



HAL
open science

FDP control in mass-univariate linear models using the residual bootstrap

Samuel Davenport, Bertrand Thirion, Pierre Neuvial

► **To cite this version:**

Samuel Davenport, Bertrand Thirion, Pierre Neuvial. FDP control in mass-univariate linear models using the residual bootstrap. *Electronic Journal of Statistics*, 2025, 19, pp.1313 - 1336. <10.1214/25-EJS2354>. <hal-05002607>

HAL Id: hal-05002607

<https://hal.science/hal-05002607v1>

Submitted on 24 Mar 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

FDP control in mass-univariate linear models using the residual bootstrap

Samuel Davenport^{*1,2}, Bertrand Thirion^{†1,2} and Pierre Neuvial^{†3}

¹*Division of Biostatistics, University of California San Diego, United States,
e-mail: sdavenport@health.ucsd.edu; bertrand.thirion@inria.fr*

²*Inria, Université Paris-Saclay, France*

³*Institut de Mathématiques de Toulouse, Université de Toulouse, France,
e-mail: pierre.neuvial@math.univ-toulouse.fr*

Abstract: In this article we develop a method for performing post hoc inference of the False Discovery Proportion (FDP) over multiple contrasts of interest in the mass-univariate linear model. To do so we use the residual bootstrap to simulate from the distribution of the null contrasts. We combine the bootstrap with the post hoc inference bounds of [11] and prove that doing so provides simultaneous asymptotic control of the FDP over all subsets of hypotheses. We demonstrate, via simulations, that our approach provides simultaneous control of the FDP over all subsets and is typically more powerful than existing, state of the art, parametric methods. We illustrate our approach on functional Magnetic Resonance Imaging data from the Human Connectome project and on a transcriptomic dataset of chronic obstructive pulmonary disease.

MSC2020 subject classifications: Primary 62J15.

Keywords and phrases: FDP control, residual bootstrap, simultaneous inference, post hoc inference.

Received October 2023.

Contents

1	Introduction	1314
2	Modelling framework	1316
	2.1 Notation	1316
	2.2 Linear model framework	1317
	2.3 Bounds on the false discovery proportion	1318
	2.4 Resampling from the null distribution	1320
3	Joint error rate control in the linear model	1321
	3.1 Joint error rate control	1322
	3.2 Bootstrap step-down procedure	1323
4	Simulation results	1324
	4.1 Simulation setup	1324
	4.2 False positive control	1325

*National Institute of Health

†Agence nationale de la recherche

4.3	Power	1325
5	Real data results	1327
5.1	Neuroimaging data application	1327
5.2	Transcriptomic data application	1328
6	Discussion	1332
	Acknowledgments	1332
	Funding	1332
	Supplementary Material	1333
	References	1333

1. Introduction

Statistical analysis of functional Magnetic Resonance Imaging data grounds the inference of associations between external conditions (such as disease status and experimental factors) and the signals recorded in brain regions, that are assumed to reflect brain activity. In particular, practitioners typically aim to uncover associations between local signals and conditions that are specific to a given area; such specificity is essential for interpretation purposes. The most standard framework is that of mass-univariate inference, in which models are fit separately at each brain location, in order to detect significant associations. This framework is simple and computationally efficient but, given mm-scale resolution reached by current imaging setups, results in a dire multiple comparison issue.

Statistical analysis of genomic data encounters a similar multiple comparison problem. In particular, this is the case in Genome-Wide Association Studies that aim to identify Single Nucleotide Polymorphisms that are associated with one or more phenotypes of interest, and in gene expression studies, the goal of which is to identify genes where activity is associated with one or more variables of biological or clinical interest. In this field, the state-of-the-art framework is also based on univariate tests that are performed for each genomic marker. While imaging data typically consist of a smooth volume-domain voxel grid, the dependence structure of genomic data is dictated by the interdependence between genomic markers, which is mediated by haplotypic blocks encountered in Genome-Wide Association Studies and by gene networks or pathways in expression studies.

In both of these scientific fields (and in others), control of the false discovery rate (FDR) has quickly become a de facto standard, as it yields image or genome-level error control together with acceptable power [24, 46]. In practice, most researchers control the FDR using the Benjamini-Hochberg procedure [5], under the assumption of positive regression dependence [6]. This assumption is generally considered reasonable given the positive correlation that typically exists between voxels or genomic markers. However, users often interpret FDR-control as a control of false discovery proportion (FDP), which is incorrect, as the FDR is only the expected value of the FDP. Overall, this approach can result in unreliable error control, especially when there is dependence within the data,

see [31] and Figure 2.1 in [37]. It is instead desirable to provide probabilistic control on the *proportion* or *number* of false discoveries.

In this work we will be considering mass-univariate models in which many separate linear models are fit at each given voxel/gene [50]. In particular genomic and brain imaging datasets frequently involve the simultaneous test of several contrasts [45, 1]. Such simultaneous tests are important because they can ground double dissociation [30], ensuring the specificity of discoveries and leading to unambiguous interpretations of the results. A difficulty arises here as the tests of the different contrasts that are considered at each feature (voxel/gene) are typically dependent and it may no longer be reasonable to assume positive regression dependence. It is thus of interest to consider controlling the FDP under the null hypothesis for each contrast, without making unwanted assumptions.

The notion of *post hoc inference* was introduced by [28], following earlier works by [25, 36] on the probabilistic control of the FDP. The idea of post hoc inference is to provide confidence bounds on the number or proportion of true/false discoveries among arbitrary and possibly data-driven subsets of variables of interest. By construction, such guarantees address the issue of circular inference [40].

Post hoc bounds can be obtained as a by-product of the control of a multiple testing risk called the joint error rate (JER) by a simple interpolation argument [11]. Using this construction, state-of-the-art post hoc bounds [28, 40] can be recovered from the Simes inequality, a classical result from the multiple testing literature [43]. The resulting bounds are valid under positive regression dependence. They have also been shown to be conservative in genomics and neuroimaging applications [10].

Since the joint error rate only depends on the joint distribution of the test statistics under the null hypothesis, joint error rate control can alternatively be obtained by randomization techniques [11, 10, 29, 2]. In particular, sharp data-driven joint error rate control and associated post hoc bounds have been obtained for one-sample tests using sign-flipping, and for two-sample tests using permutations [11, 10]. When inferring on contrasts in the linear model, on which we will focus, exchangeability does not generally hold. As such it is not typically possible to provide valid finite sample inference using permutation. In order to perform post hoc inference over contrasts we will need to be able to obtain the joint null distribution of the test-statistics of multiple contrasts within the framework of the mass-univariate linear model. To do so we use the residual bootstrap, adjusting the approach of [49] to the mass-univariate setting. Justification for bootstrapping the residuals in a one-dimensional linear model was first provided in [23] based on theory proved in [7]. These results and their proofs were extended to mass-univariate linear models in [20].

An alternative appealing strategy for the statistical control of false discoveries is the Knockoffs framework [4, 12]. Knockoffs have been used to obtain FDR control in the context of conditional variable selection in high-dimensional linear models with a large number of samples. These results have recently been extended to post hoc FDP control [33, 8]. However, knockoff-based methods

cannot be applied to the mass-univariate testing scenario considered in this paper. These approaches are instead tailored to high-dimensional linear models in which a response of interest (e.g. a given phenotype) is regressed against all features at once (voxels for imaging data, or genes for transcriptomic data). Moreover knockoffs are not typically suited to voxel level resolution inference in smooth images as the variance explained by a given voxel conditional on all surrounding voxels is usually very small, and so knockoffs require aggregation for effective application [8].

The primary contribution of our work is to theoretically extend the results of [11] to the non-exchangeable setting of the linear model, allowing us to provide asymptotically valid simultaneous FDP control in a wide range of application settings. To do so we prove that the residual bootstrap can be combined with the interpolation approach of [11] in order to provide asymptotically valid post-hoc bounds. This requires a different proof strategy to that of [11], in order to establish asymptotic control of the joint error rate, because their proofs strongly rely on exchangeability of the resampled pivotal statistics which does not hold in our setting. As part of this contribution we illustrate how these methods can be applied in practice on both brain imaging and genetics datasets and provide software tools to implement them which are available in the `pyperm` python package [14] and `StatBrainz` MATLAB toolbox [15]. As a further contribution we include an alternative proof of the consistency of the bootstrap in the mass-univariate linear model which relies on the Lindeberg CLT.

Proofs and further theoretical and simulation results are available in the supplementary material – sections of which will be denoted using the suffix S. Code to reproduce the analyses and figures for this paper is available at: https://github.com/sjdavenport/2023_glmfdp.

2. Modelling framework

2.1. Notation

Throughout we will take $(\Omega, \mathcal{F}, \mathbb{P})$ to be a probability space, write \mathbb{E} to denote expectation and will define random variables with respect to this space. We will also take \mathbb{N} to be the set of positive integers. We will primarily be working with random fields, observed at a finite number of points, as our data. These defined as follows.

Definition 2.1. Given $D, L \in \mathbb{N}$ and a finite set $\mathcal{V} \subset \mathbb{R}^D$, we define a **random field** on \mathcal{V} to be a measurable mapping $g : \Omega \rightarrow \{h : \mathcal{V} \rightarrow \mathbb{R}^L\}$. We say that g has **dimension** L .

Remark 2.2. The random field notation which we have described above is motivated by our main application in neuroimaging in which \mathcal{V} is the set of voxels and $D = 3$. However we will not assume any spatial structure on the random fields themselves meaning that the framework is fully general. In particular note that any vector $a \in \mathcal{R}^{|\mathcal{V}|}$ can also be seen as a measurable mapping from \mathcal{V} to \mathbb{R} and therefore falls under this framework.

Given $\omega \in \Omega$ and $v \in \mathcal{V}$ we will write $g(\omega, v) = g(\omega)(v)$ and will typically drop dependence on ω and simply refer to the random variable $g(v) : \Omega \rightarrow \mathbb{R}^L$ when indexing g and say that g is a random field on \mathcal{V} . We define the mean of g to be the function $\mu : \mathcal{V} \rightarrow \mathbb{R}^L$ sending $v \in \mathcal{V}$ to $\mathbb{E}[g(v)]$. To each g we associate a covariance \mathbf{c} and a correlation function ρ which map $\mathcal{V} \times \mathcal{V}$ to $\mathbb{R}^{L \times L}$ and are defined as

$$\mathbf{c}(u, v) = \text{cov}(g(u), g(v)) = \mathbb{E} \left[(g(u) - \mu(u))(g(v) - \mu(v))^T \right]$$

and $\rho(u, v) = \mathbf{c}(u, v)(\mathbf{c}(u, u)\mathbf{c}(v, v))^{-1/2}$ for all $u, v \in \mathcal{V}$.

In order for us to state results about the limiting bootstrap distribution (in Section 2.4) we define Gaussian random fields as follows.

Definition 2.3. Given functions $\mu : \mathcal{V} \rightarrow \mathbb{R}^L$ and $\mathbf{c} : \mathcal{V} \times \mathcal{V}$ we write $g \sim \mathcal{G}(\mu, \mathbf{c})$ if g is a random field with mean μ and covariance \mathbf{c} and the vector $(g_j(v) : v \in \mathcal{V}, 1 \leq j \leq L)$ has a multivariate Gaussian distribution.

2.2. Linear model framework

Let $\mathcal{V} \subset \mathbb{R}^D$ be a finite set of points corresponding to the domain of interest (this could for instance be the voxels of the brain or points representing genes). Suppose that we observe random fields $y_i : \mathcal{V} \rightarrow \mathbb{R}$, for $1 \leq i \leq n$ and some number of subjects $n \in \mathbb{N}$. At each point $v \in \mathcal{V}$, we assume that

$$Y_n(v) = X_n \beta(v) + E_n(v) \tag{1}$$

where for each $v \in \mathcal{V}$, $Y_n(v) = [y_1(v), \dots, y_n(v)]^T$ is a vector giving the observed data, $\beta(v) \in \mathbb{R}^p$ is the vector of parameters (for some $p \in \mathbb{N}$), $X_n \in \mathbb{R}^{n \times p}$ is the design matrix of the covariates (note that this may include nuisance variables) and $E_n(v) = [\epsilon_1(v), \dots, \epsilon_n(v)]^T$ is an n -dimensional random vector which represents the unobserved noise, where $(\epsilon_n)_{n \in \mathbb{N}}$ is an i.i.d sequence of 1-dimensional random fields on \mathcal{V} . Note that we do not make any distributional assumptions on the sequence of noise fields $(\epsilon_n)_{n \in \mathbb{N}}$ other than i.i.d. Above we have given the design matrix X_n a subscript n as we will allow it to grow with n . Let \mathbf{c}_ϵ and ρ_ϵ be the covariance and correlation functions of ϵ_1 respectively.

Remark 2.4. In our results in Section 5 we consider two main applications settings. The first is neuroimaging in which $D = 3$ and the domain \mathcal{V} is a subset of 3-dimensional space corresponding to the set of voxels of the brain. The second setting is a transcriptomic example in which $D = 1$ and \mathcal{V} corresponds to a set of genes. Since (1) does not impose any distributional or spatial assumptions on the structure of the errors, both of these settings fit into the described model.

Then, given contrasts $c_1, \dots, c_L \in \mathbb{R}^p$ for some number of contrasts $L \in \mathbb{N}$, we are interested in testing the null hypotheses: $H_{0,l}(v) : c_l^T \beta(v) = 0$, for $1 \leq l \leq L$ and each $v \in \mathcal{V}$. For each $v \in \mathcal{V}$ we can test the intersection null hypothesis

$$H_0(v) : c_l^T \beta(v) = 0 \text{ for } 1 \leq l \leq L$$

using an F -test at each $v \in \mathcal{V}$ given by

$$F_n(v) = \frac{(C\hat{\beta}_n(v))^T(C(X_n^T X_n)^{-1}C^T)^{-1}(C\hat{\beta}_n(v))/\text{rank}(C)}{\hat{\sigma}_n(v)^2}. \quad (2)$$

Here $\hat{\beta}_n(v) = (X_n^T X_n)^{-1}X_n^T Y_n(v)$ and $C = (c_1, \dots, c_L)^T \in \mathbb{R}^{L \times p}$ is the matrix of contrasts. $\hat{\sigma}_n^2 : \mathcal{V} \rightarrow \mathbb{R}$ is the estimate of the variance based on the residuals which sends $v \in \mathcal{V}$ to

$$\hat{\sigma}_n^2(v) = \frac{1}{n - r_n} \left\| Y_n(v) - X_n \hat{\beta}_n(v) \right\|^2.$$

where r_n is the rank of X_n . The individual null hypotheses can be tested using test statistics:

$$T_{n,l}(v) = \frac{c_l^T \hat{\beta}_n(v)}{\sqrt{\hat{\sigma}_n(v)^2 c_l^T (X_n^T X_n)^{-1} c_l}}. \quad (3)$$

Under $H_{0,l}(v)$ and assuming that the noise is Gaussian, conditional on X_n , $T_{n,l}(v)$ is distributed as a t -statistic with $n - r_n$ degrees of freedom. This allows a p -value to be calculated to test $H_{0,l}(v)$ for each contrast l at each point v , namely, $p_{n,l}(v) = 2(1 - \Phi_{n-r_n}(|T_{n,l}(v)|))$ where Φ_d is the CDF of a t -statistic with $d \in \mathbb{N}$ degrees of freedom. Dropping the Gaussianity assumption, the p -values are still asymptotically valid under reasonable assumptions (see e.g. Theorem S-3.4 [18]). Moreover, for each $1 \leq l \leq L$, $T_{n,l}$ is a 1-dimensional random field and we define $T_n = [T_{n,1}, \dots, T_{n,L}]^T$.

2.3. Bounds on the false discovery proportion

The above framework gives us $L|\mathcal{V}|$ different hypothesis tests, and results in a multiple testing problem, which can be quite severe e.g. if the size of \mathcal{V} is large. Let $\mathcal{H} = \{(l, v) : 1 \leq l \leq L \text{ and } v \in \mathcal{V}\}$ index the hypotheses. For $H \subseteq \mathcal{H}$, let $|H|$ denote the number of elements within H . Finally let $\mathcal{N} \subseteq \mathcal{H}$ index the true null hypotheses. Given $0 < \alpha < 1$ we will seek to provide a function $V : \{H : H \subseteq \mathcal{H}\} \rightarrow \mathbb{N}$ such that

$$\mathbb{P}(|H \cap \mathcal{N}| \leq V(H), \text{ for all } H \subseteq \mathcal{H}) \geq 1 - \alpha. \quad (4)$$

If (4) holds then simultaneously over all $H \subseteq \mathcal{H}$, with probability $1 - \alpha$, $V(H)$ provides an upper bound on the number of false positives within H . Suppose that for some $K \in \mathbb{N}$ we have sets $R_1, \dots, R_K \subseteq \mathcal{H}$ (which depend on the data) and constants $\zeta_1, \dots, \zeta_K \in \mathbb{N}$ and define

$$\text{JER}((R_k, \zeta_k)_{1 \leq k \leq K}) := \mathbb{P}(|R_k \cap \mathcal{N}| > \zeta_k, \text{ some } 1 \leq k \leq K) \quad (5)$$

to be the joint error rate of the collection $(R_k, \zeta_k)_{1 \leq k \leq K}$. [11] showed that if $\text{JER}((R_k, \zeta_k)_{1 \leq k \leq K}) \leq \alpha$, then the bound $\bar{V} : \{H : H \subseteq \mathcal{H}\} \rightarrow \mathbb{R}$, sending $H \subseteq \mathcal{H}$ to

$$\bar{V}(H) = \min_{1 \leq k \leq K} (|H \setminus R_k| + \zeta_k) \wedge |H|, \quad (6)$$

satisfies (4) and thus provides an α -level bound over the number of false positives within each chosen rejection set. If the sets R_1, \dots, R_K are nested then \bar{V} is in fact the optimal interpolation bound¹ among the post hoc bounds that can be derived from JER control. We will follow the approach of [11] and define the collections $(R_k, \zeta_k)_{1 \leq k \leq K}$, that we will consider, using template families. An important practical feature of the bound $\bar{V}(H)$ is that it can be computed in linear time in $|H|$, see e.g. Algorithm 2 in [21]. For simplicity we did not consider the closed testing-based post hoc bounds introduced by [28]. These types of post hoc bounds are briefly discussed in Section 6.

Definition 2.5. Given $K \in \mathbb{N}$, we say that a family of functions $(t_k)_{1 \leq k \leq K}$ is a **template family** if for each $1 \leq k \leq K$, $t_k : [0, 1] \rightarrow \mathbb{R}$, $t_k(0) = 0$ and t_k is strictly increasing and continuous. The parameter K is called the **size** of the template.

The simplest and most commonly used template family is the linear template which, for $K \in \mathbb{N}$, is given by $t_k(x) = \frac{xk}{L|\mathcal{V}|}$ for $1 \leq k \leq K$ and $x \in [0, 1]$. Existing post hoc bounds associated with this template are described in Section S-4.1 [18]. However other choices are available and the optimal choice of template may depend on the dataset under consideration: we refer to Section 6 for further details and a discussion of the choice of template as well as to [29, 11, 9, 2]. Given a template family and $\lambda \in [0, 1]$, for each $1 \leq k \leq K$ and $n \in \mathbb{N}$, we will take $R_k(\lambda) = \{(l, v) \in \mathcal{H} : p_{n,l}(v) \leq t_k(\lambda)\}$, set $\zeta_k = k - 1$, and let $p_{(k:\mathcal{N})}^n$ be the k th smallest p -value in the set $\{p_{n,l}(v) : (l, v) \in \mathcal{N}\}$ (setting $p_{(k:\mathcal{N})}^n = 1$ if $k > |\mathcal{N}|$). We will refer to the collection $(R_k(\lambda), k - 1)_{1 \leq k \leq K}$ as the canonical reference family.

Lemma 2.6. For each $\lambda \in [0, 1]$,

$$JER((R_k(\lambda), k - 1)_{1 \leq k \leq K}) = \mathbb{P}\left(\min_{1 \leq k \leq K \wedge |\mathcal{H}|} t_k^{-1}(p_{(k:\mathcal{N})}^n) \leq \lambda\right).$$

Thus for a given template family, in order to obtain an upper bound on the number of false positives we can choose a threshold $\lambda \in [0, 1]$ such that

$$\mathbb{P}\left(\min_{1 \leq k \leq K \wedge |\mathcal{H}|} t_k^{-1}(p_{(k:\mathcal{N})}^n) \leq \lambda\right) \leq \alpha. \tag{7}$$

Then the joint error rate of the family $(R_k(\lambda), k - 1)_{1 \leq k \leq K}$ is controlled to a level α and so the corresponding bound: \bar{V} , provides a $(1 - \alpha)$ -level simultaneous upper bound on the number of false positives.

[11] chose λ via permutation testing, using the fact that under an exchangeability assumption permutation allows the probability in (7) to be controlled exactly. In the linear model, permutation of the response does not satisfy the exchangeability assumption when there are multiple potentially non-zero covariates in the model (see Appendix S-4.4 [18] for a discussion of this). In what

¹Here we mean optimal in the sense that given another function $V' : \{H : H \subseteq \mathcal{H}\}$, such that for all $A \subseteq \mathcal{H}$, $|R_k \cap A| \leq \zeta_k$ for $1 \leq k \leq K$ implies that $|R \cap A| \leq V'(R)$ for all $R \subseteq \mathcal{H}$, it follows that $V'(R) \leq \bar{V}(R)$. Optimality in this sense follows by [11]’s Proposition 2.5.

follows we take a different approach that proceeds via bootstrapping the data and results in asymptotic control of the error rate.

For $\alpha \in (0, 1)$ $\overline{V}(H)$ provides an $(1 - \alpha)$ -level simultaneous upper bound on the number of false positives within H . From (4) we have

$$\mathbb{P}\left(\frac{|H \cap \mathcal{N}|}{|H|} \leq \frac{\overline{V}(H)}{|H|}, \forall H \subseteq \mathcal{H}\right) \geq 1 - \alpha. \tag{8}$$

It thus follows that for each $H \in \mathcal{H}$, $\frac{\overline{V}(H)}{|H|}$ provides an upper bound on the proportion of false positives within H also known as the **false discovery proportion** or **FDP**. Similarly $\frac{|H| - \overline{V}(H)}{|H|}$ provides a $(1 - \alpha)$ -level simultaneous lower bound on the **true discovery proportion** or **TDP**.

2.4. Resampling from the null distribution

In what follows we shall resample our data using the residual bootstrap [23]. Bootstrapping the residuals allows us to target the null distribution because any potential effect is first estimated then and subtracted before resampling. Given $n \in \mathbb{N}$, in our notation, this proceeds by calculating the residuals

$$\hat{E}_n = Y_n - X_n \hat{\beta}_n = (I_n - X_n(X_n^T X_n)^{-1} X_n^T) E_n, \tag{9}$$

where I_n is the $n \times n$ identity matrix and

$$\hat{\beta}_n = (X_n^T X_n)^{-1} X_n^T Y_n = \beta + (X_n^T X_n)^{-1} X_n^T E_n. \tag{10}$$

Given a number of bootstraps to perform: $B \in \mathbb{N}$ for each $1 \leq b \leq B$, conditional on the data, a selection: $\hat{\epsilon}_1^b, \dots, \hat{\epsilon}_n^b$ is chosen independently with replacement from $\{\hat{E}_{n,1}, \dots, \hat{E}_{n,n}\}$ resulting in a combined random field $E_n^b = [\hat{\epsilon}_1^b, \dots, \hat{\epsilon}_n^b]^T$. Given this let $Y_n^b = X_n \hat{\beta}_n + E_n^b$ and define bootstrapped parameter estimates

$$\hat{\beta}_n^b = (X_n^T X_n)^{-1} X_n^T Y_n^b. \tag{11}$$

Define the estimate of variance using the bootstrap residuals to be

$$(\hat{\sigma}_n^b)^2 = \frac{1}{n} \sum_{i=1}^n (E_{n,i}^b)^2 - \left(\frac{1}{n} \sum_{i=1}^n E_{n,i}^b\right)^2. \tag{12}$$

We can use bootstrap to infer on our desired null distribution. To do so we will require the following assumption.

Assumption 1.

a) For $n \in \mathbb{N}$, $X_n = [x_1, \dots, x_n]^T$ for a sequence of i.i.d vectors $(x_n)_{n \in \mathbb{N}}$ in \mathbb{R}^p with bounded density and such that $\mathbb{E}[\|x_1\|^{2+\delta}] < \infty$ for some $\delta > 0$.

b) $(\epsilon_n)_{n \in \mathbb{N}}$ is an i.i.d sequence of 1-dimensional random fields on \mathcal{V} , independent of $(x_n)_{n \in \mathbb{N}}$, with $\max_{v \in \mathcal{V}} \mathbb{E}[\epsilon_1(v)^4] < \infty$ and $\min_{v \in \mathcal{V}} \text{var}(\epsilon_1(v)) > 0$.

Suppose $(X_n)_{n \in \mathbb{N}}$ and $(\epsilon_n)_{n \in \mathbb{N}}$ satisfy Assumption 1. Then consistency of the multivariate residual bootstrap [23, 20] implies that, conditional on $(X_m, Y_m)_{m \in \mathbb{N}}$ for almost all sequences $(X_m, Y_m)_{m \in \mathbb{N}}$, for each $1 \leq b \leq B$, as $n \rightarrow \infty$,

$$\begin{aligned} \sqrt{n}(\hat{\beta}_n^b - \hat{\beta}_n) &\xrightarrow{d} \mathcal{G}(0, \mathbf{c}_\epsilon \Sigma_X^{-1}) \\ \text{and } \hat{\sigma}_n^b &\xrightarrow{\mathbb{P}} \sigma, \end{aligned} \tag{13}$$

where $\Sigma_X = \mathbb{E}[x_1 x_1^T]$. In particular, let $T_n^b : \mathcal{V} \rightarrow \mathbb{R}$ be the L -dimensional random field on \mathcal{V} such that, for $1 \leq l \leq L$,

$$T_{n,l}^b = \frac{c_l^T (\hat{\beta}_n^b - \hat{\beta}_n)}{\hat{\sigma}_n^b \sqrt{c_l^T (X_n^T X_n)^{-1} c_l}}. \tag{14}$$

Then conditional on $(X_m, Y_m)_{m \in \mathbb{N}}$, for almost every sequence $(X_m, Y_m)_{m \in \mathbb{N}}$, for each $1 \leq b \leq B$,

$$T_n^b \xrightarrow{d} \mathcal{G}(0, \mathbf{c}') \tag{15}$$

as $n \rightarrow \infty$. Here $\mathbf{c}' : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ takes $u, v \in \mathcal{V}$ to $\mathbf{c}'(u, v) = \rho_\epsilon(u, v) A C \Sigma_X^{-1} C^T A^T$ where $A \in \mathbb{R}^{L \times L}$ is a diagonal matrix with $A_{ll} = (c_l^T \Sigma_X^{-1} c_l)^{-1/2}$ for $1 \leq l \leq L$. Crucially the limiting distribution in this result is the same as the limiting distribution of the test-statistics (3) under the global null that $\beta = 0$, see Lemma S-3.4 [18]. It follows that the bootstrap provides consistent estimates of the quantiles of functionals of the data under the null, see Section S-2.10 [18]. Application of the residual bootstrap thus approximately samples from the limiting null distribution given a sufficiently large sample size.

We provide two proofs of these bootstrap consistency results. The first demonstrates how to translate the results of [20] into our notation, see Section S-2.4 [18]. The second instead uses the Lindeberg CLT to establish convergence and is provided in Section S-2.3, see Theorems S-2.5 and S-2.7 [18] and their proofs for further details.

Remark 2.7. In the univariate setting consistency of the residual bootstrap was originally proved in [23]. This result has been used widely in the mass-univariate setting however it was only recently that [20] formally showed that that the proof of [23] extends to multiple dimensions. The approach of [23] and [20], takes advantage of the fact that convergence in distribution is equivalent to convergence in the Mallows metric ([7]). We provide an alternative proof of the consistency of the residual bootstrap in Section S-2 [18] which instead relies on elementary probability tools such as the Lindeberg CLT and the triangular law of large numbers.

3. Joint error rate control in the linear model

In this section we will state and prove our main results. We will show that given $0 < \alpha < 1$, choosing λ to be the α -quantile of the bootstrapped distribution results in asymptotic $1 - \alpha$ level control of the joint error rate and thus results in simultaneous control of the FDP.

3.1. Joint error rate control

To set this up, given a test-statistic $T : \mathcal{V} \rightarrow \mathbb{R}^L$, a subset $H \subseteq \mathcal{H}$, and $n \in \mathbb{N}$, for $1 \leq k \leq |H|$, let $p_{(k:H)}^n(T)$ be the k th minimum value in the set

$$\{2 - 2\Phi_{n-r_n}(|T_l(v)|) : (l, v) \in H\}.$$

Using the results we have proved so far we can obtain the following theorem.

Theorem 3.1. For $H \subseteq \mathcal{H}$, let $f_{n,H} : \{g : \mathcal{V} \rightarrow \mathbb{R}^L\} \rightarrow \mathbb{R}$ send

$$T \mapsto \min_{1 \leq k \leq K \wedge |H|} t_k^{-1}(p_{(k:H)}^n(T))$$

and for $n, B \in \mathbb{N}$ and $\alpha \in (0, 1)$, let

$$\lambda_{\alpha,n,B}^*(H) = \inf \left\{ \lambda : \frac{1}{B} \sum_{b=1}^B 1[f_{n,H}(T_n^b) \leq \lambda] \geq \alpha \right\}$$

be the α -quantile of the bootstrap distribution based on $B \in \mathbb{N}$ bootstraps, of $f_{n,H}(T_n)$ conditional on the observed data. Assume that Assumption 1 holds and that $n - r_n \rightarrow \infty$ almost surely. Then for all $H \subseteq \mathcal{H}$ such that $\mathcal{N} \subseteq H$,

$$\lim_{n \rightarrow \infty} \lim_{B \rightarrow \infty} \mathbb{P}(f_{n,\mathcal{N}}(T_n) \leq \lambda_{\alpha,n,B}^*(H)) \leq \alpha.$$

This limit holds with equality if $H = \mathcal{N}$. Furthermore, taking $H = \mathcal{H}$, it follows that

$$\lim_{n \rightarrow \infty} \lim_{B \rightarrow \infty} \mathbb{P} \left(\min_{1 \leq k \leq K \wedge |\mathcal{H}|} t_k^{-1}(p_{(k:\mathcal{N})}^n(T_n)) \leq \lambda_{\alpha,n,B}^*(\mathcal{H}) \right) \leq \alpha$$

Applying this result and using Lemma 2.6 we are thus able to obtain asymptotic control of the joint error rate of the canonical reference family. Following the discussion in Section 2.3 this means that we obtain asymptotic post hoc FDP control. In particular we having the following corollary.

Corollary 3.2. Under the assumptions of Theorem 3.1, for $0 < \alpha < 1$, and $H \subseteq \mathcal{H}$, let

$$\bar{V}_{\alpha,n,B}(H) = \min_{1 \leq k \leq K} (|H \setminus R_k(\lambda_{\alpha,n,B}^*(\mathcal{H}))| + k - 1) \wedge |H|.$$

$$\text{Then } \lim_{n \rightarrow \infty} \lim_{B \rightarrow \infty} \mathbb{P}(|H \cap \mathcal{N}| \leq \bar{V}_{\alpha,n,B}(H), \forall H \subseteq \mathcal{H}) \geq 1 - \alpha.$$

Thus in order to provide FDP control, given a number of bootstraps $B \in \mathbb{N}$, we can calculate $\lambda_{\alpha,n,B}^*(H)$, the α -quantile of the bootstrap distribution of $f_{n,H}(T_n)$ conditional on the observed data. Then $\bar{V}_{\alpha,n,B}(H)$ provides a $(1 - \alpha)$ level simultaneous upper bound on the number of false positives in $H \subseteq \mathcal{H}$.

3.2. Bootstrap step-down procedure

It is possible to improve on the power of the above procedure by taking a step-down approach in the spirit of [39]. This is based on the idea that joint error rate control implies familywise error rate control, see e.g. Section S-4.2 [18]. As such it is possible to obtain an estimate of the set of null hypotheses and thereby obtain a tighter bound. The procedure, which adapts the step-down procedure of [11] to our setting, can be iterated as follows.

Algorithm 1 step-down bootstrap

- 1: Set $j \leftarrow 0$ and $H_n^{(0)} \leftarrow \mathcal{H}$
 - 2: **repeat**
 - 3: Set $j \leftarrow j + 1$, $\lambda_{n,j} \leftarrow \lambda_{\alpha,n,B}^*(H_n^{(j-1)})$ and $H_n^{(j)} \leftarrow \{(l, v) : p_{n,l}(v) \geq t_1(\lambda_{n,j})\}$
 - 4: **until** $H_n^{(j)} = H_n^{(j-1)}$
 - 5: Set $\hat{H}_n \leftarrow H_n^{(j)}$ and **return** \hat{H}_n
-

As the following theorem demonstrates, the step-down approach controls the joint error rate and therefore provides simultaneous FDP control.

Theorem 3.3. *Under the assumptions of Theorem 3.1, for $0 < \alpha < 1$, let \hat{H}_n be the set generated by applying Algorithm 1. Then*

$$\lim_{n \rightarrow \infty} \lim_{B \rightarrow \infty} \mathbb{P}\left(f_{n,\mathcal{N}}(T_n) < \lambda_{\alpha,n,B}^*(\hat{H}_n)\right) \leq \alpha.$$

Thus for $H \subseteq \mathcal{H}$ we can define the bound,

$$\bar{V}_{\alpha,n,B}(H) = \min_{1 \leq k \leq K} (|H \setminus R_k(\lambda_{\alpha,n,B}^*(\hat{H}_n))| + k - 1) \wedge |H|,$$

and it follows that

$$\lim_{n \rightarrow \infty} \lim_{B \rightarrow \infty} \mathbb{P}(|H \cap \mathcal{N}| \leq \bar{V}_{\alpha,n,B}(H), \forall H \subseteq \mathcal{H}) \geq 1 - \alpha.$$

In the definition of $\lambda_{\alpha,n,B}^*$ we require the computation of $|\mathcal{H}|$ statistics for each bootstrap each of which is based on a sample of size n . As such the complexity of these algorithms is $O(nB|\mathcal{H}|)$.

We can asymptotically characterize the nature of the bound as follows.

Remark 3.4. The results in this subsection and the one previous have been stated for two-sided p -values however they also hold for one-sided p -values, $1 - \Phi_{n-r_n}(|T_{n,l}(v)|)$ without change. All that is required to show this is to re-define $p_{(k:H)}^n(T)$ as the k th minimum value in the set

$$\{1 - \Phi_{n-r_n}(T_l(v)) : (l, v) \in H\}.$$

In this scenario we would use the one-sided p -values to test the null hypotheses that $c_l^T \beta(v) \leq 0$ at each $v \in \mathcal{V}$ and $1 \leq l \leq L$.

Remark 3.5. As the number of subjects increases the signal to noise ratio goes to infinity and so it becomes easier to correct identify the true signals. In particular we can show that the post-hoc bounds on the number of true null hypotheses within a given set H converge to the actual number of true null hypotheses with high probability. We formalize this result in Corollary S-4.2, see Appendix S-4.3 [18] for further details.

4. Simulation results

4.1. Simulation setup

In order to assess empirically that our method correctly controls the joint error rate we run numerical simulations. We create the noise in these simulations by generating 2-dimensional stationary Gaussian random fields on domains which are 25 by 25, 50 by 50 and 100 by 100 pixels. To do so we smooth white noise with a Gaussian kernel with full width at half maximum (FWHM) in $\{0, 4, 8\}$ (in pixel units), accounting for edge effects to ensure stationarity (see e.g. [16]), and scaled so that the variance is 1 everywhere. For the noise distribution we consider Gaussian as well as heavier tailed t_3 and t_5 distributions, as these can arise in fMRI datasets [38, 17].

We let the total number of subjects n range from 20 to 100. For each n , smoothness level, image size and $\pi_0 \in \{0.5, 0.8, 0.9, 1\}$, we run 5000 simulations – each with 1000 bootstraps – to test the joint error rate. For each simulation we do the following. First we generate n Gaussian random fields $\epsilon_1, \dots, \epsilon_n$ as described above and add signal to them (as detailed in the next paragraph). We then randomly divide these images into 3 disjoint groups: $G_1, G_2, G_3 \subset \{1, \dots, n\}$ – performing assignment to each group with equal probability (we eliminate assignments where a given group has no entries). We test for the difference between the first and the second group and between the second and the third group – giving us $L = 2$ contrasts to differentiate between. We thus, in total, test 5000 hypotheses for the 50 by 50 scenario and 20000 for the 100 by 100 case.

We vary the amount of signal in the datasets as follows. Given a proportion π_0 we randomly choose a subset \mathcal{N} of size $\pi_0|\mathcal{H}|$ of $\mathcal{H} = \{(l, v) : 1 \leq l \leq 2, v \in \mathcal{V}\}$ to be null (which is thus different in each simulation) and add signal to ensure that the remainder are non-null. To do so, for $1 \leq i \leq n$, and each $v \in \mathcal{V}$, we set

$$Y_i(v) = 1[i \in G_2, (1, v) \notin \mathcal{N}] + 1[i \in G_3, (1, v) \notin \mathcal{N}] + 1[i \in G_3, (2, v) \notin \mathcal{N}] + \epsilon_i(v).$$

This ensures that the power of the test to detect a difference at h is equal for any $h \in \mathcal{N}^C$. If $\pi_0 = 1$ then all hypotheses are null. An example realisation is shown in Figure S-4.

In the next subsections we compare our bootstrap procedures, in terms of false positive control and power, to two parametric alternatives: the Simes procedure [28] and its step-down: all resolutions inference (ARI, [40]) which are described formally in Section S-4.1 [18]. Here we use the term parametric to indicate that

dependency assumptions on the noise are required in order for the methods to be valid.

4.2. False positive control

In each simulation setting, for $1 \leq j \leq 5000$, we calculate a test statistic random field $T_n^{(j)}$ and obtain λ thresholds for the single-step bootstrap, step-down bootstrap, Simes and ARI methods. For each method we obtain λ -thresholds $\lambda_1, \dots, \lambda_{5000}$ allowing us to estimate the joint error rate via the statistic

$$\frac{1}{5000} \sum_{j=1}^{5000} 1[f_{n,\mathcal{N}}(T_n^{(j)}) \leq \lambda_j]$$

which we refer to as the **empirical joint error rate**. Here $1[\cdot]$ denotes the indicator function.

The results for the 50 by 50 Gaussian simulations are displayed in Figure 1 and those for the other domain sizes are shown in Figures S-9 and S-10 [18]. The results for the bootstrap methods are shown in blue whilst those for the parametric methods are shown in red. The solid lines indicate the step-down methods (i.e. ARI and the step-down bootstrap). These plots demonstrate that, given a reasonable number of subjects ($N \geq 80$) and smoothness of the underlying data (FWHM = 4, 8), the joint error rate of the bootstrap procedures converges to the nominal level, in this case 0.1. Results for the t_3 and t_5 distributions are shown in Section S-6.4 [18] and show similar behaviour indicating that the methods are robust to heavy tails.

Empirically the parametric procedures are valid in all settings considered. However, their control of the joint error rate is substantially below the nominal level when the applied smoothing is non-zero, while the bootstrap approaches demonstrate tighter control. The step-down procedures provide an improvement on their single-step counterparts. This difference increases as π_0 decreases. See Section 4.3 for further details on the effect of π_0 .

4.3. Power

In this section we compare the power of the various methods in the simulation setting described in Section 4.1 in the case where the applied FWHM is 4 pixels. We have chosen to focus on this setting because it represents a realistic level of applied smoothness and illustrates the benefits that can be achieved when using the bootstrap under dependence.

Here we shall use a notion of power originally proposed in [11] to compare the ability of joint error rate controlling procedures to detect signal. Given a set $R \subseteq \mathcal{H}$, define

$$\text{Pow}(R) := \mathbb{E} \left[\frac{|R| - \bar{V}(R)}{|R \cap (\mathcal{H} \setminus \mathcal{N})|} \mid |R \cap (\mathcal{H} \setminus \mathcal{N})| > 0 \right] \tag{16}$$

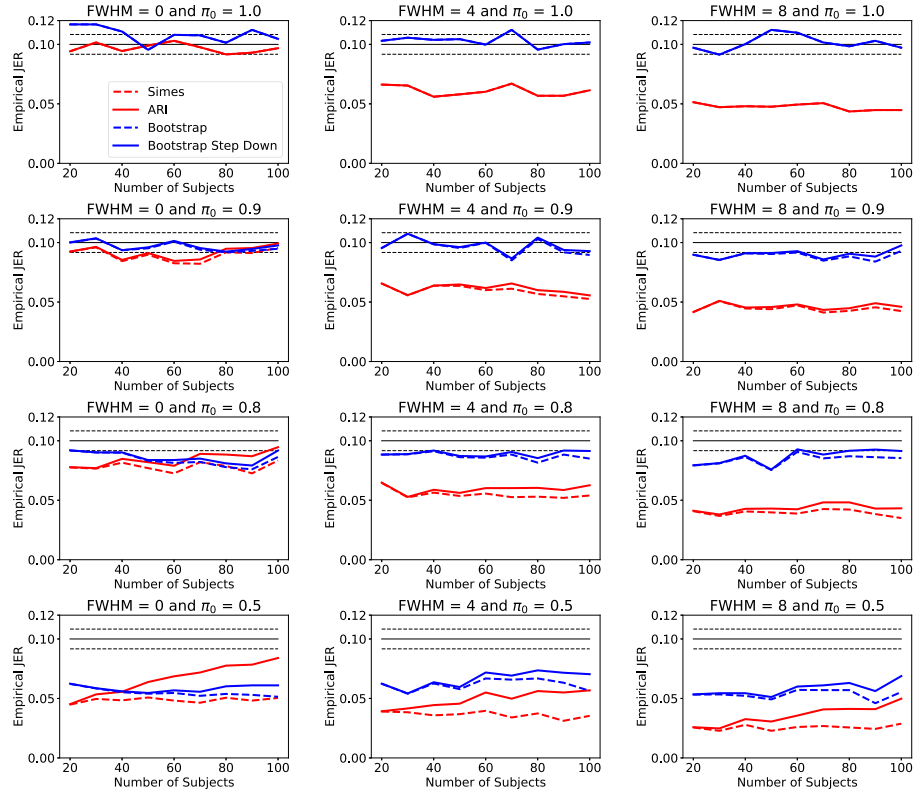


FIG 1. Comparing the empirical joint error rate across methods for the simulation setting described in Section 4.1 for $\alpha = 0.1$ on the domain of size 50 by 50 pixels. The bootstrap procedures typically provide tighter control of the joint error rate than the parametric ones, except under independence. The bootstrap methods are shown in blue whilst the parametric methods are shown in red. The solid lines indicate the step-down methods. The thin flat black dashed lines provide 95% marginal confidence bands based on the normal approximation to the binomial distribution.

where for each method \bar{V} is the corresponding post-hoc bound. Here we consider the following choices of R with which we compare the power (as in [11]). 1) $R = \mathcal{H}$ and 2) taking R to be the hypotheses of \mathcal{H} which are rejected by the Benjamini Hochberg procedure, applied to the p -values $\{p_{n,l}(v) : (l, v) \in \mathcal{H}\}$, at a level 0.05. Note that, unlike in [11], no additional level of randomness in the choice of the sets in 2) is prescribed. We also consider taking $R = \{(l, v) : p_{n,l}(v) \leq 0.05\}$, see Section S-6.7 [18], the results for which are similar in nature to scenario 1 from above. The results for cases 1) and 2) are illustrated graphically in Figure 2. These are for simulations on the 50 by 50 domain.

From these plots we can see that overall the bootstrap based approaches have a higher power than the parametric ones. The power of ARI only becomes comparable (or higher) to that of the bootstrap in the extreme scenario ($\pi_0 = 0.5$) given a large enough sample size. Additionally the bootstrap is not robust at

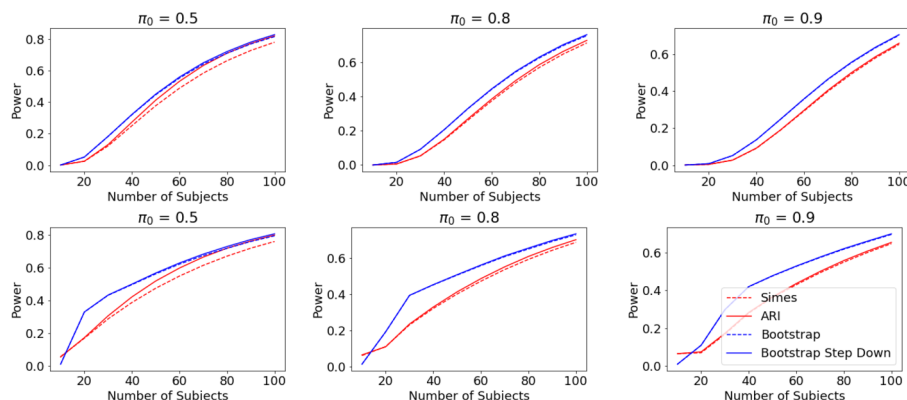


FIG 2. Plotting the power of the different methods against the number of subjects. The power for setting 1 (i.e. $R = \mathcal{H}$) is shown in the top row and the power for setting 2 (i.e. taking R to be the Benjamini-Hochberg rejection set) is shown in the bottom row.

the smallest sample size considered (i.e. $n = 10$) where it is slightly conservative. However it is important to note that in typical high-dimensional applications (neuroimaging, genetics) $\pi_0 > 0.9$ and n is often substantially greater than 10.

The lower the value of π_0 , the greater the increase in power that is obtained by using the step-down algorithms. ARI is always more powerful than Simes by construction. In the relatively sparse scenarios (i.e. $\pi_0 \geq 0.8$) they have a very similar power however for $\pi_0 = 0.5$, ARI provides a marked improvement over Simes. The bootstrap step-down always improves on the standard bootstrap approach though the difference is not particularly large: even when $\pi = 0.5$ this increase is relatively small. The similarity of the standard and step-down procedures, for both the parametric and bootstrap methods, is consistent with the results obtained on real data which are described in the next subsections.

5. Real data results

5.1. Neuroimaging data application

We have 3D functional Magnetic Resonance Imaging data from $n = 386$ unrelated subjects, who performed an m -back working memory task, from the Human Connectome Project. After pre-processing (described in Section S-5 [18]) we obtain a 3-dimensional contrast image for each subject, i.e. we observe images $y_1, \dots, y_n : \mathcal{V} \rightarrow \mathbb{R}$ where \mathcal{V} is the set of voxels making up the brain. We fit the linear model (1) at each voxel $v \in \mathcal{V}$, where the columns of the design matrix X_n consist of sex, height, weight, body mass index, two different measures of blood pressure, handedness and IQ (measured using the PMAT24_A_CR test score). We consider sex and IQ as a variables of interest, meaning that we test $L = 2$ contrasts. We obtain test-statistic contrasts for sex and IQ and a p -value at each voxel for each contrast. We form clusters using a cluster defining

threshold on the p -values of $p = 0.001$, with each cluster being a contiguous set of voxels above the threshold (clusters are defined separately for each contrast of interest).

We use our bootstrap framework, performing the resampling using 1000 bootstraps, to provide a lower bound on the proportion of active voxels within each cluster, taking $\alpha = 0.1$. This illustrates that multiple clusters, in different regions of the brain, have a relatively large proportion of active voxels for the contrast of IQ. In particular we found more informative bounds in the regions which are known to be related to IQ [41]. For the contrast of sex only a single cluster has a non-zero lower bound on the number of true positives. This is in line with previous literature which does not find substantial associations between sex and working memory data [44]. These results indicate that our method is sensitive to activation where it exists and correctly does not find activation where it does not exist. The bounds provided using the step-down bootstrap procedure are the same as the single-step version in this example.

We compare to the results that are obtained using Simes and ARI bounds (taking $\alpha = 0.1$) and see that our bootstrap approach results in higher lower bounds on the number of active voxels. In this setting the bounds obtained by the parametric procedures are very similar to each other, which is not surprising given the sparsity of the signal. For the IQ contrast the lower bounds provided by the bootstrap and ARI for the number of true positives and on the TDP within each cluster are shown graphically in the upper panel of Figure 3. The corresponding plot for the sex contrast is shown in Section S-6.2 [18]. Direct comparison of the lower bounds is shown in Figure 4.

5.2. Transcriptomic data application

In this section, we illustrate the application of our methods to a specific gene expression data set. Gene expression studies use microarray or sequencing biotechnologies in order to measure the activity (or “expression level”) of a large number of genes simultaneously. We focus on a study of chronic obstructive pulmonary disease (COPD), see [3], whose main goal was to identify genes whose expression level is significantly associated with lung function. The data consists of vectors $y_1, \dots, y_n \in \mathbb{R}^{12,531}$ with the gene expression at each gene for each subject. This setting fits into our framework, as discussed in Remarks 2.2 and 2.4. [3] fit the linear model (1) for this association at each gene, with design matrix X_n made up of the following covariates: age, sex, body mass index, parental history of COPD, and two smoking variables (smoking status and pack-years). The number of subjects is $n = 135$ while the number of genes is $V = 12,531$, leading to a large-scale multiple testing problem. A single contrast for lung function is considered, meaning that we take $L = 1$. Using the Benjamini-Hochberg method to control the FDR at the 5% level, 1,745 genes were found to be significantly associated.

We applied our framework to this data, considering the contrast for lung function. We performed 1000 bootstraps and used these to obtain $\lambda_{\alpha,135,1000}^* = 0.22$,

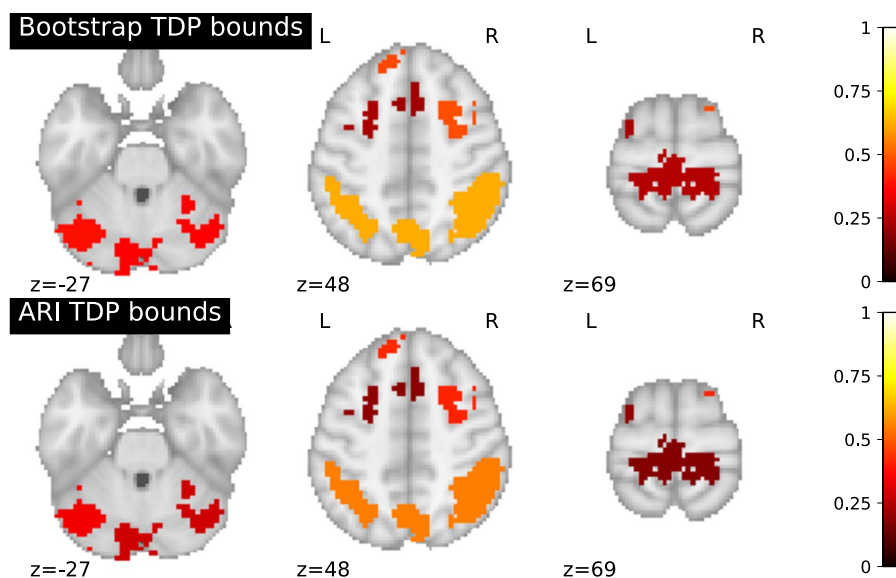


FIG 3. TDP bounds within clusters for the contrast for IQ in the linear regression model fit to the HCP data. Each cluster is shaded a single colour which is the lower bound on the TDP. The upper panel gives the TDP bounds within each cluster provided by the bootstrap procedure. The lower panel gives the bounds provided by using ARI. The bounds given by the bootstrap are larger (as indicated by the lighter colours) indicating that the method is more powerful. Note that these images are 2D slices through the 3D brain and so voxels that are part of the same cluster are not necessarily connected.

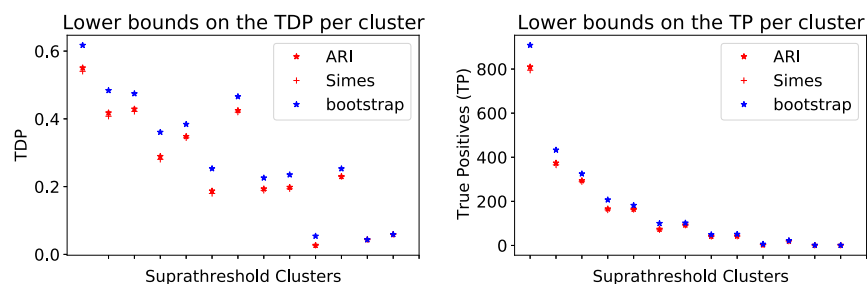


FIG 4. Comparing the TDP and true positive lower bounds across clusters for the different methods. The bootstrap lower bounds are consistently higher than the parametric methods. Clusters are organized from left to right in terms of their size. Only one cluster for the sex contrast is found: this is the 2nd smallest cluster overall with a TP lower bound of 1 voxel. The sizes and bounds of the clusters in the IQ contrast are larger. For the largest cluster we are able to conclude that it contains 908 true positives using the bootstrap approach.

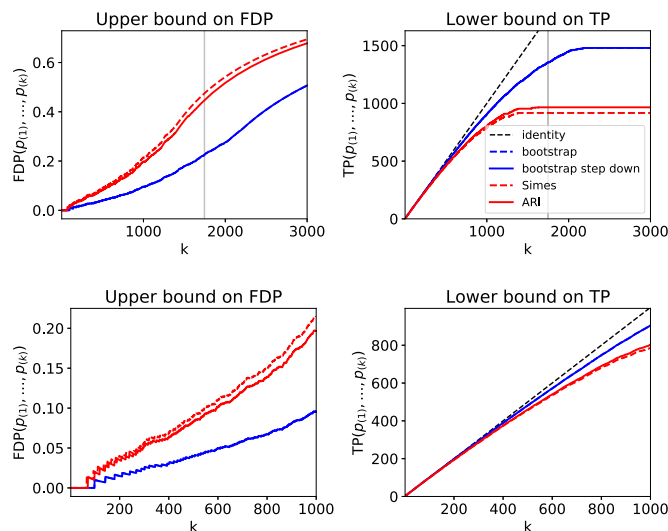


FIG 5. False discovery proportion and true positive plots for the transcriptomic dataset. In the upper panels, for $k = 1, \dots, 3000$, upper bounds on the FDP and lower bounds on the number of true positives are provided by each of the methods for the sets comprised of the hypotheses with the k smallest p -values. The silver vertical line corresponds to the location of the Benjamini-Hochberg rejection set.

The lower panels provide a zoomed in version of the same plot for for the 1000 smallest p -values. The bootstrap methods provide substantially better bounds than the parametric ones. ARI slightly improves on Simes while the step-down bootstrap is indistinguishable from the single-step bootstrap approach in this setting.

where we took $\alpha = 0.1$. This allows us to provide $(1 - \alpha)$ -level simultaneous lower bounds on the number of true positives within any specified set of genes. In particular it allows us to conclude (with 90% confidence) that at least 1,354 of the 1,745 genes within the Benjamini-Hochberg significance set are active. The stepdown bootstrap provides the same bound as the single-step version in this case. Simes and ARI provide lower bounds on the number of true positives in this set of 917 and 966 respectively, which are substantially less informative than the bootstrap bounds.

In the absence of prior information on genes, a natural idea is to rank them by decreasing statistical significance. Our post hoc methods provide upper confidence curves on the proportion of true positives among the most significant genes. Such curves are displayed in Figure 5, where the blue lines correspond to our proposed single-step and step-down bootstrap-based methods, and the red lines correspond to the parametric approaches of [28] and [40]. These results are consistent with the numerical experiments of Section 4. First, the bootstrap method yields post hoc bounds that are substantially more informative than their parametric counterpart. Second, the difference between single-step methods and their step-down counterpart is very small, which is consistent with the fact that the signal is expected to be sparse in such genomic

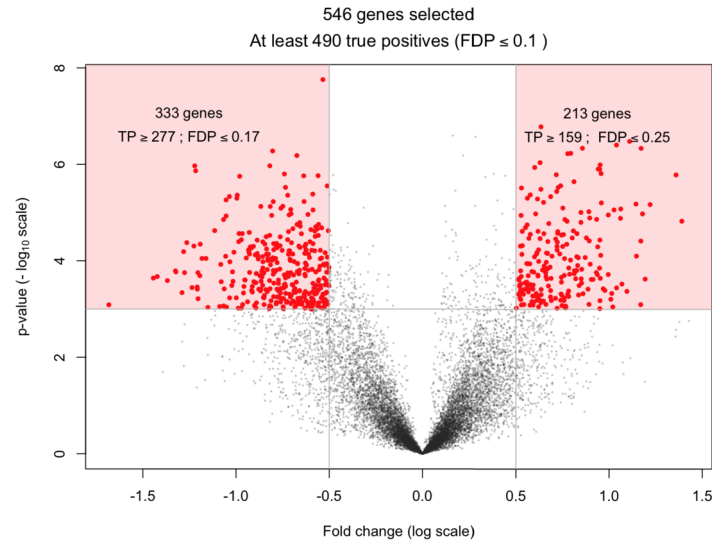


FIG 6. A volcano plot for the p -values for the transcriptomic data. For each gene this plots the estimated contrast effect size (labelled as fold change and corresponding to $c^T \hat{\beta}_{135}$ where c is the contrast vector for COPD) against the p -value, where both are measured in log scale. Two regions (shown via shading) are selected containing the genes whose p -values are less than 10^{-3} and for which the absolute fold change is greater than $10^{0.5}$. Bounds on the true positives (TP) and FDP, overall and for the shaded regions are provided on the plot.

data sets, corresponding to π_0 close to 1. For the bootstrap there is in fact no difference between the single-step and step-down approach in this example.

A widely used approach in differential expression studies is to select genes based on the conjunction of a threshold on the p -values and a threshold on its effect size [13]. [19] recently noted that this type of double selection can lead to inflated numbers of false discoveries when used in conjunction with FDR-based multiple testing corrections, whereas post hoc inference is by construction robust against this issue. The use of our proposed post hoc bounds in this context is illustrated in the volcano plot in Figure 6 [13]. In this plot, each gene is represented in two dimensions by estimates of its effect size (x axis, also known as “fold change” in genomics) and p -value (y axis), in a logarithmic scale. Figure 6 illustrates a particular selection, corresponding to the genes whose p -value is below 0.001 and whose effect size is above 0.5. Our bootstrap-based bound ensures that with probability $1 - \alpha = 90\%$, among these 546 genes, at least 490 are true positives, corresponding to a FDP below 0.1. Importantly, the p -value and effect size thresholds can be chosen post hoc, and multiple such choices can be made without compromising the statistical coverage of the associated bound. For example, the bounds associated to the gene subsets with positive and negative effect size are also displayed in Figure 6.

6. Discussion

In this paper we have introduced a bootstrap method which provides simultaneous control of the FDP over subsets of hypotheses of multiple contrasts in the linear model. We have proved the asymptotic validity of this approach and shown, via simulation, that the error rate is controlled to the correct level given a reasonable number of subjects.

From our simulations and real data examples, we can see that the bootstrap approach typically provides better bounds than existing, state of the art parametric methods (i.e. Simes and ARI). This occurs because we are able to model the dependence within the data. The parametric methods, on the other hand, rely on the Simes inequality which is only exact under independence. Moreover the Simes inequality is only valid under positive regression dependence whereas the non-parametric bootstrap makes relatively few assumptions other than finite moments of the noise and the design. In real data situations there is often relatively strong dependence within the data in which case we would expect the bootstrap to give better bounds. This is illustrated in our brain imaging and transcriptomic examples where the bootstrap bounds provided substantial improvements over the ones derived using the parametric methods.

Acknowledgments

We are grateful to Alexandre Blain at Inria for his help with demonstrating how to use the sanssouci code to get ARI to work and for checking that the Python implementation of ARI was consistent with the implementation in the ARIBrain R package. We are grateful to François Bachoc at the University of Toulouse for his help with the proof of Lemma S-4.5 [18]. SD is grateful to Fabian Telschow at Humboldt University for useful discussions on bootstrapping. We are also grateful to the two anonymous reviewers whose comments helped to improve the quality of the manuscript.

Data were provided in part by the Human Connectome Project, WU-Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Washington University.

Funding

SD was supported by NIH grant R01EB026859. SD and PN were supported by the SansSouci ANR project (ANR-16-CE40-0019). BT was supported by the KARAIB AI chair (ANR-20-CHIA-0025-01) and the FastBig ANR project (ANR-17-CE23-0011).

Supplementary Material

Supplementary material for FDP control in mass-univariate linear models using the residual bootstrap

(doi: [10.1214/25-EJS2354SUPP](https://doi.org/10.1214/25-EJS2354SUPP); .pdf). Contains proofs and additional experiments.

References

- [1] Bianca AV Alberton, Thomas E Nichols, Humberto R Gamba, and Anderson M Winkler. Multiple testing correction over contrasts for brain imaging. *NeuroImage*, 216:116760, 2020.
- [2] Angela Andreella, Jesse Hemerik, Livio Finos, Wouter Weeda, and Jelle Goeman. Permutation-based true discovery proportions for functional magnetic resonance imaging cluster analysis. *Statistics in Medicine*, 2023. [MR4596798](#)
- [3] Timothy M Bahr et al. Peripheral blood mononuclear cell gene expression in chronic obstructive pulmonary disease. *American journal of respiratory cell and molecular biology*, 49(2):316–323, 2013.
- [4] Rina Foygel Barber and Emmanuel J Candès. Controlling the false discovery rate via knockoffs. 2015. [MR3375876](#)
- [5] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995. [MR1325392](#)
- [6] Yoav Benjamini and Daniel Yekutieli. The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188, 2001. [MR1869245](#)
- [7] Peter J Bickel and David A Freedman. Some Asymptotic Theory for the Bootstrap. *Annals of Statistics*, 9(6):1196–1217, 1981. [MR0630103](#)
- [8] Alexandre Blain, Bertrand Thirion, Olivier Grisel, and Pierre Neuvial. False Discovery Proportion control for aggregated Knockoffs. *NeurIPS 2023 – 37th Conference on Neural Information Processing Systems*, December 2023.
- [9] Alexandre Blain, Bertrand Thirion, and Pierre Neuvial. Notip: Non-parametric true discovery proportion control for brain imaging. *NeuroImage*, page 119492, 2022.
- [10] Gilles Blanchard, Pierre Neuvial, and Etienne Roquain. On agnostic post hoc approaches to false positive control. In *Handbook of Multiple Comparisons*, Handbooks of Modern Statistical Methods. Chapman & Hall/CRC, November 2021.
- [11] Gilles Blanchard, Pierre Neuvial, Etienne Roquain, et al. Post hoc confidence bounds on false positives using reference families. *Annals of Statistics*, 48(3):1281–1303, 2020. [MR4124323](#)
- [12] Emmanuel Candès, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: ‘model-x’ knockoffs for high dimensional controlled variable selection.

- Journal of the Royal Statistical Society Series B: Statistical Methodology*, 80(3):551–577, 2018. [MR3798878](#)
- [13] Xiangqin Cui and Gary A Churchill. Statistical tests for differential expression in cDNA microarray experiments. *Genome Biol*, 4(4):210, 2003.
 - [14] Samuel Davenport. Pyperm python package. <https://github.com/sjdavenport/pyperm>, 2024.
 - [15] Samuel Davenport. StatBrainz matlab toolbox. <https://github.com/sjdavenport/StatBrainz>, 2024.
 - [16] Samuel Davenport and Thomas E. Nichols. The expected behaviour of random fields in high dimensions: contradictions in the results of [1]. *Magnetic Resonance Imaging*, 2022.
 - [17] Samuel Davenport, Armin Schwartzman, Thomas E Nichols, and Fabian JE Telschow. Robust fwer control in neuroimaging using random field theory: Riding the surf to continuous land part 2. *arXiv preprint arXiv:2312.10849*, 2023.
 - [18] Samuel Davenport, Bertrand Thirion, Pierre Neuvial. Supplementary material for FDP control in mass-univariate linear models using the residual bootstrap. DOI: [10.1214/25-EJS2354SUPP](https://doi.org/10.1214/25-EJS2354SUPP), 2024.
 - [19] Mitra Ebrahimipoor, Pietro Spitali, Kristina Hettne, Roula Tsonaka, and Jelle Goeman. Simultaneous enrichment analysis of all possible gene-sets: unifying self-contained and competitive methods. *Briefings in bioinformatics*, 21(4):1302–1312, 2020.
 - [20] Daniel J Eck. Bootstrapping for multivariate linear regression models. *Statistics & Probability Letters*, 134:141–149, 2018. [MR3758593](#)
 - [21] Nicolas Enjalbert-Courrech and Pierre Neuvial. Powerful and interpretable control of false discoveries in two-group differential expression studies. *Bioinformatics*, 38(23):5214–5221, 2022.
 - [22] David Freedman and David Lane. A nonstochastic interpretation of reported significance levels. *Journal of Business & Economic Statistics*, 1(4):292–298, 1983.
 - [23] David A Freedman. Bootstrapping regression models. *The Annals of Statistics*, 9(6):1218–1228, 1981. [MR0630104](#)
 - [24] Christopher R Genovese, Nicole A Lazar, and Thomas E Nichols. Thresholding of statistical maps in functional neuroimaging using the false discovery rate. *Neuroimage*, 15(4):870–878, 2002.
 - [25] Christopher R Genovese and Larry Wasserman. Exceedance control of the false discovery proportion. *Journal of the American Statistical Association*, 101(476):1408–1417, 2006. [MR2279468](#)
 - [26] Jelle J Goeman, Jesse Hemerik, and Aldo Solari. Only closed testing procedures are admissible for controlling false discovery proportions. *The Annals of Statistics*, 49(2):1218–1238, 2021. [MR4255125](#)
 - [27] Jelle J Goeman, Rosa J Meijer, Thijmen JP Krebs, and Aldo Solari. Simultaneous control of all false discovery proportions in large-scale multiple hypothesis testing. *Biometrika*, 106(4):841–856, 2019. [MR4046036](#)
 - [28] Jelle J. Goeman and Aldo Solari. Multiple Testing for Exploratory Research. *Statistical Science*, 26(4):584–597, 2011. [MR2951390](#)

- [29] Jesse Hemerik, Aldo Solari, and Jelle J Goeman. Permutation-based simultaneous confidence bounds for the false discovery proportion. *Biometrika*, 106(3):635–649, 2019. [MR3992394](#)
- [30] R. Henson. Forward inference using functional neuroimaging: Dissociations versus associations. *Trends in cognitive sciences*, 10:64–69, 2006.
- [31] Edward L Korn, James F Troendle, Lisa M McShane, and Richard Simon. Controlling the number of false discoveries: application to high-dimensional genomic data. *Journal of Statistical Planning and Inference*, 124(2):379–398, 2004. [MR2080371](#)
- [32] Michael R Kosorok. Bootstraps of sums of independent but not identically distributed stochastic processes. *Journal of Multivariate Analysis*, 84(2):299–318, 2003. [MR1965224](#)
- [33] Jinzhou Li, Marloes H Maathuis, and Jelle J Goeman. Simultaneous false discovery proportion bounds via knockoffs and closed testing. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkae012, March 2024.
- [34] Enno Mammen. Bootstrap and wild bootstrap for high dimensional linear models. *The annals of statistics*, 21(1):255–285, 1993. [MR1212176](#)
- [35] Bryan FJ Manly. Randomization and regression methods for testing for associations with geographical, environmental and biological distances between populations. *Researches on Population Ecology*, 28(2):201–218, 1986.
- [36] Nicolai Meinshausen. False discovery control for multiple tests of association under general dependence. *Scandinavian Journal of Statistics*, 33(2):227–237, 2006. [MR2279639](#)
- [37] Pierre Neuvial. *Contributions to statistical inference from genomic data*. Habilitation thesis, Université Toulouse III (France), 2020. Available from <https://tel.archives-ouvertes.fr/tel-02969229>.
- [38] Fatma Parlak et al. A robust multivariate, non-parametric outlier identification method for scrubbing in fMRI. *arXiv preprint arXiv:2304.14634*, 2023.
- [39] Joseph P Romano and Michael Wolf. Exact and Approximate Stepdown Methods for Multiple Hypothesis Testing. *Journal of the American Statistical Association*, 100(469):94–108, 2005. [MR2156821](#)
- [40] Jonathan D Rosenblatt, Livio Finos, Wouter D Weeda, Aldo Solari, and Jelle J Goeman. All-resolutions inference for brain imaging. *Neuroimage*, 181:786–796, 2018.
- [41] Emiliano Santarnecchi, Alexandra Emmendorfer, Sayedhedayatollah Tadayon, Simone Rossi, Alessandro Rossi, and Alvaro Pascual-Leone. Network connectivity correlates of variability in fluid intelligence performance. *Intelligence*, 65:35–47, 2017.
- [42] Galen R Shorack. Bootstrapping robust regression. *Communications in statistics-theory and methods*, 11(9):961–972, 1982. [MR0655465](#)
- [43] R J Simes. An improved Bonferroni procedure for multiple tests of significance. *Biometrika*, 73(3):751–754, 1986. [MR0897872](#)
- [44] Stephen M Smith, Thomas E Nichols, Diego Vidaurre, Anderson M Winkler, Timothy EJ Behrens, Matthew F Glasser, Kamil Ugurbil, Deanna M

- Barch, David C Van Essen, and Karla L Miller. A positive-negative mode of population covariation links brain connectivity, demographics and behavior. *Nature neuroscience*, 18(11):1565–1567, 2015.
- [45] Gordon K. Smyth. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Statistical Methods in Genetics and Molecular Biology*, 3(3), 2004. [MR2101454](#)
- [46] John D Storey and Robert Tibshirani. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences*, 100(16):9440–9445, 2003. [MR1994856](#)
- [47] Fabian Telschow and Samuel Davenport. Precise FWER Control for Gaussian Related Fields: Finding the SuRF to continuous land – Part 1, 2023.
- [48] A.W. van der Vaart. *Asymptotic Statistics*. 1998. [MR1652247](#)
- [49] Peter H. Westfall. On Using the Bootstrap for Multiple Comparisons. *Journal of Biopharmaceutical Statistics*, 21(6):1187–1205, 2011. [MR2861901](#)
- [50] Keith J. Worsley, Sean Marrett, Peter Neelin, Alain C Vandal, Karl J Friston, and Alan C Evans. A unified statistical approach for determining significant signals in images of cerebral activation. *Human brain mapping*, 4(1):58–73, 1996.