

<b>Statistica Sinica Preprint No: SS-2017-0468</b>	
<b>Title</b>	Optimal Rates and Tradeoffs in Multiple Testing
<b>Manuscript ID</b>	SS-2017-0468
<b>URL</b>	<a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>
<b>DOI</b>	10.5705/ss.202017.0468
<b>Complete List of Authors</b>	Maxim Rabinovich Aaditya Ramdas Michael I. Jordan and Martin J. Wainwright
<b>Corresponding Author</b>	Maxim Rabinovich
<b>E-mail</b>	rabinovich@berkeley.edu
Notice: Accepted version subject to English editing.	

## Optimal Rates and Tradeoffs in Multiple Testing

Maxim Rabinovich<sup>†</sup>   Aaditya Ramdas<sup>\*,†</sup>

Michael I. Jordan<sup>\*,†</sup>   Martin J. Wainwright<sup>\*,†</sup>

*Departments of Statistics\* and EECS<sup>†</sup>, University of California, Berkeley*

*Abstract:* Multiple hypothesis testing is a central topic in statistics, but despite abundant work on the false discovery rate (FDR) and the corresponding Type II error concept known as the false non-discovery rate (FNR), a fine-grained understanding of the fundamental limits of multiple testing has not been developed. The main contribution of this paper is to derive a precise non-asymptotic trade-off between FNR and FDR for a variant of the generalized Gaussian sequence model. Our analysis is flexible enough to permit analyses of settings where the problem parameters vary with the number of hypotheses  $n$ , including various sparse and dense regimes (with  $o(n)$  and  $\mathcal{O}(n)$  signals). Moreover, we prove that the Benjamini-Hochberg algorithm as well as the Barber-Candès algorithm are both rate-optimal up to constants across these regimes.

*Key words and phrases:* multiple testing, minimax lower bounds

---

## 1. Introduction

The problem of multiple comparisons has been a central topic in statistics ever since Tukey's influential book (Tukey, 1953). In broad terms, suppose that one observes a sequence of  $n$  independent random variables  $X_1, \dots, X_n$ , of which some unknown subset are drawn from a null distribution, corresponding to the absence of a signal or effect, whereas the remainder are drawn from a non-null distribution, corresponding to signals or effects. Within this framework, one can pose three problems of increasing hardness: the *detection* problem of testing whether or not there is at least one signal; the *localization* problem of identifying the positions of the nulls and signals; and the *estimation* problem of returning estimates of the means and/or distributions of the observations. Note that these problems form a hierarchy of difficulty: identifying the signals implies that we know whether there is at least one of them, and estimating each mean implies we know which are zero and which are not. The focus of this paper is on the problem of localization.

There are a variety of ways of measuring type I errors for the localization problem, including the *family-wise error rate*, which is the probability of incorrectly rejecting at least one null, and the *false discovery rate* (FDR), which is the expected ratio of incorrect rejections to total rejec-

---

tions. An extensive literature has developed around both of these metrics, resulting in algorithms geared towards controlling one or the other. Our focus is the FDR metric, which has been widely studied, but for which relatively little is known about the behavior of existing algorithms in terms of the corresponding type II error concept, namely the *false non-discovery rate* (FNR). Indeed, it is only very recently that Arias-Castro and Chen (2017), working within a version of the sparse generalized Gaussian sequence model, established asymptotic consistency for the FDR-FNR localization problem. Informally, in this framework, we receive  $n$  independent observations  $X_1, \dots, X_n$ , out of which  $n^{1-\beta_n}$  are non-nulls, and the remainder are nulls. The  $n - n^{1-\beta_n}$  null variables are drawn from a centered distribution with tails decaying as  $\exp(-\frac{|x|^\gamma}{\gamma})$ , whereas the non-nulls are drawn from the same distribution shifted by  $(\gamma r_n \log n)^{1/\gamma}$ . Using this notation, Arias-Castro and Chen (2017) considered the setting with fixed problem parameters  $r_n = r$  and  $\beta_n = \beta$ , and showed that when  $r < \beta < 1$ , all procedures must have risk  $\text{FDR} + \text{FNR} \rightarrow 1$ . They also showed that in the achievable regime  $r > \beta > 0$ , the Benjamini-Hochberg (BH) procedure is

---

We follow Arias-Castro and Chen (2017) in defining the FNR as the ratio of undiscovered to total non-nulls, which differs from the definition of Genovese and Wasserman (Genovese and Wasserman, 2002).

---

consistent, meaning that  $\text{FDR} + \text{FNR} \rightarrow 0$ . Finally, they proposed a new “distribution-free” method inspired by the knockoff procedure by Barber and Candès (2015), and they showed that the resulting procedure is also consistent in the achievable regime.

These existing consistency results are asymptotic. To date, there has been no study of the important non-asymptotic questions that are of interest in comparing procedures. For instance, for a given FDR level, what is the best possible achievable FNR? What is the best-possible non-asymptotic behavior of the risk  $\text{FDR} + \text{FNR}$  attainable in finite samples? And, perhaps most importantly, non-asymptotic questions, regarding whether or not procedures such as BC and BH are *rate-optimal* for the  $\text{FDR} + \text{FNR}$  risk, remain unanswered. The main contributions of this paper are to develop techniques for addressing such questions, and to essentially resolve them in the context of the sparse generalized Gaussians model.

Specifically, we establish the tradeoff between FDR and FNR in finite samples (and hence also asymptotically), and we use the tradeoff to determine the best attainable rate for the  $\text{FDR} + \text{FNR}$  risk. Our theory is sufficiently general to accommodate sequences of parameters  $(r_n, \beta_n)$ , and thereby to reveal new phenomena that arise when  $r_n - \beta_n = o(1)$ . For a fixed pair of parameters  $(r, \beta)$  in the achievable regime  $r > \beta$ , our theory

---

leads to an explicit expression for the optimal rate at which FDR+FNR can decay. In particular, defining the  $\gamma$ -“distance”  $D_\gamma(a, b) := |a^{1/\gamma} - b^{1/\gamma}|^\gamma$  between pairs of positive numbers, we show that the equation

$$\kappa = D_\gamma(\beta + \kappa, r)$$

has a unique solution  $\kappa_*$ , and moreover that the combined risk of any threshold-based multiple testing procedure  $\mathcal{I}$  is lower bounded as  $\mathcal{R}_n(\mathcal{I}) \gtrsim n^{-\kappa_*}$ . Moreover, by direct analysis, we are able to prove that both the Benjamini-Hochberg (BH) and the Barber-Candès (BC) algorithms attain this optimal rate.

At the core of our analysis is a simple comparison principle. The flexibility of the resulting proof strategy allows us to identify a new critical regime in which  $r_n - \beta_n = o(1)$ , but the problem is infeasible, meaning that if the FDR is driven to zero, then the FNR must remain bounded away from zero. Moreover, we are able to study some challenging settings in which the fraction of signals is a constant  $\pi_1 \in (0, 1)$  and not asymptotically vanishing, which corresponds to the setting  $\beta_n = \frac{\log(1/\pi_1)}{\log n}$ , so that  $\beta_n \rightarrow 0$ . Perhaps surprisingly, even in these regimes, the BH and BC algorithms continue to be optimal, though the best rate can weaken from polynomial to subpolynomial in the number of hypotheses  $n$ .

## 1.1 Related work

As noted above, our work provides a non-asymptotic generalization of recent work by Arias-Castro and Chen (2017) on asymptotic consistency in localization, using  $\text{FDR} + \text{FNR}$  as the notion of risk. It should be noted that this notion of risk is distinct from the asymptotic Bayes optimality under sparsity (ABOS) studied in past work by Bogdan et al. (2011) for Gaussian sequences, and more recently by Neuvial and Roquain (2012) for binary classification with extreme class imbalance. The ABOS results concern a risk derived from the probability of incorrectly rejecting a single null sample (false positive, or FP for short) and the probability of incorrectly failing to reject a single non-null sample (false negative, or FN for short). Concretely, one has  $\mathcal{R}_n^{\text{ABOS}} = w_1 \cdot \text{FP} + w_2 \cdot \text{FN}$  for some pair of positive weights  $(w_1, w_2)$  that need not be equal. As this risk is based on the error probability for a single sample, it is much closer to misclassification risk or single-testing risk than to the ratio-based  $\text{FDR} + \text{FNR}$  risk studied in this paper.

Using the notation of this paper, the work of Neuvial and Roquain (2012) can be understood as focusing on the particular setting  $r = \beta$ , a regime referred to as the “verge of detectability” by these authors, and with performance metric given by the Bayes classification risk, rather than

## 1.1 Related work

the combination of FDR and FNR studied here. In comparison, our results provide additional insight into models that are close to the verge of detectability, in that even when  $\beta_n = \beta$  is fixed, we can provide quantitative lower and upper bounds on the FDR/FNR ratio as  $r_n \rightarrow \beta$  from above; moreover, these bounds depend on how quickly  $r_n$  approaches  $\beta$ . These conclusions actually make it clear that a further transition in rates occurs in the case where  $r = \beta$  exactly for all  $n$ , though we do not explore the latter case in depth. We suspect that the methods developed in this paper may have sufficient precision to answer the non-asymptotic minimaxity questions posed by Neuvial and Roquain (2012) as to whether any threshold-based procedure can match the Bayes optimal classification error rate up to an additive error  $\ll \frac{1}{\log n}$ .

For the special case of  $\gamma = 2$ , Ji and Jin (2012) and Ji and Zhao (2014) prove bounds for localization that are closely related to, but distinct from, our bounds on the overall risk. Both deal with sparse high-dimensional regression: the former work proposes a new method called UPS for variable selection that has advantages over Lasso and Subset Selection in certain settings, while the latter builds on the first to prove upper and lower bounds for multiple testing using the so-called mFNR and mFDR. These metrics replace the expected ratio in the definitions of FDR and FNR (see defi-



## 1.1 Related work

---

nition (2.3) below) by a ratio of expectations—a modification that should lead to qualitatively similar behavior as  $n$  becomes large. The resulting bounds in both papers can be used to recover our bounds up to polylogarithmic factors in the special case where  $\gamma = 2$ . The main advantage of their work, relative to ours, is the handling of dependence between the p-values. Unlike our work, however, they do not establish the tradeoff between FDR and FNR when both quantities can decay to zero at different rates, and as mentioned they only consider  $\gamma = 2$ . Nor do they consider regimes where sparsity and signal strength vary with  $n$ . Our results can handle this more general setting, which encompasses dense regimes with qualitatively different behavior from the more commonly investigated sparse one.

The above line of work is complementary to the well-known asymptotic results by Donoho and Jin (2004; 2015) on phase transitions in detectability using Tukey’s higher-criticism statistic, employing the standard type I and type II errors for testing of the single global null hypothesis. Note that Donoho and Jin use the generalized Gaussian assumption directly on the PDFs, while our assumption (2.5) is on the survival function. Just as in Arias-Castro and Chen (2017), Donoho and Jin also consider the asymptotic setting with  $r_n = r$  and  $\beta_n = \beta$ , which they sometimes call the RW (rare and weak) model. We are not aware of any non-asymptotic results for

---

## 1.1 Related work

detection akin to the results that the current paper provides for localization.

Our paper is also complementary to work on estimation, the most notable result being the asymptotic minimax optimality of BH-derived thresholding for denoising an approximately-sparse high-dimensional vector (Abramovich et al., 2006; Donoho and Jin, 2006). The relevance of our results on the minimaxity of BH for approximately-sparse denoising problems lies primarily in the use of deterministic thresholds as a useful proxy for BH and other procedures that determine their threshold in a manner that has complex dependence on the input data (Donoho and Jin, 2006). Unlike the strategy of Donoho and Jin (2006), which depends on establishing concentration of the empirical threshold around the population-level value, we use a more flexible comparison principle. Deterministic approximations to optimal FDR thresholds are also studied by Chi (2007) and Genovese et al. (2006). Other related papers are discussed in Section 5, when discussing directions for future work.

The remainder of this paper is organized as follows. In Section 2, we provide background on the multiple testing problem, as well as the particular model we consider. In Section 3, we provide an overview of our main results: namely, optimal tradeoffs between FDR and FNR, which imply lower bounds on the FDR+FNR risk, and optimality guarantees for the

---

BH and BC algorithms. In Section 4, we prove our main results, focusing first on the lower bounds and then using the ideas we have developed to provide matching upper bounds for the well-known and popular Benjamini-Hochberg (BH) procedure and the recent Barber-Candès (BC) algorithm for multiple testing with FDR control. Proofs of some technical lemmas are given in the appendices.

## 2. Problem formulation

In this section, we provide background and a precise formulation of the problem under study.

### 2.1 Multiple testing and false discovery rate

Suppose that we observe a real-valued sequence  $X_1^n := \{X_1, \dots, X_n\}$  of  $n$  independent random variables. When the null hypothesis is true,  $X_i$  is assumed to have zero mean; otherwise, it is assumed that the mean of  $X_i$  is equal to some unknown number  $\mu_n > 0$ . We introduce the sequence of binary labels  $\{H_1, \dots, H_n\}$  to encode whether or not the null hypothesis holds for each observation; the setting  $H_i = 0$  indicates that the null hypothesis holds. We define

$$\mathcal{H}_0 := \{i \in [n] \mid H_i = 0\}, \quad \text{and} \quad \mathcal{H}_1 := \{i \in [n] \mid H_i = 1\}, \quad (2.1)$$

## 2.1 Multiple testing and false discovery rate

corresponding to the *nulls* and *signals*, respectively. Our task is to identify a subset of indices that contains as many signals as possible, while not containing too many nulls.

More formally, a testing rule  $\mathcal{I} : \mathbb{R}^n \rightarrow 2^{[n]}$  is a measurable mapping of the observation sequence  $X_1^n$  to a set  $\mathcal{I}(X_1^n) \subseteq [n]$  of *discoveries*, where the subset  $\mathcal{I}(X_1^n)$  contains those indices for which the procedure rejects the null hypothesis. There is no single unique measure of performance for a testing rule for the localization problem. In this paper, we study the notion of the *false discovery rate* (FDR), paired with the *false non-discovery rate* (FNR). These can be viewed as generalizations of the type I and type II errors for single hypothesis testing.

We begin by defining the false discovery proportion (FDP), and false non-discovery proportion (FNP), respectively, as

$$\text{FDP}_n(\mathcal{I}) := \frac{\text{card}(\mathcal{I}(X_1^n) \cap \mathcal{H}_0)}{\text{card}(\mathcal{I}(X_1^n)) \vee 1}, \quad \text{and} \quad \text{FNP}_n(\mathcal{I}) := \frac{\text{card}(\mathcal{I}(X_1^n)^c \cap \mathcal{H}_1)}{\text{card}(\mathcal{H}_1)}. \quad (2.2)$$

Since the output  $\mathcal{I}(X_1^n)$  of the testing procedure is random, both quantities are random variables. The FDR and FNR are given by taking the

## 2.1 Multiple testing and false discovery rate

expectations of these random quantities—that is

$$\text{FDR}_n(\mathcal{I}) := \mathbb{E} \left[ \frac{\text{card}(\mathcal{I}(X_1^n) \cap \mathcal{H}_0)}{\text{card}(\mathcal{I}(X_1^n)) \vee 1} \right], \quad \text{and} \quad \text{FNR}_n(\mathcal{I}) := \mathbb{E} \left[ \frac{\text{card}(\mathcal{I}(X_1^n)^c \cap \mathcal{H}_1)}{\text{card}(\mathcal{H}_1)} \right], \quad (2.3)$$

where the expectation is taken over the random samples  $X_1^n$ .

It is worth noting that our definition of FNP and FNR, which follows that of Arias-Castro and Chen (2017), differs from an alternative definition of  $\text{FNR}_{alt}$ , where the denominator is set to the number of non-rejections. In general, however, the number of non-rejections will be close to  $n$  for any procedure with low FDR and thus in the sparse regime, the  $\text{FNR}_{alt}$  would trivially go to zero for any procedure that controls FDR at any level strictly below one. Our definition is therefore better suited to studying transitions in difficulty in the multiple testing problem.

In this paper, we measure the overall performance of a procedure in terms of its *combined risk*

$$\mathcal{R}_n(\mathcal{I}) := \text{FDR}_n(\mathcal{I}) + \text{FNR}_n(\mathcal{I}). \quad (2.4)$$

Finally, when the testing rule  $\mathcal{I}$  under discussion is clear from the context, we frequently omit explicit reference to this dependence from all of these quantities.

## 2.2 Tail generalized Gaussians model

### 2.2 Tail generalized Gaussians model

In this paper, we describe the distribution of the observations for both nulls and non-nulls in terms of a *tail generalized Gaussians model*. Our model is a variant of the generalized Gaussian sequence model studied in past work (Arias-Castro and Chen, 2017; Donoho and Jin, 2004); the only difference is that whereas a  $\gamma$ -generalized Gaussian has a density proportional to  $\exp\left(-\frac{|x|^\gamma}{\gamma}\right)$ , we focus on distributions whose tails are proportional to  $\exp\left(-\frac{|x|^\gamma}{\gamma}\right)$ . This alteration is in line with the asymptotically generalized Gaussian (AGG) distributions studied by Arias-Castro and Chen (2017), with the important caveat that our assumptions are imposed in a non-asymptotic fashion.

For a given degree  $\gamma \geq 1$ , a  $\gamma$ -tail generalized Gaussian random variable with mean 0, written as  $G \sim \text{tGG}_\gamma(0)$ , has a survival function  $\Psi(t) := \mathbb{P}(G \geq t)$  that satisfies the bounds

$$\frac{e^{-\frac{|t|^\gamma}{\gamma}}}{Z_\ell} \leq \min\{\Psi(t), 1 - \Psi(t)\} \leq \frac{e^{-\frac{|t|^\gamma}{\gamma}}}{Z_u}, \quad t \in \mathbb{R}, \quad (2.5)$$

for some constants  $Z_\ell > Z_u > 0$ . (Note that  $t \mapsto \Psi(t)$  is a decreasing function, and becomes smaller than  $1 - \Psi(t)$  at the origin.) As a concrete example, a  $\gamma$ -tail generalized Gaussian with  $Z_\ell = Z_u = 1$  can be generated by sampling a standard exponential random variable  $E$  and

### 2.3 Threshold-based procedures

a Rademacher random variable  $\varepsilon$  and putting  $G = \varepsilon(\gamma E)^{1/\gamma}$ . We use the terminology “tail generalized Gaussian” because of the following connection: the survival function of a 2-tail Gaussian random variable is on the order of  $\exp(-|x|^2/2)$ , whereas that of a Gaussian is on the order of  $\frac{1}{\text{poly}(x)} \exp(-x^2/2)$ . In particular, this observation implies a  $\text{tGG}_2$  random variable has tails that are equivalent to a Gaussian in terms of their exponential decay rates.

In terms of this notation, we assume that each observation  $X_i$  is distributed as

$$X_i \sim \begin{cases} \text{tGG}_\gamma(0) & \text{if } i \in \mathcal{H}_0 \\ \text{tGG}_\gamma(0) + \mu_n & \text{if } i \in \mathcal{H}_1, \end{cases} \quad (2.6)$$

where our notation reflects the fact that the mean shift  $\mu_n$  is permitted to vary with the number of observations  $n$ . See Section 3.1 for further discussion of the scaling of the mean shift.

### 2.3 Threshold-based procedures

Following prior work (Arias-Castro and Chen, 2017; Donoho and Jin, 2004), we restrict attention to testing procedures of the form

$$\mathcal{I}(X_1^n) = \{i \in [n] \mid X_i \geq T_n(X_1^n)\}, \quad (2.7)$$

## 2.4 Benjamini-Hochberg (BH) and Barber-Candès (BC) procedures

where  $T_n(X_1^n) \in \mathbb{R}_+$  is a data-dependent threshold. We refer to such methods as *threshold-based procedures*. The BH and BC procedures both belong to this class. Moreover, from an intuitive standpoint, the observations are exchangeable in the absence of prior information, and we are considering testing between a single unimodal null distribution and a single positive shift of that distribution. In this setting, it is hard to conceive of reasonable procedures that would reject the hypothesis corresponding to one observation while rejecting a hypothesis with a smaller observation value.

It will be convenient to reason about the performance metrics associated with rules of the form

$$\mathcal{I}_t(X_1^n) = \{i \in [n] \mid X_i \geq t\}, \quad (2.8)$$

where  $t > 0$  is a pre-specified (fixed, non-random) threshold. In this case, we adopt the notation  $\text{FDR}_n(t)$ ,  $\text{FNR}_n(t)$  and  $\mathcal{R}_n(t)$  to denote the metrics associated with the rule  $X_1^n \mapsto \mathcal{I}_t(X_1^n)$ .

## 2.4 Benjamini-Hochberg (BH) and Barber-Candès (BC) procedures

Arguably the most popular threshold-based procedure that provably controls FDR at a user-specified level  $q_n$  is the *Benjamini-Hochberg* (BH) procedure. More recently, Arias-Castro and Chen (2017) proposed a method



## 2.4 Benjamini-Hochberg (BH) and Barber-Candès (BC) procedures

that we refer to as the *Barber-Candès* (BC) procedure. Both algorithms are based on estimating the  $\text{FDP}_n$  that would be incurred at a range of possible thresholds and choosing one that is as large as possible (maximizing discoveries) while satisfying an upper bound linked to  $q_n$  (controlling  $\text{FDR}_n$ ). Further, they both only consider thresholds that coincide with one of the values  $X_1^n$ , which we denote as a set by  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ . The data-dependent threshold for both can be written as

$$t_n(X_1, \dots, X_n) = \min \{t \in \mathcal{X}_n : \widehat{\text{FDP}}_n(t) \leq q_n\}. \quad (2.9)$$

The two algorithms differ in the estimator  $\widehat{\text{FDP}}_n(t)$  they use. The BH procedure assumes access to the true null distribution through its survival function  $\Psi$  and sets

$$\widehat{\text{FDP}}_n^{\text{BH}}(t) = \frac{\Psi(t)}{\#(X_i \geq t)/n}, \quad \text{for } t \in \mathcal{X}_n. \quad (2.10)$$

The BC procedure instead estimates the survival function  $\Psi(t)$  from the data and therefore does not even need to know the null distribution. This approach is viable when  $\#(X_i \leq -t)/n$  is a good proxy for  $\Psi(t)$ , which our upper and lower tail bounds guarantee; more typically, the BC procedure is applicable when the null distribution is (nearly) symmetric, and the signals are shifted by a positive amount (as they are in our case). Then, the BC

---

estimator is given by

$$\widehat{\text{FDP}}_n^{\text{BC}}(t) = \frac{[\#(X_i \leq -t) + 1]/n}{\#(X_i \geq t)/n}, \quad \text{for } t \in \mathcal{X}_n. \quad (2.11)$$

With these definitions in place, are now ready to describe our main results.

### 3. Main results

We now turn to a statement of our main results, along with some illustrations of their consequences. Our first main result (Theorem 1) characterizes the optimal tradeoff between FDR and FNR for any testing procedure. By optimizing this tradeoff, we obtain a lower bound on the combined FDR and FNR of any testing procedure (Corollary 1). Our second main result (Theorem 2), shows that BH achieves the optimal FDR-FNR tradeoff up to constants and that BC almost achieves it. In particular, our result implies that with the proper choice of target FDR, both BH and BC can achieve the optimal combined FDR-FNR rate (Corollary 2).

#### 3.1 Scaling of sparsity and mean shifts

We study a sparse instance of the multiple testing problem in which the number of signals is assumed to be small relative to the total number of hypotheses. In particular, motivated by related work in multiple hypothesis testing (Arias-Castro and Chen, 2017; Donoho and Jin, 2004, 2015; Jin and

### 3.1 Scaling of sparsity and mean shifts

Ke, 2014), we assume that the number of signals scales as

$$\text{card}(\mathcal{H}_1) = m_n = n^{1-\beta_n} \quad \text{for some } \beta_n \in (0, 1). \quad (3.12)$$

Note that to the best of our knowledge, all previous results in the literature assume that  $\beta_n = \beta$  is actually independent of  $n$ . In this case, the sparsity assumption (3.12) implies that all but a polynomially vanishing fraction of the hypotheses are null. In contrast, as indicated by our choice of notation, the set-up in this paper allows for a sequence of parameters  $\beta_n$  that can vary with the number of hypotheses  $n$ . In this way, our framework is flexible enough to handle relatively dense regimes (e.g., those with  $\frac{n}{\log n}$  or even  $\mathcal{O}(n)$  signals).

The non-null hypotheses are distinguished by a positively shifted mean  $\mu_n > 0$ . It is natural to parameterize this mean shift in terms of a quantity  $r_n > 0$  via the relation

$$\mu_n = (\gamma r_n \log n)^{1/\gamma}. \quad (3.13)$$

As shown by Arias-Castro and Chen (2017), when the pair  $(\beta, r)$  are fixed such that  $r < \beta$ , the problem is asymptotically infeasible, meaning that there is no procedure such that  $\mathcal{R}_n(\mathcal{I}) \rightarrow 0$  as  $n \rightarrow \infty$ . Accordingly, we focus on sequences  $(\beta_n, r_n)$  for which  $r_n > \beta_n$ . Further, even though the asymptotic consistency boundary of  $r < \beta$  versus  $r > \beta$  is apparently

### 3.2 Lower bound on any threshold-based procedure

independent of  $\gamma$ , we will see that the rate at which the risk decays to zero is determined jointly by  $r, \beta$  and  $\gamma$ .

### 3.2 Lower bound on any threshold-based procedure

In this section, we assume :

$$\beta_n \stackrel{(i)}{\geq} \frac{\log 2}{\log n} \iff n^{1-\beta_n} \leq n/2, \quad \text{and} \quad (3.14a)$$

$$\max\{\beta_n, \frac{1}{\log^{\frac{\gamma-1/2}{\gamma}} n}\} \stackrel{(ii)}{<} r_n \stackrel{(iii)}{<} r_{\max} \quad \text{for some constant } r_{\max} < 1. \quad (3.14b)$$

Condition (i) requires that the proportion  $\pi_1$  of non-nulls is at most  $1/2$ .

Condition (ii) asserts that the natural requirement of  $r_n > \beta_n$  is not enough,

but further insists that  $r_n$  cannot approach zero too fast. The constants

$\log 2$  and  $\frac{\gamma-1/2}{\gamma}$  are somewhat arbitrary and can be replaced, respectively,

by  $\log \frac{1}{\pi_{\max}}$  for any  $0 < \pi_{\max} < 1$  and  $\frac{\gamma-1+\rho}{\gamma}$  for any  $\rho > 0$ , but we fix

their values in order not to introduce unnecessary extra parameters. As

for condition (iii), although the assumption  $r_n < 1$  is imposed because the

problem becomes qualitatively easy for  $r_n \geq 1$ , the assumption that it is

bounded away from one is a technical convenience that simplifies some of

our proofs.

Our analysis shows that the FNR behaves differently depending on the

closeness of the parameter  $r_n$  to the boundary of feasibility given by  $\beta_n$ . In

### 3.2 Lower bound on any threshold-based procedure

order to characterize this closeness, we define

$$r_{\min} = r_{\min}(\kappa_n) := \begin{cases} \beta_n + \kappa_n + \frac{\log \frac{1}{6Z_\ell}}{\log n} & \text{if } \kappa_n \leq 1 - \beta_n - \frac{\log \frac{3}{\log 16}}{\log n}, \\ 1 + \frac{\log \frac{1}{24Z_\ell}}{\log n} & \text{otherwise.} \end{cases} \quad (3.15)$$

Here  $\kappa_n$  is to be interpreted as the “exponent” of a target FDR rate  $q_n$ , in the sense that  $q_n = n^{-\kappa_n}$ . The rate  $q_n$  may differ from the actual achieved  $\text{FDR}_n$ , but it is nonetheless useful for parameterizing the quantities that enter into our analysis. When we need to move between  $q_n$  and  $\kappa_n$ , we shall write  $\kappa_n = \kappa_n(q_n) = \frac{\log(1/q_n)}{\log n}$  and  $q_n = q_n(\kappa_n) = n^{-\kappa_n}$ . For mathematical convenience, we wish to have the target FDR  $q_n$  to be bounded away from one, and we therefore impose one further technical but inessential assumption in this section:

$$q_n \leq \min \left\{ \frac{1}{24}, \frac{1}{6Z_\ell} \right\} \iff \kappa_n \geq \frac{\log \max \{24, 6Z_\ell\}}{\log n}. \quad (3.16)$$

The theorem that follows will apply to all sample sizes  $n > n_{\min, \ell}$  (subscript  $\ell$  for lower), where

$$n_{\min, \ell} := \min \left\{ n \in \mathbb{N} : \exp \left( -\frac{n^{1-r_{\max}}}{24(Z_\ell \vee 1)} \right) \leq \frac{1}{4} \right\} \quad (3.17)$$

$$= \left\lceil [24(Z_\ell \vee 1) \log 4]^{\frac{1}{1-r_{\max}}} \right\rceil, \quad (3.18)$$

which is an explicit known function of the problem parameters and can therefore be computed whenever the problem setting is fixed.

### 3.2 Lower bound on any threshold-based procedure

Finally, for  $\gamma \in [1, \infty)$  and non-negative numbers  $a, b > 0$ , let us define the associated  $\gamma$ -“distance”:

$$D_\gamma(a, b) := |a^{1/\gamma} - b^{1/\gamma}|^\gamma. \quad (3.19)$$

Our first main theorem states that for  $r_n > r_{\min}(\kappa_n)$ , the FNR decays as a power of  $1/n$ , with exponent specified by the  $\gamma$ -distance.

**Theorem 1.** Consider the  $\gamma$ -tail generalized Gaussians testing problem with sparsity  $\beta_n$  and signal level  $r_n$  satisfying conditions (3.14a), and (3.14b), and with sample size  $n > n_{\min, \ell}$  from definition (3.18). Then, for any choice of exponent  $\kappa_n \in (0, 1)$  satisfying condition (3.16), there exists a minimum signal strength  $r_{\min}(\kappa_n)$  from definition (3.15), such that any threshold-based procedure  $\mathcal{I}$  that satisfies  $\text{FDR}_n(\mathcal{I}) \leq n^{-\kappa_n}$  must have its FNR lower bounded as

$$\text{FNR}_n(\mathcal{I}) \geq \begin{cases} \frac{1}{32} & \text{if } r_n \in [\beta_n, r_{\min}] \\ c(\beta_n, \gamma) n^{-D_\gamma(\beta_n + \kappa_n, r_n)} & \text{otherwise,} \end{cases} \quad (3.20)$$

where  $c(\beta_n, \gamma) := c_0 \exp(c_1 \beta_n^{\frac{1-\gamma}{\gamma}})$ , with  $(c_0, c_1)$  being positive constants depending only on  $(Z_\ell, Z_u, \gamma)$ .

The proof of this theorem is provided in Section 4.1. Note that the theorem holds for any choice of  $\kappa_n \in (0, 1)$ . In the special case of constant pairs

### 3.2 Lower bound on any threshold-based procedure

$(\beta, r)$ , this choice can be optimized to achieve the best possible lower bound on the risk  $\mathcal{R}_n(\mathcal{I}) = \text{FDR}_n(\mathcal{I}) + \text{FNR}_n(\mathcal{I})$ . **Since we obtain this lower bound by optimizing the sum of the FDR and FNR lower bounds from Theorem 1, we want to balance the contributions from these two bounds. Doing so requires us to set the FDR rate  $\kappa$  equal to the corresponding FNR rate  $D_\gamma(\beta + \kappa, r)$ , which leads to a fixed-point equation for the overall rate, as summarized below.**

**Corollary 1.** When  $r > \beta$ , let  $\kappa_* = \kappa_*(\beta, r, \gamma) > 0$  be the unique solution to the equation

$$\kappa = D_\gamma(\beta + \kappa, r). \quad (3.21)$$

Then the combined risk of any threshold-based multiple testing procedure  $\mathcal{I}$  is lower bounded as

$$\mathcal{R}_n(\mathcal{I}) \gtrsim n^{-\kappa_*}, \quad (3.22)$$

where  $\gtrsim$  denotes inequality up to a pre-factor independent of  $n$ .

The proof of this corollary is provided in Appendix S3. Figure 1 provides an illustration of the predictions in Corollary 1. In particular, panel (a) shows how the unique solution  $\kappa_*$  to equation (3.21) is determined for varying settings of the triple  $(r, \beta, \gamma)$ . Panel (b) shows how  $\kappa_*$  varies over the

### 3.3 Upper bounds for some specific procedures

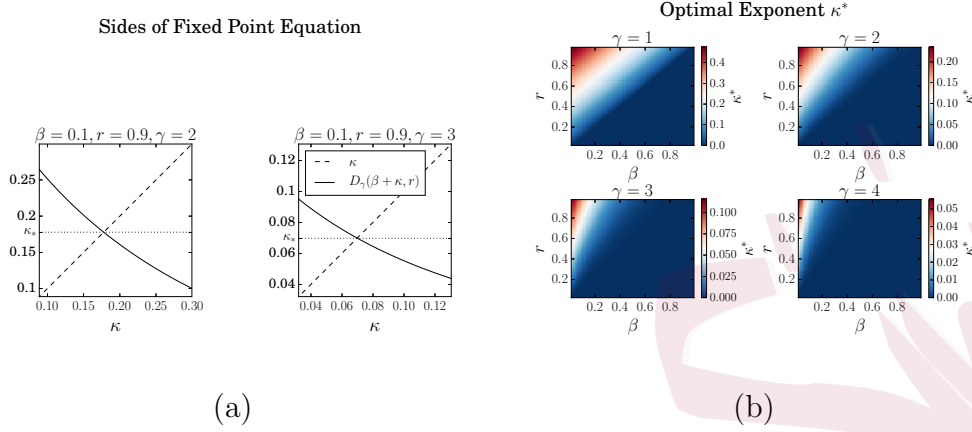


Figure 1: Visualizations of the fixed-point equation (3.21). (a) Plots comparing the left- and right-hand sides of the fixed-point equation. (b) The optimal exponent  $\kappa^*$  as a function of  $r$  and  $\beta$ .

interval  $(0, 0.5)$ , again for different settings of the triple  $(r, \beta, \gamma)$ . As would be expected, the fixed point  $\kappa_*$  increases as a function of the difference  $r - \beta > 0$ .

### 3.3 Upper bounds for some specific procedures

Thus far, we have provided general lower bounds applicable to any threshold procedure. We now turn to the complementary question—how do these lower bounds compare to the results achievable by the BH and BC algorithms introduced in Section 2.3? Remarkably, we find that up to the constants defining the prefactor, both the BH and BC procedures achieve



### 3.3 Upper bounds for some specific procedures

the minimax lower bound of Theorem 1.

We state these achievable results in terms of the fixed point  $\kappa_*$  from equation (3.21). Moreover, they apply to all problems with sample size  $n > n_{\min,u}$  (subscript  $u$  for upper), where

$$\begin{aligned} n_{\min,u} &:= \min \left\{ n \in \mathbb{N} : \exp \left( -\frac{n^{1-r_{\max}}}{24} \right) \leq \frac{1}{Z_u n} \right\} \\ &= \min \left\{ n \in \mathbb{N} : n \geq [24 \log(Z_u n)]^{\frac{1}{1-r_{\max}}} \right\}. \end{aligned} \quad (3.23)$$

Just as we did for (3.18), we note that this lower bound on  $n$  is explicitly computable from problem parameters.

In order to state our results cleanly, let us introduce the constants

$$c_{\text{BH}} := \frac{Z_u}{36Z_\ell}, \quad c_{\text{BC}} := \frac{Z_u}{48Z_\ell}, \quad \text{and} \quad \zeta := \max \left\{ 6Z_\ell, \frac{1}{6Z_\ell} \right\}, \quad (3.24)$$

and require in particular that  $r_n \geq r_{\min}(\kappa_n(c_A q_n))$  for algorithm  $A \in \{\text{BH}, \text{BC}\}$ . Note that  $c_A < 1$  since  $Z_\ell \geq Z_u$  by definition, and that the introduction of  $c_A$  into the argument of  $r_{\min}$  only changes the minimum allowed value of  $r_n$  by a conceptually negligible amount of  $\mathcal{O}(\frac{1}{\log n})$ .

Lastly, we note that BC requires an additional mild condition that the number of non-nulls  $n^{1-\beta_n}$  is large relative to the target FDR  $q_n = n^{-\kappa_n}$  (otherwise, in some sense, the problem is too hard if there are too few non-nulls and a very strict target FDR). Specifically, we need that both quantities cannot simultaneously be too small, formalized by the assump-

### 3.3 Upper bounds for some specific procedures

tion:

$$\exists n_{\min, \text{BC}} \text{ such that for all } n \geq n_{\min, \text{BC}} \text{ we have } \frac{3c_{\text{BC}}}{4} \cdot \frac{q_n}{\log \frac{1}{q_n}} \cdot n^{1-\beta_n} \geq 1. \quad (3.25)$$

We note that when  $r_n = r$  and  $\beta_n = \beta$  are constants, this decay condition is satisfied by  $q_n = n^{-\kappa_*}$ .

Our second main theorem delivers an optimality result for the BH and BC procedures, showing that under some regularity conditions, their performance achieves the lower bounds in Theorem 1 up to constant factors.

**Theorem 2.** Consider the  $\beta_n$ -sparse  $\gamma$ -tail generalized Gaussians testing problem with target FDR level  $q_n$  upper bounded as in condition (3.16).

(a) Guarantee for BH procedure: Given a signal strength  $r_n \geq r_{\min}(\kappa_n(c_{\text{BH}}q_n))$

and sample size  $n > n_{\min, u}$  as in condition (3.23), the BH procedure

satisfies the bounds

$$\text{FDR}_n \leq q_n \text{ and } \text{FNR}_n \leq \frac{2\zeta_{\text{BH}}^{\frac{1-\gamma}{2\beta_n}}}{Z_u} \cdot n^{-D_\gamma(\beta_n + \kappa_n, r_n)}, \quad \text{where } \zeta_{\text{BH}} := \frac{\zeta}{c_{\text{BH}}}. \quad (3.26)$$

(b) Guarantee for BC procedure: Given a signal strength  $r_n \geq r_{\min}(\kappa_n(c_{\text{BC}}q_n))$

and sample size  $n > \max\{n_{\min, \text{BC}}, n_{\min, u}\}$  as in condition (3.25), the

### 3.3 Upper bounds for some specific procedures

BC procedure satisfies the bounds

$$\text{FDR}_n \leq q_n \quad \text{and} \quad \text{FNR}_n \leq \frac{2\zeta_{\text{BC}}^{2\beta_n \frac{1-\gamma}{\gamma}}}{Z_u} \cdot n^{-D_\gamma(\beta_n + \kappa_n, r_n)} + q_n, \quad \text{where } \zeta_{\text{BC}} := \frac{\zeta}{c_{\text{BC}}}. \quad (3.27)$$

The proof of the theorem can be found in Section 4.2. For constant pairs  $(r, \beta)$ , Theorem 2 can be applied with a target FDR proportional to  $n^{-\kappa_*}$  to show that both BH and BC achieve the optimal decay of the combined FDR-FNR up to constant factors, as stated formally below.

**Corollary 2.** For  $\beta < r$  and  $q_* = c_* n^{-\kappa_*}$  with  $0 < c_* \leq \min\{\frac{1}{24}, \frac{1}{6Z_\ell}\}$ , the BH and BC procedures with target FDR  $q_*$  satisfy

$$\mathcal{R}_n \lesssim n^{-\kappa_*}. \quad (3.28)$$

The proof of this corollary is given in Appendix S5. To help visualize the result of the corollary, Figure 2 displays the results of some simulations of the BH procedure that show correspondence between its performance and the theoretically predicted rate of  $n^{-\kappa_*}$ .

Despite the optimality, Figures 1 and 2 paint a fairly dark picture from a practical point of view: while asymptotic consistency can be achieved when  $r > \beta$ , the convergence of the risk to zero can be extremely slow, exhibiting “nonparametric” rates far slower than  $n^{-1/2}$ . Figure 2 shows in particular that the decay to zero may be barely evident even for sample sizes as large

### 3.3 Upper bounds for some specific procedures

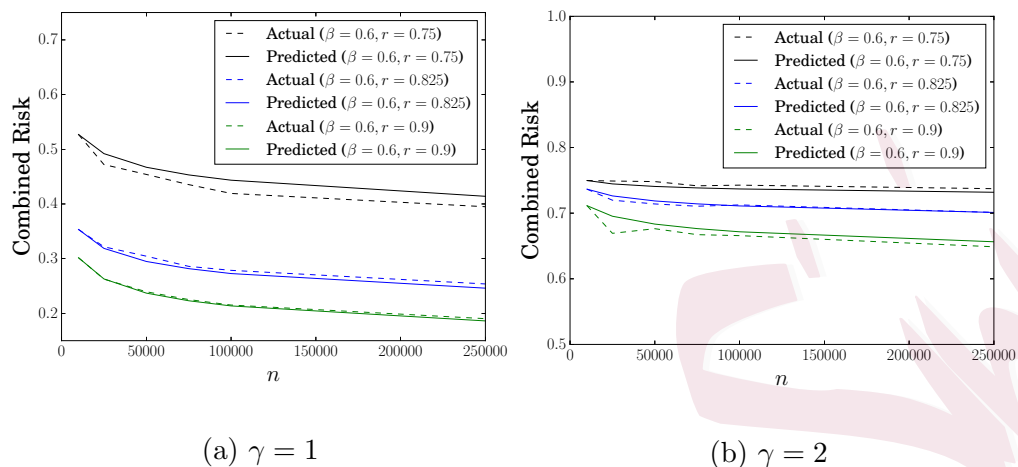


Figure 2: Results of simulations comparing the predicted combined risk with the actual experimentally-observed risk for the BH procedure. Agreement is good across the board and improves as the gap  $(r - \beta)$  increases. We believe the latter phenomenon arises because the sampling error is a smaller fraction of the risk as the separation increases.

as  $n = 250,000$ , even with comparatively strong signals. The “nonparametric” nature may arise because the dimensionality of the decision space increases linearly with sample size, and asymptotically, the upside of having increasing data seems to *just* overcome the downside of having to make an increasing number of decisions. However, non-asymptotically, one cannot hope to drive both FDR and FNR to zero at any practical sample size in this general setting, at least when the mean signal lies below the maximum

### 3.3 Upper bounds for some specific procedures

of the nulls (i.e.,  $r_n < 1$ ).

**Intuition for the  $\gamma$ -distance.** The distance  $D_\gamma$  plays a crucial role because of the scaling of order statistics under the  $\text{tGG}_\gamma$  model. If  $W_1, \dots, W_n$  are iid from a  $\text{tGG}_\gamma(0)$  model, then—ignoring constants inside the logarithm—we expect the  $i^{\text{th}}$ -largest order statistic  $W_{(i)}$  to be around  $(\gamma i \log n)^{1/\gamma}$  if  $i \ll n/2$  and around  $-(\gamma i \log n)^{1/\gamma}$  if  $n - i \ll n/2$ . If an algorithm is to achieve an FNR on the order of  $n^{-\kappa'}$ , it must successfully identify all but the smallest  $n^{-\kappa'}$  fraction of true signals. The algorithm's cutoff for rejection must thus exceed the  $m - n^{-\kappa'}m$  order statistic of the signals, which is approximately

$$\mu - \left(\gamma \log \frac{m}{n^{-\kappa'}m}\right)^{1/\gamma} = (\gamma r \log n)^{1/\gamma} - (\gamma \kappa' \log n)^{1/\gamma}. \quad (3.29)$$

If we suppose that the FDR is also vanishing at a rate  $n^{-\kappa}$ , then first of all the algorithm must identify about  $(1 \pm o(1))m$  indices as signals, since otherwise either the FDR or FNR would fail to vanish. Second, it must be that the  $n^{-\kappa}m$ -th, or equivalently  $n^{1-\beta-\kappa}$ -th, largest null is of the order of the quantity in (3.29).

### 3.3 Upper bounds for some specific procedures

Combining these insights, we obtain the relation

$$(\gamma(\beta + \kappa)) = (\gamma r \log n)^{1/\gamma} - (\gamma \kappa' \log n)^{1/\gamma},$$

which after rearranging yields the heuristic

$$\kappa' = \left( r^{1/\gamma} - (\beta + \kappa)^{1/\gamma} \right)^{1/\gamma}. \quad (3.30)$$

The theorems and corollaries in this paper together show that this intuition is exactly right.

**Regime of linear sparsity:** We turn to the regime of *linear sparsity*—that is, when the number of signals scales as  $\pi_1 n$  for some scalar  $\pi_1 \in (0, 1)$ . Recalling that we have parameterized the number of signals as  $n^{1-\beta_n}$ , some algebra leads to  $\beta_n = \frac{\log \frac{1}{\pi_1}}{\log n}$ , so both Theorem 1 and Theorem 2 predict an upper and lower bound on the risk of the form

$$c_0 \exp \left( c_1 \left[ \frac{\log n}{\log \frac{1}{\pi_1}} \right]^{\frac{\gamma-1}{\gamma}} \right) \cdot n^{-\kappa_*}. \quad (3.31)$$

Note that here we overload the exponent  $\kappa_*$  to the case when it is nonconstant. In order to interpret this result, observe that if  $r_n = r$  is constant, then  $\kappa_* = \frac{r}{2^\gamma} - o(1)$ , so the rate is  $n^{-r/2^\gamma}$  up to subpolynomial factors in  $n$ . On the other hand, if  $r_n = \frac{1}{\log^{\frac{\gamma-1/2}{\gamma}} n}$  is at the extreme lower limit permitted by the lower bound (ii) in (3.14b), then it is not hard to see

---

that  $\kappa_* \approx \log^{-\frac{\gamma-1/2}{\gamma}} n$ , which ensures that  $n^{\kappa_*} \gg \exp\left(\log^{\frac{\gamma-1}{\gamma}} n\right)$ , so that the risk (3.31) still approaches zero asymptotically, albeit subpolynomially in  $n$ .

## 4. Proofs

We now turn to the proofs of our main results, namely Theorems 1 and 2.

The proofs of the associated corollaries can be found in the appendix.

### 4.1 Proof of Theorem 1

The main idea of the proof is to reduce the problem of lower bounding the  $\text{FNR}_n$  of threshold-based procedures that use random, data-dependent thresholds  $T_n$ , to the easier problem of lower bounding the  $\text{FNR}_n$  of threshold-based procedures that use a deterministic, data-independent threshold  $t_n$ .

We refer to the latter class of procedures as *fixed threshold procedures*, and we parameterize them by their target FDR  $q_n = n^{-\kappa_n}$ . Concretely, we define the *critical threshold*, derived from the critical regime boundary  $r_{\min}$  from equation (3.15), by

$$\tau_{\min}(\kappa_n) := (\gamma r_{\min}(\kappa_n) \log n)^{1/\gamma} \quad \equiv \quad \tau_{\min}(q_n) := \left( \gamma r_{\min} \left( \frac{\log(1/q_n)}{\log n} \right) \log n \right)^{1/\gamma}. \quad (4.32)$$

#### 4.1 Proof of Theorem 1

Here and throughout the proof, we express  $\tau_{\min}$  and  $r_{\min}$  as functions of  $q_n$  rather than  $\kappa_n$ ; this formulation turns out to make certain calculations in the proof simpler to express.

**From data-dependent threshold to fixed threshold:** Our first step is to reduce the analysis from data-dependent to fixed threshold procedures. In particular, consider a threshold procedure, using a possibly random threshold  $T_n$ , that satisfies the FDR upper bound  $\text{FDR}_n(T_n) \leq q_n$ . We claim that the FNR of any such procedure must be lower bounded as

$$\mathbb{E}[\text{FNP}_n(T_n)] \geq \frac{\text{FNR}_n(\tau_{\min}(4q_n))}{16}. \quad (4.33)$$

This lower bound is crucial, as it reduces the study of random threshold procedures (LHS) to study of fixed threshold procedures (RHS). Its proof can be found in Appendix S1.

Our next step is to lower bound the FNR for choices of the threshold  $t \geq \tau_{\min}(q_n)$ :

**Lemma 1.** For any  $t \geq \tau_{\min}(q_n)$ , we have

$$\text{FNR}_n(t) \geq \begin{cases} \frac{\zeta^{\frac{1-\gamma}{2\beta_n\gamma}}}{Z_\ell} \cdot n^{-D_\gamma(\beta_n + \kappa_n, r)} & \text{if } r > r_{\min}(\kappa_n(q_n)), \\ \frac{1}{2} & \text{otherwise,} \end{cases} \quad (4.34)$$

where  $\zeta$  was previously defined (3.24).



#### 4.1 Proof of Theorem 1

The proof of this lemma can be found in Appendix S2. Armed with Lemma 1 and the lower bound (4.33), we can now complete the proof of Theorem 1. We split the argument into two cases:

**Case 1:** First, suppose that  $r \leq r_{\min}(\kappa_n(4q_n))$ . In this case, we have

$$\text{FNR}_n(T_n) \stackrel{(i)}{\geq} \frac{\text{FNR}_n(\tau_{\min}(4q_n))}{16} \stackrel{(ii)}{\geq} \frac{1}{32},$$

where step (i) follows from the lower bound (4.33), and step (ii) follows by lower bounding the FNR by  $1/2$ , as is guaranteed by Lemma 1 in the regime  $r \leq r_{\min}(\kappa_n(4q_n))$ .

**Case 2:** Otherwise, we may assume that  $r > r_{\min}(4q_n)$ . In this case, we have

$$\text{FNR}_n(T_n) \stackrel{(i)}{\geq} \frac{\text{FNR}_n(\tau_{\min}(4q_n))}{16} \stackrel{(ii)}{\geq} \frac{(4\zeta)^{2\beta_n^{\frac{1-\gamma}{\gamma}}}}{Z_\ell} \cdot n^{-D_\gamma(\beta_n+\kappa_n,r)}.$$

Here step (i) follows from the lower bound (4.33), whereas step (ii) follows from applying Lemma 1 in the regime  $r > r_{\min}(\kappa_n(4q_n))$ . With some further algebra, we find that

$$\text{FNR}_n(T_n) \geq \frac{1}{Z_\ell} \exp(2 \log(4\zeta) \cdot \beta_n^{\frac{1-\gamma}{\gamma}}) n^{-D_\gamma(\beta+\kappa_n,r)} = c_0 \exp(c_1 \beta_n^{\frac{1-\gamma}{\gamma}}) n^{-D_\gamma(\beta+\kappa_n,r)},$$

where  $c_0 := \frac{1}{Z_\ell}$  and  $c_1 := 2 \log(4\zeta)$ . Note that since  $Z_\ell > 0$  and  $\zeta \geq 1$ , both of the constants  $c_0$  and  $c_1$  are positive, as claimed in the theorem statement.

## 4.2 Proof of Theorem 2

We now sketch the proof that the Benjamini-Hochberg (BH) and Barber-Candès (BC) algorithms achieve the minimax rate (3.20) when  $r_n > r_{\min}(\kappa_n(c_A q_n))$ , where  $A \in \{\text{BH}, \text{BC}\}$  and  $c_A$  is the algorithm-dependent constant defined in (3.24). For reasons of space, the details are relegated to Appendix S4.

The proof strategy for both algorithms is essentially the same. Given a target FDR rate  $q_n$ , we apply each algorithm  $q_n$  as the target FDR level and prove that the resulting threshold satisfies  $t_A \leq \tau_{\min}(c_A q_n)$  with high probability. Letting  $\tau_{\min,A} = \tau_{\min}(c_A q_n)$ , we can formulate the specific claims we seek as:

$$\mathbb{P}(t_{\text{BH}} > \tau_{\min,\text{BH}}) \leq \exp\left(-\frac{n^{1-r_{\max}}}{24}\right) \quad (4.35)$$

and

$$\mathbb{P}(t_{\text{BC}} > \tau_{\min,\text{BC}}) \leq q_n + \exp\left(-\frac{n^{1-r_{\max}}}{24}\right). \quad (4.36)$$

The known properties of the algorithms guarantee the required FDR bounds (as studied by Arias-Castro and Chen, 2017; Foygel Barber and Candès, 2015; Benjamini and Hochberg, 1995), while the following converse to Lemma 1, coupled with the probabilistic upper bounds (4.35) and (4.36), provides the requisite upper bounds on the FNR.

---

**Lemma 2.** If  $r_n > r_{\min}(cq_n)$  and  $t \leq \tau_{\min}(cq_n)$  for some  $c > 0$ , then we have

$$\text{FNR}_n(t) \leq \frac{(\max\{c, 1/c\} \cdot \zeta)^{2\beta_n^{\frac{1-\gamma}{\gamma}}}}{Z_u} \cdot n^{-D_\gamma(\beta_n + \kappa_n, r)},$$

where constant  $\zeta$  is defined in (3.24).

## 5. Discussion

Despite considerable interest in multiple testing with false discovery rate (FDR) control, there has been relatively little understanding of the non-asymptotic trade-off between controlling FDR and the analogous measure of power known as the false non-discovery rate (FNR). In this paper, we explored this issue in the context of the sparse generalized Gaussians model, and derived the first non-asymptotic lower bounds on the sum of FDR and FNR. We complemented these lower bounds by establishing the non-asymptotic minimaxity of both the Benjamini-Hochberg (BH) and Barber-Candès (BC) procedures for FDR control. The theoretical predictions are validated in simple simulations, and our results recover recent asymptotic results (Arias-Castro and Chen, 2017) as special cases. Our work introduces a simple proof strategy based on reduction to deterministic and data-oblivious procedures. We suspect this core idea may apply to other multiple testing settings: in particular, since our arguments do not depend on CDF

---

asymptotics in the way that many classical analyses of both global null testing and FDR control procedures do, we hope they will be possible to adapt for other problems described below.

As mentioned after the statement of Theorem 2, the practical implications of our results are somewhat pessimistic. Even for rather simple problems having  $r - \beta$  of constant order, the resulting rate at which the risk tends to zero can be far slower than  $n^{-1/2}$ . (Indeed, it seems like such a parametric rate is only achievable when  $\gamma = 1, r_n \rightarrow 1, \beta_n \rightarrow 0$ .) Hence, in practice, one must carefully consider whether good FDR or good FNR is more important, as achieving both may not be possible unless most of the signals to be identified are rather large.

### **Future directions**

We have focused on establishing a non-asymptotic tradeoff between FDR and FNR in what is arguably the simplest interesting model of the problem. By way of contrast, a large part of the multiple testing literature in recent years has focused on the development of valid FDR control procedures that can gain power or precision by explicitly using prior knowledge and structure in various ways: whether through null-proportion adaptivity (Storey, 2002; Storey et al., 2004), grouping of hypotheses (Foygel Bar-

---

ber and Ramdas, 2016; Hu et al., 2010), prior or penalty weights (Benjamini and Hochberg, 1997; Genovese et al., 2006), or other forms of structure (Li and Foygel Barber, 2016; Ramdas et al., 2017).

Similarly, the issue of dependence—either positive or arbitrary—between test statistics has been an area of focus (Benjamini and Yekutieli, 2001; Blanchard and Roquain, 2008; Ramdas et al., 2017). (Dependence has already been explored for the higher criticism statistic applied to the detection problem (Hall and Jin, 2008; Jin and Ke, 2014; Hall and Jin, 2010).) Non-exchangeability of hypotheses, either in the context of multiple scales of signal strength, or in the context of online FDR procedures, has also been studied (Foster and Stine, 2008; Javanmard and Montanari, 2015).

Due to the increasing importance of the structured, dependent, and non-exchangeable settings, developing analogues of our results for those settings is a worthwhile direction for future work. It is, furthermore, far from clear that known procedures are optimal under assumptions of structure, dependence, or various kinds of non-exchangeability, so that an improved understanding of the fundamental difficulty of the multiple testing problem under such assumptions may also yield improved algorithms. Chen and Arias-Castro (2017) have made progress in this direction by providing *upper* bounds for existing procedures for the online FDR problem (Javanmard

---

and Montanari, 2015), but much still remains unknown.

Finally, a general proof technique for establishing non-asymptotic lower bounds in multiple testing remains an important direction for future work. In this work, we pursued an approach based on reduction to a class of non-adaptive procedures, and this principle could perhaps be applied to other multiple testing problems. Our arguments are, however, based on analytical calculations, and they are therefore sensitive to the specific observation model under consideration. One especially pressing problem is thus to develop approaches that depend on more intrinsic structural properties of the test statistic distributions and that are less brittle when it becomes inconvenient to reason about the analytical forms.

## Supplementary Materials

The supplementary materials contain proofs that—for space reasons—we could not accommodate in the main body. These include parts of proofs of theorems, as well as proofs of corollaries and technical lemmas.

## Acknowledgements

This work was partially supported by Office of Naval Research Grant DOD-ONR-N00014, Air Force Office of Scientific Research Grant AFOSR-FA9550-14-1-0016, and Army Research Office grant W911NF-16-1-0368. In

addition, MR was supported by an NSF Graduate Research Fellowship and a Fannie and John Hertz Foundation Google Fellowship.

## References

- F. Abramovich, Y. Benjamini, D. Donoho, and I.M. Johnstone. Adapting to unknown sparsity by controlling the false discovery rate. *Annals of Statistics*, 34(2):584–653, 2006.
- E. Arias-Castro and S. Chen. Distribution-free multiple testing. *Electronic Journal of Statistics*, 11(1):1983–2001, 2017.
- Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 57(1):289–300, 1995.
- Y. Benjamini and Y. Hochberg. Multiple hypotheses testing with weights. *Scandinavian Journal of Statistics*, 24(3):407–418, 1997.
- Y. Benjamini and D. Yekutieli. The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, 29(4):1165–1188, 2001.
- G. Blanchard and E. Roquain. Two simple sufficient conditions for FDR control. *Electronic Journal of Statistics*, 2:963–992, 2008.
- M. Bogdan, A. Chakrabarti, F. Frommlet, and J.K. Ghosh. Asymptotic Bayes-optimality under sparsity of some multiple testing procedures. *Annals of Statistics*, 39(3):1551–1579, 2011.
- S. Chen and E. Arias-Castro. Sequential multiple testing. *arXiv preprint arXiv:1705.10190*,

## REFERENCES

---

2017.

Z. Chi. On the performance of FDR control: Constraints and a partial solution. *Annals of Statistics*, 35(4):1409–1431, 2007.

D. Donoho and J. Jin. Higher criticism for detecting sparse heterogeneous mixtures. *Annals of Statistics*, 32(3):962–994, 2004.

D. Donoho and J. Jin. Asymptotic minimaxity of false discovery rate thresholding for sparse exponential data. *Annals of Statistics*, 34(6):2980–3018, 2006.

D. Donoho and J. Jin. Higher criticism for large-scale inference, especially for rare and weak effects. *Statistical Science*, 30(1):1–25, 2015.

D.P. Foster and R.A. Stine.  $\alpha$ -investing: A procedure for sequential control of expected false discoveries. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(2):429–444, 2008.

R. Foygel Barber and E.J. Candès. Controlling the false discovery rate via knockoffs. *Annals of Statistics*, 43(5):2055–2085, 2015.

R. Foygel Barber and A. Ramdas. The  $p$ -filter: Multi-layer FDR control for grouped hypotheses. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2016.

C.R. Genovese and L. Wasserman. Operating characteristics and extensions of the false discovery rate procedure. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):499–517, 2002.



## REFERENCES

---

- C.R. Genovese, K. Roeder, and L. Wasserman. False discovery control with  $p$ -value weighting. *Biometrika*, 93(3):509–524, 2006.
- P. Hall and J. Jin. Properties of higher criticism under strong dependence. *Annals of Statistics*, 36(1):381–402, 2008.
- P. Hall and J. Jin. Innovated higher criticism for detecting sparse signals in correlated noise. *Annals of Statistics*, 38(3):1686–1732, 2010.
- J. Hu, H. Zhao, and H. Zhou. False discovery rate control with groups. *Journal of the American Statistical Association*, 105(491):1215–1227, 2010.
- A. Javanmard and A. Montanari. On online control of false discovery rate. *arXiv preprint arXiv:1502.06197*, 2015.
- P. Ji and J. Jin. UPS delivers optimal phase diagram in high-dimensional variable selection. *Annals of Statistics*, 40(1):73–103, 2012.
- P. Ji and Z. Zhao. Rate optimal multiple testing procedure in high-dimensional regression. *arXiv preprint arXiv:1404.2961*, 2014.
- J. Jin and T. Ke. Rare and weak effects in large-scale inference: Methods and phase diagrams. *arXiv preprint arXiv:1410.4578*, 2014.
- A. Li and R. Foygel Barber. Multiple testing with the structure adaptive Benjamini-Hochberg algorithm. *arXiv preprint arXiv:1606.07926*, 2016.
- Pierre Neuvial and Etienne Roquain. On false discovery rate thresholding for classification

## REFERENCES

---

under sparsity. *Annals of Statistics*, 40(5):2572–2600, 2012.

A. Ramdas, R. Foygel Barber, M.J. Wainwright, and M.I. Jordan. A unified treatment of multiple testing with prior knowledge. *arXiv preprint arXiv:1703.06222*, 2017.

J. Storey. A direct approach to false discovery rates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):479–498, 2002.

J. Storey, J. Taylor, and D. Siegmund. Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: A unified approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(1):187–205, 2004.

J.W. Tukey. *The Problem of Multiple Comparisons: Introduction and Parts A, B, and C*. Princeton University, 1953.

Department of EECS, UC Berkeley

E-mail: [rabinovich@berkeley.edu](mailto:rabinovich@berkeley.edu)

Department of Statistics and EECS, UC Berkeley

E-mail: [ramdas@berkeley.edu](mailto:ramdas@berkeley.edu)

Department of Statistics and EECS, UC Berkeley

E-mail: [jordan@berkeley.edu](mailto:jordan@berkeley.edu)

Department of Statistics and EECS, UC Berkeley

E-mail: [wainwrig@berkeley.edu](mailto:wainwrig@berkeley.edu)