Bounding the k-family-wise error rate using resampling methods

Tijl De Bie and John Shawe-Taylor

University of Bristol, K.U.Leuven, University College London

May 2006

A (1) > (1) = (1) (1)

Overview

- I will discuss a methodology to bound
 - the Family-Wise Error Rate (FWER)
 - the *k*-FWER
 - the False Discovery Proportion (FDP)

by means of randomisation (resampling)

- Want to achieve:
 - quantification and understanding of conclusions of multiple testing
 - take account of dependencies between tests
- (Results build on previous work by Meinshausen and others! See references)

- 4 同 6 4 日 6 4 日 6

Overview

- Single and multiple hypothesis testing (notation / formalism / terminology)
- Randomisation testing
- Bounding the FWER, k-FWER, FDP by resampling
- Practical results

・ロト ・同ト ・ヨト ・ヨト

Testing a single hypothesis Testing a set of hypotheses

Motivating example

Assumption:

• Gene A is unrelated to tissue type (e.g. 'cancerous' / 'non-cancerous')

Observation:

- Expression of Gene A in a set of tissue samples
- Label of each of these tissue samples

labels: expression data:

< ロ > < 同 > < 回 > < 回 >

Question:

- Is the observation 'surprising', given our assumption?
- If so, something more interesting than the null hypothesis may be true...

A pattern function to quantify the state

- Aim: draw statistical conclusions about state of the world X based on a finite sample
- 'Particular state': quantified by means of a *pattern function*

$$\pi:X\to\mathbb{R}$$

• For example, $\pi=$ correlation of the expression with the label (-1/1)

- 4 回 ト 4 ヨト 4 ヨト

Testing a single hypothesis Testing a set of hypotheses

The null hypothesis

- Which statistical conclusions can we draw?
- One computes an upper bound on the probability to see a pattern as strong as $\pi(X) \ge \sigma$, assuming a null hypothesis, the 'dull' assumption
- Null hypothesis = a set Ω_0 of distributions ${\mathcal D}$ over X
- Hence, level of surprise is quantified by p:

$$\forall \mathcal{D} \in \Omega_{0} : P_{X \sim \mathcal{D}} \left(\pi \left(X \right) \geq \sigma \right) \leq p$$

• *p* is called the p-value

- 4 回 ト 4 ヨト 4 ヨト

Testing a single hypothesis Testing a set of hypotheses

Rejecting the null hypothesis

- If p is small, the level of surprise is large $ightarrow \Omega_0$ rejected
- If Ω_0 holds nevertheless: false rejection (aka: type I error)
- If Ω_0 does not hold: true rejection
- If Ω_0 does not hold but still $\pi(X) < \sigma$: type II error

test $\setminus \Omega_0$	true	false
accept	just acceptance	type II error
reject	type I error	true rejection

• Challenge: keep type I and II errors small

・ロト ・同ト ・ヨト ・ヨト

Multiple testing: the formalism

Randomisation testing A zoo of error rates In practice Testing a single hypothesis Testing a set of hypotheses

Motivating example

Back to our genes...

• Microarrays gather data for all *n* genes in genome



expression data:



A 10

Testing a single hypothesis Testing a set of hypotheses

Motivating example

Back to our genes...

- Microarrays gather data for all *n* genes in genome
- Doing *n* single hypothesis tests will cause many false rejections
- Assuming the genes are independent and Ω_0 holds, *pn* genes are expected to be rejected! (e.g. for 20,000 genes and p = 0.01, this is 200)
- Should all these genes be subjected to a closer look of the biologist??
- Need for a new formalism

イロト イポト イラト イラト

Testing a single hypothesis Testing a set of hypotheses

A set of pattern functions

• Multiple testing: a set of pattern functions, each quantifying one aspect of the state of the world:

$$\Pi = \{\pi_{\alpha} : \alpha \in A\}$$

with A and index set of size |A|

- For example: π_{α} measures correlation of gene α with the labels
- Now we have |A| different tests

$$\pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}$$

(note: X contains all expression data, for all genes)

< 口 > < 同 > < 三 > < 三

Testing a single hypothesis Testing a set of hypotheses

A set of pattern functions

- Null hypothesis: Ω_0 contains all distributions of X where the labels and the gene expression data are independent
- Each gene may provide evidence to reject null hypothesis now! If:

$$\pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}$$

- FWER: probability to reject null hypothesis assuming that it holds (i.e. the probability than any of the genes motivates rejection)
- The FWER is $\leq p$ if

$$\forall \mathcal{D} \in \Omega_{0} : P_{X \sim \mathcal{D}} \left(\exists \alpha \in A : \pi_{\alpha} \left(X \right) \geq \sigma_{\alpha} \right) \leq p$$

- 4 回 ト 4 ヨト 4 ヨト

Testing a single hypothesis Testing a set of hypotheses

Alternative null hypotheses

- Assuming that all genes are unrelated to the labels is probably too simplistic (hopefully!)
- Some will be truly unrelated
- How can we do more than merely reject the null hypothesis or not, based on some genes as a witness?
- In particular: can we say anything about these genes that show surprising behaviour?

イロト イポト イラト イラト

Testing a single hypothesis Testing a set of hypotheses

Alternative null hypotheses

- Null hypothesis Ω₀ (A
 = the sets of all distributions where all but the subset {π_α : α ∈ A} ⊆ Π of k = |A| pattern functions is jointly distributed as under a D ∈ Ω₀
- Hence: we tolerate $k = |\overline{A}|$ positive genes (i.e. genes that are *potentially* related with the labels)
- Then the probability of at least one *false* rejection is

$$\forall \mathcal{D} \in \Omega_{0}\left(\overline{A}\right) : P_{X \sim \mathcal{D}}\left(\exists \alpha \in A \backslash \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\right) \leq p$$

- 4 回 ト 4 ヨト 4 ヨト

Testing a single hypothesis Testing a set of hypotheses

Alternative null hypotheses

• Further define:

$$\overline{\Omega}_0 = igcup_{\overline{\mathcal{A}}\subseteq \mathcal{A}} \Omega_0\left(\overline{\mathcal{A}}
ight)$$

 Under any distribution from null hypothesis Ω₀, the probability of at least one *false* rejection is bounded by *p* if:

$$\forall \overline{A}, \forall \mathcal{D} \in \Omega_{0}\left(\overline{A}\right): P_{X \sim \mathcal{D}}\left(\exists \alpha \in A ackslash \overline{A}: \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\right) \leq p$$

 Hence, if p is small and there are rejections α, then we must conclude that these are true rejections and belong to A

伺 ト イ ヨ ト イ ヨ

Testing a single hypothesis Testing a set of hypotheses

Alternative null hypotheses

- Note: we never reject the null hypothesis $\overline{\Omega}_0$
- But: we would reject many $\Omega_0(\overline{A})$ with too small sets \overline{A}
- Namely: we would reject all $\Omega_0(\overline{A})$ for which

$$\exists \alpha \in A \backslash \overline{A} : \pi_{\alpha} (X) \geq \sigma_{\alpha}$$

Now, how to compute this upper bound p?

イロト イポト イラト イラト

Testing a single hypothesis Testing a set of hypotheses

Bonferroni correction

Lemma (Bonferroni correction)

Suppose $A \subseteq \mathbb{N}$ is a countable index set for a class of pattern functions $\{\pi_{\alpha} : \alpha \in A\}$. Let $q_{\alpha} \in (0, 1)$ such that

$$\sum_{lpha\in A}q_lpha\leq 1.$$

If we choose σ_{α} such that $\forall \mathcal{D} \in \Omega_0$:

$$P_{X\sim\mathcal{D}}\left(\pi_{\alpha}\left(X\right)\geq\sigma_{\alpha}\right)\leq pq_{\alpha},$$

then $\forall \overline{A} \subseteq A, \forall \mathcal{D} \in \Omega_0\left(\overline{A}\right)$:

$$P_{X\sim\mathcal{D}}\left(\exists \alpha\in A\setminus\overline{A}:\pi_{\alpha}\left(X\right)\geq\sigma_{\alpha}\right)\leq\rho.$$

< ロ > < 同 > < 回 > < 回 >

Testing a single hypothesis Testing a set of hypotheses

Bonferroni correction

- The most well-known approach to bound the FWER
- Proof is a straightforward application of the union bound
- Correct for all Ω_0 for arbitrary dependencies between genes
- Therefore very conservative :-(

- 4 同 2 4 日 2 4 日 2

Randomisation (permutation) testing

So far for the general introduction to multiple testing... Now a general introduction to randomisation testing as a technique!

- Computation intensive method to perform hypothesis testing
- I Will first discuss single hypothesis testing
- Explicitly taking finiteness of random sample into account
- Then expand ideas to multiple hypothesis testing
- Everything relies on a transformation group G, such that

$$P_{X\sim\mathcal{D}}(X) = P_{X\sim\mathcal{D}}(g(X))$$

for $g \in G$ and $\mathcal{D} \in \Omega_0$

(人間) システレ イテレ

Randomisation (permutation) testing

Theorem (The p-value using randomisation testing)

- Given a pattern function π ,
- a null hypothesis Ω_0 invariant under $g \in G$,
- a fixed $p \in (0, 1)$,
- and the set $\{g(X) : g \in G\}$.

Choose $\sigma(X)$ such that

$$p = \frac{\# \left\{ g \in G : \pi \left(g \left(X \right) \right) \ge \sigma \left(X \right) \right\}}{|G|}$$

Then, $\forall \mathcal{D} \in \Omega_0$:

$$P_{X\sim\mathcal{D}}\left(\pi\left(X\right)\geq\sigma\left(X\right)\right)=p.$$

Subsampling the transformation group

- Usually, the number of group elements is prohibitively large
- Then a subsample G_m of size $|G_m| = m$ is used to estimate p:

$$\widehat{p} = \frac{\# \{ g \in G_m : \pi \left(g \left(X \right) \right) \ge \sigma \left(X \right) \}}{m}$$

• It can be shown (Langford, see references) that with probability $\geq 1 - \delta$ over the random sample G_m ,

$$p \leq \overline{\mathsf{Bin}}(m, \widehat{p}m, \delta)$$

where $\hat{p} \leq \overline{\text{Bin}}(m, \hat{p}m, \delta) \leq \hat{p} + \sqrt{\frac{\ln(\frac{1}{\delta})}{2m}}$ approaches \hat{p} from above for increasing m

イロト イポト イラト イラト

Subsampling the transformation group

- Upper bound the bias p (probabilistically) by random sampling
- E.g.: 20 coin tosses, 2 successes



A (10) < A (10) </p>

- Upper bound Bin (20, 2, 0.1): largest bias that gives probability ≥ 0.1 to see 2 or fewer successes
- Any larger value for *p* would imply that this few observations is unlikely

Subsampling the transformation group

Theorem (The p-value using randomisation testing)

- Given a pattern function π ,
- a null hypothesis Ω_0 invariant under $g \in G$,
- a fixed $\widehat{p} \in (0, 1)$,
- and the set $\{g(X) : g \in G_m\}$ with G_m a random subset of G.

Choose $\sigma(X)$ such that

$$\widehat{p} = \frac{\# \left\{ g \in G_m : \pi \left(g \left(X \right) \right) \ge \sigma \left(X \right) \right\}}{m}$$

Then, with probability $\geq 1 - \delta$ over G_m , $\forall \mathcal{D} \in \Omega_0$:

$$P_{X \sim \mathcal{D}}\left(\pi\left(X\right) \geq \sigma\left(X\right)\right) \leq \overline{Bin}\left(m, \widehat{p}m, \delta\right).$$

I → < ▷ → < ▷ → < ⊇ → < ⊇ → < ⊇ → ○ Q ↔</p>
Bounding the k-family-wise error rate using resampling methods

The FWER The k-FWER A uniform bound The FDP

The FWER

- Thus far:
 - Multiple testing framework
 - How to use transformation invariants of the null hypothesis and randomisation testing to perform single hypothesis testing

• Now: use of randomisation testing for multiple tests, bounding

- FWER
- *k*-FWER
- *k*-FWER uniformly
- FDP

• First the FWER: the probability of at least one false rejection

$$P_{X \sim \mathcal{D}}\left(\exists \alpha \in A \setminus \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\left(X\right)\right)$$

・ロト ・同ト ・ヨト ・ヨト

The FWER The k-FWER A uniform bound The FDP

The FWER

Theorem (Bounding the FWER using randomisation testing)

- Given a pattern class ∏,
- a null hypothesis Ω_0 invariant under $g \in G$,
- a fixed $\widehat{p} \in (0, 1)$,
- and the set $\{g(X) : g \in G_m\}$ with G_m a random subset of G.

Choose $\sigma_{\alpha}\left(X
ight)$ such that

$$\widehat{p} = \frac{\# \left\{ g \in G_m : \exists \alpha : \pi_{\alpha} \left(g \left(X \right) \right) \ge \sigma_{\alpha} \left(X \right) \right\}}{m}$$

Then, with probability $\geq 1 - \delta$ over G_m , $\forall \overline{A} \subseteq A$, $\forall \mathcal{D} \in \Omega_0$ (\overline{A}) :

$$P_{X\sim\mathcal{D}}\left(\exists \alpha\in A\setminus\overline{A}:\pi_{\alpha}\left(X
ight)\geq\sigma_{\alpha}\left(X
ight)
ight)\leq\overline{Bin}\left(m,\widehat{p}m,\delta
ight).$$

Bounding the k-family-wise error rate using resampling methods

The FWER **The k-FWER** A uniform bound The FDP

The k-FWER

- In practice, we are happy to tolerate more than one false rejection
- If only one is tolerated: \widehat{p} is very large for any reasonable choice of $\sigma_{\alpha}\left(X
 ight)$
- Overly conservative very few true rejections (many type II errors!)
- Remedy: upper bound probability to observe at most k false rejections (instead of at most 1): the k-FWER
- Instead of bounding

$$P_{X \sim \mathcal{D}}\left(\exists \alpha \in A \setminus \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\left(X\right)\right)$$

now bound

$$P_{X \sim \mathcal{D}}\left(\exists k \text{ different } \alpha \in A \backslash \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\left(X\right)
ight).$$

・ロト ・同ト ・ヨト ・ヨト

The FWER **The k-FWER** A uniform bound The FDP

The k-FWER

Theorem (Bounding the k-FWER using randomisation testing)

- Given a pattern class ∏,
- a null hypothesis Ω_0 invariant under $g \in G$,
- a fixed $\widehat{p} \in (0, 1)$, and a fixed k,
- and the set $\{g(X) : g \in G_m\}$ with G_m a random subset of G.

Choose $\sigma_{\alpha}\left(X
ight)$ such that

$$\widehat{p} = \frac{\# \left\{ g \in G_m : \exists k \text{ different } \alpha : \pi_{\alpha} \left(g \left(X \right) \right) \geq \sigma_{\alpha} \left(X \right) \right\}}{m}$$

Then, with probability $\geq 1 - \delta$ over $\mathcal{G}_m, \forall \overline{A} \subseteq A, \forall \mathcal{D} \in \Omega_0$ (\overline{A}) :

 $P_{X \sim \mathcal{D}}\left(\exists k \text{ different } \alpha \in A \backslash \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\left(X\right)\right) \leq \overline{Bin}\left(m, \widehat{p}m, \delta\right)$

The FWER The k-FWER A uniform bound The FDP

A uniform bound on the k-FWER

- In practice it may be hard to prespecify the *k*, the number of false rejections tolerated
- However, the previous bound holds for any fixed k, not uniformly for all k
- Can we obtain a uniform bound, for all $k \in \mathcal{K} \subseteq \{1, 2, ..., |A|\}$?
- \bullet Here the set ${\cal K}$ should be chosen so as to include all potentially interesting values

・ロト ・同ト ・ヨト ・ヨト

The FWER The k-FWER A uniform bound The FDP

A uniform bound on the k-FWER

• Instead of considering

$$P_{X\sim\mathcal{D}}\left(\exists k \text{ different } \alpha\in Aackslash\overline{A}:\pi_{lpha}\left(X
ight)\geq\sigma_{lpha}\left(X
ight)
ight)$$

now consider

 $P_{X \sim \mathcal{D}}\left(\exists k \in \mathcal{K} : \exists k \text{ different } \alpha \in A \backslash \overline{A} : \pi_{\alpha}\left(X\right) \geq \sigma_{\alpha}\left(X, k\right)\right)$

- Note that we have different thresholds $\sigma_{\alpha}(X, k)$ for different k (we have to make this explicit now)
- For larger k, the threshold should be smaller (we tolerate more false rejections)
- Hence: the $\sigma_{\alpha}\left(X,k
 ight)$ should be non-increasing

・ロト ・同ト ・ヨト ・ヨト

The FWER The k-FWER A uniform bound The FDP

A uniform bound on the k-FWER

Theorem (Bounding the k-FWER using randomisation testing)

- Given a pattern class Π,
- a null hypothesis Ω_0 invariant under $g \in G$,
- a fixed $\widehat{p} \in (0,1)$, a set $\mathcal{K} \subseteq \{1,2,...,|\mathcal{A}|\}$,
- and the set $\{g(X) : g \in G_m\}$ with G_m a random subset of G.

Choose $\sigma_{\alpha}(X, k)$ non-increasing in k, such that

$$\widehat{p} = \frac{\# \{ g \in G_m : \exists k \in \mathcal{K} : \exists k \text{ different } \alpha : \pi_{\alpha} (g(X)) \ge \sigma_{\alpha} (X, k) \}}{m}$$

Then, with probability $\geq 1 - \delta$ over $\mathcal{G}_m, \forall \overline{A} \subseteq A, \forall \mathcal{D} \in \Omega_0$ (\overline{A}) :

 $P_{X \sim \mathcal{D}} \left(\exists k \in \mathcal{K} : \exists k \text{ different } \alpha \in A \setminus \overline{A} : \pi_{\alpha} \left(X \right) \geq \sigma_{\alpha} \left(X, k \right) \right)$ $\leq \overline{Bin} \left(m, \widehat{p}m, \delta \right).$

Tijl De Bie and John Shawe-Taylor

Bounding the k-family-wise error rate using resampling methods

Multiple testing: the formalism Randomisation testing A zoo of error rates In practice The k-FWER A uniform bound The FDP

The FDP

- We can easily check the total number of rejections for $\sigma_{\alpha}(X, k)$: #rej(k)
- The uniform k-FWER provides an immediate upper bound on the number of false rejections: k-1
- \Rightarrow a lower bound on the number of true rejections as: #trej(k) \geq #rej(k)-(l-1)
- Note however that #trej(k) cannot decrease with increasing k. Hence:

$$\underline{\#\mathsf{trej}(k)} = \max_{l \leq k} \left(\#\mathsf{rej}(l) - (l-1) \right)$$

In particular, <u>#trej(max {k ∈ K})</u> gives a lower bound on the total number of positives

イロト イポト イヨト イヨト

The FWER The k-FWER A uniform bound The FDP

The FDP

- We already have an upper bound on the false rejections of k-1
- But, we can get a possibly tighter one: $\overline{\#}$ frej(k)=#rej(k)- $\frac{\#}{trej}(k)$
- Based on these quantities, we can bound the False Discovery Proportion as:

$$\overline{\mathsf{FDP}(k)} = rac{\#\mathsf{frej}(k)}{\#\mathsf{rej}(k)}$$

イロト イポト イヨト イヨト

Technical details Experiments Conclusions

Technical details

- 2 microarray datasets: Alon and Golub
- Compare point-wise bound for the k-FWER with the uniform bound
- $\bullet\,$ Do this for different choices of ${\cal K}\,$
- One technical issue remains: how to choose $\sigma_{\alpha}(X, \cdot)$?
- In principle, this is entirely free
- For convenience, we choose $\sigma_{\alpha}\left(X,\cdot\right)=\sigma\left(X,\cdot\right)$ independent of α

(日) (同) (三) (三)

Technical details Experiments Conclusions

How to choose the thresholds?

- Compute $\{\pi_{\alpha}\left(g\left(X\right)\right):g\in G_{m}\subseteq G, \alpha\in A\}$
- For each $g \in G_m$, sort $\{\pi_{\alpha} (g (X)) : \alpha \in A\}$ in decreasing ordering
- Then, pick the k'th largest numbers from $\{\pi_{\alpha} (g(X)) : \alpha \in A\}$ (there are m of them), and put this in a set S(X, k)
- Interpretation: S(X, k) contains an empirical estimate of the distribution of the k'th strongest correlation with the labels, under Ω_0
- Choose $\sigma^{r}\left(X,k
 ight)=$ the r'th largest value of $\mathcal{S}\left(X,k
 ight)$
- Result: non-increasing functions of k that vary roughly as the sorted {π_α (g (X)) : α ∈ A}

イロン 不同 とくほう イロン

Technical details Experiments Conclusions

Alon and Golub datasets

• Statistics of data sets

# genes	# arrays tissue 1	# arrays tissue 2
2000	40	22
7129	47	25

• Pattern functions used: the Wilcoxon rank sum test statistic

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

Technical details Experiments Conclusions

Alon and Golub datasets



Tijl De Bie and John Shawe-Taylor

Bounding the k-family-wise error rate using resampling methods

Technical details Experiments Conclusions

Achievements

- Based on randomisation testing, which ensures:
 - More limited applicability (need to identify G)
 - But all dependencies are adequately taken into account
 - For that reason: practical relevance!
- Finite sample bound for randomisation testing
 - Separates randomness in g and X (δ and $\overline{Bin}(m, \hat{p}m, \delta)$)
 - This is relevant when one sample g is used for several multiple tests
- Uniform bounds over regions of k in k-FWER bounds
 - Relevant where a priori unclear how large k should be taken
 - Allows to bound FDP

イロト イポト イヨト イヨト