

Multiple Testing Procedures for Hierarchically Related Hypotheses

Przemyslaw.Biecek@gmail.com, MIM University of Warsaw

Outline

- ➊ Brief introduction to multiple hypothesis testing,
- ➋ Testing under a hierarchical relation,
- ➌ Testing procedures that control:
 - FWER,
 - PFER,
 - FDR.
- ➍ Simulation study,
- ➎ Application to Gene Ontology.

This talk will NOT be about...

I am interested in applications in biology, genetics and medicine

- QTL mapping (model selection criteria, accurate estimates for unbalanced data, QTL-QTL interaction),
- microarray analyses and other $p \gg n$ problems,
- survival regression for longitudinal studies (kidney graft function, breast cancer study).



A set of null hypotheses \mathcal{H}

- Consider a set of m random experiments $\{\mathcal{E}^{(i)} : i \in \{1, \dots, m\}\}$. For each experiment $\mathcal{E}^{(i)}$ we introduce a family of distributions $\mathcal{P}^{(i)} = \{P_{\theta^{(i)}}^{(i)} : \theta^{(i)} \in \Theta^{(i)}\}$.
- We assume that observations generated by experiment $\mathcal{E}^{(i)}$ follow distribution $P_{\theta_0^{(i)}}^{(i)}$, for some $\theta_0^{(i)} \in \Theta^{(i)}$.
- For each random experiment $\mathcal{E}^{(i)}$ we specify null and alternative hypotheses. The set of all null hypotheses is denoted by $\mathcal{H} = \{H_0^{(i)} : i = 1, \dots, m\}$.

Null hypothesis and alternative hypothesis

Null hypothesis $H_0^{(i)}$ is the supposition that $\theta_0^{(i)} \in \Theta_0^{(i)}$.

Alternative hypothesis $H_A^{(i)}$ is the supposition that $\theta_0^{(i)} \notin \Theta_0^{(i)}$.

Standard notation

The outcomes of testing the set \mathcal{H} may be described in the following form.

	#accepted nulls	#rejected nulls	
#true nulls	U	V	m_0
#false nulls	T	S	m_1
sum	$m - R$	R	$m.$

Variables V and T denote number of wrong decisions.

The error rates

Per-family error rate (PFER)

$$PFER = E(V).$$

Family-wise error rate (FWER)

$$FWER = Pr(V > 0).$$

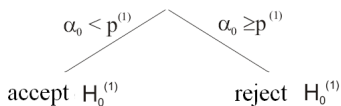
False discovery rate (FDR)

$$FDR = E(Q),$$

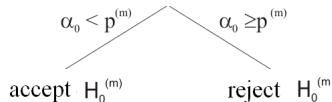
where

$$Q = \begin{cases} V/R & R > 0, \\ 0 & R = 0. \end{cases}$$

Single step procedure



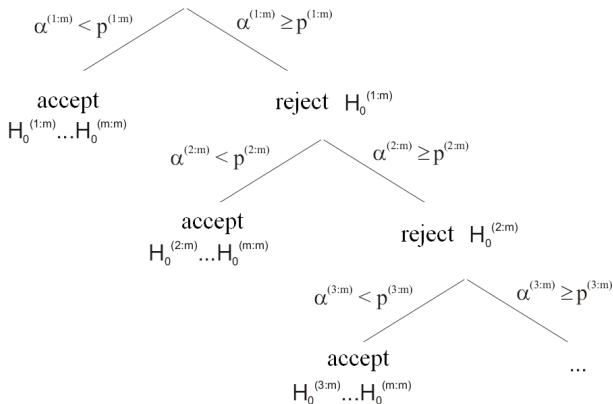
...



Bonferroni procedure (1936)

The single step testing procedure with significance level $\alpha_0 = \alpha/m$, controls the FWER at the level α .

Step down



Holm procedure (1979)

The step-down testing procedure with significance levels

$$\alpha^{(i:m)} = \alpha / (m - i + 1), \quad (1)$$

controls FWER at the level α .

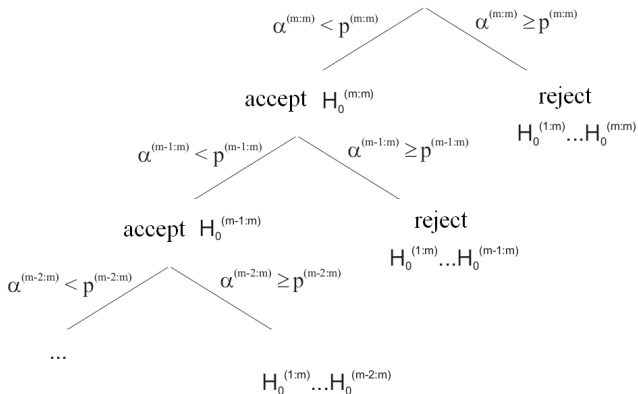
Hochberg procedure (1988)

The step-up testing procedure with significance levels

$$\alpha^{(i:m)} = \alpha / (m - i + 1), \quad (2)$$

controls FWER at the level α .

Step up



Benjamini and Hochberg procedure (1995)

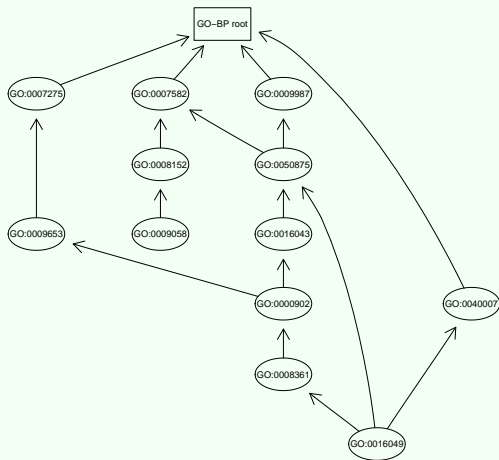
The step-up procedure with significance levels

$$\alpha^{(i:m)} = \frac{i}{m} \alpha, \quad (3)$$

controls FDR at the level α .

The control of PFER and FWER is more conservative than control of FDR (i.e. leads to a smaller number of rejections). In applications procedures that control FDR are more popular.

Gene Ontology



- There are some papers concerned with the control of FWER under such hierarchical relations (closure testing). First results were published by Gabriel (1969), R. Marcus (1976), U. Naik (1977), more recently by H. Finner (2002).
- In these papers one common parameter space Θ is introduced for all null hypotheses. The relation between hypotheses is derived from the relation between the sets $\Theta_0^{(i)}$, e.g. if $\Theta_0^{(i)} \subset \Theta_0^{(j)}$, then the rejection of $H_0^{(j)}$ requires the rejection of $H_0^{(i)}$.
This approach is very interesting, but hard to apply to the analysis with the Gene Ontology, since it is not obvious how to describe the sets $\Theta_0^{(i)}$ in terms of a common parameter space.
- We propose a different approach to incorporating hierarchical relations into the testing scheme.

A toy example

We observe $x_i \sim \mathcal{N}(\mu_i, 1)$, $i \in \{1, 2\}$, $(\mu_1, \mu_2) \in \Theta = \mathbb{R}^2$.

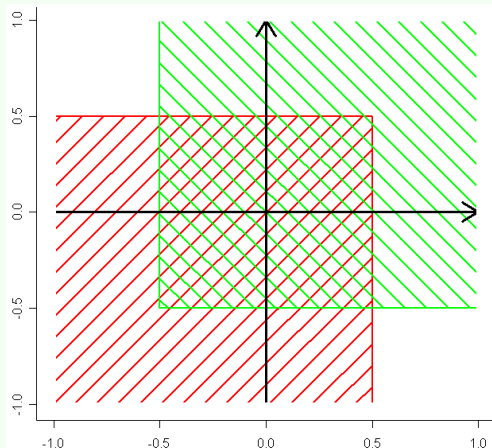
Let's consider following hypotheses:

$$H_0^{(1)} : -0.5 < \min(\mu_1; \mu_2);$$

$$H_0^{(2)} : -0.5 < \mu_1; \mu_2 < 0.5;$$

$$H_0^{(3)} : \max(\mu_1; \mu_2) < 0.5.$$

We cannot apply this approach to analyses based on Gene Ontology!!!



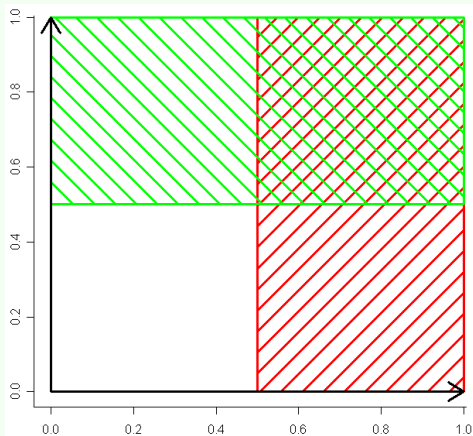
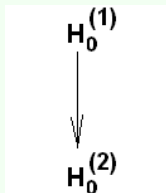
A toy example

We observe $x_i \sim \mathcal{B}(\rho^{(i)})$, $i \in \{1, 2\}$, $(\rho^{(1)}, \rho^{(2)}) \in \Theta = [0, 1]^2$.

Let's consider following hypotheses:

$$H_0^{(1)} : \rho^{(1)} \leq 0.5;$$

$$H_0^{(2)} : \rho^{(2)} \leq 0.5.$$



Set of null hypotheses

As before, consider the set of m random experiments $\{\mathcal{E}^{(i)} : i \in \{1, \dots, m\}\}$, the corresponding set $\{\Theta^{(i)} : i \in \{1, \dots, m\}\}$ and the set of null hypotheses $\mathcal{H} = \{H_0^{(i)} : i \in \{1, \dots, m\}\}$ of the form $H_0^{(i)} : \theta_0^{(i)} \in \Theta_0^{(i)}$.

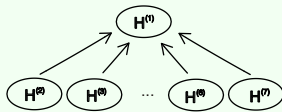
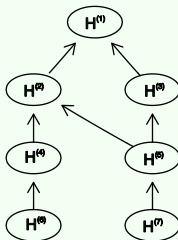
Hierarchical relation \mathcal{R}

Let \mathcal{R} denote an irreflexive (aliorelative), asymmetric and transitive binary relation on the set \mathcal{H} . Relation \mathcal{R} fulfills following conditions

- ① $R(i, i) = 0$,
- ② $(R(i, j) = R(j, k) = 1) \Rightarrow (R(i, k) = 1)$,
- ③ $R(i, j) * R(j, i) = 0$.

Relation \mathcal{R} describes possible states of null hypotheses.

$$(R(i, j) = 1) \Rightarrow (H_0^{(j)} \text{ is true} \Rightarrow H_0^{(i)} \text{ is true}).$$



Testing outcomes

Let $\psi = (\psi_1, \dots, \psi_m)$, where $\psi_i \in \{0, 1\}$, denote the outcomes of the tests as standard.

Coherency

The outcomes of the tests ψ are coherent with relation \mathcal{R} if and only if

$$(R(i, j) = 1) \Rightarrow (\psi_i \geq \psi_j).$$

We want to obtain coherent results!!!

Control of FWER

Closure of the outcomes of tests

Let $\hat{\psi}$ stand for the closure of outcomes of tests ψ , where

$$\hat{\psi}_i = \max(\psi_i, \{\psi_j : R(i, j) = 1\}).$$

In other words, this means that the rejection of a test leads to rejection of all related tests.

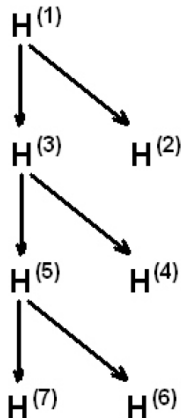
Theorem 1

Closure of the outcomes of tests does not affect FWER.

Note

Closure of the outcomes of tests results in an increase of both PFER and FDR.

Follow up testing procedure



Step 1 We test $H_0^{(7)}$

$p^{(7)} > \alpha^{(7)}$, thus we accept $H_0^{(7)}$

Step 2: We test $H_0^{(6)}$

$p^{(6)} < \alpha^{(6)}$, thus we reject $H_0^{(1)}$, $H_0^{(3)}$,
 $H_0^{(5)}$, $H_0^{(6)}$

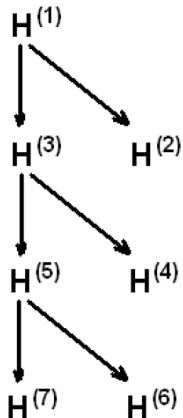
Step 3 We test $H_0^{(4)}$

$p^{(4)} > \alpha^{(4)}$, thus we accept $H_0^{(4)}$

Step 4 We test $H_0^{(2)}$

$p^{(2)} > \alpha^{(2)}$, thus we accept $H_0^{(2)}$

Follow down testing procedure



Step 1: We test $H_0^{(1)}$
 $p^{(1)} < \alpha^{(1)}$, thus we reject $H_0^{(1)}$

Step 2: We test $H_0^{(3)}$
 $p^{(3)} > \alpha^{(3)}$, thus we accept
 $H_0^{(3)}$, $H_0^{(4)}$, $H_0^{(5)}$, $H_0^{(6)}$, $H_0^{(7)}$

Step 3: We test $H_0^{(2)}$
 $p^{(2)} < \alpha^{(2)}$, thus we reject $H_0^{(2)}$

The control of PFER

Theorem 2

Let $\phi(i)$ stand for the maximum cardinality of the set of unrelated hypotheses which contains $H_0^{(i)}$

$$\phi(i) = \max_{A \in \mathcal{C}_i} \#A,$$

where

$$\mathcal{C}_i = \{C : i \in C \wedge \forall_{j,k \in C} R(j,k) = 0\}.$$

The follow down strategy with significance levels

$$\alpha^{(i)} = \alpha / [\phi(i)(1 + \alpha)] \quad (4)$$

controls PFER at the level α .

The control of PFER

Theorem 3

Let $\rho(i)$ stand for the number of hypotheses related to $H_0^{(i)}$

$$\rho(i) = 1 + \#\{j : R(j, i) = 1\}.$$

The follow up procedure with significance levels

$$\alpha^{(i)} = \alpha / [m * \rho(i)] \quad (5)$$

controls PFER at the level α .

The control of FDR

Theorem 4

Let \mathcal{C}_i stand for the family of sets

$$\mathcal{C}_i = \{C : i \in C \wedge \forall_{j,k \in C} R(j, k) = 0\}.$$

Let $\rho(B)$ stand for the number of hypotheses related to any hypothesis from the set B

$$\rho(B) = 1 + \#\{j : i \in B \wedge R(j, i) = 1\}.$$

The follow down strategy with significance levels

$$\alpha^{(i)} = \min_{B \in \mathcal{C}_i} \{\alpha_{lin}^{(\rho(B), m)} / \#B\}, \quad (6)$$

controls FDR at the level α .

Control of the FDR

Theorem 4 cont.

Significance levels $\alpha_{lin}^{(i,m)}$ are derived in the following way

$$\begin{aligned}\alpha_{lin}^{(1,m)} &= \alpha, \\ \alpha_{lin}^{(i,m)} &= \min \left(0.5, \alpha \left[\sum_{k=i}^{m-1} \frac{k-i+1}{k} (1 - \alpha_{lin}^{(k+1,m)}) \prod_{l=i+1}^k \alpha_{lin}^{(l,m)} + \right. \right. \\ &\quad \left. \left. \frac{m-i+1}{m} \prod_{l=i+1}^m \alpha_{lin}^{(l,m)} \right]^{-1} \right), \text{ for } 1 < i < m, \\ \alpha_{lin}^{(m,m)} &= \min(0.5, m\alpha).\end{aligned}\tag{7}$$

Example: Relation \mathcal{R}_1

$\mathcal{R}_1(i, j)$	j=1	j=2	j=3	j=4	j=5	j=6	j=7
i=1	0	1	1	1	1	1	1
i=2	0	0	1	1	1	1	1
i=3	0	0	0	1	1	1	1
i=4	0	0	0	0	1	1	1
i=5	0	0	0	0	0	1	1
i=6	0	0	0	0	0	0	1
i=7	0	0	0	0	0	0	0

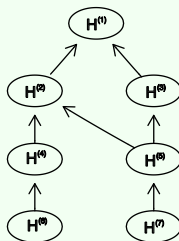


Significance levels for \mathcal{R}_1

	Relation \mathcal{R}_1				
	$\rho(i)$	$\phi(i)$	$\alpha_{FU}^{PFER}(i)$	$\alpha_{FD}^{PFER}(i)$	$\alpha_{FD}^{FDR}(i)$
i=1	1	1	0.00714	0.04761	0.05000
i=2	2	1	0.00357	0.04761	0.09523
i=3	3	1	0.00238	0.04761	0.13636
i=4	4	1	0.00178	0.04761	0.17391
i=5	5	1	0.00142	0.04761	0.20833
i=6	6	1	0.00119	0.04761	0.24000
i=7	7	1	0.00102	0.04761	0.35000

Example: Relation \mathcal{R}_2

$\mathcal{R}_2(i, j)$	j=1	j=2	j=3	j=4	j=5	j=6	j=7
i=1	0	1	1	1	1	1	1
i=2	0	0	0	1	0	1	0
i=3	0	0	0	0	1	0	1
i=4	0	0	0	0	0	1	0
i=5	0	0	0	0	0	0	1
i=6	0	0	0	0	0	0	0
i=7	0	0	0	0	0	0	0



Significance levels for \mathcal{R}_2

	Relation \mathcal{R}_2				
	$\rho(i)$	$\phi(i)$	$\alpha_{FU}^{PFER}(i)$	$\alpha_{FD}^{PFER}(i)$	$\alpha_{FD}^{FDR}(i)$
i=1	1	1	0.00714	0.04761	0.05000
i=2	2	2	0.00357	0.02380	0.04762
i=3	2	2	0.00357	0.02380	0.04762
i=4	3	2	0.00238	0.02380	0.06818
i=5	4	2	0.00178	0.02380	0.08696
i=6	4	2	0.00178	0.02380	0.10417
i=7	5	2	0.00142	0.02380	0.10417

Simulation design

Consider a set of 7 experiments. Observations are drawn from $\mathcal{N}(\mu, I_{7 \times 7})$, where $\mu = (\mu_1, \dots, \mu_7)$.

The corresponding null hypotheses are of the form $H_0^{(i)} : \mu_i = 0$.

We consider 6 different scenarios for μ

$$\mu^{(1)} = (0, 0, 0, 0, 0, 0, 0),$$

$$\mu^{(2)} = (2, 2, 2, 2, 0, 0, 0),$$

$$\mu^{(3)} = (3, 2.75, 2.5, 2.25, 0, 0, 0),$$

$$\mu^{(4)} = (2, 2, 2, 2, 2, 2, 2),$$

$$\mu^{(5)} = (3, 2.75, 2.5, 2.25, 2, 1.75, 1.5),$$

$$\mu^{(6)} = (3, 0, 0, 0, 0, 0, 0).$$

Closure of testing procedures

		closure of single-step Bonferroni procedure	closure of step-up Benjamini-Hochberg procedure
Relation	Vector μ	PFER	FDR
\mathcal{R}_1	$\mu^{(1)}$	0.201	0.050
\mathcal{R}_1	$\mu^{(2)}$	0.073	0.046
\mathcal{R}_1	$\mu^{(6)}$	0.151	0.085
\mathcal{R}_2	$\mu^{(1)}$	0.151	0.050
\mathcal{R}_2	$\mu^{(2)}$	0.045	0.035
\mathcal{R}_2	$\mu^{(6)}$	0.102	0.076
\mathcal{R}_3	$\mu^{(1)}$	0.093	0.050
\mathcal{R}_3	$\mu^{(2)}$	0.294	0.028
\mathcal{R}_3	$\mu^{(6)}$	0.044	0.042

Results based on 1 000 000 repetitions.

Results for proposed testing procedures

Relation	Vector μ	follow-up PFER	follow-down PFER	follow-down FDR	follow-up	follow-down E(S)	follow-down
\mathcal{R}_1	$\mu^{(1)}$	0.050	0.050	0.050	—	—	—
\mathcal{R}_1	$\mu^{(2)}$	0.007	0.008	0.017	1.712	1.433	1.852
\mathcal{R}_1	$\mu^{(3)}$	0.007	0.022	0.035	2.503	2.761	3.236
\mathcal{R}_1	$\mu^{(4)}$	—	—	—	3.384	1.634	2.670
\mathcal{R}_1	$\mu^{(5)}$	—	—	—	3.425	3.258	4.842
\mathcal{R}_1	$\mu^{(6)}$	0.031	0.045	0.045	0.712	0.909	0.912
\mathcal{R}_2	$\mu^{(1)}$	0.050	0.050	0.050	—	—	—
\mathcal{R}_2	$\mu^{(2)}$	0.006	0.008	0.010	1.552	1.432	1.722
\mathcal{R}_2	$\mu^{(3)}$	0.006	0.022	0.023	2.314	2.679	3.034
\mathcal{R}_2	$\mu^{(4)}$	—	—	—	2.975	1.638	2.265
\mathcal{R}_2	$\mu^{(5)}$	—	—	—	3.092	3.184	4.186
\mathcal{R}_2	$\mu^{(6)}$	0.028	0.044	0.043	0.712	0.909	0.912
\mathcal{R}_3	$\mu^{(1)}$	0.050	0.050	0.050	—	—	—
\mathcal{R}_3	$\mu^{(2)}$	0.011	0.015	0.006	1.448	1.272	1.303
\mathcal{R}_3	$\mu^{(3)}$	0.011	0.022	0.007	2.225	2.365	2.394
\mathcal{R}_3	$\mu^{(4)}$	—	—	—	2.345	1.915	1.967
\mathcal{R}_3	$\mu^{(5)}$	—	—	—	2.784	3.069	3.116
\mathcal{R}_3	$\mu^{(6)}$	0.021	0.043	0.022	0.714	0.909	0.912

The presented procedures were applied to functional enrichment analysis. We used a data set provided by Adam Zagdanski, which contains protein-protein interactions and GO BP annotations. To predict functional attribute for a protein, we use Fisher's exact test which identifies overrepresented functional attributes

$$p^{(i,j)} = 1 - \sum_{l=0}^{k_{i,j}-1} \frac{\binom{k_j}{l} \binom{n-k_j}{n_i-l}}{\binom{n}{n_i}},$$

where

- n , number of proteins,
- k_j , number of proteins with functional attribute j ,
- n_i , number of proteins interacting with protein i ,
- $k_{i,j}$, number of proteins interacting with protein i with functional attribute j .

- We selected a subset of 1169 proteins and 258 functional attributes, which yields 301 602 null hypotheses.
- For this data set we applied the closure of the Hochberg procedure, follow up and follow down testing procedures, to control PFER, FWER and FDR, respectively, at the level $\alpha = 0.05$.
- To compute empirical error rates, the one-leave out CV scheme was used.

Results

Testing procedure	#rejected nulls	empirical error rate	#true positives
Follow up PFER controlling procedure	1706	PFER = 0.0209	1683
Follow down PFER controlling procedure	1191	PFER = 0.0164	1173
Closure of the Hochberg step up procedure	2300	FWER = 0.0382	2255
Follow down FDR controlling procedure	2892	FDR = 0.0173	2823

Choosing the number of significant variables

Let's consider the linear model

$$Y = X\beta + \sigma\varepsilon, \quad (8)$$

where

- Y stands for $n \times 1$ dependent variable,
- X is a full rank matrix $n \times p$,
- β is a vector of parameters $p \times 1$,
- σ is an unknown constant,
- ε is a vector of random fluctuations $\mathcal{N}(0, 1)$.

Let M_0 stand for the number of nonzero elements in β .

The problem

Find a good estimate for M_0 .

The „forward selection” strategy

- ➊ In the first step we choose an empty model $M^{best(0)} = \{\}$. Likelihood function for this model is $L(M^{best(0)}|Y)$.
- ➋ Consider p models of the form $M_i = \{i\}$. Let $L(M^{best(1)}|Y) = \max_{i \in \{1 \dots p\}} L(M_i|Y)$, and for model $M^{best(1)}$ we have $L(M_i|Y) = L(M^{best(1)}|Y)$.
- ➌ Do a likelihood test for nested models, compare $M^{best(1)}$ with $M^{best(0)}$ on the significance level $\alpha^{(1)}$.
- ➍ If we reject this hypothesis we stop the procedure. Otherwise we choose model $M^{best(1)}$ and go to the next step.
- ➎ Best model from step i is marked as $M^{best(i)}$. Consider $p - i + 1$ models of the form $M_i = M^{best(i)} \cup \{i\}$. Note that $L(M^{best(i+1)}|Y) = \max_i L(M_i|Y)$, and model which $L(M_i|Y) = L(M^{best(i+1)}|Y)$ will be denoted as $M^{best(i+1)}$.
- ➏ Do the likelihood test and compare $M^{best(i+1)}$ with $M^{best(i)}$ on the significance level $\alpha^{(i+1)}$.
- ➐ Repeat these steps until any hypothesis is accepted.

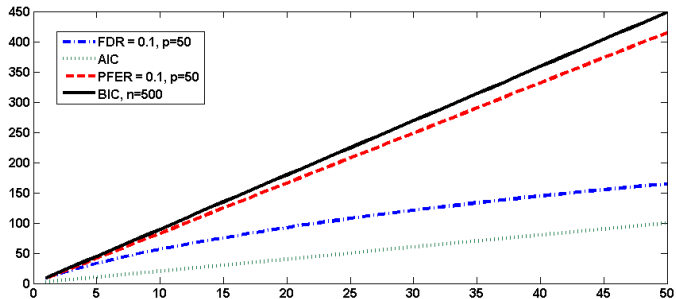
The „forward selection” procedure choose the model wich maximize $\mathcal{S}(M_i)$

$$\mathcal{S}(M_i) = -2 \log L(M_i|Y) - \sum_{j=1}^{|M_i|} c^{(j)} = -2 \log L(M_i|Y) - \sum_{j=1}^{|M_i|} (\chi_1^2)^{-1} (1 - \alpha^{(j)})$$

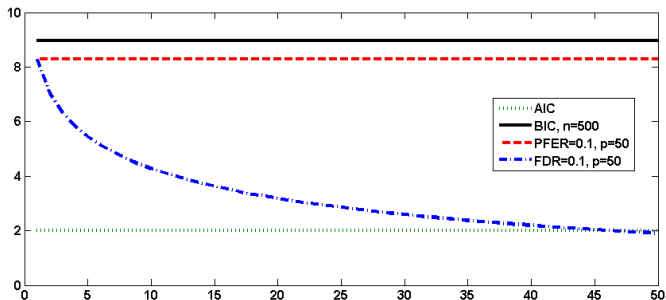
where $\sum_{i=1}^{|M_i|} c^{(i)}$ stands for penalty for the model size. If $\alpha^{(j)}$ are equal then this criteria are equal to GLC (Generalized Infortmation Criteria), in which we maximize

$$\mathcal{G}(M_i) = -2 \log L(M_i|Y) + \lambda |M_i|$$

where λ is a nonnegative constatn. For $\lambda = 2$ GLC is equal to AIC, and for $\lambda = \log_2(n)$ it becomes BIC.



Penalty for model with k components.



Penalty for one additional component for model with k components.

ρ	$ M_0 $	β	n	p	BIC	AIC	$FDR = 0.1$
0	1	0.5	100	50	0.01	3.50	0.00
0	1	0.4	200	50	0.00	2.21	0.00
0	1	0.5	200	50	0.01	2.59	0.01
0	10	0.5	100	50	-4.77	3.58	-0.49
0	10	0.5	200	50	-0.03	2.48	0.21
0	10	0.4	200	50	-1.12	2.52	0.26
0	30	0.4	200	50	-26.98	1.76	1.07
0	30	0.3	500	50	-0.59	1.10	0.60
0.9	1	0.4	500	50	0.00	1.55	0.01
0.9	10	0.4	500	50	-2.18	1.75	0.29
0.9	10	0.5	200	50	-3.20	1.90	0.11
0.9	10	0.6	200	50	-2.22	3.12	1.09

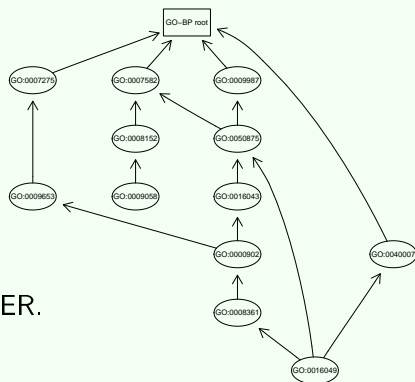
Average differences between number of nonzero elements in the linear model and our prediction. We consider 50 regressors. ρ is equal to correlation among test statistics ($\sigma = 1$).

Short summary

This talk was about testing strategies under hierarchical relation among hypotheses.

We proposed testing procedures that control FWER, FDR and PFER.

There is a lot of applications for such procedures.
And still there is a lot to do.



Acknowledgments

- Many thanks for your attention.
- Many thanks to Adam Zagdanski for arousing my interest in Gene Ontology.
- This research was financially supported by the Polish State Committee for Scientific Research, grant 1 P03A 017 29, so many thanks to the Polish Government ;-).

Literature



Marcus R., Peritz E., Gabriel K. R. 1976. On closed testing procedures with special reference to ordered analysis of variance. *Biometrika* 1976 63(3):655-660;



Finner H., Strassburger K. 2002. The partitioning principle: A powerful tool in multiple decision theory. *The Annals of Statistics* 2002, Vol. 30, No. 4, 1194-1213.



Benjamini Y., Hochberg Y. "Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple Testing". *J. R. Statist. Soc. B* (1995) 57.



Benjamini, Y., and Yekutieli, D. (2001). "The Control of the False Discovery Rate in Multiple Testing under Dependency,". *The Annals of Statistics* 29 (4).



Dudoit S., van der Laan M. J., Pollard K. S. (2004), Multiple Testing. Part I. Single-Step Procedures for Control of General Type I Error Rates, *Statistical Applications in Genetics and Molecular Biology*, 3(1).



Storey, J. D. (2002). „A direct approach to false discovery rates." *Journal of the Royal Statistical Society, Series B* 64, 479-498



Hochberg, Y. and Tamhane, A. C. *Multiple Comparison Procedures*, (1987) , New York John Wiley and Sons Inc.