

Improved post hoc bounds for localized signal

Guillermo Durand

Joint work with Gilles Blanchard, Pierre Neuvial and Etienne Roquain
Young Researchers' Meeting in Mathematical Statistics 2018



23/09/2018

Table of contents

1. Motivation

2. Background

3. New methods

4. Simulations

5. Conclusion

Replication crisis

Toy example

GWAS study with 10^6 genetic variants, apply a multiple testing procedure (MTP) over the 1000 smallest p -values

p -hacking

- ▶ Pre-selecting variables that seem significant, exclude others from experiment
- ▶ Theoretical results no longer hold
- ▶ Results poorly interpretable and non reproducible

Selective inference

Develop new methods that account for the selection step to provide statistical guarantees.

- 1 Conditionally to the selection event [Fithian, Sun, et al. (2017), Lee, Sun, et al. (2016), and Tibshirani, Taylor, et al. (2016)]
- 2 Simultaneously over all possible selection events [Goeman and Solari (2011), Berk, Brown, et al. (2013), and Blanchard, Neuvial, et al. (2018)]

Table of contents

1. Motivation

2. Background

3. New methods

4. Simulations

5. Conclusion

Multiple testing setting

- ▶ Data $X \in (\mathcal{X}, \mathfrak{X})$ with $X \sim P \in \mathcal{P}$ a collection of distributions, P unknown
- ▶ m null hypotheses $H_{0,i}$ on P which are subsets of \mathcal{P}
- ▶ m is large!
- ▶ $\mathcal{H}_0 = \{i : P \in H_{0,i}\}$: $i \in \mathcal{H}_0 \Leftrightarrow H_{0,i}$ is true
- ▶ m p -values $p_i = p_i(X)$ such that $p_i \succeq \mathcal{U}([0, 1])$ if $i \in \mathcal{H}_0$
 - ▶ Each p_i provides an α level test : $\mathbb{P}_{P \in H_{0,i}}(p_i \leq \alpha) \leq \alpha$
- ▶ Classic theory: form a rejection set R with a guarantee on $V(R) = |R \cap \mathcal{H}_0|$

Our goal

A simultaneous inference problem

Confidence bounds on any set of selected variables

A confidence bound is a (random) function \hat{V} such that

$$\mathbb{P} \left(\forall S \subset \mathbb{N}_m, V(S) \leq \hat{V}(S) \right) \geq 1 - \alpha$$

- ▶ Hence for any selected \hat{S} , $\mathbb{P} \left(V(\hat{S}) \leq \hat{V}(\hat{S}) \right) \geq 1 - \alpha$ holds
- ▶ Originates from [\[Genovese and Wasserman \(2006\) and Meinshausen \(2006\)\]](#)
- ▶ Not a classic MTP: a guarantee over any selected set instead of a rejected set
- ▶ Also named “post hoc inference”

BNR methodology

[Blanchard, Neuvial, et al. (2018)]

Key concept: reference family

- ▶ $\mathfrak{R} = (R_k, \zeta_k)$ (random) such that Joint Error Rate (JER) control:

$$\mathbb{P}(\forall k, |R_k \cap \mathcal{H}_0| \leq \zeta_k) \geq 1 - \alpha$$

- ▶ Confidence bound only on the members of \mathfrak{R}
- ▶ \implies Derivation of a global confidence bound

Two different bounds

- ▶ $V_{\mathfrak{R}}^*(S) = \max \{|S \cap A|, \forall k, |R_k \cap A| \leq \zeta_k\}$ difficult to compute
- ▶ $\overline{V}_{\mathfrak{R}}(S) = \min_k (\zeta_k + |S \setminus R_k|) \wedge |S|$ worst but easy to compute

BNR methodology

A flexible and unified approach

- ▶ Compatible with previous works, like the closed testing approach of [Goeman and Solari (2011)]
- ▶ BNR approach: $\zeta_k = k - 1$ and find $R_k = \{i : p_i < t_k\}$ such that JER control. Example: $t_k = \alpha k/m$ (Simes inequality)
 - ▶ JER control becomes “simultaneous k -FWER control”
- ▶ Property: if R_k nested, then $\overline{V}_{\mathfrak{R}} = V_{\mathfrak{R}}^*$

In this talk

- ▶ Consider the realistic case where the signal is localized in some regions (e.g. active genes in same chromosomes)
- ▶ Accordingly, find adapted new structures where $V_{\mathfrak{R}}^*$ is also easy to compute
- ▶ Fix the regions (i.e. R_k non random) and estimate the true nulls inside them (i.e. ζ_k random): opposite of [Blanchard, Neuvial, et al. (2018)]

Table of contents

1. Motivation

2. Background

3. New methods

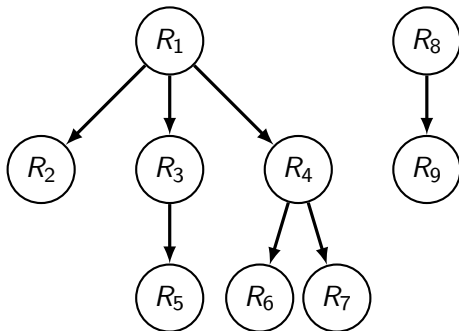
4. Simulations

5. Conclusion

Forest structure

New structure

- ▶ $\forall k, k' \in \mathcal{K}, R_k \cap R_{k'} \in \{R_k, R_{k'}, \emptyset\}$
- ▶ Includes nested families or totally disjoint families
- ▶ Accommodates to signal localization through the “leaves”



Forest structure

Important property

Property

There is a partition $(L_n)_{1 \leq n \leq N}$ of \mathbb{N}_m (the leaves) such that for each $k \in \mathcal{K}$, there exists some (i, j) with $1 \leq i \leq j \leq N$ and $R_k = L_{i:j}$, where we denote

$$L_{i:j} = \bigcup_{i \leq n \leq j} L_n$$

Identification:

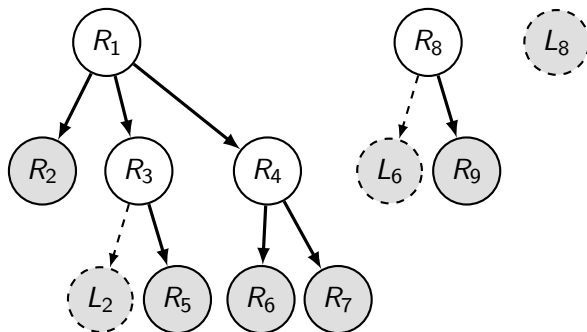
$$\mathfrak{R} = (R_k, \zeta_k)_{k \in \mathcal{K}} \quad \text{or} \quad \mathfrak{R} = (L_{i:j}, \zeta_{i:j})_{(i,j) \in \mathcal{K}}$$

\implies leaves represent the thinnest division possible of the structure

Forest structure

Other important property

Each forest structure can be completed to includes all leaves



New post hoc bounds

Goal: compute $V_{\mathfrak{R}}^*$ easily with forest structure

Definition

For any $q \leq K$,

$$\tilde{V}_{\mathfrak{R}}^q(S) = \min_{Q \subset \mathcal{K}, |Q| \leq q} \left(\sum_{k \in Q} \zeta_k \wedge |S \cap R_k| + \left| S \setminus \bigcup_{k \in Q} R_k \right| \right),$$

and

$$\tilde{V}_{\mathfrak{R}}(S) = \tilde{V}_{\mathfrak{R}}^K(S).$$

Property

$$V_{\mathfrak{R}}^*(S) \leq \tilde{V}_{\mathfrak{R}}(S) \leq \tilde{V}_{\mathfrak{R}}^{K-1}(S) \leq \dots \leq \tilde{V}_{\mathfrak{R}}^2(S) \leq \tilde{V}_{\mathfrak{R}}^1(S) = \overline{V}_{\mathfrak{R}}(S)$$

Main results

Compute $V_{\mathfrak{H}}^*$ easily with forest structure

Theorem

$$V_{\mathfrak{H}}^*(S) = \tilde{V}_{\mathfrak{H}}(S)$$

More precisely,

$$V_{\mathfrak{H}}^*(S) = \tilde{V}_{\mathfrak{H}}^{\ell}(S),$$

with ℓ = number of leaves (without completion).

Corollary

$\ell = 1$ for nested families and BNR property is recovered

Main results

Compute $V_{\mathfrak{R}}^*$ easily with forest structure

Lemma

There is a simple algorithm to compute $\tilde{V}_{\mathfrak{R}}$ if \mathfrak{R} is complete.

Lemma

Completing the family does not change $V_{\mathfrak{R}}^*$ and $\tilde{V}_{\mathfrak{R}}$.

Corollary

There is a simple algorithm to compute $V_{\mathfrak{R}}^*(S)$ by:

- 1 Completing the family
- 2 Travel across the forest from the leaves

True nulls estimation

Under independence

- ▶ K deterministic regions, let $C = \sqrt{\frac{1}{2} \log \left(\frac{K}{\alpha} \right)}$
- ▶ $\zeta_k = |R_k| \wedge \min_{t \in [0,1)} \left[\frac{C}{2(1-t)} + \left(\frac{C^2}{4(1-t)^2} + \frac{\sum_{i \in R_k} \mathbf{1}\{p_i > t\}}{1-t} \right)^{1/2} \right]^2$
- ▶ Comes from carefully handling the DKWM inequality [Dvoretzky, Kiefer, et al. (1956) and Massart (1990)]
- ▶ Replace $\min_{t \in [0,1)}$ and t above by $\min_{0 \leq \ell \leq s}$ and $p_{(\ell)}$ for practical computation
- ▶ α/K instead of α in C : union bound
- ▶ Dependence on α only through a log
- ▶ $\zeta_k > 0$ (entry cost)

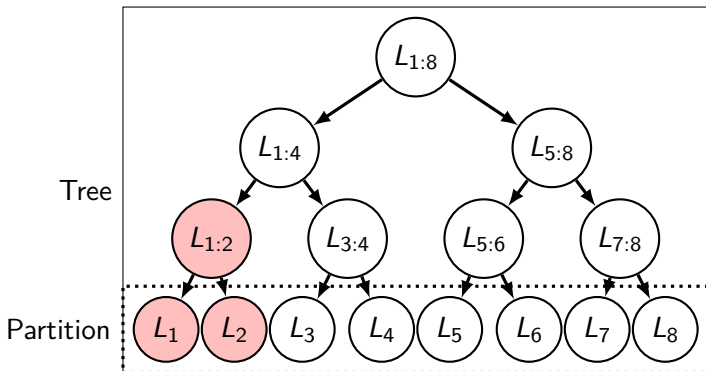
Table of contents

1. Motivation
2. Background
3. New methods
4. Simulations
5. Conclusion

Comparison of 3 bounds

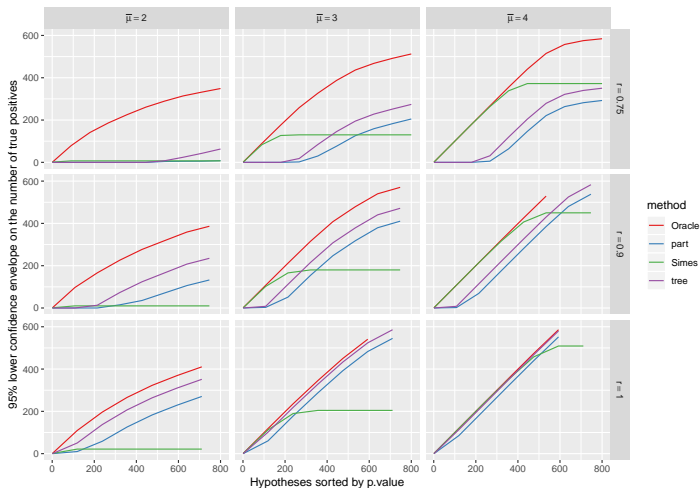
Simes bound of BNR, and 2 new

- ▶ V_{tree} and V_{part} , from a complete binary tree or only the partition of leaves
- ▶ Signal in adjacent leaves, good performance of V_{tree} expected
- ▶ Parameters: signal $\bar{\mu}$ and signal proportion in active leaves r



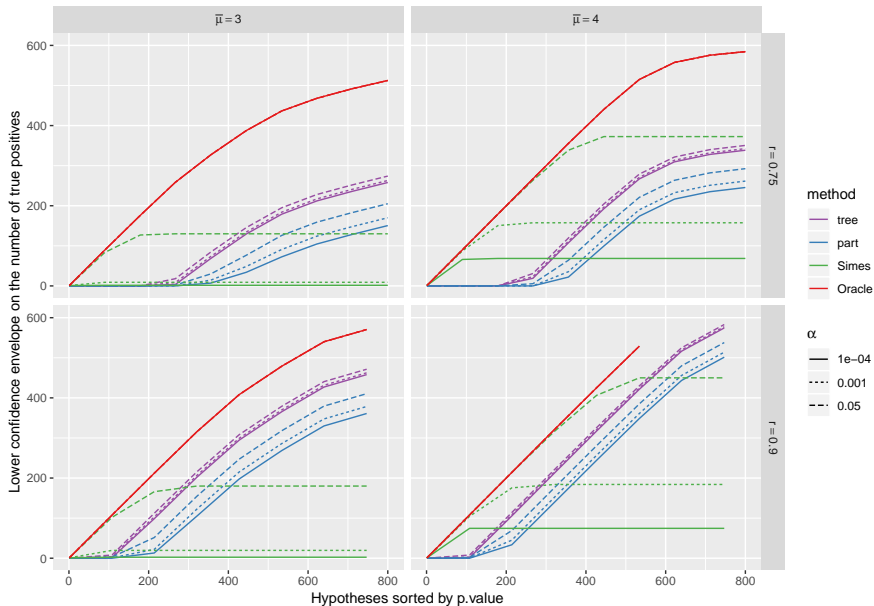
Comparison of 3 bounds

- ▶ The choice of S favors the Simes bound of BNR
- ▶ But for large r , new bounds better
- ▶ V_{tree} better than V_{part} as expected, despite a worst union bound correction



Comparison of 3 bounds

Influence of α



New hybrid bound suggested by the simulations

- ▶ $V_{\text{hybrid}}^{\gamma}(\alpha, S) = \min(V_{\text{Simes}}((1 - \gamma)\alpha, S), V_{\text{tree}}(\gamma\alpha, S))$
- ▶ $\gamma = 0.02$: favors Simes, not a problem because V_{tree} is little sensitive to small α

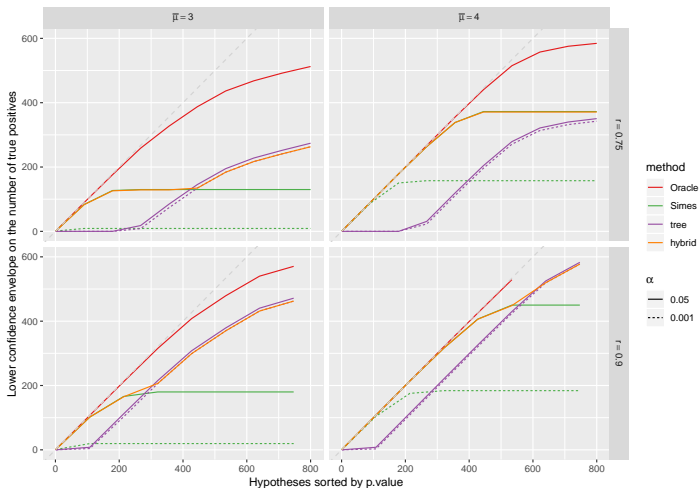


Table of contents

1. Motivation
2. Background
3. New methods
4. Simulations
5. Conclusion

Conclusion

- ▶ DKWM inequality involves independence
- ▶ Other over-estimators of true nulls ? [Blanchard, Neuvial, et al. (2018) and Hemerik and Goeman (2018)]
- ▶ Other families combining BNR approach and deterministic regions ?
 - ▶ $\mathfrak{R} = (R_{k,i_k}, \zeta_{k,i_k})_{\substack{k \in \mathcal{K} \\ 1 \leq i_k \leq |R_k|}}$

Preprint available: [arXiv:1807.01470](https://arxiv.org/abs/1807.01470)

R package available from github: [sansSouci](https://github.com/sansSouci)

Bibliography I

- Berk, Richard et al. (2013). "Valid post-selection inference". In: *The Annals of Statistics* 41.2, pp. 802–837.
- Blanchard, Gilles, Pierre Neuvial, and Etienne Roquain (2018). "A unified approach to post hoc false positive control". In: *arXiv preprint arXiv:1703.02307*.
- Dvoretzky, Aryeh, Jack Kiefer, and Jacob Wolfowitz (1956). "Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator". In: *The Annals of Mathematical Statistics*, pp. 642–669.
- Fithian, William, Dennis Sun, and Jonathan Taylor (2017). "Optimal inference after model selection". In: *arXiv preprint arXiv:1410.2597*.
- Genovese, Christopher R and Larry Wasserman (2006). "Exceedance control of the false discovery proportion". In: *Journal of the American Statistical Association* 101.476, pp. 1408–1417.
- Goeman, Jelle J and Aldo Solari (2011). "Multiple testing for exploratory research". In: *Statistical Science*, pp. 584–597.
- Hemerik, Jesse and Jelle J Goeman (2018). "False discovery proportion estimation by permutations: confidence for significance analysis of microarrays". In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80.1, pp. 137–155.
- Lee, Jason D et al. (2016). "Exact post-selection inference, with application to the lasso". In: *The Annals of Statistics* 44.3, pp. 907–927.

Bibliography II

- Marcus, Ruth, Peritz Eric, and K Ruben Gabriel (1976). "On closed testing procedures with special reference to ordered analysis of variance". In: *Biometrika* 63.3, pp. 655–660.
- Massart, Pascal (1990). "The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality". In: *The Annals of Probability*, pp. 1269–1283.
- Meinshausen, Nicolai (2006). "False discovery control for multiple tests of association under general dependence". In: *Scandinavian Journal of Statistics* 33.2, pp. 227–237.
- Tibshirani, Ryan J et al. (2016). "Exact post-selection inference for sequential regression procedures". In: *Journal of the American Statistical Association* 111.514, pp. 600–620.

Classical theory

Family-Wise Error Rate (FWER)

- ▶ $\text{FWER}(R) = \mathbb{P}(V(R) > 0)$
- ▶ Bonferroni method: reject all $p_i \leq \frac{\alpha}{m}$ (union bound)
- ▶ Variant: $k\text{-FWER}(R) = \mathbb{P}(V(R) \geq k)$

False Discovery Rate (FDR)

- ▶ $\text{FDR}(R) = \mathbb{E} \left[\frac{V(R)}{|R| \vee 1} \right]$
- ▶ Benjamini-Hochberg method for positive dependence

Closed testing for post hoc inference

Designed for FWER control [Marcus, Eric, et al. (1976)]

- ▶ Form $H_{0,I} = \bigcap_{i \in I} H_{0,i}$ all intersection hypotheses
- ▶ Have a collection of α level local test ϕ_I
- ▶ Examples:
 - ▶ Bonferroni test $\phi_I = 1$ if $\exists i \in I : p_i \leq \alpha/|I|$
 - ▶ Simes test $\phi_I = 1$ if $\exists i \in I : p_{(i:I)} \leq \alpha i/|I|$ (under PRDS)
- ▶ Test $H_{0,I}$ only if all $H_{0,J}$, $J \supseteq I$, are rejected
- ▶ Reject the individual hypotheses $H_{0,i}$ such that $H_{0,\{i\}}$ has been rejected that way
- ▶ Then $\text{FWER}(\text{Closed testing}) \leq \alpha$

Closed testing for post hoc inference

[Goeman and Solari (2011)]

Main idea

The closed testing provides more information than just the individual rejects:

- ▶ Let \mathcal{X} the set of all I such that we rejected $H_{0,I}$
- ▶ Simultaneous guarantee over all $H_{0,I}$, $I \in \mathcal{X}$:

$$\mathbb{P}(\forall I \in \mathcal{X}, H_{0,I} \text{ is false}) \geq 1 - \alpha$$

Confidence bound derivation:

- ▶ $V_{\text{GS}}(S) = \max_{\substack{I \subseteq S \\ I \notin \mathcal{X}}} |I|$ is a confidence bound because

$$\begin{aligned} \exists S, |S \cap \mathcal{H}_0| > V_{\text{GS}}(S) &\implies \exists S, S \cap \mathcal{H}_0 \in \mathcal{X} \\ &\implies \exists I \in \mathcal{X}, H_{0,I} \text{ is true} \end{aligned}$$

- ▶ $V_{\text{GS}}(S) = V_{\mathfrak{R}}^*(S)$ with $\mathfrak{R} = (I, |I| - 1)_{I \in \mathcal{X}}$

DKWM use

- ▶ Let $S \subset \mathbb{N}_m$
- ▶ $N_t(S) = \sum_{i \in S} \mathbf{1}\{p_i(X) > t\}$
- ▶ $v = |S \cap \mathcal{H}_0|$

$$v \leq \min_{t \in [0,1)} \left(\frac{\sqrt{\log(1/\lambda)/2}}{2(1-t)} + \left\{ \frac{\log(1/\lambda)/2}{4(1-t)^2} + \frac{N_t(S)}{1-t} \right\}^{1/2} \right)^2$$

comes from

$$v^{-1} \sum_{i=1}^v \mathbf{1}\{U_i > t\} - (1-t) \geq -\sqrt{\log(1/\lambda)/(2v)}, \quad \forall t \in [0, 1],$$

with probability at least $1 - \lambda$ (U_1, \dots, U_v i.i.d. uniform, $N_t(S)$ dominates $\sum_{i=1}^v \mathbf{1}\{U_i > t\}$ by independence)

- ▶ $S = R_k$ and $\lambda = \alpha/K$ (union bound)

Forest algorithm

