


A fast algorithm to compute a curve of confidence upper bounds for the False Discovery Proportion using a reference family with a forest structure

Guillermo Durand ¹ Laboratoire de Mathématiques d'Orsay, Université Paris-Saclay

Date published: 2024-07-05 Last modified: 2024-07-05

Abstract

This paper presents a new algorithm (and an additional trick) that allows to compute fastly an entire curve of post hoc bounds for the False Discovery Proportion when the underlying bound $V_{\mathfrak{R}}^*$ construction is based on a reference family \mathfrak{R} with a forest structure à la [Durand et al. \(2020\)](#). By an entire curve, we mean the values $V_{\mathfrak{R}}^*(S_1), \dots, V_{\mathfrak{R}}^*(S_m)$ computed on a path of increasing selection sets $S_1 \subsetneq \dots \subsetneq S_m, |S_t| = t$. The new algorithm leverages the fact that going from S_t to S_{t+1} is done by adding only one hypothesis.

Keywords: multiple testing, algorithmic, post hoc inference, false discovery proportion, confidence bound

1	Contents	
2	1 Introduction	2
3	2 Notation and reference family methodology	3
4	2.1 Multiple testing notation	3
5	2.2 Post hoc bounds with reference families	3
6	2.3 Deterministic regions with a forest structure	3
7	3 New algorithms	3
8	3.1 Pruning the forest	3
9	3.2 Fast algorithm to compute a curve of confidence bounds on a path of selection sets	4
10	4 Numerical experiments	4
11	5 Conclusion	4
12	6 Acknowledgments	4
13	References	5
14	Session information	6

¹Corresponding author: guillermo.durand@universite-paris-saclay.fr

1 Introduction

Multiple testing theory is often used for exploratory analysis, like Genome-Wide Association Studies, where multiple features are tested to find promising ones. Classical multiple testing theory like Family-Wise Error Rate (FWER) control or False Discovery Rate (FDR) control (Benjamini and Hochberg, 1995) can be used, but a more recent trend consists in the computation of post hoc bounds, also named post selection bounds or confidence envelopes, for the number of false positives, or, equivalently, for the False Discovery Proportion (FDP). This approach is notably advocated for in the context of exploratory research by (Goeman and Solari, 2011, Section 1).

Mathematically speaking, a confidence upper bound (we prefer to say upper bound instead of envelope for obvious reasons) is a function $\hat{V} : \mathcal{P}(\mathbb{N}_m^*) \rightarrow \mathbb{N}_m$, where $\mathbb{N}_m = \{0, \dots, m\}$, $\mathbb{N}_m^* = \{1, \dots, m\}$ and m is the number of hypotheses, such that

$$\forall \alpha \in]0, 1[, \mathbb{P} \left(\forall S \subseteq \mathbb{N}_m^*, |S \cap \mathcal{H}_0| \leq \hat{V}(S) \right) \geq 1 - \alpha. \quad (1)$$

Here, α is a target error rate and \mathcal{H}_0 is the set of hypotheses indices that are true null hypotheses. Note that the construction of \hat{V} depends on α and on the random data X and the dependence is omitted to lighten notation and because there is no ambiguity. The meaning of Equation 1 is that \hat{V} provides an upper bound of the number of null hypotheses in S for any selection set $S \subseteq \mathbb{N}_m^*$, which allows the user to perform post hoc selection on their data without breaching the statistical guarantee. Also note that by dividing by $|S| \wedge 1$ in Equation 1 we also get a confidence bound for the FDP:

$$\forall \alpha \in]0, 1[, \mathbb{P} \left(\forall S \subseteq \mathbb{N}_m^*, \text{FDP}(S) \leq \frac{\hat{V}(S)}{|S| \wedge 1} \right) \geq 1 - \alpha. \quad (2)$$

So post hoc bounds provide ways to construct FDP-controlling sets instead of FDR-controlling sets, which is much more desirable given the nature of the FDR as an expected value. See for example (Bogdan et al., 2015, Figure 4) for a credible example where the FDR is controlled but the FDP has a highly undesirable behavior (either 0 because no discoveries at all are made, either higher than the target level).

The first confidence bounds are found in (Genovese and Wasserman, 2006) and (Meinshausen, 2006), although, in the latter, only for selection sets of the form $\{i \in \mathbb{N}_m : P_i \leq t\}$ where P_i is the p -value associated to the null hypothesis $H_{0,i}$. In (Goeman and Solari, 2011) the authors re-wrote the generic construction of (Genovese and Wasserman, 2006) in terms of closed testing (Marcus et al., 1976), proposed several practical constructions and sparked a new interest in multiple testing procedures based on confidence envelopes. This work was followed by a prolific series of works like (Meijer et al., 2015) and (Vesely et al., 2023). In (Blanchard et al., 2020), the authors introduce the new point of view of reference families (see Section 2.2) to construct post hoc bounds, and show the links between this meta-technique and the closed testing one, along with new bounds.

Following the reference family trail, in (Durand et al., 2020) the authors introduce new reference families with a special set-theoretic constraint that allows an efficient computation of the bound denoted by $V_{\mathfrak{R}}^*$ on a single selection set S . The problem is that one often wants to compute $V_{\mathfrak{R}}^*$ on a whole path of selection sets $(S_t)_{t \in \mathbb{N}_m^*}$, for example the hypotheses attached to the t smallest p -values. Whereas the algorithm provided the aforementioned work (Durand et al., 2020, Algorithm 1) is fast for a single evaluation, it is slow and inefficient to repeatedly call it to compute each $V_{\mathfrak{R}}^*(S_t)$. If the S_t 's are nested, and growing by one, that is $S_1 \subsetneq \dots \subsetneq S_m$ and $|S_t| = t$, there is a way to efficiently compute $(V_{\mathfrak{R}}^*(S_t))_{t \in \mathbb{N}_m}$ by leveraging the nested structure.

This is the main contribution of the present paper: a new and fast algorithm computing the curve $(V_{\mathfrak{R}}^*(S_t))_{t \in \mathbb{N}_m}$ for a nested path of selection sets, that is presented in Section 3.2. An additional

algorithm that can speed up computations both for the single-evaluation algorithm and the new curve-evaluation algorithm is also presented, in Section 3.1. In Section 2.1, all necessary notation and vocabulary is re-introduced, most of it being the same as in (Durand et al., 2020). Finally, a few numerical experiments are presented in Section 4 to demonstrate the computation time gain.

2 Notation and reference family methodology

2.1 Multiple testing notation

2.2 Post hoc bounds with reference families

2.3 Deterministic regions with a forest structure

3 New algorithms

3.1 Pruning the forest

Algorithm 1 Pruning of \mathfrak{R}

```

1: procedure PRUNING( $\mathfrak{R} = (R_k, \zeta_k)_{k \in \mathcal{K}}$  with  $\mathfrak{R}$  complete)
2:    $\mathcal{K}^{\text{pr}} \leftarrow \mathcal{K}$ 
3:    $H \leftarrow \max_{k \in \mathcal{K}} \phi(k)$  ▷ maximum depth
4:   for  $h = H - 1, \dots, 1$  do
5:      $\mathcal{K}^h \leftarrow \{k \in \mathcal{K} : \phi(k) = h\}$ 
6:      $\text{newVec} \leftarrow (0)_{k \in \mathcal{K}^h}$ 
7:     for  $k \in \mathcal{K}^h$  do
8:        $\text{Succ}_k \leftarrow \{k' \in \mathcal{K}^{h+1} : R_{k'} \subseteq R_k\}$ 
9:       if  $\text{Succ}_k = \emptyset$  then
10:         $\text{newVec}_k \leftarrow \zeta_k$ 
11:       else
12:        if  $\zeta_k \geq \sum_{k' \in \text{Succ}_k} \text{Vec}_{k'}$  then
13:          $\mathcal{K}^{\text{pr}} \leftarrow \mathcal{K}^{\text{pr}} \setminus \{k\}$ 
14:        end if
15:         $\text{newVec}_k \leftarrow \min(\zeta_k, \sum_{k' \in \text{Succ}_k} \text{Vec}_{k'})$ 
16:       end if
17:     end for
18:      $\text{Vec} \leftarrow \text{newVec}$ 
19:   end for
20:   return  $(\mathcal{K}^{\text{pr}}, \sum_{k \in \mathcal{K}^1} \text{Vec}_k)$ 
21: end procedure

```

Algorithm 2 Formal computation of $(V_{\mathfrak{R}}^*(S_t))_{0 \leq t \leq m}$

```
1: procedure CURVE( $\mathfrak{R} = (R_k, \zeta_k)_{k \in \mathcal{K}}$  with  $\mathfrak{R}$  complete, path  $(S_t)_{1 \leq t \leq m}$  with  $S_t = \{i_1, \dots, i_t\}$ )
2:    $\mathcal{P}^0 \leftarrow \{(i, i) : 1 \leq i \leq n\}$  ▷ the set of all atoms indices
3:    $\mathcal{K}_0^- \leftarrow \{k \in \mathcal{K} : \zeta_k = 0\}$ 
4:    $\eta_k^0 \leftarrow 0$  for all  $k \in \mathcal{K}$ 
5:   for  $t = 1, \dots, m$  do
6:     if  $i_t \in \bigcup_{k \in \mathcal{K}_{t-1}^-} R_k$  then
7:        $\mathcal{P}^t \leftarrow \mathcal{P}^{t-1}$ 
8:        $\mathcal{K}_t^- \leftarrow \mathcal{K}_{t-1}^-$ 
9:        $\eta_k^t \leftarrow \eta_k^{t-1}$  for all  $k \in \mathcal{K}$ 
10:    else
11:      for  $h = 1, \dots, h_{\max}(t)$  do
12:         $\eta_{k^{(t,h)}}^t \leftarrow \eta_{k^{(t,h)}}^{t-1} + 1$ 
13:        if  $\eta_{k^{(t,h)}}^t < \zeta_k$  then
14:          Pass
15:        else
16:           $h_t^f \leftarrow h.$ 
17:           $\mathcal{P}^t \leftarrow \left( \mathcal{P}^{t-1} \setminus \{k \in \mathcal{P}^{t-1} : R_k \subseteq R_{k^{(t,h_t^f)}}\} \right) \cup \{k^{(t,h_t^f)}\}$ 
18:           $\mathcal{K}_t^- \leftarrow \mathcal{K}_{t-1}^- \cup \{k^{(t,h_t^f)}\}$ 
19:          Break the loop
20:        end if
21:      end for
22:      if the loop has been broken then
23:         $\eta_k^t \leftarrow \eta_k^{t-1}$  for all  $k \in \mathcal{K}$  not visited during the loop, that is all  $k \notin \{k^{(t,h)}, 1 \leq h \leq h_t^f\}$ 
24:      else
25:         $\mathcal{P}^t \leftarrow \mathcal{P}^{t-1}$ 
26:         $\mathcal{K}_t^- \leftarrow \mathcal{K}_{t-1}^-$ 
27:         $\eta_k^t \leftarrow \eta_k^{t-1}$  for all  $k \in \mathcal{K}$  not visited during the loop, that is all  $k \notin \{k^{(t,h)}, 1 \leq h \leq$ 
28:           $h_{\max}(t)\}$ 
29:        end if
30:      end for
31:    return  $\mathcal{P}^t, \eta_k^t$  for all  $t = 1, \dots, m$  and  $k \in \mathcal{K}$ 
32: end procedure
```

66 **3.2 Fast algorithm to compute a curve of confidence bounds on a path of selection**
67 **sets**

68 **4 Numerical experiments**

69 **5 Conclusion**

70 **6 Acknowledments**

71 This work has been supported by ANR-20-IDEES-0002 (PIA), ANR-19-CHIA-0021 (BISCOTTE),
72 ANR-23-CE40-0018 (BACKUP) and ANR-21-CE23-0035 (ASCAI).

Algorithm 3 Implementation of $(V_{\mathfrak{R}}^*(S_t))_{0 \leq t \leq m}$

```
1: procedure CURVE( $\mathfrak{R} = (R_k, \zeta_k)_{k \in \mathcal{K}}$  with  $\mathfrak{R}$  complete, path  $(S_t)_{1 \leq t \leq m}$  with  $S_t = \{i_1, \dots, i_t\}$ )
2:    $V_0 \leftarrow 0$ 
3:    $\mathcal{K}^- \leftarrow \{k \in \mathcal{K} : \zeta_k = 0\}$ 
4:    $\eta_k \leftarrow 0$  for all  $k \in \mathcal{K}$ 
5:   for  $t = 1, \dots, m$  do
6:     if  $i_t \in \bigcup_{k \in \mathcal{K}^-} R_k$  then
7:        $V_t \leftarrow V_{t-1}$ 
8:     else
9:       for  $h = 1, \dots, h_{\max}(t)$  do
10:        find  $k^{(t,h)} \in \mathcal{K}^h$  such that  $i_t \in R_{k^{(t,h)}}$ 
11:         $\eta_{k^{(t,h)}} \leftarrow \eta_{k^{(t,h)}} + 1$ 
12:        if  $\eta_{k^{(t,h)}} < \zeta_k$  then
13:          pass
14:        else
15:           $\mathcal{K}^- \leftarrow \mathcal{K}^- \cup \{k^{(t,h)}\}$ 
16:          break the loop
17:        end if
18:      end for
19:       $V_t \leftarrow V_{t-1} + 1$ 
20:    end if
21:  end for
22:  return  $(V_t)_{1 \leq t \leq m}$ 
23: end procedure
```

References

- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. Roy. Statist. Soc. Ser. B*, 57(1):289–300, 1995. ISSN 0035-9246. URL <https://www.jstor.org/stable/2346101>.
- Gilles Blanchard, Pierre Neuvial, and Etienne Roquain. Post hoc confidence bounds on false positives using reference families. *Ann. Statist.*, 48(3):1281–1303, 2020. ISSN 0090-5364. doi: 10.1214/19-AOS1847. URL <https://doi.org/10.1214/19-AOS1847>.
- Małgorzata Bogdan, Ewout van den Berg, Chiara Sabatti, Weijie Su, and Emmanuel J. Candès. SLOPE—adaptive variable selection via convex optimization. *Ann. Appl. Stat.*, 9(3):1103–1140, 2015. ISSN 1932-6157, 1941-7330. doi: 10.1214/15-AOAS842. URL <https://doi.org/10.1214/15-AOAS842>.
- Guillermo Durand, Gilles Blanchard, Pierre Neuvial, and Etienne Roquain. Post hoc false positive control for structured hypotheses. *Scand. J. Stat.*, 47(4):1114–1148, 2020. ISSN 0303-6898. doi: 10.1111/sjos.12453. URL <https://doi.org/10.1111/sjos.12453>.
- Christopher R. Genovese and Larry Wasserman. Exceedance control of the false discovery proportion. *J. Amer. Statist. Assoc.*, 101(476):1408–1417, 2006. ISSN 0162-1459. doi: 10.1198/016214506000000339. URL <https://doi.org/10.1198/016214506000000339>.
- Jelle J. Goeman and Aldo Solari. Multiple testing for exploratory research. *Statist. Sci.*, 26(4):584–597, 2011. ISSN 0883-4237. doi: 10.1214/11-STS356. URL <https://doi.org/10.1214/11-STS356>.
- Ruth Marcus, Eric Peritz, and K. R. Gabriel. On closed testing procedures with special reference

to ordered analysis of variance. *Biometrika*, 63(3):655–660, 1976. ISSN 0006-3444. doi: 10.1093/biomet/63.3.655. URL <https://doi.org/10.1093/biomet/63.3.655>.

Rosa J. Meijer, Thijmen J. P. Krebs, and Jelle J. Goeman. A region-based multiple testing method for hypotheses ordered in space or time. *Stat. Appl. Genet. Mol. Biol.*, 14(1):1–19, 2015. ISSN 2194-6302. doi: 10.1515/sagmb-2013-0075. URL <https://doi.org/10.1515/sagmb-2013-0075>.

Nicolai Meinshausen. False discovery control for multiple tests of association under general dependence. *Scand. J. Statist.*, 33(2):227–237, 2006. ISSN 0303-6898. doi: 10.1111/j.1467-9469.2005.00488.x. URL <https://doi.org/10.1111/j.1467-9469.2005.00488.x>.

Anna Vesely, Livio Finos, and Jelle J. Goeman. Permutation-based true discovery guarantee by sum tests. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 85(3):664–683, 2023. ISSN 1369-7412,1467-9868. doi: 10.1093/jrsssb/qkad019. URL <https://doi.org/10.1093/jrsssb/qkad019>.

Session information

```
R version 4.4.0 (2024-04-24)
Platform: x86_64-pc-linux-gnu
Running under: Ubuntu 22.04.4 LTS

Matrix products: default
BLAS:   /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/libopenblas-p0.3.20.so; LAPACK version 3.10.0

locale:
 [1] LC_CTYPE=C.UTF-8      LC_NUMERIC=C           LC_TIME=C.UTF-8
 [4] LC_COLLATE=C.UTF-8    LC_MONETARY=C.UTF-8    LC_MESSAGES=C.UTF-8
 [7] LC_PAPER=C.UTF-8      LC_NAME=C               LC_ADDRESS=C
[10] LC_TELEPHONE=C        LC_MEASUREMENT=C.UTF-8 LC_IDENTIFICATION=C

time zone: UTC
tzcode source: system (glibc)

attached base packages:
[1] stats      graphics  grDevices datasets  utils      methods    base

loaded via a namespace (and not attached):
 [1] compiler_4.4.0    fastmap_1.1.1      cli_3.6.2          htmltools_0.5.8.1
 [5] tools_4.4.0       yaml_2.3.8         rmarkdown_2.26     knitr_1.46
 [9] jsonlite_1.8.8    xfun_0.43          digest_0.6.35      rlang_1.1.3
[13] renv_1.0.7        evaluate_0.23
```