters (m, m_0, m_1) describes this scenario:

"ground truth"	non-rejected	rejected	total
	hypotheses	hypotheses	
null hypothesis is true	U	V	m_0
null hypothesis is false	Т	S	m_1
total	m — R	R	т

Fraction of false discoveries = $\frac{V}{R}$

In 1995 Benjamini and Hochberg proposed a procedure (BH95) to control the FDR for a given set of p-values:

- Sort and re-label the p-values, $p_{(1)} \leq p_{(2)} \leq ... \leq p_{(m)}$.
- 2 Choose $0 \le q \le 1$ the desired FDR level.
- **③** Define the set of constants $\alpha_i = \frac{i}{m}q$, i = 1, 2, ..., m.

• Identify
$$R = max\{i : p_{(i)} \leq \alpha_i\}$$
.

- If R ≥ 1 reject all hypotheses (i) = 1, 2, ., R, else no hypothesis is rejected.
- BH proved that $FDR \leq \frac{m_0q}{m} \leq q$.

The bound is not tight, there is room for improvement.

- Control (q ⇒ R) significance is preset at a desired level q. The procedure yields a set of rejected hypotheses with FDR≤ q.
- Estimation (R ⇒ q) the threshold is preset at a level that yields a desired number of rejections R. The corresponding FDR is estimated.

There were many attempts to produce tighter bounds on the FDR, using an estimator for m_0 .

Difficulty: the estimator \hat{m}_0 is a fluctuating random variable.

Produce an improved BH procedure using an estimator \hat{m}_0 for m_0 , the number of true null hypotheses.

An estimator for m_0 is a family of functions $\hat{m}_0 \equiv \hat{m}_0^{(m)} : [0, 1]^m \to \mathbb{R}$, $\hat{m}_0 \equiv \hat{m}_0(p_1, .., p_m)$. \hat{m}_0 is a monotonic estimator if it satisfies:

Given *m* hypotheses of which m_0 are null, let $p_1, ..., p_m$ be the respective p-values. The modified BH procedure with estimator \hat{m}_0 is:

• Compute
$$\hat{m}_0 \equiv \hat{m}_0(p_1,..,p_m)$$
.

- Sort and relabel the p-values $p_{(1)} \leq ... \leq p_{(m)}$.
- Solution Define the set of constants $q_k = \frac{qk}{\hat{m}_0}$ k = 1, 2..., m.

• Let
$$R = \max\{k : p_{(k)} \le q_k\}$$
.

So If $R \ge 1$ reject $p_{(1)}, ..., p_{(R)}$ else don't reject any hypothesis.

Let $\hat{m}_0 \equiv \hat{m}_0(p_1, ..., p_m)$ be a monotonic estimator for m_0 . Let $\hat{m}_0^{(1)}(p_1, ..., p_m) \equiv \hat{m}_0(p_2, ..., p_m)$ be the same estimator, but disregarding the first p-value p_1 . Assume that the null p-values are i.i.d. U[0, 1]. Then the modified BH procedure satisfies:

$$FDR = E\left[\frac{V}{R^{+}}\right] \le m_{0}qE\left[\frac{1}{\hat{m}_{0}^{(1)}}\right]$$
(1)
Note: if $E\left[\frac{1}{\hat{m}_{0}^{(1)}}\right] \le \frac{1}{m_{0}}$, then FDR $\le q$.

June 15, 2009 9 / 17

The two procedures are based on the estimators:

• IBHsum: $\hat{m}_0 = C(m) \cdot \min\left[m, \max(s(m), 2\sum_{j=1}^m p_j)\right]$. C(m), s(m) are universal correction factors.



2 IBHlog: $\tilde{m}_0 = 2 - \sum_{i=1}^m \log(1 - p_i)$.

Both procedures satisfy $E(V/R^+) \le q$

Performance: simulations (IBHsum)

 ρ is the correlation between test statistics:



Compare to other methods

Results from simulations, $m = 500, \mu_1 = 3.5$:



Applying to 33 gene expression datasets



Amit Zeisel, Or Zuk, Eytan Domany (W.I.S.)Improving the Performance of the FDR Proce

June 15, 2009 13 / 17

q	BKY	STS	IBHsum	IBHlog
a. Two tailed, large number of discoveries (10 studies)				
0.05	1.110	1.239	1.200	1.222
	(0.043)	(0.138)	(0.110)	(0.130)
b. Two tailed, small number of discoveries (10 studies)				
0.05	1.003	1.316	1.231	1.291
	(0.003)	(0.197)	(0.140)	(0.179)
c. One tailed, large number of discoveries (8 studies)				
0.05	1.049	1.011	1.014	0.108
	(0.019)	(0.033)	(0.026)	(0.306)
d. One tailed, small number of discoveries (5 studies)				
0.05	0.998	1.027	1.025	0.882
	(0.020)	(0.052)	(0.017)	(0.123)

June 15, 2009 14 / 17

◆□ > ◆□ > ◆臣 > ◆臣 > ○ 臣 = ∽ 9 Q (?)

- We proved a theorem that provides a bound on the FDR for any improved procedure based on an monotonic estimator of *m*₀.
- We proposed two improved procedures based on the estimators $2\sum_{j=1}^{m} p_j$ and $\sum_{i=1}^{m} \log(1-p_i)$.
- For the case of independent statistics all improved procedures provide similar results: saturation of the bound and more power than BH95.
- We showed by simulations that even in the case of dependent statistics our procedures provide a reliable bound and improved power.
- For real gene expression data, where dependencies are expected, our methods improve, in general, over existing ones.

Let $\vec{p} = (p_{1..m})$ be a set of independent p-values. Assume that f, the marginal probability density function of the alternatives, is monotonically non-increasing and differentiable. Let $B^{(i)}$ be two threshold FDR procedures rejecting $R^{(i)}(\vec{p})$ hypotheses and each having $FDR^{(i)}$, i = 1, 2. Assume that for any q, $R^{(1)}(\vec{p}) \leq R^{(2)}(\vec{p})$, $\forall \vec{p}$. Then it also holds that $FDR^{(1)} \leq FDR^{(2)}$.

THANKS

Amit Zeisel, Or Zuk, Eytan Domany (W.I.S.)Improving the Performance of the FDR Proce

크