

Zheng Gao  
Stilian Stoev

# Concentration of Maxima and Fundamental Limits in High-Dimensional Testing and Inference

January 15, 2021

Springer



Stilian: I suggest to:

- Create a chapter called Supplement, which will become an “online supplement”.
- Begin moving some lengthy background material and the exercises to the Supplement.
- Respond to R2’s question about the role of dependence on the rate of growth of maxima with a sub-section and/or examples in the Supplement. This will be an illustration of the “Folklore Theorem” *Dependence can only slow down the rate of growth of the maxima of identically distributed random variables.*
- Add index. **Should we have an Index?**
- Add more examples and references to the planted matrix problems. **If we have time, it would be great to broaden the scope of applications.**



# Contents

<b>1</b>	<b>Introduction and Guiding Examples</b> . . . . .	1
1.1	The additive error model . . . . .	2
1.2	Genome-wide association studies and the chi-square model . . . . .	3
1.3	Contents . . . . .	5
<b>2</b>	<b>Risks, Procedures, and Error Models</b> . . . . .	7
2.1	Statistical risks . . . . .	7
2.2	Statistical procedures . . . . .	10
2.3	Related literature and our contributions . . . . .	13
2.4	Relationships between the asymptotic risks . . . . .	16
2.5	The asymptotic generalized Gaussian (AGG) models . . . . .	18
2.6	Rapid variation and relative stability . . . . .	19
<b>3</b>	<b>The Phase Transition Phenomena in Independent Gaussian Error Models</b> . . . . .	23
3.1	Sparse signal detection problems . . . . .	24
3.2	Sparse signal support recovery problems . . . . .	27
3.3	Asymptotic power analysis: A discussion . . . . .	31
<b>4</b>	<b>Exact Support Recovery Under Dependence</b> . . . . .	33
4.1	Contributions and related work . . . . .	33
4.2	Sufficient conditions for exact support recovery . . . . .	35
4.3	Dependence and uniform relative stability . . . . .	39
4.4	Necessary conditions for exact support recovery . . . . .	39
4.5	Dense signals . . . . .	42
4.6	Numerical illustrations for independent errors . . . . .	43
4.7	Additional proofs . . . . .	44
<b>5</b>	<b>Bayes and minimax optimality</b> . . . . .	49
5.1	Bayes optimality in support recovery problems . . . . .	49
5.2	Minimax optimality in exact support recovery . . . . .	57

5.3	A discussion on optimality . . . . .	60
<b>6</b>	<b>Uniform Relative Stability for Gaussian Arrays</b> . . . . .	<b>61</b>
6.1	Ramsey's Coloring theorem and structure of correlation matrices . .	64
6.2	URS implies UDD ('only if' part of Theorem 6.1) . . . . .	66
6.3	UDD implies URS ('if' part of Theorem 6.1) . . . . .	68
6.4	Numerical illustrations of exact support recovery under dependence	75
<b>7</b>	<b>Fundamental Statistical Limits in Genome-wide Association Studies</b>	<b>81</b>
7.1	Support recovery problems in chi-squared models . . . . .	82
7.2	Odds ratios and statistical power . . . . .	86
7.3	Optimal study designs and rare variants . . . . .	88
7.4	Phase transitions in large-scale association screening studies . . . . .	89
7.5	Numerical illustrations of the phase transitions in chi-square models	92
7.6	Proofs . . . . .	95
	References . . . . .	113
	<b>Proofs for Chapter 3</b> . . . . .	<b>119</b>
A.1	Auxiliary facts of Gaussian distributions . . . . .	119
A.2	Monotonicity of the Benjamini-Hochberg procedure . . . . .	120
A.3	Proof of Theorem 3.1 . . . . .	121
A.4	Proof of Theorem 3.3 . . . . .	124
A.5	Proof of Theorems 3.4 and 3.5 . . . . .	128
	<b>Exact support recovery in non AGG models</b> . . . . .	<b>131</b>
B.1	Strong classification boundaries in other light-tailed error models . .	131
B.2	Thresholding procedures under heavy-tailed errors . . . . .	134

## Acronyms

AGG	asymptotically generalized Gaussian
BH	Benjamini-Hochberg
CDF	cumulative distribution function
FDR	false discovery rate
FNR	false non-discovery rate
FWER	family-wise error rate
FWNR	family-wise non-discovery rate
GWAS	genome-wide association studies
HC	higher-criticism
iid	independent and identically distributed
LR	likelihood-ratio
RS	relatively stable/relative stability
SNP	single-nucleotide polymorphisms
URS	uniform relatively stable/uniform relative stability





## Preface

This text presents a collection of new results and recent developments on the phase transition phenomena in sparse signal problems. The main theme is the study of the fundamental limits in high-dimensional testing and inference. Since the seminal works of Ingster (1998) and Donoho and Jin (2004), the subject has received a lot of attention in the literature with important contributions from Ji and Jin (2012); Genovese et al. (2012); Jin et al. (2014); Arias-Castro and Chen (2017); Butucea et al. (2018). These works, among many others, have discovered some fundamental limits in the so-called *needle in haystack* problems, where a sparse signal is observed with high-dimensional additive noise. In this setting, two archetypal problems arise – the *signal detection* and *signal support recovery*. The signal detection refers to a global hypothesis testing problem that amounts to determining the presence of non-zero signal in any of its dimensions. The support recovery, on the other hand, can be seen either as a multiple testing problem where the presence of non-zero signal is tested for each signal location of interest, or alternatively, as an inference problem that aims to estimate the signal support, i.e., the locations of the non-zero signal components. The fundamental limits of these problems are studied in the so-called high-dimensional asymptotic regime where the dimension  $p$  of the underlying signal grows to infinity, and the sample size  $n$  is either bounded or grows slowly relative to  $p$ .

From a probabilistic perspective, these aforementioned fundamental limits are stated as asymptotic zero-one type laws, as dimensionality diverges. Namely, consider a sparse signal with *support size* on the order of  $p^{1-\beta}$  for some parameter  $\beta \in (0, 1)$ . Parameterize the non-zero *signal amplitude* by  $\Delta(p^r)$ , for some  $r > 0$  and a suitable monotone non-decreasing function  $\Delta(\cdot)$ . Then, for a broad range of error distributions and statistical problems, one encounters a sharp transition between the regimes where the problem is solvable and unsolvable depending on the signal magnitude  $r$  and signal sparsity  $\beta$ . More precisely, there exists a boundary function  $f(\beta)$  such that if the signal magnitudes are *above* the boundary,  $r > f(\beta)$ , then the problem can be solved with vanishing loss as  $p \rightarrow \infty$ , with a suitable statistical procedure. On the other hand, if the signal is below that same boundary, i.e.,  $r < f(\beta)$ , all statistical procedures fail to provide a solution with a vanishing loss, as  $p \rightarrow \infty$ .

Of course, depending on whether one considers the detection (testing) or support recovery (inference) problems, different loss functions quantify success and failure. The choice of the loss functions is often guided by the applications, resulting in a rich picture of phase-transitions (see e.g. Figure 3.2).

**The contributions of this work.** The fundamental limits of the classic detection problem hinge of the analysis of the discrepancy between the *null* and *alternative* hypotheses e.g., via Hellinger distance. Thus, perhaps for technical reasons, much of the analysis in the existing literature has been done under the assumption that the additive errors are independent and/or Gaussian, or using loss functions unaffected by the dependence such as the Hamming loss. In this work, we demonstrate that the support recovery problems, especially *exact support recovery*, are best understood from the novel perspective of the *concentration of maxima* phenomenon in extreme value theory. It turns out that under a very broad range of light-tailed error distributions and under a *very* broad range of error dependence structures, the maxima of the errors, when rescaled (but not centered!) converge in probability to a positive constant. This concentration property leads to a complete solution of the exact support recovery problem for the broad family of thresholding procedures. Most if not all existing support estimation procedures are types thresholding procedures (see Section 2.2). That is, the signal support estimate comprises of all components exceeding a suitable (potentially data-dependent) threshold. We show, by exploiting concentration of maxima, that thresholding procedures obey a phase-transition, where if the signal is above a certain boundary, asymptotically exact recovery is possible while below the boundary all thresholding procedures fail, as  $p \rightarrow \infty$ . Remarkably, light-tailed maxima concentrate under very broad and strong dependence. This is exemplified by our characterization of the concentration of maxima phenomenon for Gaussian triangular arrays. For example, in the special case of stationary Gaussian time series, vanishing auto-covariance is necessary and sufficient for the maxima to concentrate in the same way as independent standard normal random variables. Stilian: This is in stark contrast with the behavior of sums, commonly studied under short- and long-range dependence conditions (see e.g. Dedecker et al., 2007; Pipiras and Taqqu, 2017). Simply put, the notion of weak dependence that entails that the maxima of dependent variables concentrate at the same rate as in the case of independence is fundamentally weaker than the conventional mixing conditions widely used in the study of sums. added

Stilian: Our ~~These~~ probabilistic contributions may be of independent interest and extend classic work of Berman (1964). Concentration of maxima is a type of super-concentration phenomenon studied also in Chatterjee (2014) and Tanguy (2015a). The robustness of the concentration of maxima phenomenon to dependence can perhaps explain the universality of phase transitions in support recovery problems.

The use of concentration of maxima phenomenon highlights one core idea in our work, which allows for a first of its kind comprehensive treatment of thresholding procedures under very broad error-dependence conditions. The text involves also a full spectrum of related results such as minimax-optimality and finite-sample Bayes optimality in support estimation. Using different type of loss functions and type I error controls, we obtain a rich picture of the exact and approximate support

recovery problems in high dimensions. Many of these phase transition results have not appeared in previously published literature.

High-dimensional support recovery problems arise in many modern applications ranging from cybersecurity, theoretical computer science, to statistical genetics. Genome-wide association studies (GWAS) in genetics are particularly natural applications, where the asymptotic phase-transition results help explain and quantify a previously observed empirical phenomenon of the so-called *steep part of the power curve*. In the last chapter of this work, we detail this application and highlight future theoretical and practical consequences of our work.

**Target audience.** The original research presented in this text originates from the doctoral dissertation of the first author in the Statistics department at the University of Michigan, Ann Arbor. The main goal of this text is to provide a comprehensive treatment of the exact and approximate support recovery problems by utilizing existing and newly developed probabilistic tools on concentration of maxima. The text also provides a quick introduction to the state-of-the-art in the dynamic area of phase-transitions in high-dimensional testing and inference. It is accessible to doctoral students in Statistics with background in measure-theoretic probability and statistics as well as to researchers in applied fields working with high-dimensional data sets. The text can be used as a reference and a supplement to a special topics course on high-dimensional inference.

**Acknowledgements.** The authors gratefully acknowledge the support of their families and all colleagues from the Statistics Department at the University of Michigan, Ann Arbor. Special thanks (in alphabetical order) go to Xuming He, Tailen Hsing, Michalis Kallitsis, Liza Levina, Yuezhi Li, Rodderick Little, Ya'acov Ritov, Kerby Shedden, Jinqi Shen, Jonathan Terhorst, Gongjun Xu. The authors were partially supported by the NSF program *Algorithms for Threat Detection*.

Chicago  
Ann Arbor, January 15, 2021

Zheng Gao  
Stilian Stoev



## Chapter 1

# Introduction and Guiding Examples

The proliferation of information technology has enabled us to collect and consume huge volumes of data at unprecedented speeds and at very low costs. This convenient access to data gave rise to a fundamentally different way of pursuing scientific questions. In contrast with the traditional hypothesis–experiment–analysis cycle where data are collected from the experiments, nowadays abundant data are often available before specific questions are even formulated. Such data can be used for not just evaluating hypotheses, but also for *generating*, and *selecting* the hypotheses to pursue. As a result, multiple testing — where a large number of hypotheses are formulated and screened for their plausibility simultaneously — has become a staple of modern data-driven studies.

An archetypal example of multiple testing problems is genetic association studies (Bush and Moore, 2012). In these studies, scientists test hypotheses relating each of the hundreds of thousands of genetic marker locations to phenotypic traits of interest. For a phenotypic trait on which we have little prior knowledge, we cannot simply test for association on one or a few specific genetic locations, as there are often not enough empirical evidence or biological theory to pin point these genetic locations in the first place. Rather, the goal here is to select the set of most promising genetic markers from a large number of candidate locations for subsequent investigation.

Another example of multiple testing problems arise in cybersecurity, where millions of IP addresses are monitored in real time. In this engineering application, statistics are collected and tests are performed for each IP address, in an attempt to locate the IP addresses with anomalous network activities, so that and malicious traffic and volumetric attacks can be filtered to protect end users of network services (Kallitsis et al., 2016). Similar to the genetic application above, we use data to search over candidate IP addresses and identify locations of interest.

We are motivated very much by these examples to study high-dimensional multiple testing problems where a large number of hypotheses are tested simultaneously. In the rest of the introduction, we shall more review the main objectives of high-dimensional multiple testing, and elaborate on these objectives with two classes of data models in the context of various applications.

## 1.1 The additive error model

Consider the canonical signal-plus-noise model where the observation  $x$  is a high-dimensional vector in  $\mathbb{R}^p$ ,

$$x(i) = \mu(i) + \epsilon(i), \quad i = 1, \dots, p. \quad (1.1)$$

The signal,  $\mu = (\mu(i))_{i=1}^p$ , is a vector with  $s$  non-zero components supported on the set  $S = \{i : \mu(i) \neq 0\}$ ; the second term  $\epsilon$  is a random error vector. The goal of high-dimensional statistics is usually two-fold:

- I. *Signal detection*: to detect the presence of non-zero components in  $\mu$ . That is, to test the global hypothesis  $\mu = 0$ .
- II. *Support recovery*: to estimate the support set  $S$ . This is also sometimes referred to as the *support estimation* or *signal identification* problem.

To illustrate, in the engineering application of cybersecurity, Internet service providers (ISP) routinely monitor a large number of network traffic streams to determine if there are abnormal surges, blackouts, or other types of anomalies. The data vector  $x$  could represent, for example, incoming traffic volumes to each server node, internet protocol (IP) address, or port that the ISP monitors. In this case, the vector  $\mu$  represents the average traffic volumes in each of the streams under normal operating conditions, and  $\epsilon$ 's – the fluctuations around these normal levels of traffic. The signal detection problem in this context is then equivalent to determining if there are *any* anomalies among all data streams, and the support recovery problem is equivalent to *identifying* the streams experiencing anomalies. Similar questions of signal detection and support recovery are pursued in large-scale microarray experiments (Dudoit et al., 2003), brain imaging and fMRI analysis (Nichols and Hayasaka, 2003), and numerous other anomaly detection applications.

A common theme in such applications is that the errors are *correlated*, and that the signal vectors are believed to be *sparse*: the number of non-zero (or large) components in  $\mu$  is small compared to the number of tests performed. In the cybersecurity context, while a very large number of data streams are monitored, typically only just a few of them will be experiencing problems at any time, barring large-scale outages or distributed denial of service attacks. Under such sparsity assumptions, it is natural to ask if and when one can reliably (1) detect the signals, and (2) recover the support set  $S$ . In this text, we explore both the *detection* and the *support recovery* problems. More precisely, we are interested in the theoretical feasibility of both problems, and seek minimal conditions under which these problems can be consistently solved in large dimensions.

Model (1.1) is simple yet ubiquitous. Consider the linear model

$$Y = X\mu + \xi,$$

where  $\mu$  is a  $p$ -dimensional vector of regression coefficients of interest to be inferred from observations of  $X$  and  $Y$ . If the design matrix  $X$  is of full column rank<sup>1</sup>, then the ordinary least squares (OLS) estimator of  $\mu$  can be formed

$$\hat{\mu} = (X'X)^{-1} X'Y = \mu + \epsilon, \quad (1.2)$$

where  $\epsilon := (X'X)^{-1} X'\xi$ . Hence we recover the generic problem (1.1). Signal detection is therefore equivalent to the problem of testing the global null model, and support recovery problem corresponds to the fundamental problem of variable selection.

Note that the components of the observation vector  $x$  (and equivalently, the noise  $\epsilon$ ) in (1.1) need not be independent. In the linear regression example, even when the components of the noise term  $\xi$  are independent, those of the OLS estimator (1.2) need not be, except in the case of orthogonal designs. Indeed, in practice, independence is the exception rather than the rule. Therefore, a general theory of feasibility must address the role of the *error dependence* structure in such testing and support estimation problems. It is also important to identify practical and/or optimal procedures that attain the performance limits in independent as well as dependent cases, as soon as the problems become theoretically feasible. We address both themes in this text.

## 1.2 Genome-wide association studies and the chi-square model

The second data model we analyze is the high-dimensional chi-square model,

$$x(i) \sim \chi_\nu^2(\lambda(i)), \quad i = 1, \dots, p, \quad (1.3)$$

where the data  $x(i)$ 's follow independent (non-central) chi-square distributions with  $\nu$  degrees of freedom and non-centrality parameter  $\lambda(i)$ .

Model (1.3) is motivated by large-scale categorical variable screening problems, typified by GWAS where millions of genetic factors are examined for their potential influence on phenotypic traits.

In a GWAS with a case-control design, a total of  $n$  subjects are recruited, consisting of  $n_1$  subjects possessing some defined traits, and  $n_2$  subjects without the traits serving as controls. The genetic compositions of the subjects are then examined for variations known as SNP at an array of  $p$  genomic marker locations, and compared between the case and the control group. These physical traits are commonly referred to as *phenotypes*, and the genetic variations are known as *genotypes*.

Focusing on one specific genomic location, the counts of observed genotypes, if two variants are present, can be tabulated as follows.

---

<sup>1</sup> This, of course, requires that we have more samples than dimensions, i.e.,  $n > p$ . Nevertheless, multiplicity of tests is still present when  $p$  itself is large – the multiple testing problem is by no means exclusive to situations where  $p \gg n$ .

# Observations	Genotype		Total by phenotype
	Variant 1	Variant 2	
Cases	$O_{11}$	$O_{12}$	$n_1$
Controls	$O_{21}$	$O_{22}$	$n_2$

Researchers test for associations between the genotypes and phenotypes using, for example, the Pearson chi-square test with statistic

$$x = \sum_{j=1}^2 \sum_{k=1}^2 \frac{(O_{jk} - E_{jk})^2}{E_{jk}}, \quad (1.4)$$

where  $E_{jk} = (O_{j1} + O_{j2})(O_{1k} + O_{2k})/n$ .

Under the mild assumption that the counts  $O_{jk}$ 's follow a multinomial distribution (or a product-binomial distribution, if we decide to condition on one of the marginals), the statistic  $x$  in (1.4) can be shown to have an approximate  $\chi^2(\lambda)$  distribution with  $\nu = 1$  degree of freedom at large sample sizes (see, e.g., classical results in Ferguson (2017) and Agresti (2018)). Independence between the genotypes and phenotypes would imply a non-centrality parameter  $\lambda$  value of zero; if dependence exists, we would have a non-zero  $\lambda$  where its value depends on the underlying multinomial probabilities. More generally, if we have a  $J$  phenotypes and  $K$  genetic variants, assuming a  $J \times K$  multinomial distribution, the statistic will follow approximately a  $\chi^2_\nu(\lambda)$  distribution with  $\nu = (J - 1)(K - 1)$  degrees of freedom, when sample sizes are large.

The same asymptotic distributional approximations also apply to the likelihood ratio statistic, and many other statistics under slightly different modeling assumptions (Gao et al., 2019). These association tests are performed at each of the  $p$  SNP marker locations throughout the whole genome, and we arrive at  $p$  statistics having approximately (non-central) chi-square distributions,  $\chi^2_{\nu(i)}(\lambda(i))$ , for  $i = 1, \dots, p$ , where  $\lambda = (\lambda(i))_{i=1}^p$  is the  $p$ -dimensional non-centrality parameter.

Although the number of tested genomic locations  $p$  can sometimes exceed  $10^5$  or even  $10^6$ , it is often believed that only a small set of genetic locations have tangible influences on the outcome of the disease or the trait of interest. Under the stylized assumption of sparsity,  $\lambda$  is assumed to have  $s$  non-zero components, with  $s$  being much smaller than the problem dimension  $p$ . The goal of researchers is again two-fold: (1) to test if  $\lambda(i) = 0$  for all  $i$ , and (2) to estimate the set  $S = \{i : \lambda(i) \neq 0\}$ . In other words, we look to first determine if there are *any* genetic variations associated with the disease; and if there are associations, we want to locate them.

The chi-square model (1.3) also plays an important role in analyzing variable screening problems under omnidirectional alternatives. A primary example is multiple testing under two-sided alternatives in the additive error model (1.1) where the errors  $\epsilon$  are assumed to have standard normal distributions.

Under two-sided alternatives, unbiased test procedures call for rejecting the hypothesis  $\mu(i) = 0$  at locations where observations have large absolute values, or equivalently, large squared values. Taking squares on both sides of (1.1), and we arrive at Model (1.3) with non-centrality parameters  $\lambda(i) = \mu^2(i)$  and degree-of-



freedom parameter  $\nu = 1$ . In this case, the support recovery problem is equivalent to locating the set of observations with mean shifts,  $S = \{i : \mu(i) \neq 0\}$ , where the mean shifts could take place in both directions.

Therefore, a theory for the chi-square model (1.3) naturally lends itself to the study of two-sided alternatives in the Gaussian additive error model (1.1). In comparing such results with existing theory on one-sided alternatives, we will be able to quantify if, and how much of a price has to be paid for the additional uncertainty when we have no prior knowledge on the direction of the signals.

### 1.3 Contents

Important notions and definitions in high-dimensional testing problems are recalled in Chapter 2. We review related literature as well as key concepts and technical results used in our subsequent analyses.

In Chapter 3 we study the sparse signal detection and support recovery problems for the additive error model (1.1) when components of the noise term  $\epsilon$  are independent standard Gaussian random variables. In particular, we point out several new *phase transitions* in signal detection problems, and provide a unified account of recently discovered phase transitions in support recovery problems. These results show that as the dimension  $p \rightarrow \infty$ , the tasks of detecting the existence of signals, or identifying the support set  $S$  are either doable or impossible depending on the sparsity and signal sizes of the problems. We also identify commonly used procedures that attain the performance limits in both detection and support recovery problems.

Both the Gaussianity assumption and the independence assumption are relaxed in Chapter 4. Established are the necessary and sufficient conditions for exact support recovery in the high-dimensional asymptotic regime. This is a major theoretical contribution of our approach, which solves and expands on open problems in the recent literature (see Butucea et al. (2018); Gao and Stoev (2020)). The analysis of support recovery problem is intimately related to a *concentration of maxima* phenomena in the analysis of extremes. The latter concept is key to understanding the role played by dependence in the phase transition phenomena of high-dimensional testing problems. Using this probabilistic concept, we establish minimax optimality results that hold for a very large class of dependence structures.

The dependence condition defined by the concentration of maxima concepts is further demystified in Chapter 6 for Gaussian errors. We offer a complete characterization of the concentration of maxima phenomenon, known as uniform relative stability, in terms of the covariance structures of the Gaussian arrays. This result may be of independent interest since it relates to the so-called *superconcentration* phenomenon coined by Chatterjee (2014). See also, Gao and Stoev (2020) and Kartsioukas et al. (2019).

Chapter 7 returns to high-dimensional multiple testing problems, and study the chi-square model (1.3) inspired by the marginal association screening problems. We establish four new phase-transition-type results in the chi-square model, and

illustrate their practical implications in the GWAS application. Our theory enables us to explain the long-standing empirical observation that small perturbations in the frequency and penetrance of genetic variations lead to drastic changes in the discoverability in genetic association studies.

## Chapter 2

# Risks, Procedures, and Error Models

We establish the background necessary for the study of sparse signal detection and support recovery problems in this chapter. Sections 2.1 and 2.2 provide a refresher on the definitions of statistical risks and some commonly used statistical procedures. Section 2.3 describes the asymptotic regime under which we analyze these procedures, and reviews the related literature in high-dimensional statistics. We discuss in Section 2.4 the connections among the risk metrics, and point out some common fallacies. The remaining sections collect the technical preparations for this text. Section 2.5 defines an important class of error distributions which will be analyzed in detail in later chapters. And finally, Section 2.6 introduces the concepts of concentration of maxima, which plays a crucial role in the analysis of high-dimensional support recovery problems.

### 2.1 Statistical risks

We define the statistical risk metrics for signal detection and signal support recovery problems in this section. Formally, we denote a statistical procedure, i.e., measurable function of the data, as  $\mathcal{R} = \mathcal{R}(x)$ . In the testing context, a procedure  $\mathcal{R}$  produces a binary decision  $T$  that represents our judgment on the presence or absence of a signal. In the support recovery problem, a procedure  $\mathcal{R}$  produces an index set  $\hat{S}$  that represents our estimate of the signal support. The statistical risks are then suitable functionals of  $T$  and  $\hat{S}$  in respective contexts.

**Signal detection.** Recall that in sparse signal detection problems, our goal is to come up with a procedure,  $\mathcal{R}(x)$ , such that the null hypothesis is rejected if the data  $x$  is deemed incompatible with the null. In the additive error models context (1.1), we wish to tell apart two hypotheses

$$\mathcal{H}_0 : \mu(i) = 0, i = 1, \dots, p, \quad \text{v.s.} \quad \mathcal{H}_1 : \mu(i) \neq 0, \text{ for some } i \in \{1, \dots, p\}, \quad (2.1)$$

based on the  $p$ -dimensional observation  $x$ . Similarly, in the chi-square model (1.3), we look to test

$$\mathcal{H}_0 : \lambda(i) = 0, i = 1, \dots, p, \quad \text{v.s.} \quad \mathcal{H}_1 : \lambda(i) \neq 0, \text{ for some } i \in \{1, \dots, p\}. \quad (2.2)$$

Since the decision is binary, we may write the outcome of the procedure in the form of an indicator function,  $T(\mathcal{R}(x)) \in \{0, 1\}$ , where  $T = 1$  if the null is to be rejected in favor of the alternative, and 0 if we fail to reject the null. The Type I and Type II errors of the procedure, i.e., the probability of wrong decisions under the null hypothesis  $\mathcal{H}_0$  and alternative hypothesis  $\mathcal{H}_1$ , respectively, are defined as

$$\alpha(\mathcal{R}) := \mathbb{P}_{\mathcal{H}_0}(T(\mathcal{R}(x)) = 1) \quad \text{and} \quad \beta(\mathcal{R}) := \mathbb{P}_{\mathcal{H}_1}(T(\mathcal{R}(x)) = 0). \quad (2.3)$$

The Neyman-Pearson framework of hypothesis testing then seeks tests that minimize the Type II error of the test, while controlling the Type I error of the test at low levels. We are particularly interested in the sum of the two errors,

$$\text{risk}^D(\mathcal{R}) := \alpha(\mathcal{R}) + \beta(\mathcal{R}), \quad (2.4)$$

which shall be referred to as the risk of signal detection (of the procedure  $\mathcal{R}$ ). It is trivial that a small  $\text{risk}^D$  would imply both small Type I and Type II errors of the procedure.

**Signal support recovery.** Turning to support recovery problems, our goal is to design a procedure that produces a set estimate  $\widehat{S}(\mathcal{R}(x))$  of the true index set of relevant variables  $S$ . For example, in the sparse additive error model (1.1) we aim to estimate  $S = \{i : \mu(i) \neq 0\}$ , while in the sparse chi-square model (1.3) the goal is to estimate  $S = \{i : \lambda(i) \neq 0\}$ . For simplicity of notation, we shall write  $\widehat{S}$  for the support estimator  $\widehat{S}(\mathcal{R}(x))$ .

For a given procedure  $\mathcal{R}$ , its false discovery rate (FDR) and false non-discovery rate (FNR) are defined, respectively, as

$$\text{FDR}(\mathcal{R}) := \mathbb{E} \left[ \frac{|\widehat{S} \setminus S|}{\max\{|\widehat{S}|, 1\}} \right] \quad \text{and} \quad \text{FNR}(\mathcal{R}) := \mathbb{E} \left[ \frac{|S \setminus \widehat{S}|}{\max\{|S|, 1\}} \right], \quad (2.5)$$

where the maxima in the denominators resolve the possible division-by-0 problem. Roughly speaking, FDR measures the expected fraction of false findings, while FNR describes the proportion of Type II errors among the true signals, and reflects the average marginal power of the procedure.

A more stringent criterion for false discovery is the family-wise error rate (FWER), defined to be the probability of reporting at least one finding not contained in the true index set. Correspondingly, a more stringent criterion for false non-discovery is the family-wise non-discovery rate (FWNR), i.e., the probability of missing at least one signal in the true index set. That is,

$$\text{FWER}(\mathcal{R}) := 1 - \mathbb{P}[\widehat{S} \subseteq S] \quad \text{and} \quad \text{FWNR}(\mathcal{R}) := 1 - \mathbb{P}[S \subseteq \widehat{S}]. \quad (2.6)$$

We introduce five different statistical risk metrics, each having different asymptotic limits in the support recovery problems as we will see in Chapter 3. Following Arias-Castro and Chen (2017), we define the risk for *approximate* support recovery as

$$\text{risk}^A(\mathcal{R}) := \text{FDR}(\mathcal{R}) + \text{FNR}(\mathcal{R}). \quad (2.7)$$

Analogously, we define the risk for *exact* support recovery as

$$\text{risk}^E(\mathcal{R}) := \text{FWER}(\mathcal{R}) + \text{FWNR}(\mathcal{R}). \quad (2.8)$$

Two closely related measures of success in the exact support recovery risk are the probability of exact recovery,

$$\mathbb{P}[\widehat{S} = S] = 1 - \mathbb{P}[\widehat{S} \neq S], \quad (2.9)$$

and the Hamming loss

$$H(\widehat{S}, S) := |\widehat{S} \Delta S| = \sum_{i=1}^p |\mathbb{1}_{\widehat{S}}(i) - \mathbb{1}_S(i)|. \quad (2.10)$$

which counts the number of mismatches between the estimated and true support sets.

The relationship between probability of support recovery  $\mathbb{P}[\widehat{S} = S]$ , exact support recovery risk  $\text{risk}^E$ , and the expected Hamming loss  $\mathbb{E}[H(\widehat{S}, S)]$  will be discussed in Section 2.4 below.

Notice that all risk metrics introduced so far penalize false discoveries and missed signals somewhat symmetrically — the approximate support recovery risk combines proportions of errors, the exact support recovery risk combines probabilities of errors, and the Hamming loss increments the risk by one regardless of the types of errors made. In applications, however, attitudes towards Type I and Type II errors are often different. In the example of GWAS, where the number of candidate locations  $p$  could be in the millions, researchers are typically interested in the marginal (location-wise) power of discovery, while exercising stringent (family-wise) false discovery control. Such types of asymmetric risks, while important in applications, have not been studied theoretically. This asymmetric consideration, and in particular the GWAS application, prompts us to consider risks that weigh both the family-wise error rate and the marginal power of discovery. One such risk metric is what shall be referred to as the *exact-approximate* support recovery risk

$$\text{risk}^{EA}(\mathcal{R}) := \text{FWER}(\mathcal{R}) + \text{FNR}(\mathcal{R}). \quad (2.11)$$

The somewhat cumbersome name and notation are chosen to reflect the asymmetry in dealing with the two types of errors in support recovery. Namely, when the risk metric (2.11) vanishes, we have “exact false discovery control, and approximate false non-discovery control” asymptotically.

Analogously, we consider the *approximate-exact* support recovery risk

$$\text{risk}^{\text{AE}}(\mathcal{R}) := \text{FDR}(\mathcal{R}) + \text{FWNR}(\mathcal{R}), \quad (2.12)$$

which places more emphasis on non-discovery control over false discovery.

Theoretical limits and performance of procedures in support recovery problems will be studied in terms of the five risk metrics (2.7), (2.8), (2.9), (2.11) and (2.12), in Chapters 3, 4, and 7. We are particularly interested in fundamental limits of signal detection and support recovery problems in terms of these metrics, as well as the optimality of commonly used procedures in high dimensional settings.

## 2.2 Statistical procedures

We review some popular procedures for signal detection and signal support recovery tasks in this section.

**Signal detection.** One of the commonly used statistics in sparse signal detection problems such as (2.1) and (2.2) are the  $L_q$  norms of the observations  $x$ ,

$$L_q(x) = \left( \sum_{i=1}^p |x(i)|^q \right)^{1/q}. \quad (2.13)$$

Typical choices of  $q$  include  $q = 1, 2$  and  $\infty$ , where  $L_\infty(x)$  is interpreted as the limit of  $L_q(x)$  norms as  $q \rightarrow \infty$ , and is equivalent to  $\max_i |x(i)|$ . Test procedures based on (2.13) may then be written as  $T(\mathcal{R}(x)) = \mathbb{1}_{(t, +\infty)}(L_q(x))$ , where the cutoff  $t$  can be chosen to control the Type I error at desired levels.

While (2.13) measures the deviation of the data from the origin in an omnidirectional manner, statistics that are tailored to the alternatives can be used in the hopes of power improvement if the directions of the alternatives are known. For example, in the additive error model (1.1), suppose we want to test for positive mean shifts, i.e., one-sided alternative

$$\mathcal{H}_1 : \mu(i) > 0, \text{ for some } i \in \{1, \dots, p\}. \quad (2.14)$$

Then, one might consider monitoring the sum (or equivalently, the arithmetic average) of the observations,

$$S(x) := \sum_{i=1}^p x(i), \quad (2.15)$$

or the maximum of the observations,

$$M(x) := \max_{i=1, \dots, p} x(i). \quad (2.16)$$

Other tests based on the empirical CDF are also available. Assuming the same one-sided alternative, let

$$q(i) = 1 - \sup\{F_i(y) : y < x(i)\}, \quad i = 1, \dots, p, \quad (2.17)$$

be the p-values of the individual observations, where  $F_i$  is the CDF of the  $i$ -th component  $x(i)$  under  $\mathcal{H}_0$ . We define empirical CDF of the p-values as

$$\widehat{F}_p(t) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}_{[0,t]}(q(i)). \quad (2.18)$$

Viewed as random elements in the space of càdlàg functions with the Skorohod  $J_1$  topology, the centered and scaled CDFs converge weakly to a Brownian bridge,

$$\left\{ \sqrt{p} \left( \widehat{F}_p(t) - t \right) \right\}_{t \in [0,1]} \Longrightarrow \{ \mathbb{B}(t) \}_{t \in [0,1]}, \quad \text{as } p \rightarrow \infty,$$

under the global null  $\mathcal{H}_0$  and mild continuity assumptions on the  $F_i$ 's (Skorokhod, 1956). Therefore, goodness-of-fit statistics such as Kolmogorov-Smirnov distance (Smirnov, 1948), Cramer-von Mises-type statistics (Cramér, 1928; Anderson and Darling, 1952) that measure the departure from this limiting behavior can be used for testing  $\mathcal{H}_0$  against  $\mathcal{H}_1$ . Of particular interest is the higher criticism (HC) statistic, first proposed by Tukey (1976),

$$HC(x) = \max_{0 \leq t \leq \alpha_0} \frac{\widehat{F}_p(t) - t}{\sqrt{t(1-t)/p}}. \quad (2.19)$$

Each of the above statistics  $L_q$ ,  $S$ ,  $M$ , or  $HC$ , gives rise to a decision rule, whereby the null hypothesis is rejected if the statistic exceeds a suitably calibrated threshold. The choice of the threshold is typically determined based on large-sample limit theorems. For example, as shown in Theorem 1.1 of Donoho and Jin (2004), under the null hypothesis

$$\frac{HC(x)}{\sqrt{2 \log \log(p)}} \longrightarrow 1, \quad \text{in probability,}$$

as  $p \rightarrow \infty$ . Thus, one decision rule is to reject  $\mathcal{H}_0$ , if  $HC(x) > t(p, \alpha_p)$ , where  $t(p, \alpha_p) = \sqrt{2 \log \log(p)}(1 + o(1))$ . As we will see, this yields an optimal signal detection procedure (see also Theorem 1.2 in Donoho and Jin, 2004). The performance of these statistics in high-dimensional sparse signal detection problems will be reviewed in Section 2.3, and analyzed in Chapter 3.

**Signal support recovery.** In signal support recovery tasks, we shall study the performance of five procedures, all of which belong to the broad class of thresholding procedures.

**Definition 2.1 (Thresholding procedures).** A thresholding procedure for estimating the support  $S := \{i : \lambda(i) \neq 0\}$  is one that takes on the form

$$\widehat{S} = \{i \mid x(i) \geq t(x)\}, \quad (2.20)$$

where the threshold  $t(x)$  may depend on the data  $x$ .

Examples of thresholding procedures include ones that aim to control FWER (2.6) — Bonferroni's (Dunn, 1961), Sidák's (Šidák, 1967), Holm's (Holm, 1979), and Hochberg's procedure (Hochberg, 1988) — as well as procedures that target FDR (2.5), such as the Benjamini-Hochberg Benjamini and Hochberg (1995) and the Barber-Candès procedure (Barber and Candès, 2015). Indeed, the class of thresholding procedures (2.20) is so general that it contains most (but not all) of the statistical procedures in the multiple testing literature.

Under the assumption that the data  $x(i)$ 's under the null have a common marginal distribution  $F$ , we review five thresholding procedures for support recovery, starting with the well-known Bonferroni's procedure which aims at controlling family-wise error rates.

**Definition 2.2 (Bonferroni's procedure).** Bonferroni's procedure with level  $\alpha$  is the thresholding procedure that uses the threshold

$$t_p = F^{\leftarrow}(1 - \alpha/p). \quad (2.21)$$

where  $F^{\leftarrow}(u) = \inf \{x : F(x) \geq u\}$  is the generalized inverse function.

The Bonferroni procedure is deterministic, i.e. non data-dependent, and only depends on the dimension of the problem and the null distribution. A closely related procedure is Sidák's procedure (Šidák, 1967), which is a more aggressive (and also deterministic) thresholding procedure that uses the threshold

$$t_p = F^{\leftarrow}((1 - \alpha)^{1/p}). \quad (2.22)$$

The third procedure, strictly more powerful than Bonferroni's, is the so-called Holm's procedure (Holm, 1979). On observing the data  $x$ , its coordinates can be ordered from largest to smallest  $x(i_1) \geq x(i_2) \geq \dots \geq x(i_p)$ , where  $(i_1, \dots, i_p)$  is a permutation of  $\{1, \dots, p\}$ . Denote these order statistics as  $x_{[1]}, x_{[2]}, \dots, x_{[p]}$ .

**Definition 2.3 (Holm's procedure).** Let  $i^*$  be the largest index such that

$$\bar{F}(x_{[i]}) \leq \alpha/(p - i + 1), \quad \text{for all } i \leq i^*.$$

Holm's procedure with level  $\alpha$  is the thresholding procedure with threshold

$$t_p(x) = x_{[i^*]}. \quad (2.23)$$

In contrast to the Bonferroni procedure, Holm's procedure is data-dependent. A closely related, more aggressive (and also data-dependent) thresholding procedure is Hochberg's procedure (Hochberg, 1988). It replaces the index  $i^*$  in Holm's procedure with the largest index such that

$$\bar{F}(x_{[i]}) \leq \alpha/(p - i + 1).$$

Notice that both Holm's and Hochberg's procedures compare p-values to the same thresholds  $\alpha/(p - i + 1)$ . However, Holm's procedure only rejects the set of hypotheses



whose p-values are all smaller than their respective thresholds. On the other hand, Hochberg's procedure rejects the set of hypotheses as long as the largest of their p-values fall below its threshold, and therefore, can be more powerful than Holm's procedure.

It can be shown that both Bonferroni's and Holm's procedures control FWER at their nominal levels, regardless of dependence in the data (Holm, 1979). In contrast, Sidák's and Hochberg's procedures control FWER at nominal levels when data are independent (Šidák, 1967; Hochberg, 1988). [\[Stilian: DISCUSS Explain in the Supplement?\]](#)

Last but not least, we review the BH procedure, which aims at controlling FDR in (2.5), proposed by Benjamini and Hochberg (1995).

Recall the order statistics of our observations  $x_{[1]} \geq x_{[2]} \geq \dots \geq x_{[p]}$ .

**Definition 2.4 (Benjamini-Hochberg's procedure).** Let  $i^*$  be the largest index such that

$$\bar{F}(x_{[i]}) \leq \alpha i/p.$$

The Benjamini-Hochberg (BH) procedure with level  $\alpha$  is the thresholding procedure with threshold

$$t_p(x) = x_{[i^*]}, \quad (2.24)$$

The BH procedure is shown to control the FDR at level  $\alpha$  when the  $x(i)$ 's are independent (Benjamini and Hochberg, 1995). [\[Stilian: DISCUSS Again, explain in the Supplement?\]](#) Variations of this procedure have been proposed to control the FDR under certain models of dependent observations (Benjamini and Yekutieli, 2001).

The performance of these procedures in high-dimensional sparse signal support recovery problems will be reviewed in Section 2.3, and analyzed in Chapters 3, 4, and 7.

## 2.3 Related literature and our contributions

We look to derive useful asymptotic approximations for high-dimensional problems, and analyze the afore-mentioned procedures in the regime where the dimensionality of the observations diverge. Throughout this text, we consider triangular arrays of observations as described in Models (1.1) and (1.3), and study the performance of various procedures in the signal detection and support recovery tasks when

$$p \rightarrow \infty.$$

The criteria for success and failure in support recovery problems under this high-dimensional asymptotic regime are defined as follows.

**Definition 2.5.** We say a sequence of procedures  $\mathcal{R} = \mathcal{R}_p$  succeeds asymptotically in the detection problem (and respectively, exact, exact-approximate, approximate-exact, and approximate support recovery problem) if

$$\text{risk}^P(\mathcal{R}) \rightarrow 0, \quad \text{as } p \rightarrow \infty, \quad (2.25)$$

where  $P = D$  (respectively,  $E, EA, AE, A$ ).

Conversely, we say the exact support recovery fails asymptotically in the detection problem (and respectively, exact, exact-approximate, approximate-exact, and approximate support recovery problem) if

$$\liminf \text{risk}^P(\mathcal{R}) \geq 1, \quad \text{as } p \rightarrow \infty, \quad (2.26)$$

where  $P = D$  (respectively,  $E, EA, AE, A$ ).

The choice of the constant 1 in the definition (2.26) allows us to declare failure for trivial testing procedures. For example, trivial deterministic procedures that always reject, and ones that always fail to reject, both have statistical risks 1 in either the detection or the support recovery problem. Similarly, a trivial randomized procedure that reject the nulls uniformly at random also has risk of 1, and is declared as a failure in both problems.

**Signal detection.** The asymptotic behavior of the statistical risk in signal detection problems (2.4) in high dimensions was first studied by [Stilian: Yurii Ismailovich Ingster](#) where did you find the full name? in the context of sparse additive models (1.1) with independent and Gaussian components. Specifically, Ingster (1998) considered the behavior of the theoretically optimal likelihood ratio (LR) test in the high-dimensional asymptotic regime, where the dimension  $p$  grows to infinity. Then, under certain parameterization of the size and sparsity of the signal  $\mu$ , a dichotomy exists: either  $\text{risk}^D(\mathcal{R})$  vanishes as  $p \rightarrow \infty$  where  $\mathcal{R}$  is the LR test, or  $\liminf_{p \rightarrow \infty} \text{risk}^D(\mathcal{R}) = 1$  for any procedure. The precise signal size and sparsity parameterizations as well as the so-called *detection boundary* discovered by Ingster are described in Chapter 3.

The LR test, unfortunately, relies on the knowledge of the signal sparsity and signal sizes which are not available in practice. The sparsity-and-signal-size-agnostic statistic HC in (2.19) was identified to attain such optimal performance limits in sparse Gaussian models in Donoho and Jin (2004). A modified goodness-of-fit test statistic in Zhang (2002), and two statistics based on thresholded- $L_1$  and  $L_2$  norms in Zhong et al. (2013) were also shown to be asymptotically optimal in the detection problem. Recent studies have also focused on the behavior of detection risk (2.4) in dense and scale mixture models Cai et al. (2011), under general distributional assumptions (Cai and Wu, 2014; Arias-Castro and Wang, 2017), as well as when the errors are dependent (Hall and Jin, 2010). A comprehensive review focusing on the role of HC in detection problems can be found in Donoho and Jin (2015). Notwithstanding the extensive literature on the detection problem, the performances of simple statistics such as  $L_q$  norms (2.13) and sums (2.15), to the best of our knowledge, have only been sparingly documented. We gather relevant results in Chapter 3, and make several new contributions on the performance of several statistics commonly used in practice.

**Exact support recovery.** There is a wealth of literature on the so-called sparsistency (i.e.,  $\mathbb{P}[\widehat{S} = S] \rightarrow 1$  as  $p \rightarrow \infty$ ) problem in the regression context. Sparsistency problems were pursued, among many others, by Zhao and Yu (2006) and Wasserman and Roeder (2009) in the high-dimensional regression setting (where the number of samples  $n \ll p$ ), and by Meinshausen and Bühlmann (2006) in graphical models. Although there have been numerous studies on the sufficient conditions for sparsistency, efforts on necessary conditions have been scarce. Notable exceptions include Wainwright (2009a,b) and Comminges and Dalalyan (2012) in regression problems. We refer the reader to the recent book by Wainwright (2019) (and in particular, the bibliographical sections of Chapters 7 and 15 therein) for a comprehensive review.

Elaborate asymptotic minimax optimality results under the Hamming loss were derived for methods proposed in Ji and Jin (2012) and Jin et al. (2014) for regression problems. More recently, Butucea et al. (2018) also obtained similar minimax optimality results for a specific procedure in the Gaussian additive error model (1.1) in terms of the expected Hamming loss.

Nevertheless, two important questions remained unanswered. Namely, precise phase-transition-type results for the exact support recovery risk (2.8) and for the support recovery probability (2.9) have not been established. And secondly, performance of commonly used statistical procedures reviewed in Section 2.2 in terms of these risk metrics have not been studied. Some of our main contributions in this text are the solutions to these problems, presented in Chapters 3 and 4 below. Specifically, we show that the Bonferroni thresholding procedure (among others) is asymptotically optimal for the exact support recovery problem in (1.1) under broad classes of error distributions. Furthermore, a phase-transition in the exact support recovery problem for thresholding procedures is established under broad dependence conditions on the errors using the concentration of maxima phenomenon (Chapter 4). We also establish finite-sample Bayes optimality and sub-optimality results for these procedures under independence, and by extension, arrive at minimax optimality results for the exact support recovery problem.

The landscape of the fundamental statistical limits in support estimation is yet to be fully charted. We conjecture, however, that the general concentration of maxima phenomenon will lead to its complete solution under very broad error dependence scenarios.

**Approximate support recovery.** The performance limits of FDR-controlling procedures in the support recovery problem have been actively studied in recent years. The asymptotic optimality of the Benjamini-Hochberg procedure was analyzed under decision theoretic frameworks in Genovese and Wasserman (2002); Bogdan et al. (2011); Neuvial and Roquain (2012), with main focus on location/scale models. In particular, these papers show that the statistical risks of the procedures come close to that of the oracle procedures under suitable asymptotic regimes. Strategies for dealing with multiple testing under general distributional assumptions can be found in, e.g., Efron (2004), Storey (2007), and Sun and Cai (2007). The two-sided alternative in the additive error model was featured as the primary example in Sun and Cai (2007).

In the additive error model (1.1) under independent Gaussian errors and one-sided alternatives (2.14), Arias-Castro and Chen (2017) showed that a phase transition exists for the approximate support recovery risk (2.7). The BH procedure Benjamini and Hochberg (1995), and the Barber-Candès procedure (Barber and Candès, 2015) was identified to be asymptotically optimal. However, Arias-Castro and Chen (2017), as all related work so far, assumed the non-nulls to follow a common alternative distribution. We derive a new phase transition result that relaxes this assumption on the alternatives in Chapter 3.

**Asymmetric statistical risks.** Although weighted sums of false discovery and non-discovery have been studied in the literature mentioned above, the case of simultaneous family-wise error control and marginal, location-wise power requirements, has not been previously considered. As a result, asymmetric statistical risks such as (2.11) and (2.12) have not been investigated. As argued in Section 2.1, the properties of these asymmetric risks are of important practical concern in applications such as GWAS. We study the asymptotic behavior of these risks in Chapters 3 and 7 of this text.

**Chi-square models and GWAS.** The high-dimensional chi-square model (1.3) seemed to have received little attention in the literature. While the sparse signal detection problem in the chi-square model has been studied Donoho and Jin (2004), to the best of our knowledge, asymptotic limits of the support recovery problems have not been studied. The chi-squared model and the motivating GWAS application are analyzed in Chapter 7. The results obtained therein help us explain a phenomenon in GWAS where statistical power decays sharply as function of sample size when the latter is in a small region known as the *steep part of the power curve*. This empirical fact has long been observed by statistical geneticists but has not been mathematically quantified. Gao et al. (2019) provide further details on the power and design in GWAS as well as an accompanying interactive statistical software (Gao, 2019).

## 2.4 Relationships between the asymptotic risks

We now elaborate on the relationship between statistical risks, as promised in Section 2.1. The first lemma concerns the asymptotic relationship between the probability of exact recovery (2.9) and the risk of exact support recovery (2.8).

**Lemma 2.1.** *For any sequence of procedures for support recovery  $\mathcal{R} = \mathcal{R}_p$ , we have*

$$\mathbb{P}[\widehat{S} = S] \rightarrow 1 \iff \text{risk}^E(\mathcal{R}) \rightarrow 0, \quad (2.27)$$

and

$$\mathbb{P}[\widehat{S} = S] \rightarrow 0 \implies \liminf \text{risk}^E(\mathcal{R}) \geq 1, \quad (2.28)$$

as  $p \rightarrow \infty$ . Dependence on  $p$  and  $\mathcal{R}$  was suppressed for notational convenience.

*Proof (Proof of Lemma 2.1).* Notice that  $\{\widehat{S} = S\}$  implies  $\{\widehat{S} \subseteq S\} \cap \{\widehat{S} \supseteq S\}$ , therefore we have for every fixed  $p$ ,

$$\text{risk}^E = 2 - \mathbb{P}[\widehat{S} \subseteq S] - \mathbb{P}[S \subseteq \widehat{S}] \leq 2 - 2\mathbb{P}[\widehat{S} = S]. \quad (2.29)$$

On the other hand, since  $\{\widehat{S} \neq S\}$  implies  $\{\widehat{S} \not\subseteq S\} \cup \{\widehat{S} \not\supseteq S\}$ , we have for every fixed  $p$ ,

$$1 - \mathbb{P}[\widehat{S} = S] = \mathbb{P}[\widehat{S} \neq S] \leq 2 - \mathbb{P}[\widehat{S} \subseteq S] - \mathbb{P}[S \subseteq \widehat{S}] = \text{risk}^E. \quad (2.30)$$

Relation (2.27) follows from (2.29) and (2.30), and Relation (2.28) from (2.30).  $\square$

By virtue of Lemma 2.1, it is sufficient to study the probability of exact support recovery  $\mathbb{P}[\widehat{S} = S]$  in place of  $\text{risk}^E$ , if we are interested in the asymptotic properties of the risk in the sense of (2.25) and (2.26).

Keen readers must have noticed the asymmetry in Relation (2.28) when we discussed the relationship between the exact support recovery risk (2.8) and the probability of exact support recovery (2.9). While a trivial procedure that never rejects and a procedure that always rejects both have  $\text{risk}^E$  equal to 1, the converse is not true. For example, it is possible that a procedure selects the true index set  $S$  with probability  $1/2$ , but otherwise makes one false inclusion *and* one false omission simultaneously. In this case the procedure will have

$$\text{risk}^E = 1, \quad \text{and} \quad \mathbb{P}[\widehat{S} = S] = 1/2,$$

showing that the converse of Relation (2.28) is in fact false.

The same argument applies to  $\text{risk}^A$ : a procedure may select the true index set  $S$  with probability  $1/2$ , but makes enough false inclusions and omissions the rest of the time, so that  $\text{risk}^A$  is equal to one. Therefore, although the class of methods with risks equal to or exceeding 1 certainly contains the trivial procedures that we mentioned, they are not necessarily “useless” as some researchers have claimed (cf. Remark 2 in Arias-Castro and Chen, 2017).

Upper and lower bounds for FDR and FNR can be immediately derived by replacing the numerators in (2.5) with the Hamming loss,

$$\mathbb{E} \left[ \frac{H(\widehat{S}, S)}{\max\{|\widehat{S}|, |S|, 1\}} \right] \leq \text{FDR} + \text{FNR} \leq \mathbb{E} \left[ \frac{H(\widehat{S}, S)}{\max\{\min\{|\widehat{S}|, |S|\}, 1\}} \right]. \quad (2.31)$$

Therefore, it is sufficient, but not necessary, that the Hamming loss vanish in order to have vanishing approximate support recovery risks (2.7).

Turning to the relationship between the probability of exact support recovery (2.9) and Hamming loss (2.10), we point out a natural lower bound of the former using the expectation of the latter,

$$\mathbb{P}[\widehat{S} = S] \geq 1 - \mathbb{E}[H(\widehat{S}, S)] = 1 - \sum_{i=1}^p \mathbb{E} |\mathbb{1}_{\widehat{S}}(i) - \mathbb{1}_S(i)|. \quad (2.32)$$

A key observation in Relation (2.32) is that the expected Hamming loss decouples into  $p$  terms, and the dependence of the estimates  $\mathbb{1}_{\widehat{S}}(i)$  among the  $p$  locations no longer plays any role in the sum. Therefore, studying support recovery problems via the expected Hamming loss is not very informative especially under severe dependence, as the bound (2.32) may become very loose. Vanishing Hamming loss is again sufficient, but not necessary for  $\mathbb{P}[\widehat{S} = S]$  or the exact support recovery risk to fo to zero.

## 2.5 The asymptotic generalized Gaussian (AGG) models

We introduce a fairly general class of distributions known as asymptotic generalized Gaussians AGG. We also state some of their tail properties which play important roles in the analysis of phase transitions of high-dimensional testing problems.

**Definition 2.6.** A distribution  $F$  is called asymptotic generalized Gaussian with parameter  $\nu > 0$  (denoted  $\text{AGG}(\nu)$ ) if

1.  $F(x) \in (0, 1)$  for all  $x \in \mathbb{R}$ , and
2.  $\log \bar{F}(x) \sim -\frac{1}{\nu}x^\nu$  and  $\log F(-x) \sim -\frac{1}{\nu}(-x)^\nu$ ,

where  $\bar{F}(x) = 1 - F(x)$  is the survival function, and  $a(x) \sim b(x)$  is taken to mean  $\lim_{x \rightarrow \infty} a(x)/b(x) = 1$ .

The AGG models include, for example, the standard Gaussian distribution ( $\nu = 2$ ) and the Laplace distribution ( $\nu = 1$ ) as special cases. Since the requirement is only placed on the tail behavior, this class encompasses a large variety of light-tailed models. This class is commonly used in the literature on high-dimensional testing (Cai et al., 2007; Arias-Castro and Chen, 2017).

**Proposition 2.1.** *The  $(1/p)$ -th upper quantile of  $\text{AGG}(\nu)$  is*

$$u_p := F^{\leftarrow}(1 - 1/p) \sim (\nu \log p)^{1/\nu}, \quad \text{as } p \rightarrow \infty, \quad (2.33)$$

where  $F^{\leftarrow}(q) = \inf_x \{x : F(x) \geq q\}$ ,  $q \in (0, 1)$ .

*Proof (Proof of Proposition 2.1).* By the definition of AGG, for any  $\epsilon > 0$ , there is a constant  $C = C(\epsilon)$  such that for all  $x \geq C$ , we have

$$-\frac{1}{\nu}x^\nu(1 + \epsilon) \leq \log \bar{F}(x) \leq -\frac{1}{\nu}x^\nu(1 - \epsilon).$$

Therefore, for all  $x < x_l := ((1 + \epsilon)^{-1} \nu \log p)^{1/\nu}$ , we have

$$-\log p = -\frac{1}{\nu} x_l^\nu (1 + \epsilon) \leq \log \bar{F}(x_l) \leq \log \bar{F}(x), \quad (2.34)$$

and for all  $x > x_u := ((1 - \epsilon)^{-1} \nu \log p)^{1/\nu}$ , we have

$$\log \bar{F}(x) \leq \log \bar{F}(x_u) \leq -\frac{1}{\nu} x_u^\nu (1 - \epsilon) = -\log p. \quad (2.35)$$

By definition of generalized inverse,

$$u_p := F^{\leftarrow}(1 - 1/p) = \inf\{x : \bar{F}(x) \leq 1/p\} = \inf\{x : \log \bar{F}(x) \leq -\log p\}.$$

We know from relations (2.34) and (2.35) that

$$[x_u, +\infty) \subseteq \{x : \log \bar{F}(x) \leq -\log p\} \subseteq [x_l, +\infty),$$

and so  $x_l \leq u_p \leq x_u$ , and the expression for the quantiles follows.  $\square$

## 2.6 Rapid variation and relative stability

The behavior of the maxima of identically distributed random variables has been studied extensively in the extreme value theory literature (see, e.g., Leadbetter et al., 1983; Resnick, 2013; Embrechts et al., 2013; De Haan and Ferreira, 2007, and the references therein). The concept of rapid variation plays an important role in the light-tailed case.

**Definition 2.7 (Rapid variation).** The survival function of a distribution,  $\bar{F}(x) = 1 - F(x)$ , is said to be rapidly varying if

$$\lim_{x \rightarrow \infty} \frac{\bar{F}(tx)}{\bar{F}(x)} = \begin{cases} 0, & t > 1 \\ 1, & t = 1 \\ \infty, & 0 < t < 1 \end{cases}. \quad (2.36)$$

When  $F(x) < 1$  for all finite  $x$ , Gnedenko (1943) showed that the distribution  $F$  has rapidly varying tails if and only if the maxima of independent observations from  $F$  are *relatively stable* in the following sense.

**Definition 2.8 (Relative stability).** Let  $\epsilon_p = (\epsilon_p(i))_{i=1}^p$  be a sequence of random variables with common marginal distribution  $F$ . Define the sequence  $(u_p)_{p=1}^\infty$  to be the  $(1 - 1/p)$ -th generalized quantile of  $F$ , i.e.,

$$u_p = F^{\leftarrow}(1 - 1/p). \quad (2.37)$$

The triangular array  $\mathcal{E} = \{\epsilon_p, p \in \mathbb{N}\}$  is said to have relatively stable (RS) maxima if

$$\frac{1}{u_p} M_p := \frac{1}{u_p} \max_{i=1, \dots, p} \epsilon_p(i) \xrightarrow{\mathbb{P}} 1, \quad \text{as } p \rightarrow \infty. \quad (2.38)$$

In the case of independent and identically distributed  $\epsilon_p(i)$ 's, Barndorff-Nielsen (1963) and Resnick and Tomkins (1973) obtained necessary and sufficient conditions for the *almost sure stability* of maxima, where the convergence in (2.38) holds almost surely. See also Klass (1984) for further sharp results on almost sure stability, and Naveau (2003) for almost sure stability in stationary sequences. Here, we will only need the weaker notion in (2.38) but extend our inquiry to the case of dependent  $\epsilon_p(i)$ 's.

While relative stability (and almost sure stability) is well-understood in the independent case, the role of dependence has not been fully explored. We start this investigation with a small refinement of Theorem 2 in Gnedenko (1943) valid under *arbitrary dependence*.

**Proposition 2.2 (Rapid variation and relative stability).** *Assume that the array  $\mathcal{E}$  consists of identically distributed and possibly dependent random variables with cumulative distribution function  $F$ , where  $F(x) < 1$  for all finite  $x > 0$ .*

1. *If  $F$  has rapidly varying right tail in the sense of (2.36), then for all  $\delta > 0$ ,*

$$\mathbb{P} \left[ \frac{1}{u_p} M_p \leq 1 + \delta \right] \geq 1 - \frac{\bar{F}((1 + \delta)u_p)}{\bar{F}(u_p)} \rightarrow 1. \quad (2.39)$$

2. *If the array  $\mathcal{E}$  has independent entries, then it is relatively stable if and only if  $F$  has rapidly varying tail, i.e., (2.36) holds.*

*Proof (Proof of Proposition 2.2).* By the union bound and the fact that  $p\bar{F}(u_p) \leq 1$ , we have

$$\mathbb{P}[M_p > (1 + \delta)u_p] \leq p\bar{F}((1 + \delta)u_p) \leq \frac{\bar{F}((1 + \delta)u_p)}{\bar{F}(u_p)}. \quad (2.40)$$

In view of (2.36) (rapid variation) and the fact that  $u_p \rightarrow \infty$ , as  $p \rightarrow \infty$ , the right-hand side of (2.40) vanishes as  $p \rightarrow \infty$ , for all  $\delta > 0$ . This completes the proof of (2.39). Part 2 is a re-statement of a classic result dating back to Gnedenko (1943).  $\square$

*Remark 2.1.* Part (1) of Proposition 2.2 is equivalent to

$$\mathbb{P} \left[ \frac{1}{u_p} M_p > 1 + \delta_p \right] \rightarrow 0, \quad \text{as } p \rightarrow \infty, \quad (2.41)$$

for some positive sequence  $\delta_p \rightarrow 0$ . Notice on the other hand that, if  $M_p^*$  is the maximum of  $p$  iid variables with distribution  $F$ , the relative stability property entails  $M_p^*/u_p \rightarrow 1$ , in probability, as  $p \rightarrow \infty$ . Since the sequence  $1 + \delta_p \rightarrow 1$ , Relation (2.41) means that the rate of growth of the maxima  $M_n$  in  $\mathcal{E}$  cannot be faster than that of the independent maxima  $M_p^*$ . This somewhat surprising fact holds regardless



of the dependence structure of  $\mathcal{E}$  and is solely a consequence of the rapid variation of  $F$ .

We demonstrate next that the Gaussian, Exponential, Laplace, and Gamma distributions all have rapidly varying tails.

*Example 2.1 (Generalized AGG).* A distribution is said to have *Generalized AGG* right tail, if  $\log \bar{F}$  is regularly varying,

$$\log \bar{F}(x) = -x^\nu L(x), \quad (2.42)$$

where  $\nu > 0$  and  $L : (0, +\infty) \rightarrow (0, +\infty)$  is a slowly varying function. (A function is said to be slowly varying if  $\lim_{x \rightarrow \infty} L(tx)/L(x) = 1$  for all  $t > 0$ .) Note that the AGG( $\nu$ ) model corresponds to the special case where  $L(x) \rightarrow 1/\nu$ , as  $x \rightarrow \infty$ .

Relation (2.39) holds for all arrays  $\mathcal{E}$  with *generalized* AGG marginals; if the entries are independent, the maxima are relatively stable. This follows directly from Proposition 2.2, once we show that  $F$  has rapidly varying tail. Indeed, by (2.42), we have

$$\log \left( \bar{F}(tx) / \bar{F}(x) \right) = -L(x)x^\nu \left( t^\nu \frac{L(tx)}{L(x)} - 1 \right),$$

which converges to  $-\infty$ , 0, and  $+\infty$ , as  $x \rightarrow \infty$ , when  $t > 1$ ,  $t = 1$ , and  $t < 1$ , respectively, since  $x^\nu L(x) \rightarrow \infty$  as  $x \rightarrow \infty$  by definition of  $L$ .

The AGG class encompasses a wide variety of rapidly varying tail models such as Gaussian and double exponential distributions. The larger class (2.42) is needed, however, for the Gamma distribution.

More generally, distributions with heavier tails (e.g., log-normal) and lighter tails (e.g., Gompertz) outside the generalized AGG class (2.42) may also possess rapidly varying tails; heavy-tailed distributions like the Pareto and t-distributions, on the other hand, do not. [Stilian: More details on these alternative classes of](#)

[models can be found in the Supplement.](#) [Move these to the online Supplement?](#)

These alternative classes of models will be introduced when we study the phase transitions in Chapter 4.



## Chapter 3

# The Phase Transition Phenomena in Independent Gaussian Error Models

**Stilian:** The purpose of this chapter is to provide a unified review of the fundamental statistical limits in the sparse signal detection and support recovery problems. Our goal is to convey the main ideas and thus we shall focus on the simple but important setting of independent Gaussian errors. ~~We study the fundamental limits of the signal detection and support recovery problems in the vector models with independent Gaussian errors in this chapter.~~ Specifically, we derive the conditions under which the detection and support recovery problems succeed and fail in the sense of (2.25) and (2.26), in the additive error model

$$x(i) = \mu(i) + \epsilon(i), \quad i = 1, \dots, p, \quad (3.1)$$

where the errors  $\epsilon(i)$ 's are iid standard Gaussians random variables. Once again, we restrict our analysis to models with independent and identically distributed Gaussian errors for the moment. Both the distributional assumption and the independence assumption will be relaxed substantially in the following chapters.

As laid out in Section 2.3, we work under the asymptotic regime where the problem dimension  $p$  diverges to infinity. The set of non-zero entries of the signal vector  $\mu = \mu_p$  will be referred to as its *support* and denoted by

$$S_p := \{i : \mu(i) \neq 0\}.$$

We shall assume that the size of the support is

$$|S_p| = \lfloor p^{1-\beta} \rfloor, \quad \beta \in (0, 1], \quad (3.2)$$

where  $\beta$  parametrizes the problem sparsity. The closer  $\beta$  is to 1, the sparser the support  $S_p$ . Conversely, when  $\beta$  is close to 0, the support is dense with many non-null signals. We consider one-sided alternatives (2.14), and parametrize the range of the non-zero (and perhaps unequal) signals with

$$\underline{\Delta} = \sqrt{2r \log p} \leq \mu(i) \leq \bar{\Delta} = \sqrt{2\bar{r} \log p}, \quad \text{for all } i \in S_p, \quad (3.3)$$

for some constants  $0 < \underline{r} \leq \bar{r} \leq +\infty$ .

The parametrization of signal sparsity (3.2) and signal sizes (3.3) in the Gaussian model was first introduced in Ingster (1998), and later adopted by Hall and Jin (2010), Cai et al. (2011), Zhong et al. (2013), Cai and Wu (2014), Arias-Castro and Wang (2017), and numerous others for studying the signal detection problem in Gaussian location-scale models. Similar scalings of sparsity and signal size are also used in, e.g., Ji and Jin (2012), Jin et al. (2014), Butucea et al. (2018) to study the phase transitions of the support recovery problems under Gaussianity assumptions.

### 3.1 Sparse signal detection problems

The optimality of sparse signal detection was first studied by Ingster (1998), who showed that a phase transition in the  $r$ - $\beta$  plane exists for the signal detection problem. Specifically, consider the so-called *detection boundary* function:

$$f(\beta) = \begin{cases} \max\{0, \beta - 1/2\} & 0 < \beta \leq 3/4 \\ (1 - \sqrt{1 - \beta})^2 & 3/4 < \beta \leq 1. \end{cases} \quad \beta \in (0, 1]. \quad (3.4)$$

Assume that the non-zero signal sizes are all equal and parameterized as  $\sqrt{2r \log p}$ . If the signal size parameter  $r$  is *above* the detection boundary, i.e.,  $r > f(\beta)$ , then the global null hypothesis  $\mu(i) = 0$  for all  $i = 1, \dots, p$  can be distinguished from the alternative as  $p \rightarrow \infty$  in the sense of (2.25) using the likelihood ratio test. Otherwise, when the signal sizes fall below the boundary, i.e.,  $r < f(\beta)$ , no test can do better than a random guess. We visualize the detection boundary in the upper panel of Figure 3.1.

Adaptive tests such as Tukey's HC in (2.19) (Donoho and Jin, 2004) and a modified goodness-of-fit test statistic of Zhang (2002) have been identified to attain this performance limit without knowledge of the sparsity and signal sizes. It is also known that the max-statistic (2.16) is only efficient when  $r > (1 + \sqrt{1 - \beta})^2$ , and is therefore sub-optimal for denser signals where  $1/2 \leq \beta \leq 3/4$ ; see Cai et al. (2011). In contrast, the sum-of-square-type statistics such as  $L_2$  was shown in Fan (1996) to be asymptotically powerless when the  $L_2$ -norm of the signal  $\|\mu\|_2^2$  is  $o(\sqrt{p})$ , or equivalently, when  $\beta > 1/2$  in our parametrization.

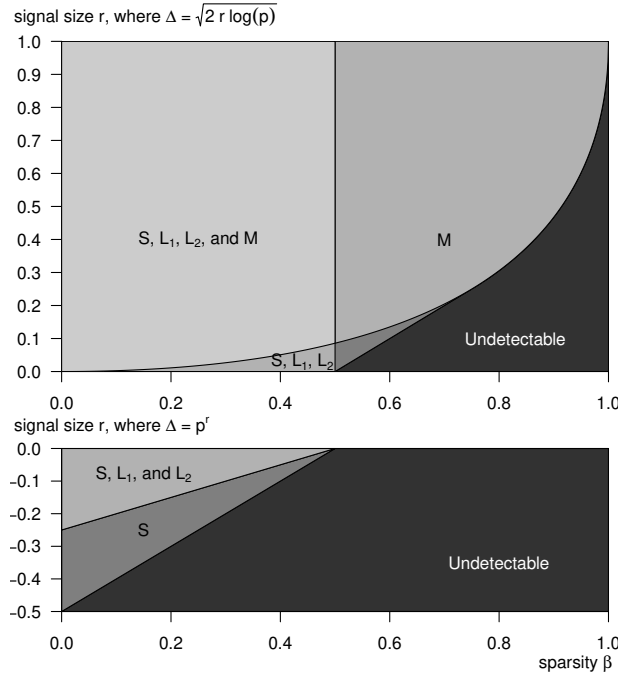
Notice that the scaling for the signal magnitude  $\Delta = \sqrt{2r \log p}$  is useful for studying very sparse signals ( $\beta > 1/2$ ), but fails to reveal the difficulties of the detection problems when the signals are relatively dense ( $\beta < 1/2$ ). This is because  $f(\beta) = 0$ ,  $\beta \in (0, 1/2]$ . Thus, a different scaling is needed to study the regime of small but dense signals. In this case, with slight overloading of notation, we parametrize signal sizes as

$$\underline{\Delta} = p^{\underline{r}} \leq \mu(i) \leq \bar{\Delta} = p^{\bar{r}}, \quad \text{for all } i \in S_p, \quad (3.5)$$

where  $\underline{r}$  and  $\bar{r}$  are negative constants and the signal magnitude vanishes, as  $p \rightarrow \infty$ . In this scaling, for the so-called faint signal regime, Cai et al. (2011) established a phase transition result characterized by the following boundary,

$$f_s(\beta) = \beta - 1/2, \quad 0 < \beta \leq 1/2. \quad (3.6)$$

Specifically, if  $\bar{r} < f_s(\beta)$ , the signal detection fails in the sense of (2.26) regardless of the procedures, while the HC statistic continues to attain asymptotically perfect detection when  $\underline{r} > f_s(\beta)$ . We visualize this boundary in the lower panel of Figure 3.1.



**Fig. 3.1** The phase diagrams of the sparse signal detection problem. Signal size and sparsity are parametrized by  $r$  and  $\beta$ , respectively. The diagrams illustrate the regions where the signal detection problem can be solved asymptotically by some of the commonly used statistics: the maximum ( $M$ ), the sum-of-squares ( $L_2$ ), the sum-of-absolute values ( $L_1$ ), and the sum ( $T$ ). In each region of the diagram, the annotated statistics can make the detection risk (2.4) vanish, as dimension  $p$  diverges. Conversely, the risks has  $\liminf$  at least one. The detection problem is unsolvable for very sparse and weak signals in the undetectable regions. Notice that the  $L_1$  and  $L_2$  statistics are in fact sub-optimal for all sparsity levels. On the other hand, the max-statistic remains powerful for sparse signals ( $\beta > 1/2$ ), and is fully efficient when the problem is very sparse ( $\beta \geq 3/4$ ). The HC statistic can detect signals in all configurations in the detectable regions. See text and Theorem 3.1.

To the best of our knowledge, performance of simple statistics such as  $L_1$ ,  $L_2$  norms, and [Stilian:](#) the sum statistic  $T$  in (2.15) R2, Typo 7 in this weak signal

setting have not been reported in the literature. Our first theorem establishes the performance of these simple but popular statistics for detecting sparse signals in high-dimensions, and summarizes the known results.

**Theorem 3.1.** *Consider the signal detection problem in the triangular array of Gaussian error models (3.1) where the sparsity is parametrized as in (3.2).*

- *For signals whose sizes are parametrized as in (3.3), the detection problem can be asymptotically solved in the sense of (2.25) with  $L_2$ ,  $L_1$ , or  $S$  statistic when  $\beta \leq 1/2$ ; on the other hand, these statistics are asymptotically powerless in the sense of (2.26) when  $\beta > 1/2$ .*
- *For small and dense signals whose signal sizes are parametrized as in (3.5), the detection problem can be asymptotically solved in the sense of (2.25) with  $L_2$  or  $L_1$  statistic when  $\underline{r} > \beta/2 - 1/4$ ; on the other hand, these statistics are asymptotically powerless in the sense of (2.26) when  $\bar{r} < \beta/2 - 1/4$ . Further, tests based on the  $S$  statistic can succeed asymptotically in the sense of (2.25) when  $\underline{r} > \beta - 1/2$ , hence attaining the boundary of detectability in (3.6).*

Theorem 3.1 is proved in Section ?? . We visualize the results in Theorem in Figure 3.1. It is worth noting that the  $\beta$ - $r$  parameter regions where  $L_1$  and  $L_2$  statistics are asymptotically powerful coincide, and these statistics are theoretically suboptimal for both sparse regimes ( $\beta > 1/2$ ) and relatively dense regimes ( $\beta \leq 1/2$ ).

Ideas have been proposed to combine statistics that are powerful for different alternatives to create adaptive tests that maintain high power for at all sparsity levels. Such adaptive tests can be constructed, for example, by leveraging the asymptotic independence of the sum- and supremum-type statistics (Hsing, 1995). Recently, Xu et al. (2016) showed that for dependent observations under mixing and moment conditions, the sum-of-power-type statistics

$$\tilde{L}_q(x) = \sum_{i=1}^p x^q(i) \quad (3.7)$$

with distinct positive integer powers (i.e.,  $q = 1, 2, \dots$ ) are asymptotically jointly independent, and proposed an adaptive test that monitors the minimum p-value of tests constructed with  $\tilde{L}_q$ 's. This idea is further developed in Wu et al. (2019) for generalized linear models and in He et al. (2018) with U-statistics.

Optimality properties of such adaptive tests and the optimal choice of the  $q$ -combinations, however, remain open problems. Xu et al. (2016) suggested combine  $q = 1, 2, 3, \dots, 6$ , and  $q = \infty$ , based empirical evidence from numerical experiments. Theorem 3.1 here implies that, at least for detecting one-sided alternatives, the  $\tilde{L}_2$  statistic (i.e.,  $L_2$  norm) and the  $L_1$  norm are asymptotically dominated by the  $\tilde{L}_1$  statistic (or equivalently, the sum  $S$ ). Therefore it is sufficient to include only the latter in the construction of the adaptive test.

### 3.2 Sparse signal support recovery problems

Turning to support recovery problems in the Gaussian error model (3.1), we will analyze the asymptotic performance limits in terms of the risk metrics for exact, exact-approximate, approximate-exact support recovery problems (i.e., (2.8), (2.11), and (2.12), respectively), as well as the probability of support recovery (2.9). We will also review the recent result for exact support recovery risk (2.7) by Arias-Castro and Chen (2017), to reveal a rather complete landscape of support recovery problems in high-dimensional Gaussian error models.

We restrict our attention to the class of thresholding procedures in this section. Specifically, the lower bounds that we develop in Theorems 3.2 through 3.5 below are only meant to apply to thresholding procedures. Although it is intuitively appealing to consider only data-thresholding procedures in multiple testing problems, such procedures are not always optimal in more general settings. The optimality of thresholding procedures and the consequences of this restriction will be treated in Chapter 4.

A technical ingredient is needed in order to state our main results. We define a rate at which the nominal levels of FWER or FDR go to zero.

**Definition 3.1.** We say the nominal level of errors  $\alpha = \alpha_p$  vanishes slowly, if

$$\alpha \rightarrow 0, \quad \text{and} \quad \alpha p^\delta \rightarrow \infty \text{ for any } \delta > 0. \quad (3.8)$$

As an example, the sequence of nominal levels  $\alpha_p = 1/\log(p)$  is slowly vanishing, while the sequence  $\alpha_p = 1/\sqrt{p}$  is not.

#### 3.2.1 The exact support recovery problem

Our study of the exact support recovery risk (2.8) begins with a brief review of existing results for the Hamming loss (2.10). Indeed, as discussions in Section 2.3 suggest, the latter can be informative of the exact support recovery problems for models with independent components.

Inspired by the phase transition results for the signal detection problem, Ji and Jin (2012), Genovese et al. (2012), and Jin et al. (2014) derived interesting sharp results on support recovery problems in linear models under the Hamming loss  $H(\hat{S}, S)$ . Specifically, these papers establish minimax-type phase transition results in their respective settings. Under the sparsity parametrization in (3.2) and assuming equal signal sizes of  $(2r \log p)^{1/2}$ , Hamming losses were shown to diverge to  $+\infty$  when  $r$  falls below the threshold

$$g(\beta) = (1 + (1 - \beta)^{1/2})^2, \quad (3.9)$$

for any method of support estimation. Conversely, under orthogonal, or near-orthogonal random designs, if  $r > g(\beta)$ , they showed that the methods they proposed achieve vanishing Hamming loss.

Very recently, Butucea et al. (2018) studied both asymptotics and non-asymptotics of support recovery problems in the additive noise model (3.1) under the assumption of equal signal sizes, using the Hamming loss. Again, the analysis of asymptotic optimality focused on a newly proposed procedure which is very specific to the Gaussian model. It is not at all clear if the optimality properties are a consequence of its mysterious construction.

We now show that commonly used and computationally efficient procedures can also be asymptotically optimal in the exact support recovery problem.

**Theorem 3.2.** *Consider the high-dimensional additive error model (3.1) under independent standard Gaussian errors, with signal sparsity and size as described in (3.2) and (3.3). The function (3.9) characterizes the phase transition of the exact support recovery problem. Specifically, if  $\bar{r} > g(\beta)$ , then Bonferroni's, Sidák's, Holm's, and Hochberg's procedures with slowly vanishing nominal FWER levels (as defined in Definition 3.1) all achieve asymptotically exact support recovery in the sense of (2.25).*

*Conversely, if  $\bar{r} < g(\beta)$ , then for any thresholding procedure  $\hat{S}_p$ , we have  $\mathbb{P}[\hat{S}_p = S_p] \rightarrow 0$ . Therefore, in view of Lemma 2.1, exact support recovery asymptotically fails for all thresholding procedures in the sense of (2.26).*

We visualize the result in a  $\beta$ - $r$  phase diagram in Figure 3.2.

Theorem 3.2 is in fact a special case of the more general Theorem 4.1 which covers dependent and non-Gaussian errors. We will study the exact support recovery problem in greater detail, and prove the more general version of the Theorem in Chapter 4.

### 3.2.2 The approximate support recovery problem

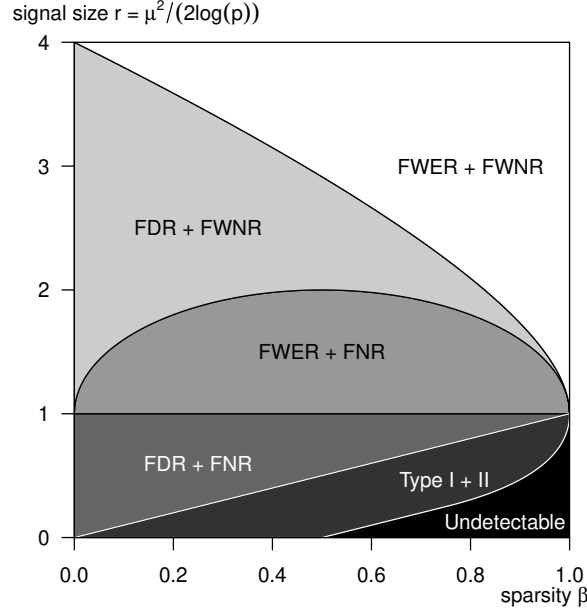
Arias-Castro and Chen (2017) studied the performance of the Benjamini-Hochberg procedure (Benjamini and Hochberg, 1995) and a stripped-down version of the Candès-Barber procedure (Barber and Candès, 2015) in approximate support recovery problems when the components of the noise term  $\epsilon$  in (3.1) have independent and symmetric distributions. A phase transition phenomenon for the approximate support recovery risk (2.7) was established in the Gaussian additive error model, where the two aforementioned methods are both shown to be asymptotically optimal.

The analysis therein, however, assumed equal signal sizes for the alternatives. We generalize the main results of Arias-Castro and Chen (2017) to allow for unequal signal sizes.

**Theorem 3.3.** *In the context of Theorem 3.2, the function*

$$h(\beta) = \beta \tag{3.10}$$





**Fig. 3.2** The phase diagram of support recovery problems for the high-dimensional chi-square model (3.1), illustrating the boundaries of the exact support recovery (FWER + FWNR; top curve; Theorem 3.2), the approximate-exact support recovery (FDR + FWNR; second curve from top; Theorem 3.5), the exact-approximate support recovery (FWER + FNR; horizontal line  $r = 1$ ; Theorem 3.4), and the approximate support recovery problems (FDR + FNR; tilted line  $r = \beta$ ; Theorem 3.3). The signal detection problem (Type I + Type II errors of the global test; lower curve) was studied in Donoho and Jin (2004). In each region of the diagram and above, the annotated statistical risk can be made to vanish, as dimension  $p$  diverges. Conversely, the risks has  $\liminf$  at least one.

characterizes the phase transition of approximate support recovery problem. Specifically, if  $\underline{r} > h(\beta)$ , then the Benjamini-Hochberg procedure (defined in Section 2.2) with slowly vanishing nominal FDR levels (as defined in Definition 3.1) achieves asymptotically approximate support recovery in the sense of (2.25).

Conversely, if  $\bar{r} < h(\beta)$ , then approximate support recovery asymptotically fails in the sense of (2.26) for all thresholding procedures.

Proof of Theorem 3.3 is presented in Section ???. The key to proving this generalization is a monotonicity property of the BH procedure. Namely, the power of the BH procedure in terms of FNR monotonically increases for stochastically larger alternatives. This fact will be formalized in Lemma A.1, and may be of independent interest.

### 3.2.3 The exact-approximate support recovery problem

We now derive two new asymptotic phase transition results for the *asymmetric* statistical risks, (2.11) and (2.12), in the Gaussian error models. The next theorem describes the phase transition in the exact-approximate support recovery problem.

**Theorem 3.4.** *In the context of Theorem 3.2, the function*

$$\tilde{g}(\beta) = 1 \quad (3.11)$$

*characterizes the phase transition of exact-approximate support recovery problem. Specifically, if  $\underline{r} > \tilde{g}(\beta)$ , then the procedures listed in Theorem 3.2 with slowly vanishing nominal FWER levels (as defined in Definition 3.1) achieve asymptotically exact-approximate support recovery in the sense of (2.25).*

*Conversely, if  $\bar{r} < \tilde{g}(\beta)$ , then for any thresholding procedure  $\hat{S}$ , the exact-approximate support recovery fails in the sense of (2.26).*

Theorem 3.4 is proved in Section ?? . The phase transition boundary (3.11) is visualized in Figure 3.2.

*Remark 3.1.* Boundary (3.11) was briefly suggested by Arias-Castro and Chen (2017). Unfortunately, it was falsely claimed that the boundary characterized the phase transition of the *exact* support recovery problem, and the alleged proof was left as an “exercise to the reader”. This exercise was completed in Chapter 4, where the correct boundary (7.4) was identified.

Theorem 3.4 here shows that the boundary (3.11) *does* exist, though for the slightly different *exact-approximate* support recovery problem. As we will see in Section 7.1, the boundary (3.11) also applies to the exact-approximate support recovery problem in chi-square models (1.3).

### 3.2.4 The approximate-exact support recovery problem

The last phase transition is in terms of the approximate-exact support recovery risk (2.12).

**Theorem 3.5.** *In the context of Theorem 3.2, the function*

$$\tilde{h}(\beta) = \left( \sqrt{\beta} + \sqrt{1 - \beta} \right)^2 \quad (3.12)$$

*characterizes the phase transition of approximate-exact support recovery problem. Specifically, if  $\underline{r} > \tilde{h}(\beta)$ , then the Benjamini-Hochberg procedure with slowly vanishing nominal FDR levels (as defined in Definition 3.1) achieves asymptotically approximate-exact support recovery in the sense of (2.25).*

*Conversely, if  $\bar{r} < \tilde{h}(\beta)$ , then for any thresholding procedure  $\hat{S}$ , the approximate-exact support recovery fails in the sense of (2.26).*

Theorem 3.5 is proved in Section ?? . The phase transition boundary (3.12) is visualized in Figure 3.2.

### 3.3 Asymptotic power analysis: A discussion

Stilian: We end this chapter with a general discussion on the effect of the signal sparsity on the interplay between signal size and sparsity in the context of the detection and support recovery problems. Specifically, different support recovery criteria are compared. added, check

Theorems 3.2 through 3.5 allow us to asymptotically quantify the required signals sizes in support recovery problems, as well as in the global hypothesis testing problem in the Gaussian additive error model (3.1). Specifically, these results indicate that at all sparsity levels  $\beta \in (0, 1)$ , the difficulties of the problems in terms of the required signal sizes have the following ordering

$$f(\beta) < h(\beta) < \tilde{g}(\beta) < \tilde{h}(\beta) < g(\beta),$$

as previewed in Figure 3.2. The ordering aligns with our intuition that the required signal sizes must increase as we move from detection to support recovery problems. Similarly, more stringent criteria for error control (e.g., FWER compared to FDR) require larger signals. We can now also compare  $\tilde{g}(\beta)$  and  $\tilde{h}(\beta)$ , whose ordering may not be clear from this line of reasoning.

Our last comment is on the gap between FDR and FWER under sparsity assumptions. Although it is believed that FWER control is sometimes too stringent compared to, say, FDR control in support recovery problems, the fact that all three thresholds (detection, weak, and strong classification) involve the same scaling indicates that the difficulties of the three problems (signal detection, approximate, and exact support recovery) are comparable when signals are very sparse, i.e., when  $\beta$  is close to 1. This is illustrated with the next example.

*Example 3.1 (Power analysis for variable selection).* For Gaussian errors (AGG with  $\nu = 2$ ), when  $\beta = 3/4$ , the signal detection boundary (3.4) says that signals will have to be at least of magnitude  $\sqrt{(\log p)/2}$ , while approximate support recovery (3.10) requires signal sizes of at least  $\sqrt{3(\log p)/2}$ , and exact support recovery (3.9) calls for signal sizes of at least  $\sqrt{9(\log p)/2}$ . The required signal sizes increases, but are within the same order of magnitude.

If  $m$  independent copies  $x_1, \dots, x_m$  of the observations were made on the same set of  $p$  locations, then by taking location-wise averages,  $\bar{x}_m(j) = \frac{1}{m} \sum_{i=1}^m x_i(j)$ , we can reduce error standard deviation, and hence boost the signal-to-noise ratio, by a factor of  $\sqrt{m}$ . By the simple calculations above, if  $m$  samples are needed to detect (sparse) signals of a certain magnitude, then  $3m$  samples will enable approximate support recovery with FDR control, and in fact,  $9m$  samples would enable exact support recovery with FWER control.

On the other hand, the gap between FDR and FWER is much larger when signals are dense. For example, if the signals are only *approximately* sparse, i.e., having a few components above (3.9) but many smaller components above (3.10), then FDR-controlling procedures will discover substantially larger proportion of signals than FWER-controlling procedures.

Indeed, as  $\beta \rightarrow 0$ , the required signal size for approximate support recovery (3.10) tends to 0, while the required signal size for exact support recovery (3.9) tends to 4 in the Gaussian error models. While Example 3.1 indicates that the exact support recovery is not much more stringent than approximate support recovery when signals are sparse, the gap between required signal sizes widens when signals are dense.

## Chapter 4

# Exact Support Recovery Under Dependence

We focus on exact support recovery problems in this chapter, and generalize the results we obtained in Chapter 3 to additive error models with much relaxed distributional and dependence assumptions. Sufficient conditions for exact support recovery (4.4) are given in Section 4.2. A very general class of dependence structures characterized by the uniform relative stability concept will be introduced in Section 4.3, to prepare us for the necessary condition in Section 4.4. Section 4.5 discusses the dense signal regime, and Section 4.6 illustrates the phase transition phenomena with numerical examples. We begin with an overview of our results in the context of the existing literature.

### 4.1 Contributions and related work

Consider the additive error model (1.1) with the triangular array of errors,

$$\mathcal{E} = \{(\epsilon_p(i))_{i=1}^p, p = 1, 2, \dots\}, \quad (4.1)$$

where the  $\epsilon_p(i)$ 's have common cumulative distribution function  $F(x) = \mathbb{P}[\epsilon_p(i) \leq x]$ . In contrast to the assumptions in Chapter 3, we only require the errors to have common marginal distributions, and allow them to have potentially arbitrary dependence.

Although our method of analysis applies to all light-tailed error distributions with rapidly varying tails (see Definition 2.7), to be concrete and better convey the main ideas, we will focus on the class of AGG( $\nu$ ) laws (see Definition 2.6). Extensions of the results to other error models are presented in Section B.1.

As in Chapter 3, we assume the signals in model (1.1) to be a sparse vector  $\mu = (\mu(i))_{i=1}^p$  where the sparsity, with a few exceptions which will be explicitly stated, is parametrized as

$$s = |S_p| = \lfloor p^{1-\beta} \rfloor, \quad (4.2)$$

with  $0 < \beta \leq 1$  fixed.

We assume that the non-zero entries of  $\mu$  are positive and take values in the interval  $[\underline{\Delta}, \bar{\Delta}] \subset (0, \infty)$ . That is,  $0 < \underline{\Delta} \leq \mu(i) < \bar{\Delta} \leq +\infty$ , for all  $i \in S_p$ . The lower and upper bound on the signal sizes  $\underline{\Delta}$  and  $\bar{\Delta}$  are parametrized as

$$\underline{\Delta} = \underline{\Delta}(p) = (\nu \underline{r} \log p)^{1/\nu} \quad \text{and} \quad \bar{\Delta} = \bar{\Delta}(p) = (\nu \bar{r} \log p)^{1/\nu}, \quad (4.3)$$

with parameters  $0 < \underline{r} \leq \bar{r} \leq +\infty$ . Notice that the parametrization now depends on the shape of the assumed error distributions  $\text{AGG}(\nu)$  through the parameter  $\nu$ .

According to Lemma 2.1, in order to study the asymptotic behaviors of  $\text{risk}^E$ , it is sufficient to establish minimal conditions under which the support sets can be consistently estimated, i.e.,

$$\mathbb{P}[\widehat{S}_p = S_p] \longrightarrow 1 \quad \text{as } p \rightarrow \infty, \quad (4.4)$$

where  $\widehat{S}_p$  is an estimate of the true signal support set  $S_p$ .

Several authors have studied the support recovery problem in terms of the Hamming loss and obtained minimax optimality results (see, e.g., Ji and Jin (2012); Genovese et al. (2012); Jin et al. (2014); Butucea et al. (2018)). In the special case of Gaussian marginals, Butucea et al. (2018) showed that the boundary (3.9) exists in a minimax sense. That is, when the errors are *independent* Gaussians, the Hamming loss cannot be made to vanish if the signal sizes fall below the boundary (3.9) by any procedure. Conversely, if signal size falls below, the Hamming loss can be made to vanish with a specific thresholding procedure.

Unfortunately, as pointed out in Section 2.4, vanishing Hamming loss is only sufficient, not necessary for support recovery (4.4), and the latter sharp results do not carry over directly to the study of the probability of support recovery or exact support recovery risk. More importantly, despite their elegance, these Hamming loss-minimax studies naturally reduce to the analysis of the elementary case of iid data, and is by design blind to non-trivial error-dependence structures. This prevents us from fully exploring of the phase transition phenomena under other dependence conditions.

So far in the literature, the role of dependence, and that of the distributional assumptions in model (1.1) have remained largely unexplored. This chapter offers advances in both directions, and provides a close-to-complete solution of the exact support recovery problem. (See also, Chapter 6.) We briefly summarize our contributions next.

In Section 4.2, we study exact support recovery in the sense of (4.4) directly, under general distributional and dependence assumptions. In particular, we describe the phase transition phenomena in the dependent AGG model, under the scaling described in (4.2) and (4.3). Consider the function

$$g(\beta) = g_\nu(\beta) = (1 + (1 - \beta)^{1/\nu})^\nu, \quad \nu > 0, \quad (4.5)$$

which we refer to as the *strong classification boundary*. In Theorem 4.1 we show that, if the signal sizes are above the boundary (i.e.,  $r > g(\beta)$ ), the FWER-controlling procedures described in Section 2.2 with appropriately calibrated levels achieve *exact support recovery* as in (4.4).

Conversely, we show in Theorem 4.2, that for a surprisingly large class of dependence structures characterized by the concept of *uniform relative stability* (URS, see Definition 4.1 below), when the signal size is below the boundary (i.e.,  $r < g(\beta)$ ), no thresholding procedure can achieve the asymptotically perfect support recovery (4.4). In fact,

$$\mathbb{P} \left[ \widehat{S}_p = S_p \right] \longrightarrow 0, \quad \text{as } p \rightarrow \infty, \quad (4.6)$$

for all  $\widehat{S}_p$  in the form of (2.20).

These two results show that the thresholding procedures obey a phase transition phenomenon in a strong, *point-wise* sense over the class of URS dependence structures, and over the class of AGG( $\nu$ ),  $\nu > 0$  error distributions. The conclusions are fundamentally stronger and more informative than the minimax statements in the literature (see, e.g., Butucea et al. (2018)). The techniques developed in this here are also entirely different from those in Ji and Jin (2012) or Butucea et al. (2018), and transparent characterizations of the dependence conditions under which the phase transition type result holds will be established in the Chapter 6 later.

The phase transition phenomena for two additional classes of error distributions with either heavier or lighter tails than the AGG distributions will be described in Section B.1.

To conclude this summary, we emphasize that the sharp phase transition results established in this chapter apply only to the general class of thresholding procedures. In Chapter 5, we characterize the finite-sample Bayes optimal support estimation procedures. It will turn out that in many cases the optimal procedures are in fact thresholding procedures. This will lead to complete phase transition results valid for all types of support estimators, for certain classes of error models. In general, however, thresholding procedures can be sub-optimal. This has only recently been noticed by the statistical community in the case when the errors have heavy (regularly-varying) tails. Arias-Castro and Ying (2019) discussed the phenomenon in approximate support recovery problems. In this case, we also demonstrate the absence of a phase transition phenomenon in exact support recovery by thresholding, in Section B.2.

## 4.2 Sufficient conditions for exact support recovery

Following Butucea et al. (2018), we define the parameter space for the signals  $\mu$  as

$$\Theta_p^+(\beta, r) = \{\mu \in \mathbb{R}^p : \text{there exists a set } S_p \subseteq \{1, \dots, p\} \text{ such that } |S_p| \leq \lfloor p^{1-\beta} \rfloor, \\ \mu(i) \geq (\nu r \log p)^{1/\nu} \text{ for all } i \in S_p, \text{ and } \mu(i) = 0 \text{ for all } i \notin S_p\}. \quad (4.7)$$

Our first result states that, when  $F \in \text{AGG}(\nu)$  with  $\nu > 0$ , regardless of the error dependence structure, (asymptotic) perfect support recovery is achieved by applying Bonferroni's procedure with appropriately calibrated FWER, as long as the minimum signal size  $\underline{r}$  is above the strong classification boundary (4.5).

**Theorem 4.1.** *Let the errors have common marginal distribution  $F \in \text{AGG}(\nu)$  with  $\nu > 0$ . Let  $\widehat{S}_p$  be the Bonferroni's procedure (2.21) with vanishing FWER  $\alpha = \alpha(p) \rightarrow 0$ , such that  $\alpha p^\delta \rightarrow \infty$  for every  $\delta > 0$ . If*

$$\underline{r} > g(\beta) = (1 + (1 - \beta)^{1/\nu})^\nu, \quad (4.8)$$

then we have

$$\lim_{p \rightarrow \infty} \sup_{\mu \in \Theta_p^+(\beta, \underline{r})} \mathbb{P}[\widehat{S}_p \neq S_p] = 0. \quad (4.9)$$

**Corollary 4.1 (Classes of procedures attaining the boundary).** *Relation (4.9) holds for any FWER-controlling procedure that is strictly more powerful than Bonferroni's procedure. This includes Holm's procedure (Holm, 1979), and in the case of independent errors, Hochberg's procedure (Hochberg, 1988), and the Šidák procedure (Šidák, 1967).*

*Example 4.1.* Under Gaussian errors, the particular choice of the thresholding at  $t_p = \sqrt{2 \log p}$  in (2.21) corresponds to a Bonferroni's procedure with FWER decreasing at a rate of  $(\log p)^{-1/2}$ , and hence Theorem 4.1 applies. By Corollary 4.1, Holm's procedure — and when the errors are independent, the Šidák, and Hochberg procedures — with FWER controlled at  $(\log p)^{-1/2}$  all achieve perfect support recovery provided that  $\underline{r} > g(\beta)$ .

The claims in Example 4.1 are verified in Section 4.7.

We now turn to the proof of Theorem 4.1.

*Proof (Proof of Theorem 4.1).* Throughout the proof, dependence on  $p$  will be suppressed to simplify notations when such omissions do not lead to ambiguity.

Under the  $\text{AGG}(\nu)$  model, it is easy to see from equation (2.33) that the thresholds in Bonferroni's procedure are

$$t_p = F^{\leftarrow}(1 - \alpha/p) = (\nu \log(p/\alpha))^{1/\nu} (1 + o(1)). \quad (4.10)$$

It is known that Bonferroni's procedure  $\widehat{S}_p = \{i : x(i) > t_p\}$  controls the FWER. Indeed,

$$\begin{aligned} \mathbb{P}[\widehat{S} \subseteq S] &= 1 - \mathbb{P}\left[\max_{i \in S^c} x(i) > t_p\right] = 1 - \mathbb{P}\left[\max_{i \in S^c} \epsilon(i) > t_p\right] \\ &\geq 1 - \sum_{i=1}^p \mathbb{P}[\epsilon(i) > t_p] \geq 1 - \alpha(p) \rightarrow 1, \end{aligned} \quad (4.11)$$



where we used the union bound in the first inequality. Notice that the lower bound (4.11) is independent of the parameter  $\mu$  (as well as the dependence structures), and hence holds uniformly over the parameter space, i.e.,

$$\lim_{p \rightarrow \infty} \inf_{\mu \in \Theta_p^+(\beta, r_-)} P[\widehat{S}_p \subseteq S_p] = 1. \quad (4.12)$$

On the other hand, for the probability of no missed detection, we have:

$$\mathbb{P}[\widehat{S} \supseteq S] = \mathbb{P}\left[\min_{i \in S} x(i) > t_p\right] = \mathbb{P}\left[\min_{i \in S} x(i) - (\nu r_- \log p)^{1/\nu} > t_p - (\nu r_- \log p)^{1/\nu}\right].$$

Since the signal sizes are no smaller than  $(\nu r_- \log p)^{1/\nu}$ , we have

$$x(i) - (\nu r_- \log p)^{1/\nu} \geq \epsilon(i), \quad \text{for all } i \in S,$$

and hence we obtain

$$\mathbb{P}[\widehat{S} \supseteq S] \geq \mathbb{P}\left[\min_{i \in S} \epsilon(i) > (\nu \log(p/\alpha))^{1/\nu}(1 + o(1)) - (\nu r_- \log p)^{1/\nu}\right], \quad (4.13)$$

where we plugged in the expression for  $t_p$  in (4.10). Now, since the minimum signal size is bounded below by  $r_- > (1 + (1 - \beta)^{1/\nu})^\nu$ , we have  $r_-^{1/\nu} - (1 - \beta)^{1/\nu} > 1$ , and so we can pick a  $\delta > 0$  such that

$$\delta < \left(r_-^{1/\nu} - (1 - \beta)^{1/\nu}\right)^\nu - 1. \quad (4.14)$$

Since by assumption, for all  $\delta > 0$ , we have  $p^{-\delta} = o(\alpha(p))$ , there is an  $M = M(\delta)$  such that  $p/\alpha(p) < p^{1+\delta}$  for all  $p \geq M$ . Thus, from (4.13), we further conclude that for  $p \geq M$  we have

$$\begin{aligned} \mathbb{P}[\widehat{S} \supseteq S] &\geq \mathbb{P}\left[\min_{i \in S} \epsilon(i) > ((1 + \delta)\nu \log p)^{1/\nu}(1 + o(1)) - (\nu r_- \log p)^{1/\nu}\right] \\ &= \mathbb{P}\left[\max_{i \in S} (-\epsilon(i)) < \underbrace{\left(r_-^{1/\nu} - (1 + \delta)^{1/\nu}\right)(\nu \log p)^{1/\nu}(1 + o(1))}_{=: A}\right] \\ &\geq 1 - \lfloor p^{1-\beta} \rfloor \times \overline{F}_-(A), \end{aligned} \quad (4.15)$$

where  $\overline{F}_-(x) = \mathbb{P}[-\epsilon(i) > x]$  is the survival function of the  $(-\epsilon(i))$ 's. Notice that (4.15) follows from the union bound and the assumption that  $|S_p| \leq \lfloor p^{1-\beta} \rfloor$ . Therefore, the lower bound does not depend on  $\mu$  (nor on the error dependence structure), and holds uniformly in the parameter space. In turn, we obtain

$$\inf_{\mu \in \Theta_p^+(\beta, r_-)} \mathbb{P}[\widehat{S}_p \supseteq S_p] \geq 1 - \lfloor p^{1-\beta} \rfloor \times \overline{F}_-(A). \quad (4.16)$$

If  $\beta = 1$ , we conclude that the right-hand-side of (4.16) converges to 1, since  $A \rightarrow +\infty$ .

Let now  $\beta \in (0, 1)$  and  $u_p^- := F_-^{\leftarrow}(1 - 1/p)$ . The fact that  $p\bar{F}_-(u_p^-) \leq 1$ , implies

$$\lfloor p^{1-\beta} \rfloor \times \bar{F}_-(A) \leq \frac{\bar{F}_-(B \times u_{\lfloor p^{1-\beta} \rfloor}^-)}{\bar{F}_-(u_{\lfloor p^{1-\beta} \rfloor}^-)} \quad (4.17)$$

where  $B := A/u_{\lfloor p^{1-\beta} \rfloor}^-$ .

Notice that by assumption, the  $-\epsilon(i)$ 's are also  $\text{AGG}(\nu)$  distributed and by Proposition 2.1,  $u_p^- := F_-^{\leftarrow}(1 - 1/p) \sim (\nu \log(p))^{1/\nu}$ , as  $p \rightarrow \infty$ . Therefore, we have

$$u_{\lfloor p^{1-\beta} \rfloor}^- \sim (\nu(1 - \beta) \log p)^{1/\nu} \quad (4.18)$$

and

$$B = \frac{A}{u_{\lfloor p^{1-\beta} \rfloor}^-} = \frac{r^{1/\nu} - (1 + \delta)^{1/\nu}}{(1 - \beta)^{1/\nu}} (1 + o(1)) \rightarrow c > 1$$

as  $p \rightarrow \infty$ , by our choice of  $\delta$  in (4.14).

Finally, since the distribution  $F_-$  has *rapidly varying* tails (by Definition 2.7 and Example 2.1), applying Proposition 2.2, we conclude that (4.17) vanishes. Consequently, the lower bound on the right-hand-side of (4.16) converges to 1. This, combined with (4.12), entails  $\lim_{p \rightarrow \infty} \inf_{\mu \in \Theta_p^+(\beta, \underline{r})} \mathbb{P}[\hat{S}_p = S_p] = 1$ , and hence the desired conclusion (4.9), which completes the proof.  $\square$

The statements in Theorem 4.1 can be strengthened, to prepare us for a minimax result given in Section ?? below.

*Remark 4.1.* In the proof of Theorem 4.1, both (4.11) and (4.15) hold uniformly over all error dependence structures. Therefore, (4.12) and (4.16) may be strengthened to yield

$$\lim_{p \rightarrow \infty} \sup_{\substack{\mu \in \Theta_p^+(\beta, \underline{r}) \\ \mathcal{E} \in D(F)}} P[\hat{S}_p \neq S_p] = 0, \quad (4.19)$$

for  $\underline{r} > g(\beta)$ , where  $D(F)$  is the collection of all arrays with common marginal  $F$ , i.e.,

$$D(F) = \{\mathcal{E} = (\epsilon_p(i))_p : \epsilon_p(i) \sim F \text{ for all } i = 1, \dots, p, \text{ and } p = 1, 2, \dots\}. \quad (4.20)$$

*Remark 4.2.* We emphasize that Theorem 4.1 holds for errors with *arbitrary* dependence structures. Intuitively, this is because the maxima of the errors grow at their fastest in the case of independence. Formally, the light-tailed nature of the error distribution allowed us to obtain sharp tail estimates via simple union bounds, valid under arbitrary dependence.

### 4.3 Dependence and uniform relative stability

An important ingredient needed for a converse of Theorem 4.1 is an appropriate characterization of the error dependence structure under which the strong classification boundary (4.5) is tight. The notion of *uniform relative stability* turns out to be the key.

**Definition 4.1 (Uniform Relative Stability).** Under the notations established in Definition 2.8, the triangular array  $\mathcal{E}$  is said to have uniform relatively stable (URS) maxima if for *every* sequence of subsets  $S_p \subseteq \{1, \dots, p\}$  such that  $|S_p| \rightarrow \infty$ , we have

$$\frac{1}{u_{|S_p|}} M_{S_p} := \frac{1}{u_{|S_p|}} \max_{i \in S_p} \epsilon_p(i) \xrightarrow{\mathbb{P}} 1, \quad (4.21)$$

as  $p \rightarrow \infty$ , where  $u_q$ ,  $q \in \{1, \dots, p\}$  is the generalized quantile in (2.37). The collection of arrays  $\mathcal{E} = \{\epsilon_p(i)\}$  with URS maxima is denoted  $U(F)$ .

Uniform relative stability is, as its name suggests, a stronger requirement on dependence than relative stability (recall Definition 2.8). Proposition 2.2 states that an array with iid components sharing a marginal distribution  $F$  with rapidly varying tails (Definition 2.7) has relatively stable maxima; it is easy to see that URS also follows, by independence of the entries.

**Corollary 4.2.** *An independent array  $\mathcal{E}$  with common marginals  $F \in AGG(\nu)$ ,  $\nu > 0$ , is URS; in this case, URS holds with  $u_{|S_p|} \sim (\nu \log |S_p|)^{1/\nu}$ .*

On the other hand, RS and URS hold under much broader dependence structures than just independent errors. These conditions are extremely mild and can be shown to hold for many classes of error models. In Chapter 6, we will focus extensively on the Gaussian case, which is of great interest in applications and is rather challenging. We will provide simple necessary and sufficient condition for uniform relative stability in terms of the covariance structures.

The relative stability concepts are important because they characterize the dependence structures under which the maxima of error sequences *concentrate* around the quantiles (2.37) in the sense of (2.38). This concentration of maxima phenomena, in turn, is the key to establishing the necessary conditions of the phase transition results in support recovery problems.

### 4.4 Necessary conditions for exact support recovery

With the preparations from Section 4.3, we are ready to state the necessary conditions for exact support recovery (4.4) by thresholding procedures. It turns out that the strong classification boundary (4.5) is tight, under the general dependence structure characterized by URS (Definition 4.1).

Formally, we define the parameter space for the signals  $\mu$  to be

$$\begin{aligned} \Theta_p^-(\beta, \bar{r}) &= \{\mu \in \mathbb{R}^p : \text{there exists a set } S_p \subseteq \{1, \dots, p\} \text{ such that } |S_p| = \lfloor p^{1-\beta} \rfloor, \\ &\quad 0 < \mu(i) \leq (\nu \bar{r} \log p)^{1/\nu} \text{ for all } i \in S_p, \text{ and } \mu(i) = 0 \text{ for all } i \notin S_p\}. \end{aligned} \quad (4.22)$$

**Theorem 4.2.** *Let  $\mathcal{E}$  be a triangular array with common  $AGG(\nu)$  marginal  $F$ ,  $\nu > 0$ . Assume further that the errors  $\mathcal{E}$  have uniform relatively stable maxima and minima, i.e.,  $\mathcal{E} \in U(F)$ , and  $(-\mathcal{E}) = \{-\epsilon_p(i)\} \in U(F)$ . If*

$$\bar{r} < g(\beta) = \left(1 + (1 - \beta)^{1/\nu}\right)^\nu, \quad (4.23)$$

then

$$\lim_{p \rightarrow \infty} \inf_{\hat{S}_p \in \mathcal{T}} \inf_{\mu \in \Theta_p^-(\beta, \bar{r})} \mathbb{P}[\hat{S}_p \neq S_p] = 1, \quad (4.24)$$

where  $\mathcal{T}$  is the class of all thresholding procedures (2.20).

Our first comment is on the signal sizes, and in particular, the gap between the sufficient conditions (Theorem 4.1) and the necessary conditions (Theorem 4.2).

*Remark 4.3.* The sufficient condition in Theorem 4.1 requires that *all* signals be larger than the strong classification boundary  $g(\beta)$  in order to achieve exact support recovery (4.4), while Theorem 4.2 states that exact support recovery fails (in the sense of (4.6)) when *all* signal sizes are below the boundary — the two conditions are *not* complements of each other. This gap between the sufficient and necessary conditions on signal sizes, however, may be difficult to bridge. Indeed, in general, when signal sizes straddle the boundary  $g(\beta)$ , either outcome is possible, as we demonstrate in Example 4.2 below.

*Example 4.2 (Signals straddling the boundary).* Let the signal  $\mu$  have  $|S_p| = \lfloor p^{(1-\beta)} \rfloor$  non-zero entries, composed of two disjoint sets  $S_p = S_p^{(1)} \cup S_p^{(2)}$ . Let also the magnitude of the signals be equal within the two sets, i.e.,  $\mu(i) = \sqrt{2r^{(k)} \log p}$  if  $i \in S_p^{(k)}$  for some constants  $r^{(k)} > 0$  for  $k = 1, 2$ . For simplicity, assume that the errors are iid standard Gaussians.

Consider two scenarios

1.  $r^{(1)} = (1 + \delta)g(\beta)$ ,  $r^{(2)} = (1 + \delta)$  with  $|S_p^{(1)}| = |S_p| - 1$ ,  $|S_p^{(2)}| = 1$ ,
2.  $r^{(1)} = (1 + \delta)g(\beta)$ ,  $r^{(2)} = (1 - \delta)g(\beta)$  with  $|S_p^{(1)}| = \lfloor |S_p|/2 \rfloor$ ,  $|S_p^{(2)}| = |S_p| - |S_p^{(1)}|$ .

for some constants  $0 < \delta < 1 - \beta < 1$ . In both cases, signals in  $S_p^{(1)}$  (respectively,  $S_p^{(2)}$ ) are above (respectively, below) the strong classification boundary (4.5). However, in the first scenario, we have  $\mathbb{P}[\hat{S}_p^{\text{Bonf}} = S_p] \rightarrow 1$  where  $\hat{S}_p^{\text{Bonf}}$  is the Bonferroni's procedure described in Theorem 4.1, while in the second scenario, we have  $\mathbb{P}[\hat{S}_p = S_p] \rightarrow 0$  for *all* thresholding procedures  $\hat{S}_p$ .

The claims in Example 4.2 are verified in Section 4.7.

Our second comment is on the interplay between thresholding procedures and the dependence class characterized by URS.

*Remark 4.4.* Paraphrasing Theorems 4.1 and 4.2: if we consider only thresholding procedures, then for a very large class of dependence structures, we cannot improve upon the Bonferroni procedure  $\widehat{S}_p^{\text{Bonf}}$ . Specifically, for all  $\mathcal{E} \in U(F)$  and  $-\mathcal{E} \in U(F)$ , and for all  $S_p \in \mathcal{S}$ , where  $\mathcal{S} = \{S \subseteq \{1, \dots, p\}; |S| = \lfloor p^{1-\beta} \rfloor\}$ , we have

$$\lim_{p \rightarrow \infty} \mathbb{P}[\widehat{S}_p^{\text{Bonf}} \neq S_p] = \begin{cases} \limsup_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \mathbb{P}[\widehat{S}_p \neq S_p] = 0, & \text{if } \underline{r} > g(\beta), \\ \liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \mathbb{P}[\widehat{S}_p \neq S_p] = 1, & \text{if } \bar{r} < g(\beta) \end{cases} \quad (4.25)$$

where  $\mathcal{T}$  is the set of all thresholding procedures (2.20).

Theorem 4.2 also yields an answer the question raised in Butucea et al. (2018). In particular, the authors of (Butucea et al., 2018) commented that independent error is the ‘least favorable model’ in the problem of support recovery, and conjectured that the support recovery problem may be easier to solve under dependence, similar to how the problem of signal detection is easier under dependent errors (see Hall and Jin (2010)). Surprisingly, our results here state that asymptotically, *all* error dependence structures in the large URS class are equally difficult for *thresholding procedures*. Therefore, the phase transition behavior is universal in the class of dependence structures characterized by URS.

To facilitate comparison with results in existing literature, we will formulate explicit minimax statements in Section ??.

We must emphasize that the restriction to the URS dependence class is *not* an assumption of convenience. The condition on dependence characterized by uniform relative stability is, in fact, the weakest of its kind in the literature. We will characterize the class URS dependence class in Chapter 6 below.

We conclude with the proof of Theorem 4.2.

*Proof (Proof of Theorem 4.2).* To avoid cumbersome double subscript notations, we will sometimes suppress dependence on  $p$  of the set sequences  $\widehat{S}_p$  and  $S_p$  in the proof.

Since the estimator  $\widehat{S}_p = \{x(i) \geq t_p(x)\}$  is thresholding, exact support recovery takes place if and only if the threshold separates the signals and null part, i.e.,

$$\mathbb{P}[\widehat{S}_p = S_p] = \mathbb{P}\left[\max_{i \in S^c} x(i) < t_p(x) \leq \min_{i \in S} x(i)\right] \leq \mathbb{P}\left[\max_{i \in S^c} x(i) < \min_{i \in S} x(i)\right].$$

Since the right-hand-side does not depend on the procedure  $\widehat{S}_p$ , we also have

$$\sup_{\widehat{S}_p \in \mathcal{T}} \mathbb{P}[\widehat{S}_p = S_p] \leq \mathbb{P}\left[\max_{i \in S^c} x(i) < \min_{i \in S} x(i)\right] \leq \mathbb{P}\left[\max_{i \in S^c} \epsilon(i) < \bar{\Delta} + \min_{i \in S} \epsilon(i)\right], \quad (4.26)$$

where we used the assumption that the signal sizes are no greater than  $\bar{\Delta}$ . Let  $S^* = S_p^*$  be a sequence of support sets that maximize the right-hand-side of (4.26), i.e., let

$$S_p^* = \arg \max_{S \subseteq \{1, \dots, p\}; |S| = \lfloor p^{1-\beta} \rfloor} \mathbb{P}\left[\max_{i \in S^c} \epsilon(i) < \bar{\Delta} + \min_{i \in S} \epsilon(i)\right],$$

where ties can be broken lexicographically if multiple maximizers exist. Then, we obtain the following bound which only depends on  $\bar{r}$  and the distribution of  $\mathcal{E}$ ,

$$\begin{aligned} \sup_{\hat{S}_p \in \mathcal{T}} \sup_{\mu \in \Theta_p^-(\beta, \bar{r})} \mathbb{P}[\hat{S}_p = S_p] &\leq \mathbb{P} \left[ \max_{i \in S^{*c}} \epsilon(i) < \bar{\Delta} + \min_{i \in S^*} \epsilon(i) \right] \\ &= \mathbb{P} \left[ \frac{M_{S^{*c}}}{u_p} < \frac{\bar{\Delta} - m_{S^*}}{u_p} \right], \end{aligned} \quad (4.27)$$

where  $M_{S^{*c}} = \max_{i \in S^{*c}} \epsilon(i)$  and  $m_{S^*} = \max_{i \in S^*} (-\epsilon(i))$ . Since the error arrays  $\mathcal{E}$  and  $(-\mathcal{E})$  are URS by assumption, using the expression for the AGG quantiles (2.33), we have

$$\frac{M_{S^{*c}}}{u_p} = \frac{M_{S^{*c}}}{u_{|S^{*c}|}} \frac{u_{|S^{*c}|}}{u_p} \xrightarrow{\mathbb{P}} 1, \quad \text{and} \quad \frac{m_{S^*}}{u_p} = \frac{m_{S^*}}{u_{|S^*|}} \frac{u_{|S^*|}}{u_p} \xrightarrow{\mathbb{P}} (1 - \beta)^{1/\nu}, \quad (4.28)$$

so that the two random terms in probability (4.27) converge to constants. Notice that the second relation in (4.28) holds by URS for any  $\beta \in (0, 1)$ ; when  $\beta = 1$ , the relation holds since  $u_{|S^*|}/u_p$  vanishes while  $\{m_{S^*}/u_{|S^*|}\}$  is tight.

Since signal sizes are bounded above by  $\bar{r} < (1 + (1 - \beta)^{1/\nu})^\nu$ , we can write  $\bar{r}^{1/\nu} = 1 + (1 - \beta)^{1/\nu} - d$  for some  $d > 0$ . By our parametrization of  $\bar{\Delta}$ , we have

$$\frac{\bar{\Delta}}{u_p} = \left(1 + (1 - \beta)^{1/\nu} - d\right) (1 + o(1)). \quad (4.29)$$

Combining (4.28) and (4.29), we conclude that the right-hand-side of the probability (4.27) converges in probability to a constant strictly less than 1, that is,

$$\frac{\bar{\Delta} - m_{S^*}}{u_p} \xrightarrow{\mathbb{P}} 1 - d, \quad (4.30)$$

while  $M_{S^{*c}}/u_p \xrightarrow{\mathbb{P}} 1$ . Therefore, the probability in (4.27) must go to 0.  $\square$

## 4.5 Dense signals

We treat briefly the case of dense signals, where the size of the support set is proportional to the problem dimension, i.e.  $s \sim cp$  for some constant  $c \in (0, 1)$ . We show that in this case, a phase-transition-type result still holds, independently of the value of  $c$ . Analogous to the set-up of Theorems 4.1 and 4.2, let

$$\begin{aligned} \Theta_p^{\text{d}+}(c, \underline{r}) &= \{\mu \in \mathbb{R}^p : \text{there exists a set } S_p \subseteq \{1, \dots, p\} \text{ such that } |S_p| \leq \lfloor cp \rfloor, \\ &\quad \mu(i) \geq (\nu \underline{r} \log p)^{1/\nu} \text{ for all } i \in S_p, \text{ and } \mu(i) = 0 \text{ for all } i \notin S_p\}, \end{aligned} \quad (4.31)$$

where “d” in the notation  $\Theta_p^{d+}$  stands for “dense”. Similarly, define

$$\begin{aligned} \Theta_p^{d-}(c, \bar{r}) = \{ \mu \in \mathbb{R}^p : \text{there exists a set } S_p \subseteq \{1, \dots, p\} \text{ such that } |S_p| = \lfloor cp \rfloor, \\ 0 < \mu(i) \leq (\nu \bar{r} \log p)^{1/\nu} \text{ for all } i \in S_p, \text{ and } \mu(i) = 0 \text{ for all } i \notin S_p \}. \end{aligned} \quad (4.32)$$

**Theorem 4.3.** *Let  $c \in (0, 1)$  be a fixed constant, and let  $\widehat{S} = \widehat{S}_p^{\text{Bonf}}$  denote the Bonferroni’s procedure as described in Theorem 4.1. In the context of Theorem 4.1, if  $\underline{r} > 1$ , then we have*

$$\lim_{p \rightarrow \infty} \sup_{\mu \in \Theta_p^{d+}(c, \underline{r})} \mathbb{P}[\widehat{S}_p \neq S_p] = 0. \quad (4.33)$$

While in the context of Theorem 4.2, if  $\bar{r} < 1$ , then

$$\lim_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \inf_{\mu \in \Theta_p^{d-}(c, \bar{r})} \mathbb{P}[\widehat{S}_p \neq S_p] = 1, \quad (4.34)$$

where  $\mathcal{T}$  is the class of all thresholding procedures (2.20).

*Remark 4.5.* Notice that the boundary for the signal size parameter is identically 1 in this dense regime. Therefore, if we interpret  $\beta = 0$  of the parametrization (4.2) as  $s \sim cp$ , where  $c \in (0, 1)$ , then the strong classification boundary (4.5) may be continuously extended to the left-end point where  $g(0) = 1$ .

*Proof (Proof of Theorem 4.3).* The proof is entirely analogous to that of Theorems 4.1 and 4.2. Specifically, (4.33) follows by replacing  $\lfloor p^{1-\beta} \rfloor$  with  $\lfloor cp \rfloor$  in Relation (4.15) onward, and replacing (4.18) with

$$u_s^- \sim (\nu \log cp)^{1/\nu} \sim (\nu \log p)^{1/\nu}.$$

in the proof of Theorem 4.1. Similarly, (4.34) follows the proof of Theorem 4.2. Indeed, by using the fact that

$$\frac{u_{|S^{*c}|}}{u_p} \sim \frac{(\nu \log(1-c)p)^{1/\nu}}{(\nu \log p)^{1/\nu}} \rightarrow 1$$

and  $u_{|S^*|}/u_p \rightarrow 1$  for all  $c \in (0, 1)$ , we see that Relation (4.28) holds with  $\beta = 0$ , and the rest of Theorem 4.2 applies.  $\square$

## 4.6 Numerical illustrations for independent errors

We examine numerically the boundaries (4.5) under several error tail assumptions for independence errors in this section. Numerical experiments for dependent errors will be deferred until we characterize the URS conditions in Chapter 6.

To demonstrate the phase transition phenomenon under different error tail densities, we simulate from the additive error model (1.1) with

- Gaussian errors, where the density is given by  $f(x) = \frac{1}{\sqrt{2\pi}} \exp\{-x^2/2\}$ .
- Laplace errors, where the density is given by  $f(x) = \frac{1}{2} \exp\{-|x|\}$ .
- Generalized Gaussian  $\nu = 1/2$ , with density  $f(x) = \frac{1}{2} \exp\{-2|x|^{1/2}\}$ .

The sparsity and signal size of the sparse mean vector are parametrized as in equations (4.2) and (4.3), respectively. The support set  $S$  is estimated with  $\tilde{S} = \{i : x(i) > \sqrt{2 \log p}\}$  under the Gaussian errors,  $\tilde{S} = \{i : x(i) > \log p + (\log \log p)/2\}$  under the Laplace errors, and with  $\tilde{S} = \{i : x(i) > \frac{1}{4} (W(-c/(ep \log p)) + 1)^2\}$  under the generalized Gaussian ( $\nu = 1/2$ ) errors. Here  $W$  is the Lambert  $W$  function, i.e.,  $W = f^{-1}$  where  $f(x) = x \exp(x)$ . The choices of thresholds correspond to Bonferroni's procedures with FWER decreasing at a rate of  $1/\sqrt{\log p}$ , therefore satisfying the assumptions in Theorem 4.1. Experiments were repeated 1000 times under each sparsity-and-signal-size combination.

The results of the numerical experiments are shown in Figure 4.1. The numerical results illustrate that the predicted boundaries are not only accurate in high-dimensions ( $p = 10000$ , right panels of Figure 4.1), but also practically meaningful even at moderate dimensions ( $p = 100$ , left panels of Figure 4.1).

We establish in this section minimax versions of our results from Section ???. Specifically, if we restrict ourselves to *the class of thresholding procedures*  $\mathcal{T}$  (defined in (2.20)), then Bonferroni's procedure is minimax optimal, for *any* fixed dependence structures in the URS class. This is formalized in Corollary 5.2 in Section 5.2.1. We refer to this result as *point-wise* minimax, to emphasize the fact that this optimality holds for every *fixed* URS array.

Meanwhile, if we search over *all procedures*, but expand the space of models to include *all* dependence structures, then a different minimax optimality statement holds for Bonferroni's procedure. This result, formally stated in Section 5.2.2, is a consequence of our characterization of the finite-sample Bayes optimality of thresholding procedures in Sections 5.1 and 5.1.1.

Finally, we offer some insights into the support recovery problem in the case when errors have heavier-than-exponential tails in Section 5.1.2.

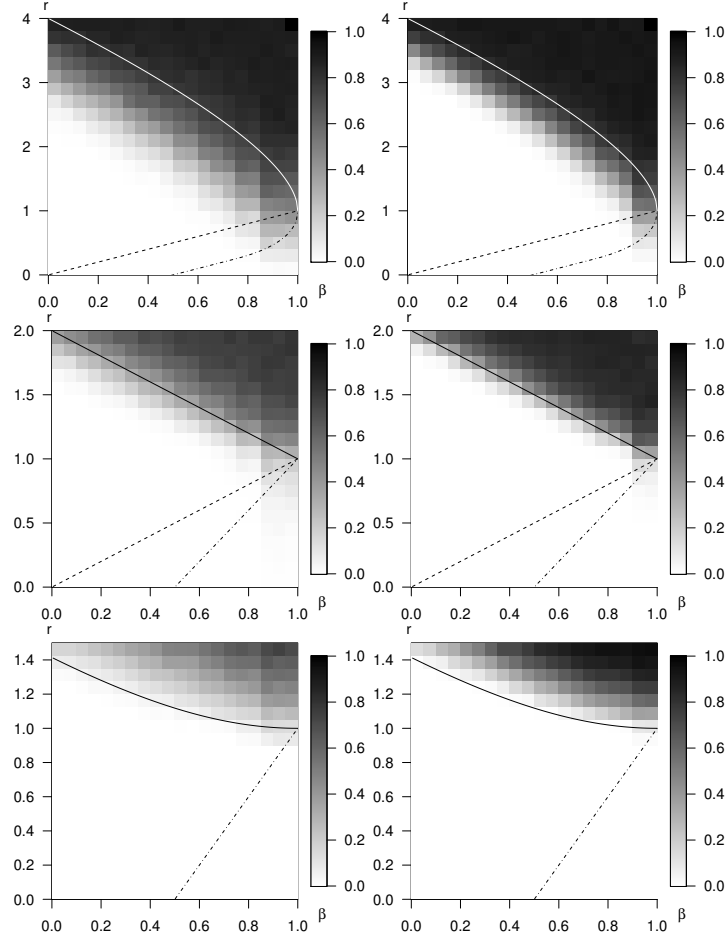
## 4.7 Additional proofs

### 4.7.1 Proof of the claims in Examples 4.1 and 4.2

*Proof (Proof of claims in Example 4.1).* By the Mill's ratio for the standard Gaussian distribution,

$$\frac{t_p \mathbb{P}[Z > t_p]}{\phi(t_p)} \rightarrow 1, \quad \text{as } t_p \rightarrow \infty,$$





**Fig. 4.1** The empirical probability of exact support recovery from numerical experiments, as a function of sparsity level  $\beta$  and signal sizes  $r$ , from Gaussian error models (upper panels), Laplace error models (middle panels), and generalized Gaussian with  $\nu = 1/2$  (lower panels); darker color indicates higher probability of exact support recovery. The experiments were repeated 1000 times for each sparsity-signal size combination, and for dimensions  $p = 100$  (left panels) and  $p = 10000$  (right panels). Numerical results agree with the boundaries described in Theorem 4.1; convergence is noticeably slower for under generalized Gaussian ( $\nu = 1/2$ ) errors. For reference, the dashed and dash-dotted lines represent the weak classification and detection boundaries (see Chapter 3).

where  $Z \sim N(0, 1)$ . Using the expression for  $t_p = \sqrt{2 \log p}$ , we have

$$p \mathbb{P}[Z > t_p] \sim \sqrt{2\pi}^{-1} (2 \log p)^{-1/2} \rightarrow 0,$$

as desired. The rest of the claims follow from Corollary 4.1.  $\square$

*Proof (Proof of claims in Example 4.2).* In the first scenario, signal sizes in  $S_p^{(1)}$  are by definition above the strong classification boundary (4.5). The signal in  $S_p^{(2)}$  has size parameter  $1 + \delta < 2 - \beta < (1 + \sqrt{1 - \beta})^2$ , and therefore falls below the boundary.

It remains to show that  $\mathbb{P}[\widehat{S}_p^{\text{Bonf}} = S_p] \rightarrow 1$ . To do so, we define two new arrays

$$\mathcal{Y}^{(k)} = \{y_p^{(k)}(j), j = 1, 2, \dots, p\}, \quad k \in \{1, 2\}_p,$$

where  $y_p^{(k)}(j) = x_p(j)$  if  $j \notin S_p^{(k)}$ , and  $y_p^{(k)}(j) = \tilde{\epsilon}_p(j)$  if  $j \in S_p^{(k)}$ , using an independent error array  $\{\tilde{\epsilon}_p(j), j = 1, \dots, p\}$  with iid standard Gaussian elements. That is, we replace the elements in  $S_p^{(1)}$  and  $S_p^{(2)}$  with iid standard Gaussian noise. Notice both arrays  $\mathcal{Y}^{(1)}$  and  $\mathcal{Y}^{(2)}$  satisfy the conditions in Theorem 4.1 (with sparsity parameter equal to  $\beta$  and 1, respectively). Hence, we have

$$\mathbb{P}[\widehat{S}_p^{\text{Bonf}} \subseteq S_p] = \mathbb{P}\left[\max_{j \in S^c} x(j) \leq t_p\right] \leq \mathbb{P}\left[\max_{j \in S^c} y^{(1)}(j) \leq t_p\right] \rightarrow 0,$$

and

$$\begin{aligned} \mathbb{P}[\widehat{S}_p^{\text{Bonf}} \supseteq S_p] &= \mathbb{P}\left[\min_{j \in S} x(j) > t_p\right] \geq 1 - \mathbb{P}\left[\min_{j \in S^{(1)}} x(j) \leq t_p\right] - \mathbb{P}\left[\min_{j \in S^{(2)}} x(j) \leq t_p\right] \\ &\geq 1 - \mathbb{P}\left[\min_{j \in S^{(1)}} y_p^{(2)}(j) \leq t_p\right] - \mathbb{P}\left[\min_{j \in S^{(2)}} y_p^{(1)}(j) \leq t_p\right] \rightarrow 1, \end{aligned}$$

where  $t_p$  is the threshold in Bonferroni's procedure. The conclusion follows.

In the second scenario, the signal sizes in  $S^{(2)}$  by definition falls below the strong classification boundary (4.5). To see that no thresholding procedure succeeds, we adapt the proof of Theorem 4.2. In particular, we obtain

$$\mathbb{P}[\widehat{S}_p = S_p] \leq \mathbb{P}\left[\max_{j \in S^c} x(j) \leq t_p < \min_{j \in S} x(j)\right] \leq \mathbb{P}\left[\max_{j \in S^c} x(j) < \min_{j \in S^{(2)}} x(j)\right].$$

By the assumption that signals in  $S^{(2)}$  have size parameter  $(1 - \delta)g(\beta)$ , we have

$$\mathbb{P}\left[\max_{j \in S^c} x(j) < \min_{j \in S^{(2)}} x(j)\right] = \mathbb{P}\left[\frac{M_{S^c}}{u_p} < \frac{\sqrt{2(1 - \delta)g(\beta) \log p} - m_{S^{(2)}}}{u_p}\right], \quad (4.35)$$

where  $M_{S^c} = \max_{j \in S^c} \epsilon(j)$  and  $m_{S^{(2)}} = \max_{j \in S^{(2)}} (-\epsilon(j))$ . The ratio on the left-hand-side of the inequality converges to 1 as in (4.28) in the main text, whereas the term on the right-hand-side

$$\begin{aligned} \frac{\sqrt{2(1-\delta)g(\beta)\log p - m_{S^{(2)}}}}{u_p} &= \sqrt{(1-\delta)g(\beta)} - \frac{m_{S^{(2)}}}{u_{|S^{(2)}|}} \frac{u_{|S^{(2)}|}}{u_p} \\ &\xrightarrow{\mathbb{P}} \sqrt{(1-\delta)} + \sqrt{1-\beta}(\sqrt{(1-\delta)} - 1) < 1. \end{aligned}$$

where we used the URS of the error arrays, and that

$$u_{|S^{(2)}|} \sim \sqrt{2 \log(p^{1-\beta}/2)} = \sqrt{2((1-\beta)\log p - \log 2)} \sim \sqrt{2(1-\beta)\log p}.$$

to conclude the convergence in probability.  $\square$



## Chapter 5

### Bayes and minimax optimality

Finally, we offer some insights into the support recovery problem in the case when errors have heavier-than-exponential tails in Section 5.1.2.

#### 5.1 Bayes optimality in support recovery problems

In studying support recovery problems (e.g., Arias-Castro and Chen (2017)), restrictions to the thresholding procedures are sometimes justified by arguing that such procedures are the “reasonable” choice for estimating the support set. We show in this section that, perhaps surprisingly, for general error models, thresholding procedures are not always optimal, even when the observations are independent.

We shall identify the optimal procedure for support recovery problems under a Bayesian setting with general distributional assumptions (including but not limited to additive models (1.1)). Specifically, we assume that there is an ordered set  $P = (i_1, \dots, i_s)$ ,  $i_i \in \{1, \dots, p\}$ , and  $s$  (not necessarily equal) densities  $f_1, \dots, f_s$ , such that the observations indexed by  $P$  have corresponding densities. That is,

$$x(i_j) \sim f_j, \quad j = 1, \dots, s. \quad (5.1)$$

Let also the rest  $(p - s)$  observations have common density  $f_0$ , i.e.,  $x(i) \sim f_0$  for  $i \notin S$ . We further assume that the observations  $x$  are mutually independent.

We adopt here a Bayesian framework to measure statistical risks. Let the ordered support  $P = (i_1, \dots, i_s)$  have prior

$$\pi((i_1, \dots, i_s)) = (p - s)!/p!, \quad (5.2)$$

for all distinct  $1 \leq i_1, \dots, i_s \leq p$ . Consequently, the unordered support  $S = \{i_1, \dots, i_s\}$  is distributed uniformly in the collection of all set of size  $s$ , with the unordered uniform distribution  $\pi^u$ . That is, for all  $S \in \mathcal{S} := \{S \subseteq \{1, \dots, p\}; |S| = s\}$ , we have

$$\pi^u(\{i_1, \dots, i_s\}) = \sum_{\sigma} \pi((i_{\sigma(1)}, \dots, i_{\sigma(s)})) = (p-s)!s!/p!, \quad (5.3)$$

where the sum is taken over all permutations of  $\{1, 2, \dots, s\}$ .

For any fixed configuration  $P$ , consider the loss function,

$$\ell(\widehat{S}, S) := \mathbb{P}[\widehat{S} \neq S] = \mathbb{P}_P[\widehat{S} \neq S],$$

where the probability is taken over the randomness in the observations only. The Bayes optimal procedures should minimize

$$\mathbb{E}_{\pi} \mathbb{P}[\widehat{S} \neq S], \quad (5.4)$$

where the expectation is taken over the random configurations  $P$ , with a uniform distribution  $\pi$  as specified in (5.2).

If, however, the sparsity  $s = |S|$  of the problem is known, then a natural estimator for  $S$  would be based on the set of top  $s$  order statistics.

Stilian: COMMENT Just like for the likelihood thresholding procedure, we should be more careful about the ties and how to define the oracle thresholding procedure so that  $\widehat{S}^*$  has exactly  $s$  elements. I propose the following change

Stilian:

For any collection of numbers  $\{a_i, i = 1, \dots, s\}$ , let

$$\langle a_1, \dots, a_s \rangle := (a_{[1]}, \dots, a_{[s]})$$

denote the vector of  $a_i$ 's arranged in a non-increasing order.

**Definition 5.1 (Oracle data thresholding).** Let  $x_{[1]} \geq \dots \geq x_{[p]}$  be the order statistics of the data vector  $x$ . Any estimator  $\widehat{S}^* := \{i_1, \dots, i_s\}$ , where

$$\langle x(i_1), \dots, x(i_s) \rangle = (x_{[1]}, \dots, x_{[s]})$$

will be referred to as an *oracle thresholding estimator*.

Simply put, the oracle thresholding estimators are comprised of the indices corresponding to the  $s$  largest values in the data. Note that, in the absence of ties among the largest  $s+1$  data values, the oracle thresholding estimator is unique. For concreteness, one can break possible ties lexicographically. In many cases, the oracle thresholding estimators will be almost surely unique. IMPORTANT, replacing the following

**Definition 5.2 (Oracle data thresholding).** We call  $\widehat{S}^* = \{i \mid x(i) \geq x_{[s]}\}$  the oracle data thresholding procedure, where  $x_{[1]} \geq \dots \geq x_{[p]}$  are the order statistics of  $x$ .

### 5.1.1 When is oracle thresholding Bayes optimal?

An elegant answer of this question can be given in terms of the *monotone likelihood ratio* (MLR) property.

**Definition 5.3 (Monotone Likelihood Ratio).** A family of positive densities on  $\mathbb{R}$ ,  $\{f_\delta, \delta \in U\}$ , is said to have the MLR property if, for all  $\delta_0, \delta_1 \in U \subseteq \mathbb{R}$  such that  $\delta_0 < \delta_1$ , the likelihood ratio  $(f_{\delta_1}(x)/f_{\delta_0}(x))$  is an increasing function of  $x$ .

The following result provides a general criterion for the finite-sample optimality of the oracle thresholding procedure  $\widehat{S}^*$ .

**Proposition 5.1.** *Let the observations  $x(i)$ ,  $i = 1, \dots, p$  be as prescribed as in (5.1) through (5.2). If each of the pairs  $\{f_0, f_1\}, \dots, \{f_0, f_s\}$  forms an MLR family, then every oracle data thresholding procedure  $\widehat{S}^*$  is finite-sample optimal in terms of Bayes risk  $\mathbb{E}_\pi \mathbb{P}[\widehat{S} \neq S]$ . That is,*

$$\widehat{S}^* \in \arg \min_{\widehat{S}} \mathbb{E}_\pi \mathbb{P}[\widehat{S} \neq S]. \quad (5.5)$$

for all  $s$  and  $p$ .

*Proof.* The problem of support recovery can be equivalently stated as a classification problem, where the discrete parameter space is  $\mathcal{S} = \{S \subseteq \{1, \dots, p\} : |S| = s\}$ , and the observation  $x \in \mathbb{R}^p$  has likelihood  $f(x|S)$  indexed by the support set  $S$ .

By the optimality of the Bayes classifier (see, e.g., Domingos and Pazzani, 1997), a set estimator that maximizes the probability of support recovery is such that

$$\widehat{S} \in \arg \max_{S \in \mathcal{S}} f(x|S)\pi(S).$$

Since we know from (5.3) that  $\pi(\cdot)$  is uniform, the problem in our context reduces to showing that  $f(x|\widehat{S}^*) = f(x|\widehat{S})$ , where  $f(x|S)$  is the conditional distribution of data given the unordered support  $S$ ,

$$f(x|S) = \sum_{P \in \sigma(S)} f(x|P)\pi^{\text{ord}}(P|S) = \frac{1}{s!} \left( \sum_{P \in \sigma(S)} \prod_{i=1}^s f_i(x(P(i))) \right) \prod_{k \notin S} f_0(x(k)),$$

where  $\sigma(S)$  is the set of all permutations of the indices in the support set  $S$ .

Suppose that  $\widehat{S}$  is *not* an oracle thresholding estimator, then there must be indices  $j \in \widehat{S}$  and  $j' \notin \widehat{S}$  such that  $x(j) < x(j')$ . We exchange the labels of  $x(j)$  and  $x(j')$ , and form a new estimate  $\widehat{S}' = (\widehat{S} \setminus \{j\}) \cup \{j'\}$ . Comparing the likelihoods under  $\widehat{S}$  and  $\widehat{S}'$ , we have

$$\begin{aligned}
f(x|\widehat{S}) - f(x|\widehat{S}') &= \frac{1}{s!} \sum_{P \in \sigma(\widehat{S})} \prod_{i=1}^s f_i(x(P(i))) f_0(x(j')) \prod_{k \notin \widehat{S} \cup \{j'\}} f_0(x(k)) - \\
&\quad - \frac{1}{s!} \sum_{P' \in \sigma(\widehat{S}')} \prod_{i=1}^s f_i(x(P'(i))) f_0(x(j)) \prod_{k \notin \widehat{S}' \cup \{j\}} f_0(x(k)) \\
&= \frac{1}{s!} \left( \sum_{i=1}^s a_i \left( f_i(x(j)) f_0(x(j')) - f_i(x(j')) f_0(x(j)) \right) \right) \prod_{k \notin \widehat{S} \cup \{j'\}} f_0(x(k)),
\end{aligned} \tag{5.6}$$

where the last equality follows by first summing over all permutations fixing  $P(i) = j$  and  $P'(i) = j'$ , and setting  $a_i = \sum_{P \in \sigma(\widehat{S} \setminus \{j\})} \prod_{i' \neq i} f_{i'}(x(P(i')))$ . Notice that the  $a_i$ 's are non-negative.

Since  $x(j) < x(j')$ , and since each of  $\{f_0, f_i\}$  is an MLR family, we have

$$\frac{f_i(x(j))}{f_0(x(j))} - \frac{f_i(x(j'))}{f_0(x(j'))} \leq 0 \implies f_i(x(j)) f_0(x(j')) - f_i(x(j')) f_0(x(j)) \leq 0.$$

Using Relation (5.6), we conclude that  $f(x|\widehat{S}) \leq f(x|\widehat{S}')$ . [Stilian: Continuing this way, we can successively improve the likelihood of every estimator until we arrive at an oracle thresholding estimator, proving the desired optimality. Note that with the same argument, any two oracle thresholding estimators have the same likelihood. replacing This implies that any estimator that is not  $\widehat{S}^*$  may be improved, and the optimality follows. □

We emphasize that under the MLR conditions in Proposition 5.1, the oracle thresholding procedures are in fact *finite-sample optimal* in the above Bayesian context. Further, our setup allows for different alternative distributions, and relaxes the assumptions of Butucea et al. (2018) when studying distributional generalizations, where the alternatives are assumed to be identically distributed.

It remains to understand when the key MLR property holds. We elaborate on this question next. Returning to the more concrete signal-plus-noise model (1.1), it turns out that the error tail behavior is what determines the optimality of data thresholding procedures. In this setting, log-concavity of the error densities is *equivalent* to the MLR property (Lemma 5.1). This, in turn, yields the finite-sample optimality of data thresholding procedures (Proposition 5.2).

**Lemma 5.1.** *Let  $\delta$  be the magnitude of the non-zero signals in the signal-plus-noise model (1.1) with positive error density  $f_0$ , and let  $f_\delta(x) = f_0(x - \delta)$ . The family  $\{f_\delta, \delta \in \mathbb{R}\}$  has the MLR property if and only if the error density  $f_0$  is log-concave.*

*Proof.* Suppose MLR holds, we will show that  $f_0(t) = \exp\{\phi(t)\}$  for some concave function  $\phi$ . By the assumption of MLR, for any  $x_1 < x_2$ , setting  $\delta_0 = 0$ , and  $\delta_1 = (x_2 - x_1)/2 > 0$ , we have

$$\log \frac{f_{\delta_1}(x_2)}{f_{\delta_0}(x_2)} = \phi\left(\frac{(x_1 + x_2)}{2}\right) - \phi(x_2) \geq \phi(x_1) - \phi\left(\frac{(x_1 + x_2)}{2}\right) = \log \frac{f_{\delta_1}(x_1)}{f_{\delta_0}(x_1)}.$$



This implies that the log-density  $\phi(t)$  is midpoint-concave, i.e., for all  $x_1$  and  $x_2$ , we have,

$$\phi\left(\frac{x_1 + x_2}{2}\right) \geq \frac{1}{2}\phi(x_1) + \frac{1}{2}\phi(x_2). \quad (5.7)$$

For Lebesgue measurable functions, midpoint concavity is equivalent to concavity by the Sierpinski Theorem (see, e.g., Sec I.3 of Donoghue, 2014). This proves the ‘only-if’ part.

For the ‘if’ part, when  $\phi(t) = \log(f_0(t))$  is log-concave, then for any  $\delta_0 < \delta_1$ , and any  $x < y$ , we have

$$\log \frac{f_{\delta_1}(y)}{f_{\delta_0}(y)} - \log \frac{f_{\delta_1}(x)}{f_{\delta_0}(x)} = \phi(y - \delta_1) - \phi(y - \delta_0) - \phi(x - \delta_1) + \phi(x - \delta_0) \geq 0, \quad (5.8)$$

where the last inequality is a simple consequence of concavity (see Lemma 5.2 below). This proves the ‘if’ part.  $\square$

**Lemma 5.2.** *Let  $\phi$  be any concave function on  $\mathbb{R}$ . For any  $x < y \in \mathbb{R}$ , and  $\delta > 0$  we have*

$$\phi(x) + \phi(y + \delta) \leq \phi(y) + \phi(x + \delta).$$

*Proof.* Pick  $\lambda = \delta/(y - x + \delta)$ , by concavity of  $f$  we have

$$\lambda\phi(x) + (1 - \lambda)\phi(y + \delta) \leq \phi(\lambda x + (1 - \lambda)(y + \delta)) = \phi(y), \quad (5.9)$$

and

$$(1 - \lambda)\phi(x) + \lambda\phi(y + \delta) \leq \phi((1 - \lambda)x + \lambda(y + \delta)) = \phi(x + \delta). \quad (5.10)$$

Summing up (5.9) and (5.10) and we arrive at the conclusion as desired.  $\square$

Proposition 5.1 and Lemma 5.1 yield immediately the following.

**Proposition 5.2.** *Consider the additive error model (1.1), where the  $\epsilon(i)$ ’s are independent with common distribution  $F$ . Let the signal  $\mu$  have  $s$  positive entries with magnitudes  $0 < \delta_1 \leq \dots \leq \delta_s$ , located on  $\{1, \dots, p\}$  as prescribed in (5.2).*

*If  $F$  has a positive, log-concave density  $f$ , then the support estimator*

$$\widehat{S}^* := \{i : x(i) \geq x_{[s]}\}$$

*is finite-sample optimal in terms of Bayes risk in the sense of (5.5).*

*Proof.* The independence and the fact that the observations have densities implies the absence of ties among the order statistics  $\{x_{[i]}\}$ , with probability one. Thus, the oracle thresholding procedure is a.s. unique and given by  $\widehat{S}^* = \{i : x(i) \geq x_{[s]}\}$ . The result then follows from Proposition 5.1 and Lemma 5.1.  $\square$

**Remark 5.1.** Propositions 5.1 and 5.2 show that under MLR (or equivalently, log-concavity of the errors in additive models), the oracle thresholding procedures are

finite-sample optimal even in the case where the signals have different (positive) sizes. This fascinating property perhaps explains the success of the thresholding estimators.

The assumption of log-concavity of the densities is compatible with the AGG model when  $\nu \geq 1$ , as demonstrated in the next example.

*Example 5.1.* The generalized Gaussian density  $f(x) \propto \exp\{-|x|^\nu/\nu\}$  is log-concave for all  $\nu \geq 1$ . Therefore in the additive error model (1.1), according to Proposition 5.2, the oracle thresholding procedure is Bayes optimal in the sense of (5.5).

Stilian: COMMENT Consider moving the following discussion & theorem to the section on Minimax below...

Consider the asymptotic Bayes risk as defined in (5.4). The statement for the necessary condition of support recovery, with the help of Proposition 5.2, can be strengthened to include all procedures (in the Bayesian context), regardless of whether they are thresholding.

**Theorem 5.1.** *Consider the additive model (1.1) where the  $\epsilon_p(i)$ 's are independent and identically distributed with log-concave densities in the AGG class. Let the signals be as prescribed in Proposition 5.2. If the signal sizes fall below the strong classification boundary (4.5), i.e.  $\bar{r} < g(\beta)$ , then we have*

$$\liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p} \mathbb{E}_\pi \mathbb{P}[\widehat{S}_p \neq S_p] = 1, \quad (5.11)$$

where the infimum on  $\widehat{S}_p$  is taken over all procedures.

*Proof.* When the errors are independent with log-concave density, the oracle thresholding procedure  $\widehat{S}_p^*$ , by Proposition 5.2, minimizes the Bayes risk (5.4) among all procedures. That is,

$$\liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p} \mathbb{E}_\pi \mathbb{P}[\widehat{S}_p \neq S_p] \geq \liminf_{p \rightarrow \infty} \mathbb{E}_\pi \mathbb{P}[\widehat{S}_p^* \neq S_p].$$

Since  $\widehat{S}_p^*$  belongs to the class of all thresholding procedures, we have

$$\begin{aligned} \liminf_{p \rightarrow \infty} \mathbb{E}_\pi \mathbb{P}[\widehat{S}_p^* \neq S_p] &\geq \liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \mathbb{E}_\pi \mathbb{P}[\widehat{S}_p \neq S_p] \\ &\geq \liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \inf_{S_p} \mathbb{P}[\widehat{S}_p \neq S_p] = 1, \end{aligned}$$

when  $\bar{r} < g(\beta)$ , where the last line follows from Theorem 4.2.  $\square$

### 5.1.2 Bayes optimality of likelihood ratio thresholding

When the MLR condition in Proposition 5.1 is violated, the oracle thresholding procedures can in fact be sub-optimal (see Example 5.2 and Section 5.1.3, below).

**Stilian:** In this section, we demonstrate that thresholding the *likelihood ratio* rather than signal values yields the finite-sample Bayes optimal procedures. We consider a special but sufficiently general case of signal models with equal densities.

Namely, let the observations  $x(i)$ ,  $i = 1, \dots, p$  have  $s$  signals as prescribed in (5.2) with *common* “signal” density  $f_a$ , and let the remaining  $(p - s)$  locations have common “error” density  $f_0$ . Define the likelihood ratios

$$L(i) := f_a(x(i)) / f_0(x(i)),$$

and let  $L_{[1]} \geq L_{[2]} \geq \dots \geq L_{[p]}$  be the order statistics of the  $L(i)$ ’s.

**Definition 5.4 (Oracle likelihood ratio thresholding).** Recall that  $\langle a_1, \dots, a_s \rangle$  denotes the vector of  $a_i$ ’s arranged in a non-increasing order. Any estimator  $\hat{S} = \{i_1, \dots, i_s\}$  such that

$$\langle L(i_1), \dots, L(i_s) \rangle = (L_{[1]}, \dots, L_{[s]}),$$

will be referred to as an *oracle likelihood thresholding* estimator of the support  $S$ .

**Proposition 5.3.** Any oracle likelihood ratio thresholding procedure  $\hat{S}_{LRT}$  is finite-sample optimal in terms of Bayes risk. That is,

$$\hat{S}_{LRT} \in \arg \min_{\hat{S} \in \mathcal{S}} \mathbb{E}_\pi \mathbb{P}[\hat{S} \neq S]. \quad (5.12)$$

for all  $s$  and  $p$ , where the infimum on  $\hat{S}$  is taken over all support estimators of size  $s$ .

Introducing part of the proposition as a definition....

**Stilian:** NOTE This change & definition was motivated by the annoying fact that we could have ties among the likelihood ratios. Also, to my understanding we must LIMIT OURSELVES to estimators of size  $s$ , so we should limit  $|\hat{S}_{LRT}|$  to  $s$ . CORRECT? I am SORRY if I made a MESS.

*Proof.* The proof is analogous to that of Proposition 5.1. We need to show that  $\hat{S}_{LRT} \in \arg \max_{S \in \mathcal{S}} f(x|S)\pi(S)$ . Since the distribution  $\pi$  of the support  $S$  is uniform (recall (5.3)), it is equivalent to prove that

$$f(x|\hat{S}_{LRT}) = \max_{S \in \mathcal{S}} f(x|S),$$

where  $f(x|S)$  is the conditional distribution of the data given the unordered support  $S$ ,

$$f(x|S) = \sum_P f(x|P)\pi^{\text{ord}}(P|S) = \prod_{j \in S} f_a(x(j)) \prod_{j \notin S} f_0(x(j)). \quad (5.13)$$

Suppose  $\widehat{S} \in \mathcal{S}$  is *not* an oracle likelihood thresholding estimator. Then from the definition of the likelihood ratio thresholding procedure, there must be indices  $j \in \widehat{S}$  and  $j' \notin \widehat{S}$  such that  $L(j) < L(j')$ . If we exchange the labels of  $L(j)$  and  $L(j')$ , that is, we form a new estimate  $\widehat{S}' = (\widehat{S} \setminus \{j\}) \cup \{j'\}$ , comparing the log-likelihoods under  $\widehat{S}$  and  $\widehat{S}'$ , we have

$$\log f(x|\widehat{S}) - \log f(x|\widehat{S}') = \log f_a(x(j)) + \log f_0(x(j')) - \log f_a(x(j')) - \log f_0(x(j)).$$

By the definition of  $L(j)$ 's, and the order relations, we obtain

$$\log f(x|\widehat{S}) - \log f(x|\widehat{S}') = \log L(j) - \log L(j') > 0$$

Stilian: This shows that  $\widehat{S}$  cannot be Bayes optimal unless it is a likelihood thresholding estimator. Note that with the same argument for every two likelihood thresholding estimators  $\widehat{S}'$  and  $\widehat{S}''$ , we have  $f(x|\widehat{S}') = f(x|\widehat{S}'')$ , proving the desired optimality.

Please check. See also the change in the oracle thresholding argument.  $\square$

The characterization of optimal likelihood ratio thresholding procedures in Proposition 5.3 may not always yield practical estimators, as the density of the alternatives, and the number of signals  $s$  are typically unknown. Still, some insights can be gained by virtue of Proposition 5.3. In particular, when MLR fails (or equivalently, when the errors in model (1.1) do not have log-concave densities), data thresholding is sub-optimal.

*Example 5.2 (Sub-optimality of data thresholding).* Let the errors have iid generalized Gaussian density with  $\nu = 1/2$ , i.e.,  $\log f_0(x) \propto -x^{1/2}$ . Let dimension  $p = 2$ , sparsity  $s = 1$  with uniform prior, and signal size  $\delta = 1$ . That is,  $\mathbb{P}[\mu = (0, 1)^T] = \mathbb{P}[\mu = (1, 0)^T] = 1/2$ . If the observations take on values  $x = (x_1, x_2)^T = (1, 2)^T$ , we see from a comparison of the likelihoods (and hence, the posteriors),

$$\log \frac{f(x|\{1\})}{f(x|\{2\})} = 2x_1^{1/2} + 2(x_2 - 1)^{1/2} - 2x_2^{1/2} - 2(x_1 - 1)^{1/2} = 4 - 2\sqrt{2} > 0,$$

that even though  $x_1 < x_2$ , the set  $\{1\}$  is a better estimate of support than  $\{2\}$ , i.e.,  $\mathbb{P}[S = \{1\} | x] > \mathbb{P}[S = \{2\} | x]$ .

This simple example shows that, in the case when the errors have super-exponential tails, the optimal procedures are in general *not* data thresholding. A slightly more general conclusion can be found in Corollary 5.1.

### 5.1.3 Sub-optimality of data thresholding procedures

We provide a slightly more general result on the sub-optimality of data thresholding procedures.

**Corollary 5.1.** *Consider the additive error model (1.1). Let the errors  $\epsilon$  be independent with common distribution  $F$ . Let each of the  $s$  signals be located on  $\{1, \dots, p\}$  uniformly at random with equal magnitude  $0 < \delta < \infty$ . Assume the errors  $\epsilon(i)$ 's are iid with density  $f$  that is log-convex on  $[K, +\infty)$ , for some  $K > 0$ .*

*If  $\hat{S}_{\text{opt}}$  is the Bayes optimal (i.e., the oracle likelihood thresholding estimator), then, whenever  $j \in \hat{S}_{\text{opt}}$  for some  $x(j) > K + \delta$ , we must necessarily have  $j' \in \hat{S}_{\text{opt}}$  for all  $j'$  such that  $K + \delta \leq x(j') < x(j)$ .*

Specifically, if there are  $m$  observations exceeding  $K + \delta$ , with  $m > s$ , then the top  $m - s$  observations will *not* be included in the optimal estimator  $\hat{S}_{\text{opt}}$ . This shows that, in the case when the errors have super-exponential tails, the optimal procedures are in general *not* data thresholding.

*Proof (Proof of Corollary 5.1).* Since the density of the alternatives  $f_a(t) = f(t - \delta)$  is log-convex on  $[K + \delta, \infty)$ , by Relation (5.8) in the proof of Lemma 5.1 and appealing to log-convexity (rather than log-concavity), the likelihood ratio

$$L(j) := \frac{f_a(x(j))}{f_0(x(j))}$$

is decreasing in  $x(j)$  on  $[K + \delta, \infty)$ . The claim follows from Proposition 5.3.  $\square$

*Remark 5.2 (An open problem).* As we have seen, the thresholding estimators are no longer optimal in the additive model with error-densities heavier than exponential. Thanks to Proposition 5.3, the oracle likelihood thresholding procedures are promising alternatives that can lead us to practical support estimators.

In the case where the signals have difference sizes, however, the argument in the proof of Proposition 5.3 breaks down since the signal densities need to be identical for Relation (5.13) to hold. In such cases, the characterization of the optimal procedure is an open problem.

## 5.2 Minimax optimality in exact support recovery

We establish in this section minimax versions of our results from Chapter 4. Specifically, if we restrict ourselves to *the class of thresholding procedures*  $\mathcal{T}$  (defined in (2.20)), then Bonferroni's procedure is minimax optimal, for *any* fixed dependence structures in the URS class. This is formalized in Corollary 5.2 below. We refer to this result as *point-wise* minimax, to emphasize the fact that this optimality holds for every *fixed* URS array.

Meanwhile, if we search over *all procedures*, but expand the model space to include *all* dependence structures, then a different minimax optimality statement holds for Bonferroni's procedure. This result, formally stated in Section 5.2.2, is a consequence of our characterization of the finite-sample Bayes optimality of thresholding procedures in Section 5.1.1.

### 5.2.1 Point-wise minimax optimality for thresholding procedures

Theorems 4.1 and 4.2 can be cast in the form of an asymptotic minimax statement.

**Corollary 5.2 (Point-wise minimax).** *Let  $\widehat{S}_p^{\text{Bonf}}$  be the sequence of Bonferroni's procedure described in Theorem 4.1. Let also the errors have common AGG( $\nu$ ) distribution  $F$  with parameter  $\nu > 0$ , and  $\Theta_p^+$  be as defined in (4.7). If  $\underline{r} > g(\beta)$ , then we have*

$$\limsup_{p \rightarrow \infty} \sup_{\mu \in \Theta_p^+(\beta, \underline{r})} \mathbb{P}(\widehat{S}_p^{\text{Bonf}} \neq S_p) = 0, \quad (5.14)$$

for arbitrary dependence structure of the error array  $\mathcal{E} = \{\epsilon_p(i)\}_p$ . Let  $\mathcal{T}$  be the class of thresholding procedures (2.20). If  $\underline{r} < g(\beta)$ , then we have

$$\liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \sup_{\mu \in \Theta_p^+(\beta, \underline{r})} \mathbb{P}(\widehat{S}_p \neq S_p) = 1, \quad (5.15)$$

for any error dependence structure such that  $\mathcal{E} \in U(F)$  and  $(-\mathcal{E}) \in U(F)$ .

*Proof.* The first conclusion (5.14) is a restatement of Theorem 4.1.

For the second statement (5.15), since  $\underline{r} < g(\beta)$ , we can pick a sequence  $\mu^* \in \Theta_p^+(\beta, \underline{r})$  such that  $|S_p| = \lfloor p^{1-\beta} \rfloor$ , with signals having the same signal size  $\mu(i) = (2r \log p)^{1/\nu}$  for all  $i \in S_p$ , where  $\underline{r} < r < g(\beta)$ . For this particular choice of  $\mu^*$  we have  $\mu^* \in \Theta_p^-(\beta, \bar{r})$  where  $r < \bar{r} < g(\beta)$ , and according to Theorem 4.2, we obtain  $\lim_{p \rightarrow \infty} \inf_{\widehat{S}_p \in \mathcal{T}} \mathbb{P}[\widehat{S}_p \neq S_p] = 1$ , for all dependence structures in the URS class.  $\square$

*Remark 5.3.* Theorem 4.2 is a stronger result than the traditional minimax claim in Relation (5.15). Indeed, (4.24) involves an infimum (over the class  $\Theta_p^-$ ) while (5.15) has a supremum (over the class  $\Theta_p^+$ ).

On the other hand, Corollary 5.2 is more informative than many minimax-type statements, since it applies “point-wise” to any fixed error dependence structure in the URS class.

*Remark 5.4.* Corollary 5.2 echoes Remark 4.4: for a very large class of dependence structures, we cannot improve upon Bonferroni's procedure in exact support recovery problems (asymptotically), unless we look beyond thresholding procedures.

### 5.2.2 Minimax optimality over all procedures

Theorem 5.1 allows us to state another minimax conclusion — one in which we search over *all procedures*, by allowing the supremum in the minimax statement to be taken over the dependence structures.

**Corollary 5.3.** *Let  $D(F)$  be the collection of error arrays with common marginal  $F$  as defined in (4.20) where  $F$  is an  $\text{AGG}(\nu)$  distribution. Let also  $\widehat{S}_p^{\text{Bonf}}$  be Bonferroni's procedure as described in Theorem 4.1. If  $\underline{r} > g(\beta)$ , then we have*

$$\limsup_{p \rightarrow \infty} \sup_{\substack{\mu \in \Theta_p^+(\beta, \underline{r}) \\ \mathcal{E} \in D(F)}} \mathbb{P}(\widehat{S}_p^{\text{Bonf}} \neq S_p) = 0. \quad (5.16)$$

Further, when  $\underline{r} < g(\beta)$ , and  $F$  has a positive log-concave density  $f$ , we have

$$\liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p} \sup_{\substack{\mu \in \Theta_p^+(\beta, \underline{r}) \\ \mathcal{E} \in D(F)}} \mathbb{P}(\widehat{S}_p \neq S_p) = 1, \quad (5.17)$$

where the infimum on  $\widehat{S}_p$  is taken over all procedures.

*Proof.* Relation (5.16) is a re-statement of Remark 4.1.

For any distribution  $\pi$  (with a slight abuse of notation) over the parameter space  $\Theta_p^+ \times D(F)$ , we have

$$\liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p} \sup_{\substack{\mu \in \Theta_p^+(\beta, \underline{r}) \\ \mathcal{E} \in D(F)}} \mathbb{P}(\widehat{S}_p \neq S_p) \geq \liminf_{p \rightarrow \infty} \inf_{\widehat{S}_p} \mathbb{E}_\pi \mathbb{P}(\widehat{S}_p \neq S_p), \quad (5.18)$$

since the supremum is bounded from below by expectations. In particular, define  $\pi$  to be the uniform distribution over the configurations  $\Theta_p^* \times I(f)$ , where

$$\begin{aligned} \Theta_p^* &= \{\mu \in \mathbb{R}^d : |S_p| = \lfloor p^{1-\beta} \rfloor, \mu(i) = 0 \text{ for all } i \notin S, \text{ and} \\ &\quad \mu(i) = (\nu r \log p)^{1/\nu} \text{ for all } i \in S, \text{ where } \underline{r} < r < g(\beta)\}, \end{aligned}$$

and

$$I(f) = \{\mathcal{E} = (\epsilon_p(i))_p : \epsilon_p(i) \text{ iid with density } f(x) \propto \exp\{-|x|^\nu/\nu\}\}.$$

Since the density  $f$  of  $F$  is log-concave, the distribution of the signal configurations satisfies the conditions in Theorem 5.1. Thus, the desired conclusion (5.17) follows from Theorem 5.1 and (5.18).  $\square$

*Remark 5.5.* Since the class  $\text{AGG}(\nu)$ ,  $\nu \geq 1$  contains distributions with log-concave densities (Example 5.1), the minimax statement (5.17) continues to hold if the supremum is taken over the entire class  $F \in \text{AGG}(\nu)$ ,  $\nu \geq 1$ . We opted for a more informative formulation which emphasizes the log-concavity condition on the density of  $F$ .

*Remark 5.6.* Corollary 5.3 is no stronger than Corollary 5.2. In Corollary 5.2 we search over only the class of thresholding procedures, but offer a tight, point-wise lower bound on the asymptotic risk over the class of URS dependence structures. On the other hand, Corollary 5.3 provides a uniform lower bound for the asymptotic

risk over all dependence structures, which may not be tight except in the case of independent errors.

### 5.3 A discussion on optimality

Stilian: Here we briefly discuss the optimality and sub-optimality of the thresholding procedures. Needs work

Consider the model (1.1) with independent errors, Proposition 5.3, and indeed, Example 5.2 demonstrate that thresholding procedures are in fact *sub-optimal* for  $\text{AGG}(\nu)$  models with  $\nu < 1$ . Therefore, the optimality of thresholding procedures (specifically, Bonferroni's procedure) only applies to  $\text{AGG}(\nu)$  models with  $\nu \geq 1$ .

If we restrict the space of methods to only thresholding procedures, then results in Section 5.2.1 state that the phase transition phenomenon — the 0-1 law in the sense of Corollary 5.2 — is universal in all error models with rapidly varying tails. This includes  $\text{AGG}(\nu)$  models *for all*  $\nu > 0$ . In contrast, models with heavy (regularly varying) tailed errors do not exhibit this phenomenon (see Theorem B.3). We summarize the properties of thresholding procedures in Table 5.1.

**Table 5.1** Properties of thresholding procedures under different error distributions when errors are independent. Properties of the error distributions are listed in brackets.

Thresholding procedure (Error distributions)	Bayes optimality (Log-concave density)	Phase transition (Rapidly-varying tails)
$\text{AGG}(\nu)$ , $\nu \geq 1$	Yes (Yes)	Yes (Yes)
$\text{AGG}(\nu)$ , $0 < \nu < 1$	No (No)	Yes (Yes)
Power laws	No (No)	No (No)



## Chapter 6

### Uniform Relative Stability for Gaussian Arrays

In this chapter, we establish a complete characterization of URS for Gaussian arrays in terms of a simple condition on the covariance structures. The condition is as follows.

**Definition 6.1 (Uniformly decreasing dependence (UDD)).** Consider a triangular array of jointly Gaussian distributed errors  $\mathcal{E} = \{(\epsilon_p(i))_{i=1}^p, p = 1, 2, \dots\}$  with unit variances,

$$\epsilon_p \sim N(0, \Sigma_p), \quad p = 1, 2, \dots$$

The array  $\mathcal{E}$  is said to be uniform decreasingly dependent (UDD) if for every  $\delta > 0$  there exists a finite  $N(\delta) < \infty$ , such that for every  $i \in \{1, \dots, p\}$ , and  $p \in \mathbb{N}$ , we have

$$\left| \{k \in \{1, \dots, p\} : \Sigma_p(i, k) > \delta\} \right| \leq N(\delta) \quad \text{for all } \delta > 0. \quad (6.1)$$

That is, for every coordinate  $i$ , the number of elements which are more than  $\delta$ -correlated with  $\epsilon_p(i)$  does not exceed  $N(\delta)$ .

Note that the bound in (6.1) holds uniformly in  $i$  and  $p$ , and only depends on  $\delta$ . Also observe that on the left-hand side of (6.1), we merely count in each row of  $\Sigma_p$  the number of exceedances of covariances (not their absolute values!) over level  $\delta$ .

*Remark 6.1.* Without loss of generality, we may require that  $N(\delta)$  be a monotone non-increasing function of  $\delta$ , for we can take

$$N(\delta) = \sup_{p,i} \left| \{k : \Sigma_p(i, k) > \delta\} \right|,$$

which is non-increasing in  $\delta$ . Definition 6.1 therefore states that the array is UDD when  $N(\delta) < \infty$  for all  $\delta > 0$ .

Observe that the UDD condition does not depend on the order of the coordinates in the error vector  $\epsilon_p = (\epsilon_p(i))_{i=1}^p$ . Often times, however, the errors are thought of coming from a stochastic process indexed by time or space. To illustrate the generality of the UDD condition, we formulate next a simple sufficient condition (UDD') that is easier to check in a time-series context.

**Definition 6.2 (UDD').** For  $\epsilon_p \sim N(0, \Sigma_p)$  with unit variances, an array  $\mathcal{E} = (\epsilon_p(i))_{i=1}^p$  is said to satisfy the UDD' condition if there exist:

- (i) permutations  $l = l_p$  of  $\{1, \dots, p\}$ , for all  $p \in \mathbb{N}$ , and
- (ii) a non-negative sequence  $(r_n)_{n=1}^\infty$  converging to zero  $r_n \rightarrow 0$ , as  $n \rightarrow \infty$ ,

such that

$$\sup_{p \in \mathbb{N}} |\Sigma_p(i', j')| \leq r_{|i-j|}. \quad (6.2)$$

where  $i' = l(i)$ ,  $j' = l(j)$ , for all  $i, j \in \{1, \dots, p\}$ .

*Remark 6.2.* Without loss of generality, we may also require that  $r_n$  be non-increasing in  $n$ , for we can replace  $r_n$  with the non-increasing sequence  $r'_n = \sup_{m \geq n} r_m$ .

**Proposition 6.1.** *UDD' implies UDD.*

*Proof.* Since  $r_n \rightarrow 0$ , for any  $\delta > 0$ , there exists an integer  $M = M(\delta) < \infty$  such that  $r_n \leq \delta$ , for all  $n \geq M$ . Thus, by (6.2), for every fixed  $j' \in \{1, \dots, p\}$ , we can have  $|\text{Cov}(\epsilon_p(k'), \epsilon_p(j'))| > \delta$ , only if  $k'$  belongs to the set:

$$\{k' \in \{1, \dots, p\} : j - M \leq k := l_p^{-1}(k') \leq j + M\},$$

where  $j := l_p^{-1}(j')$ . That is, there are at most  $2M + 1 < \infty$  indices  $k' \in \{1, \dots, p\}$ , whose covariances with  $\epsilon(j')$  may exceed  $\delta$ . Since this holds uniformly in  $j' \in \{1, \dots, p\}$ , Condition UDD follows with  $N(\delta) = 2M + 1$ .  $\square$

We now state the main result of this section: a Gaussian sequence is URS if and only if it is UDD. The URS condition essentially requires that the dependencies decay in a uniform fashion, the rate at which dependence decay does *not* matter.

**Theorem 6.1.** *Let  $\mathcal{E}$  be a Gaussian triangular array with standard normal marginals. The array  $\mathcal{E}$  has uniformly relatively stable (URS) maxima if and only if it is uniformly decreasing dependent (UDD).*

Specifically, for stationary Gaussian arrays, we have the following corollary.

**Corollary 6.1.** *Let  $\mathcal{E} = \{\epsilon_p(i) = Z(i)\}$  for a stationary Gaussian time series  $\mathcal{Z} = \{Z(i)\}$ . Then  $\mathcal{E}$  is URS if and only if the autocovariance function  $\text{Cov}(Z(k), Z(0)) \rightarrow 0$ , as  $k \rightarrow \infty$ .*

Corollary 6.1 follows by Theorem 6.1 and the observation that UDD is equivalent to vanishing autocovariance of  $\mathcal{Z}$ . A slightly weaker form of the “if” part was established in Theorem 3 of Berman (1964).

Returning again to the study of support recovery problems, Theorem 6.1 and the necessary condition for exact support recovery in Theorem 4.2 yields the following result.

**Corollary 6.2.** *For UDD Gaussian errors, the result in Theorem 4.2 holds.*

As a counterpart to Remark 4.4, we demonstrate the tightness of the dependence conditions in Theorem 4.2. Specifically, we demonstrate that if the URS dependence condition is violated, then it may be possible to recover the support of weaker signals below the boundary.

*Example 6.1.* Suppose  $\mathcal{E} = (\epsilon_p(i))_{i=1}^p$  is Gaussian, and is comprised of  $\lfloor p^{1-\beta} \rfloor$  blocks, each of size at least  $\lfloor p^\beta \rfloor$ ; let the elements of each block have correlation 1, and let elements from different blocks be independent. If  $\underline{r} \geq 4(1-\beta)$ , then the procedure  $\widehat{S} = \{i : x(i) > \sqrt{2(1-\beta)\log p}\}$  yields  $\mathbb{P}[\widehat{S} = S] \rightarrow 1$ . This requirement on signal size is strictly weaker than that of the strong classification boundary, since  $4(1-\beta) < (1 + \sqrt{1-\beta})^2$  on  $\beta \in (0, 1)$ .

The above example shows that if the correlations of the Gaussian errors do not decay in a uniform fashion (UDD fails), then we can do substantially better in terms of support recovery. The claims in the example are verified in Section 4.7, while numerical simulations of this example can be found in Section 4.6.

We conclude with a brief discussion on the relationships between UDD and other dependence conditions in the context of extreme value theory, before proceeding to the proof of Theorem 6.1.

Suppose that the array of errors  $\mathcal{E}$  comes from a stationary Gaussian time series  $\epsilon(i)$ ,  $i \in \mathbb{N}$ , with auto-covariance  $r_p = \text{Cov}(\epsilon(i+p), \epsilon(i))$ . One is interested in the asymptotic behavior of the maxima  $M_p := \max_{i=1, \dots, p} \epsilon(i)$ .

In this setting, the Berman's condition, introduced in Berman (1964), requires that

$$r_p \log p \rightarrow 0, \quad \text{as } p \rightarrow \infty. \quad (6.3)$$

This condition entails that

$$a_p(M_p - b_p) \xrightarrow{d} Z, \quad \text{as } p \rightarrow \infty, \quad (6.4)$$

with the Gumbel limit distribution  $\mathbb{P}[Z \leq x] = \exp\{-e^{-x}\}$ ,  $x \in \mathbb{R}$ , where

$$a_p = \sqrt{2 \log p}, \quad b_p = \sqrt{2 \log p} - \frac{1}{2} \left( \sqrt{2 \log p} \right)^{-1} (\log \log(p) + \log(4\pi)),$$

are *the same* centering and normalization sequences as in the case of iid  $\epsilon(i)$ 's. Berman's condition is one of the weakest dependence conditions in the literature for which this result holds. See, e.g., Theorem 4.4.8 in Embrechts et al. (2013), where (6.3) is described as “very weak”.

For dependence conditions weaker than (6.3), the sequences of normalizing and centering constants in (6.4) are *different* from the iid case, and the corresponding limit is no longer Gumbel; see, for example, Theorems 6.5.1 and 6.6.4 in Leadbetter et al. (1983), and McCormick and Mittal (1976).

On the other hand, in our high dimensional support estimation context, the notion of relative stability is sufficient and more natural than the finer notions of distributional convergence. If one is merely interested in the asymptotic relative stability of the Gaussian maxima, then Berman's condition can be relaxed significantly (see

also, Theorem 4.1 of Berman (1964)). Observe that by Proposition 6.1, the Berman condition (6.3) implies UDD and hence relative stability (Theorem 6.1), i.e.,

$$\frac{1}{b_p} M_p \xrightarrow{\mathbb{P}} 1, \quad \text{as } p \rightarrow \infty. \quad (6.5)$$

This *concentration of maxima* property can be readily deduced from (6.4), since  $a_p b_p \sim 2 \log(p) \rightarrow \infty$  as  $p \rightarrow \infty$ . Theorem 6.1 shows that (6.5) holds if the much weaker uniform dependence condition UDD holds. Note that our condition is coordinate free — neither monotonicity of the sequence  $r_p$  nor stationarity of the underlying array is required.

The rest of this chapter is devoted to the proof of the main result, i.e., Theorem 6.1. We first introduce a key lemma regarding the structure of correlation matrix of high-dimensional random variables. The proof uses a surprising, yet elegant application of Ramsey’s Theorem from the study of combinatorics. The ‘only if’ part of Theorem 6.1 follows from this lemma, in Section 6.2.

The proof of the ‘if’ part is detailed in Section 6.3. The arguments there were recently extended to establish bounds on the rate of concentration of maxima in Kartsioukas et al. (2019); see also, Tanguy (2015b) and references therein for related work on this topic.

## 6.1 Ramsey’s Coloring theorem and structure of correlation matrices

Given any integer  $k \geq 1$ , there is always an integer  $R(k, k)$  called the *Ramsey number*:

$$k \leq R(k, k) \leq \binom{2k-2}{k-1} \quad (6.6)$$

such that the following property holds: every undirected graph with at least  $R(k, k)$  vertices will contain *either* a clique of size  $k$ , or an *independent set* of  $k$  nodes. Recall that a clique is a complete sub-graph where all pairs of nodes are connected, and an independent set is a set of nodes where no two nodes are connected.

This result is a consequence of the celebrated work of Ramsey (2009), which gave birth to Ramsey Theory (see e.g., Conlon et al. (2015)). The Ramsey Theorem and the upper bound (6.6) (established first in Erdős and Szekeres (1935)) are at the heart of the proof of the following result.

**Proposition 6.2.** Fix  $\gamma \in (0, 1)$  and let  $P = (\rho(i, j))_{n \times n}$  be an arbitrary correlation matrix. If

$$k := \lfloor \log_2(n)/2 \rfloor \geq \lceil 1/\gamma \rceil + 1, \quad (6.7)$$

then there is a set of  $k$  indices  $K = \{l_1, \dots, l_k\} \subseteq \{1, \dots, n\}$  such that

$$\rho(i, j) \geq -\gamma, \text{ for all } i, j \in K. \quad (6.8)$$

*Proof (Proof of Proposition 6.2).* By using (6.6) and a refinement of the Stirling's formula, we will show at the end of the proof that for  $k \leq \log_2(n)/2$ , we have

$$R(k, k) \leq n, \quad (6.9)$$

where  $R(k, k)$  is the Ramsey number.

Now, construct a graph with vertices  $\{1, \dots, n\}$  such that there is an edge between nodes  $i$  and  $j$  if and only if  $\rho(i, j) > -\gamma$ . In view of (6.9) and Ramsey's theorem (see e.g., Theorem 1 in Fox (2009) or Conlon et al. (2015) for a recent survey on Ramsey theory), there is a subset of  $k$  nodes  $K = \{l_1, \dots, l_k\}$ , which is either a *complete graph* or an *independent set*.

If  $K$  is a complete graph, then by our construction of the graph, Relation (6.8) holds.

Now, suppose that  $K$  is a set of independent nodes. This means, again by the construction of our graph, that

$$\rho(i, j) < -\gamma, \quad \text{for all } i \neq j \in K.$$

Let  $Z_i$ ,  $i \in K$  be zero-mean random variables such that  $\rho(i, j) = \mathbb{E}[Z_i Z_j]$ . Observe that

$$\text{Var}\left(\sum_{i \in K} Z_i\right) = \sum_{i \in K} \text{Var}(Z_i) + \sum_{\substack{i \neq j \\ i, j \in K}} \text{Cov}(Z_i, Z_j) < k - k(k-1)\gamma, \quad (6.10)$$

since  $\text{Var}(Z_i) = 1$  and  $\rho(i, j) < -\gamma$  for  $i \neq j$ . By our assumption,  $k \geq (\lceil 1/\gamma \rceil + 1)$ , or equivalently,  $(k-1) \geq 1/\gamma$ , the variance in (6.10) is negative. This is a contradiction showing that there are no independent sets  $K$  with cardinality  $k$ .

To complete the proof, it remains to show that Relation (6.9) holds. In view of the upper bound on the Ramsey numbers (6.6), it is enough to show that  $k \leq \log_2(\sqrt{n})$  implies

$$\binom{2k-2}{k-1} \leq n.$$

This follows from a refinement of the Stirling formula, due to Robbins (1955):

$$\sqrt{2\pi m}^{m+1/2} e^{-m} e^{\frac{1}{(12m+1)}} \leq m! \leq \sqrt{2\pi m}^{m+1/2} e^{-m} e^{\frac{1}{12m}}.$$

Indeed, letting  $\tilde{k} := k-1$ , and applying the above upper and lower bounds to the terms  $(2\tilde{k})!$  and  $\tilde{k}!$ , respectively, we obtain:

$$\binom{2k-2}{k-1} \equiv \frac{(2\tilde{k})!}{(\tilde{k}!)^2} \leq \frac{2^{2\tilde{k}}}{\sqrt{\pi \tilde{k}}} \exp\left\{\frac{1}{24\tilde{k}} - \frac{2}{12\tilde{k}+1}\right\} < 2^{2k}$$

where the last two inequalities follow by simply dropping positive factors less than 1. Since  $2k \leq \log_2(n)$ , the above bound implies Relation (6.9) and the proof is complete.  $\square$

Using Proposition 6.2, we establish the key lemma used in the proof of Theorem 6.1.

**Lemma 6.1.** *Let  $c \in (0, 1)$ , and  $P = (\rho(i, j))_{(n+1) \times (n+1)}$  be a correlation matrix such that*

$$\rho(1, j) > c \quad \text{for all } j = 1, \dots, n+1. \quad (6.11)$$

*If  $n \geq 2^{2\lceil 2/c^2 \rceil + 4}$ , then there is a set of indices  $K = \{l_1, \dots, l_k\} \subseteq \{2, \dots, n+1\}$  of cardinality  $k = |K| = \lfloor \log_2 \sqrt{n} \rfloor$ , such that*

$$\rho(i, j) > \frac{c^2}{2} \quad \text{for all } i, j \in K. \quad (6.12)$$

*That is, all entries of the  $k \times k$  sub-correlation matrix  $P_K := (\rho(i, j))_{i, j \in K}$  are larger than  $c^2/2$ .*

*Proof (Proof of Lemma 6.1).* Let  $Z_1, \dots, Z_{n+1}$  be random variables with covariance matrix  $P$ . Denote  $\rho_j = \rho(1, j)$  and define

$$R(j) = \begin{cases} \frac{1}{\sqrt{1-\rho_j^2}} (Z(j) - \rho_j Z(1)), & \text{if } \rho_j < 1, \\ R^* & \text{if } \rho_j = 1, \end{cases} \quad (6.13)$$

where  $R^*$  is an arbitrary zero-mean, unit-variance random variable. It is easy to see that  $\text{Var}(R(j)) = 1$ , and

$$\begin{aligned} \text{Cov}(Z(i), Z(j)) &= \text{Cov} \left( \rho_i Z(1) + \sqrt{1-\rho_i^2} R(i), \rho_j Z(1) + \sqrt{1-\rho_j^2} R(j) \right) \\ &= \rho_i \rho_j + \sqrt{1-\rho_i^2} \sqrt{1-\rho_j^2} \text{Cov}(R(i), R(j)) \\ &\geq c^2 + \min \{ \text{Cov}(R(i), R(j)), 0 \}. \end{aligned}$$

Therefore, Relation (6.12) would hold if we can find a set of indices  $K = \{l_1, \dots, l_k\}$  such that  $\text{Cov}(R(i), R(j)) > -c^2/2$  for all  $i, j \in K$ , where  $k = |K| = \lfloor \log_2 \sqrt{n} \rfloor$ . This, however, follows from Proposition 6.2 applied to  $(R(j))_{j=2}^{n+1}$  with  $\gamma = c^2/2$ , provided that

$$k = \lfloor \log_2 \sqrt{n} \rfloor \geq \lceil 2/c^2 \rceil + 1.$$

The last inequality indeed follows from the assumption that  $n \geq 2^{2\lceil 2/c^2 \rceil + 4}$ .  $\square$

## 6.2 URS implies UDD ('only if' part of Theorem 6.1)

In view of Remark 6.1, UDD is equivalent to the requirement that  $N(\delta) := 1 + \sup_p N_p(\delta) < \infty$  for all  $\delta \in (0, 1)$ , where

$$N_p(\delta) := \max_{j \in \{1, \dots, p\}} \left| \{i : i \neq j, \Sigma_p(j, i) > \delta\} \right|. \quad (6.14)$$

Therefore, if  $\mathcal{E}$  is not UDD, then there must exist a constant  $c \in (0, 1)$  for which  $N(c)$  is infinite, i.e., there is a subsequence  $\tilde{p} \rightarrow \infty$  such that  $N_{\tilde{p}}(c) \rightarrow \infty$ . Without loss of generality, we may assume that  $\tilde{p} = p$ .

Let  $j_p(c)$  be the maximizers of (6.14), and let

$$S_p(c) := \{i \in \{1, \dots, p\} : \Sigma_p(j_p(c), i) > c\}. \quad (6.15)$$

Observe that  $|S_p(c)| = N_p(c) + 1 \rightarrow \infty$ , as  $p \rightarrow \infty$  (note  $j_p(c) \in S_p(c)$ ).

Applying Lemma 6.1 to the set of random variables indexed by  $S_p(c)$ , we conclude, for  $N_p(c) \geq 2^{\lceil 2/c^2 \rceil + 4}$ , there must be a further subset

$$K_p(c) \subseteq S_p(c), \quad (6.16)$$

of cardinality

$$k_p(c) := |K_p(c)| \geq \log_2 \sqrt{N_p(c)}, \quad (6.17)$$

such that all pairwise correlations of the random variables indexed by  $K_p(c)$  are greater than  $c^2/2$ . Since the sequence  $N_p(c) \rightarrow \infty$ , by (6.17), we have  $k_p(c) \rightarrow \infty$  as  $p \rightarrow \infty$ .

Therefore, we have identified a sequence of subsets  $K_p(c) \subseteq \{1, \dots, p\}$  with the following two properties:

1.  $k_p(c) := |K_p(c)| \rightarrow \infty$ , as  $p \rightarrow \infty$ , and
2. For all  $i, j \in K_p(c)$ , we have

$$\Sigma_p(i, j) > c^2/2. \quad (6.18)$$

Without loss of generality, we may assume  $K_p(c) = \{1, \dots, k_p(c)\} \subseteq \{1, \dots, p\}$ , upon re-labeling of the coordinates.

Now consider a Gaussian sequence  $\epsilon^* = \{\epsilon^*(j), j = 1, 2, \dots\}$ , independent of  $\mathcal{E}$ , defined as follows:

$$\epsilon^*(j) := Z \left( c/\sqrt{2} \right) + Z(j) \sqrt{1 - c^2/2}, \quad j = 1, 2, \dots,$$

where  $Z$  and  $Z(j), j = 1, 2, \dots$  are independent standard normal random variables. Hence,

$$\text{Var}(\epsilon^*(j)) = 1 = \text{Var}(\epsilon_p(j)), \quad (6.19)$$

and

$$\text{Cov}(\epsilon^*(i), \epsilon^*(j)) = \frac{c^2}{2} \leq \text{Cov}(\epsilon_p(i), \epsilon_p(j)), \quad (6.20)$$

for all  $p$ , and all  $i \neq j, i, j \in K_p(c)$ . Thus we have, as  $p \rightarrow \infty$ ,

$$\frac{1}{u_{k_p(c)}} \max_{j \in K_p(c)} \epsilon^*(j) = \frac{c/\sqrt{2}}{u_{k_p(c)}} Z + \frac{\sqrt{1-c^2/2}}{u_{k_p(c)}} \max_{j \in K_p(c)} Z(j) \xrightarrow{\mathbb{P}} \sqrt{1-\frac{c^2}{2}}, \quad (6.21)$$

where the convergence in probability follows from Proposition 2.2 part 2.

Relations (6.19) and (6.20), by Slepian's Lemma Slepian (1962), also imply,

$$\frac{1}{u_{k_p(c)}} \max_{j \in K_p(c)} \epsilon^*(j) \stackrel{d}{\geq} \frac{1}{u_{k_p(c)}} \max_{j \in K_p(c)} \epsilon_p(j). \quad (6.22)$$

Therefore, by (6.22) and (6.21), for all  $\sqrt{1-c^2/2} \leq \delta < 1$ , we have,

$$\mathbb{P} \left[ \frac{1}{u_{k_p(c)}} \max_{j \in K_p(c)} \epsilon_p(j) < \delta \right] \rightarrow 1 \quad \text{as } p \rightarrow \infty.$$

This contradicts the definition of URS (with the particular choice of  $S_p := K_p(c)$ ), and the proof of the ‘only if’ part is complete.

### 6.3 UDD implies URS (‘if’ part of Theorem 6.1)

Recall that our objective is to show (4.21). We will do so in two stages; namely, we will prove that for all  $\delta > 0$ , we have

$$\mathbb{P} \left[ \frac{M_{S_p}}{u_{|S_p|}} > 1 + \delta \right] \rightarrow 0, \quad (6.23)$$

and

$$\mathbb{P} \left[ \frac{M_{S_p}}{u_{|S_p|}} < 1 - \delta \right] \rightarrow 0, \quad (6.24)$$

for any sequence of subsets  $S_p$  such that  $|S_p| \rightarrow \infty$ . Although the first step (6.23) was already shown in Proposition 2.2, regardless of the dependence structure, we provide in this section a more refined result. Specifically, the following result states that for the AGG model, the constant  $\delta$  in Proposition 2.2 can be replaced by a vanishing sequence  $c_p \rightarrow 0$ .

**Lemma 6.2 (Upper tails of AGG maxima).** *Let  $\mathcal{E}$  be an array with marginal distribution  $F \in \text{AGG}(\nu)$ ,  $\nu > 0$ . If we pick*

$$c_p = \frac{u_p \log p}{u_p} - 1, \quad (6.25)$$

where  $u_p = F^{\leftarrow}(1 - 1/p)$ , then we have  $c_p > 0$ ,  $c_p \rightarrow 0$ , and

$$\mathbb{P} \left[ \frac{M_p}{u_p} - (1 + c_p) > 0 \right] \rightarrow 0. \quad (6.26)$$



The proof can be found in Section 6.3.1 below.

Since Lemma 6.2 holds regardless of the dependence structure, the same conclusions hold if one replaces  $M_p$  by  $M_{S_p} = \max_{j \in S_p} \epsilon(j)$  and  $p$  by  $q = q(p) = |S_p|$ , where  $S_p$  is any sequence of sets such that  $q \equiv |S_p| \rightarrow \infty$ . This entails (6.23).

On the other hand, the proof of (6.24) uses a more elaborate argument based on the Sudakov-Fernique bound. We proceed by first bounding the probability by an expectation. For all  $\delta > 0$ , we have

$$\begin{aligned} \mathbb{P} \left[ \frac{M_{S_p}}{u_q} < 1 - \delta \right] &= \mathbb{P} \left[ - \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right) > \delta + c_q \right] \\ &\leq \mathbb{P} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_- > \delta + c_q \right] \\ &\leq \frac{1}{\delta + c_q} \mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_- \right], \end{aligned} \quad (6.27)$$

where  $(x)_- := \max\{-x, 0\}$  and the last line follows from the Markov inequality. The next result shows that the upper bound in (6.27) vanishes.

**Lemma 6.3.** *Let  $\mathcal{E}$  be a Gaussian UDD array and  $S_p \subseteq \{1, \dots, p\}$  be an arbitrary sequence of sets such that  $q = q(p) = |S_p| \rightarrow \infty$ . Then, for  $M_{S_p} := \max_{j \in S_p} \epsilon_p(j)$  and  $c_q$  as in (6.25), we have*

$$\mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_- \right] \rightarrow 0, \quad \text{as } p \rightarrow \infty. \quad (6.28)$$

The proof of the lemma is given in Section 6.3.2 below.

Going back to the proof of Theorem 6.1, we observe that Relations (6.27) and (6.28) imply (6.24), which completes the proof of the 'if' part.  $\square$

*Remark 6.3.* Only the Sudakov-Fernique minorization argument used in the proof of Lemma 6.3, relies on the Gaussian assumption. We expect the techniques and results here to be useful in extending Theorem 6.1 to more general class of distributions, say, the AGG model.

### 6.3.1 Bounding the upper tails of AGG maxima

*Proof (Proof of Lemma 6.2).* Recall by (2.33) that

$$u_q \sim (\nu \log q)^{1/\nu}, \quad q \rightarrow \infty,$$

so that

$$c_p = \frac{u_{p \log p}}{u_p} - 1 = \left( \frac{\log p + \log \log p}{\log p} \right)^{1/\nu} (1 + o(1)) - 1 \rightarrow 0 \quad \text{as } p \rightarrow \infty. \quad (6.29)$$

By the union bound, we have

$$\begin{aligned} \mathbb{P} \left[ \frac{M_p}{u_p} > 1 + c_p \right] &\leq \sum_{j=1}^p \mathbb{P} \left[ \frac{\epsilon_p(j)}{u_p} > 1 + c_p \right] = p \bar{F}(u_p \log p) \\ &= p \bar{F} \left( F^{\leftarrow} \left( 1 - \frac{1}{p \log p} \right) \right) \leq \frac{1}{\log p} \rightarrow 0. \end{aligned} \quad (6.30)$$

where the last inequality follows from the fact that  $F(F^{\leftarrow}(u)) \geq u$  for all  $u \in [0, 1]$ .

□

In addition to Lemma 6.2, which says the upper tail vanishes in probability, we will also prepare a result which states that the upper tail also vanishes in expectation.

**Lemma 6.4.** *Let  $M_p$  and  $c_p$  be as in Lemma 6.2, and denote*

$$\xi_p := \frac{M_p}{(1 + c_p)u_p}.$$

*Then there exists  $p_0, t_0 > 0$ , and absolute constant  $C > 0$  such that*

$$\mathbb{P} [\xi_p > t] \leq \exp \{-Ct^\nu\}, \quad \text{for all } p > p_0, t > t_0. \quad (6.31)$$

*In particular, the set of random variables  $\{(\xi_p)_+, p \in \mathbb{N}\}$  is uniformly integrable.*

*Proof (Proof of Lemma 6.4).* Recalling that  $(1 + c_p)u_p = u_{p \log p}$ , and by applying the union bound as in (6.30), we have

$$\begin{aligned} \log \mathbb{P} [\xi_p > t] &\leq \log p + \log \bar{F}(u_{p \log p} t) \\ &\leq \log p - \frac{1}{\nu} (u_{p \log p} t)^\nu (1 - \delta). \end{aligned} \quad (6.32)$$

for  $t > t_0(\delta) > 0$ , where  $\delta \in (0, 1)$  is an arbitrarily small number fixed in advance. This follows from the assumption that  $F \in \text{AGG}(\nu)$  and the Definition 2.6 of AGG tails. Using in (6.32) the explicit expressions for quantiles in (2.33), we obtain

$$\log \mathbb{P} [\xi_p > t] \leq \underbrace{\log p - (1 + o(1))(1 - \delta)t^\nu \log p}_{\text{greater than 1 for large } t} - \underbrace{t^\nu \log \log p (1 + o(1))(1 - \delta)}_{\text{greater than } C \text{ for large } p}. \quad (6.33)$$

For large  $t$ , we have  $(1 + o(1))(1 - \delta)t^\nu > 1$ , so that sum of the first two terms on the right-hand side of (6.33) is negative. Also, for  $p$  larger than some constant  $p_0(\delta)$ , we have  $\log \log p (1 + o(1))(1 - \delta) > C$  for some constant  $C$  that does not depend on  $p$ . Therefore (6.31) holds for  $t > t_0(\delta)$  and  $p > p_0(\delta)$ , and the proof is complete. □

**Corollary 6.3.** *The upper tails of AGG maxima vanish in expectation, i.e.,*

$$\mathbb{E} \left[ \left( \frac{M_p}{u_p} - (1 + c_p) \right)_+ \right] \rightarrow 0 \quad \text{as } p \rightarrow \infty, \quad (6.34)$$

where  $(a)_+ := \max\{a, 0\}$ .

*Proof (Proof of Corollary 6.3).* Since  $c_p \geq 0$  is a sequence converging to 0, we have  $c_p < 1$  for  $p \geq p_0$ . Hence for any  $t > 0$ , we have

$$\begin{aligned} \mathbb{P} \left[ \left( \frac{M_p}{u_p} - (1 + c_p) \right)_+ > t \right] &= \mathbb{P} \left[ (1 + c_p) (\xi_p - 1)_+ > t \right] \\ &\leq \mathbb{P} \left[ (\xi_p - 1)_+ > t/2 \right] \leq \mathbb{P} \left[ \xi_p > t/2 \right]. \end{aligned} \quad (6.35)$$

By Lemma 6.4,  $\{(\xi_p)_+\}$  is u.i., therefore by Relation (6.35),  $\{(M_p/u_p - (1 + c_p))_+, p \in \mathbb{N}\}$  is u.i. as well. Since by Lemma 6.2,  $(M_p/u_p - (1 + c_p))_+ \rightarrow 0$  in probability, Relation (6.34) follows from the established uniform integrability.  $\square$

### 6.3.2 Bounding the lower tails of Gaussian maxima

The main goal of this section is to establish the following result.

**Proposition 6.3.** *For every UDD Gaussian array  $\mathcal{E}$ , and any sequence of subsets  $S_p \subseteq \{1, \dots, p\}$  such that  $q = q(p) = |S_p| \rightarrow \infty$ , we have*

$$\liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{M_{S_p}}{u_q} \right] \geq 1, \quad (6.36)$$

where  $M_S = \max_{j \in S} \epsilon(j)$ .

Lemma 6.3, which is the key to the proof of the 'if' part of Theorem 6.1, follows immediately from this proposition.

*Proof (Proof of Lemma 6.3).* We start with the identity

$$\mathbb{E} \left[ \frac{M_{S_p}}{u_q} - (1 + c_q) \right] = \mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_+ \right] - \mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_- \right].$$

By re-arranging terms and taking limsup/liminf, we obtain

$$\begin{aligned} 0 &\leq \limsup_{p \rightarrow \infty} \mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_- \right] \\ &\leq \limsup_{p \rightarrow \infty} \mathbb{E} \left[ \left( \frac{M_{S_p}}{u_q} - (1 + c_q) \right)_+ \right] - \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{M_{S_p}}{u_q} - (1 + c_q) \right] \end{aligned} \quad (6.37)$$

$$= - \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{M_{S_p}}{u_q} - (1 + c_q) \right], \quad (6.38)$$

where the last equality follows from the fact that the lim-sup in (6.37) vanishes by Corollary 6.3. On the other hand, since  $c_q \rightarrow 0$ , we have

$$\liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{M_{S_p}}{u_q} - (1 + c_q) \right] = \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{M_{S_p}}{u_q} - 1 \right] \geq 0,$$

where the last inequality follows from Proposition 6.3. This shows that the right-hand side of (6.38) is non-positive and hence (6.28) holds.  $\square$

A interesting fact on the relationship between the upper quantiles and the expectation of iid maxima will be needed for the proof of Proposition 6.3. The following lemma may be of independent interest.

**Lemma 6.5.** *Let  $(X_i)_{i=1}^p$  be  $p$  iid random variables with distribution  $F$  such that  $\mathbb{E}[(X_i)_-]$  exists, i.e.,*

$$\mathbb{E}[\max\{-X_i, 0\}] < \infty.$$

*Let  $M_p = \max_{i=1, \dots, p} X_i$ . Assume that  $F$  has a density  $f$ , which is eventually decreasing. More precisely, we suppose there exists a  $C_0$  such that  $0 < F(C_0) < 1$ , and  $f(x_1) \geq f(x_2)$  whenever  $C_0 < x_1 \leq x_2$ . Under these assumptions, we have,*

$$\liminf_{p \rightarrow \infty} \frac{\mathbb{E}M_p}{u_{p+1}} \geq 1,$$

where  $u_{p+1} = F^{\leftarrow}(1 - 1/(p+1))$ .

*Proof (Proof of Lemma 6.5).* Write

$$X_i = F^{\leftarrow}(U_i)$$

where  $U_i$  are iid uniform random variables on  $(0, 1)$ . Denote  $M_p^U$  as the maximum of the  $U_i$ 's, we have  $\mathbb{E}M_p = \mathbb{E}[F^{\leftarrow}(M_p^U)]$ , and by conditioning, we obtain

$$\begin{aligned} \mathbb{E}M_p &= \mathbb{E}[F^{\leftarrow}(M_p^U) \mid M_p^U \geq F(C_0)] \mathbb{P}[M_p^U \geq F(C_0)] + \\ &\quad + \mathbb{E}[F^{\leftarrow}(M_p^U) \mid M_p^U < F(C_0)] \mathbb{P}[M_p^U < F(C_0)]. \end{aligned} \quad (6.39)$$

We first handle the first term in the summation. Since  $f$  is decreasing beyond  $C_0$ ,  $F$  is concave on  $(C_0, \infty)$ , and  $F^{\leftarrow}$  is convex on  $(F(C_0), 1)$ . By Jensen's inequality, we have

$$\mathbb{E}[F^{\leftarrow}(M_p^U) \mid M_p^U \geq F(C_0)] \geq F^{\leftarrow}(\mathbb{E}[M_p^U \mid M_p^U \geq F(C_0)]).$$

With a direct calculation, one can show that

$$F^{\leftarrow}(\mathbb{E}[M_p^U \mid M_p^U \geq F(C_0)]) = F^{\leftarrow}\left(\left(1 - \frac{1}{p+1}\right)\left(\frac{1 - F(C_0)^{p+1}}{1 - F(C_0)^p}\right)\right),$$

and hence

$$\begin{aligned} \mathbb{E}[F^{\leftarrow}(M_p^U) \mid M_p^U \geq F(C_0)] &\geq F^{\leftarrow}\left(\left(1 - \frac{1}{p+1}\right)\left(\frac{1 - F(C_0)^{p+1}}{1 - F(C_0)^p}\right)\right) \\ &\geq F^{\leftarrow}\left(1 - \frac{1}{p+1}\right) = u_{p+1}. \end{aligned}$$

Since  $\mathbb{P}[M_p^U \leq m \mid M_p^U < F(C_0)] = (m/F(C_0))^p \leq m/F(C_0)$  for  $m \leq F(C_0)$ , we have

$$(M_p^U \mid M_p^U < F(C_0)) \stackrel{d}{\geq} (U_1 \mid U_1 < F(C_0)),$$

where and the latter is the uniform distribution on  $(0, F(C_0))$ . Therefore, for the second term of the sum in (6.39), by monotonicity of  $F^{\leftarrow}$ , we obtain

$$\begin{aligned} \mathbb{E}[F^{\leftarrow}(M_p^U) \mid M_p^U < F(C_0)] &\geq \mathbb{E}[F^{\leftarrow}(U_1) \mid U_1 < F(C_0)] \\ &= \mathbb{E}[X_1 \mid X_1 < C_0]. \end{aligned}$$

Finally, since  $\mathbb{P}[M_p^U < F(C_0)] = F(C_0)^p = 1 - \mathbb{P}[M_p^U \geq F(C_0)]$ , by (2.33), we have

$$\frac{\mathbb{E}M_p}{u_{p+1}} \geq (1 - F(C_0)^p) + \frac{\mathbb{E}[X_1 \mid X_1 < C_0]}{u_{p+1}} F(C_0)^p.$$

The conclusion follows since the right-hand-side of the last inequality converges to 1.  $\square$

We now ready to prove Proposition 6.3. This is where the UDD dependence assumption is used.

*Proof (Proof of Proposition 6.3).* Define the canonical (pseudo) metric on  $S_p$ ,

$$d(i, j) = \sqrt{\mathbb{E}[(\epsilon(i) - \epsilon(j))^2]}.$$

It can be easily checked that the canonical metric takes values between 0 and 2. For arbitrary  $\delta \in (0, 1)$ , take  $\gamma = \sqrt{2(1 - \delta)}$ , and let  $\mathcal{N}$  be a  $\gamma$ -packing of  $S_p$ . That is, let  $\mathcal{N}$  be a subset of  $S_p$ , such that for any  $i, j \in \mathcal{N}$ ,  $i \neq j$ , we have  $d(i, j) \geq \gamma$ , i.e.,

$$d(i, j) = \sqrt{2(1 - \Sigma_p(i, j))} \geq \gamma = \sqrt{2(1 - \delta)}, \quad (6.40)$$

or equivalently,  $\Sigma_p(i, j) \leq \delta$ . We claim that we can find a  $\gamma$ -packing  $\mathcal{N}$  whose number of elements is at least

$$|\mathcal{N}| \geq q/N(\delta). \quad (6.41)$$

Indeed,  $\mathcal{N}$  can be constructed iteratively as follows:

- 1: Set  $S_p^{(1)} := S_p$  and  $\mathcal{N} := \{j_1\}$ , where  $j_1 \in S_p^{(1)}$  is an arbitrary element. Set  $k := 1$ .
- 2: Set  $S_p^{(k+1)} := S_p^{(k)} \setminus B_\gamma(j_k)$ , where

$$B_\gamma(j_k) := \{i \in S_p : d(i, j_k) < \gamma \equiv \sqrt{2(1 - \delta)}\}.$$

- 3: If  $S_p^{(k)} \neq \emptyset$ , pick an arbitrary  $j_{k+1} \in S_p^{(k)}$ , set  $\mathcal{N} := \mathcal{N} \cup \{j_{k+1}\}$ , and  $k := k + 1$ , go to step 2; otherwise, stop.

By the definition of UDD (see Definition 6.1), there are at most  $N(\delta)$  coordinates whose covariance with  $\epsilon(j)$  exceed  $\delta$ . Therefore at each iteration,  $|B_\gamma(j_k)| \leq N(\delta)$ ,

and hence

$$\left| S_p^{(k+1)} \right| \geq \left| S_p^{(k)} \right| - |B_\gamma(j_k)| \geq q - kN(\delta).$$

The construction can continue for at least  $q/N(\delta)$  iterations, and we have  $|\mathcal{N}| \geq \lfloor q/N(\delta) \rfloor$  as desired.

Now we define on this  $\gamma$ -packing  $\mathcal{N}$  an independent Gaussian process  $(\eta(j))_{j \in \mathcal{N}}$ ,

$$\eta(j) = \frac{\gamma}{\sqrt{2}} Z(j) \quad j \in \mathcal{N},$$

where  $Z(j)$ 's are iid standard normal random variables. Observe that by the definition of  $\gamma$ -packing in (6.40), the increments of the new process are smaller than those of the original process in the following sense,

$$\mathbb{E} [(\eta(i) - \eta(j))^2] = \gamma^2 \leq d^2(i, j) = \mathbb{E} [(\epsilon(i) - \epsilon(j))^2]$$

for all  $i \neq j, i, j \in \mathcal{N}$ . Applying the Sudakov-Fernique inequality (see, e.g., Theorem 2.2.3 in (Adler and Taylor, 2009)) to  $(\eta(j))_{j \in \mathcal{N}}$  and  $(\epsilon(j))_{j \in \mathcal{N}}$ , we have

$$\mathbb{E} \left[ \max_{j \in \mathcal{N}} \eta(j) \right] \leq \mathbb{E} \left[ \max_{j \in \mathcal{N}} \epsilon(j) \right] \leq \mathbb{E} \left[ \max_{j \in S_p} \epsilon(j) \right]. \quad (6.42)$$

Since the  $(\eta(j))_{j \in \mathcal{N}}$  are independent Gaussians, Lemma 6.5 yields the lower bound,

$$\liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{\max_{j \in \mathcal{N}} \eta(j)}{u_{|\mathcal{N}|}} \right] \geq \frac{\gamma}{\sqrt{2}} = \sqrt{1 - \delta}. \quad (6.43)$$

Using the expressions (2.33) for the quantiles of AGG models (with  $\nu = 2$  here), we have

$$\frac{u_{|\mathcal{N}|}}{u_q} \geq \left( \frac{\log q - \log N(\delta)}{\log q} \right)^{1/2} (1 + o(1)) \rightarrow 1, \quad (6.44)$$

since  $N(\delta)$  does not depend on  $q = q(p) \rightarrow \infty$ , and that  $|\mathcal{N}| \geq q/N(\delta)$ .

By combining (6.42), (6.43) and (6.44), we conclude that

$$\begin{aligned} \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{\max_{j \in S_p} \epsilon(j)}{u_q} \right] &\geq \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{\max_{j \in \mathcal{N}} \eta(j)}{u_q} \right] && \text{by (6.42)} \\ &\geq \liminf_{p \rightarrow \infty} \mathbb{E} \left[ \frac{\max_{j \in \mathcal{N}} \eta(j)}{u_{|\mathcal{N}|}} \right] && \text{by (6.44)} \\ &\geq \sqrt{1 - \delta}. && \text{by (6.43)} \end{aligned}$$

Since  $\delta > 0$  is arbitrary, (6.36) follows as desired.  $\square$

## 6.4 Numerical illustrations of exact support recovery under dependence

The characterization of URS with the UDD condition allows us to simulate Gaussian errors and illustrate the effect of dependence on the phase transition behavior in finite dimensions. We shall compare the performance of the Bonferroni's procedure, which is agnostic to both sparsity and signal size, with the oracle procedure which picks the top- $s$  observations.

The first set of experiments explores short-range dependent errors from an autoregressive (AR) models.

- AR(1) Gaussian errors with parameter  $\rho = -0.5$ ,  $\rho = 0.5$ , and  $\rho = 0.9$ , where the autocovariance functions decay exponentially,  $\rho_k = \rho^k$ .

We again apply both the sparsity- and signal-size agnostic Bonferroni's procedure, i.e.,  $\tilde{S} = \{i : x(i) > \sqrt{2 \log p}\}$ , as well as the oracle procedure  $\hat{S}^* = \{i : x(i) \geq x_{[s]}\}$ ,  $s = |S|$ , to all settings. Results of the numerical experiments for the AR models are shown in Figure 6.1.

As was commented in the main text, for dependent errors the oracle procedures is able to recover support of signals with higher probability than the Bonferroni procedures in finite dimensions; compare left and right columns of Figure 6.1. Short range dependent observations, however, there is not a pronounced difference. The results of the experiments are very similar to that of the independent Gaussian case.

The second set of experiments explores exact support recovery in additive error models in the cases of long-range dependent but UDD, as well as non-UDD errors. In particular we simulate

- Fractional Gaussian noise (fGn) with Hurst parameter  $H = 0.75$  and  $H = 0.9$ . The autocovariance functions are

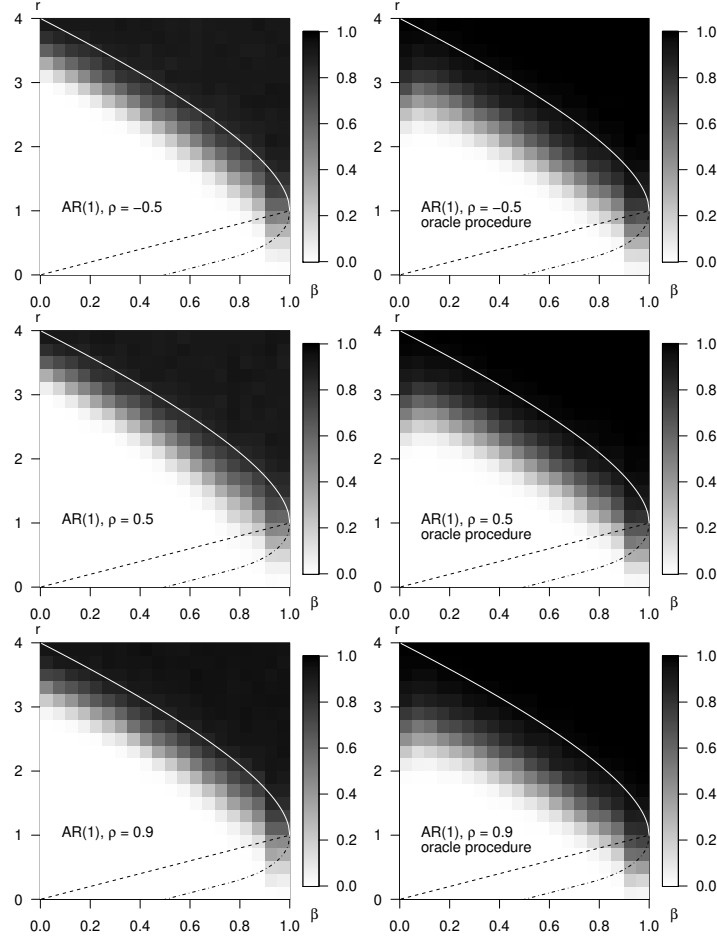
$$\rho_k \sim 0.75k^{-0.6} \quad \text{and} \quad \rho_k \sim 1.44k^{-0.2},$$

as  $k \rightarrow \infty$ . Both fGn models represent the regime of long-range dependence, where covariances decay very slowly to zero, so that  $\sum |\rho_k| = \infty$ ; see, e.g., (Taqq, 2003). Observe that every stationary Gaussian process with vanishing autocovariance gives rise to an UDD array as concluded in Corollary 6.1.

- The non-UDD Gaussian errors described in Example 6.1.

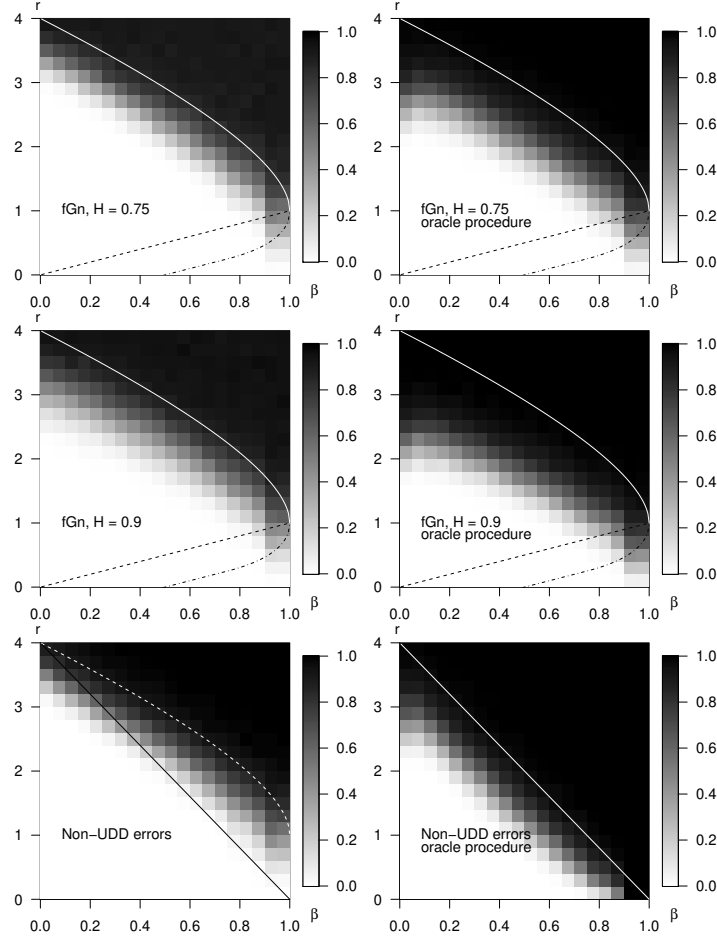
We will apply both the sparsity-and-signal-size-agnostic Bonferroni's procedure, i.e.,  $\tilde{S} = \{i : x(i) > \sqrt{2 \log p}\}$ , as well as the oracle procedure  $\hat{S}^* = \{i : x(i) \geq x_{[s]}\}$ ,  $s = |S|$ , to all settings. Results of the numerical experiments for the fGn and non-UDD models are shown in Figure 6.2.

Notice that the oracle procedure sets its thresholds more aggressively (at roughly  $\sqrt{2 \log s}$ ) than the Bonferroni procedure (at  $\sqrt{2 \log p}$ ). Although this difference vanishes as  $p \rightarrow \infty$ , in finite dimensions ( $p = 10\,000$ ) the advantage can be felt. Indeed, in all our experiments the oracle procedure is able to recover support of signals with higher probability than the Bonferroni procedures; compare left and right columns



**Fig. 6.1** The empirical probability of exact support recovery from numerical experiments, as a function of sparsity level  $\beta$  and signal sizes  $r$ . Darker colors indicate higher probability of exact support recovery. Three AR(1) models with autocorrelation functions  $(-0.5)^k$  (upper),  $0.5^k$  (middle), and  $0.9^k$  (lower) are simulated. The experiments were repeated 1000 times for each sparsity-signal size combination. In finite dimensions ( $p = 10000$ ), the Bonferroni procedures (left) suffers small loss of power compared to the oracle procedures (right). A phase transition in agreement with the predicted boundary (4.5) can be seen in the AR models. The boundaries (solid, dashed, and dash-dotted lines) are as in Fig 4.1.





**Fig. 6.2** The empirical probability of exact support recovery from numerical experiments, as a function of sparsity level  $\beta$  and signal sizes  $r$ . Darker colors indicate higher probability of exact support recovery. Two fGn models with Hurst parameter  $H = 0.75$  (upper),  $H = 0.9$  (middle), and the non-UDD errors in Example 6.1 (lower) are simulated. The experiments were repeated 1000 times for each sparsity-signal size combination. In finite dimensions ( $p = 10000$ ), the oracle procedures (right) is able to recover support for weaker signals than the Bonferroni procedures (left) when errors are heavily dependent, although they have the same phase transition limit. The non-UDD errors demonstrate qualitatively different behavior, enabling support recovery for strictly weaker signals. The boundaries (solid, dashed, and dash-dotted lines) are as in Fig 4.1. In the non-UDD example, dashed lines represent the limit attained by Bonferroni's procedures. See text for additional comments.

of Figure 6.2. Notice also that there is an increase in probability of recovery near  $\beta = 0$  for oracle procedures. This is an artifact in finite dimensions due to the fact that  $s = \lfloor p^{1-\beta} \rfloor < p/2$ , and there are more signals than nulls. The oracle procedures are able to adjust to this reversal by lowering its threshold accordingly.

For UDD errors, Theorem 4.2 predicts that exact recovery of the support is impossible when signal sizes are below the boundary (4.5), even with oracle procedures. However, the rate of this convergence (i.e.,  $\mathbb{P}[\widehat{S}^* = S] \rightarrow 0$  or 1) can be very slow when the errors are heavily dependent, even though all AR and fGn models demonstrate qualitatively the same behavior in line with the predicted boundary (4.5). In finite dimensions ( $p = 10\,000$ ), as dependence in the errors increases (fGN(H=0.75) to fGN(H=0.9)), the oracle procedure becomes more powerful at recovering signal support with high probability for weaker signals.

On the other hand, as demonstrated in Example 6.1, non-UDD errors yield qualitatively different behavior; exact support recovery is possible for signal sizes strictly weaker than that in the UDD case. Lower-right panel of Figure 6.2 demonstrates in this example that the signal support can be recovered as long as the signal sizes are larger than  $4(1 - \beta)$ .

For completeness, we verify the claims in Example 6.1.

*Proof (Proof of claims in Example 6.1).* Recall that  $\widehat{S}^* = \{j : x(j) > t_p^*\}$ , where  $t_p^* = \sqrt{2(1 - \beta) \log p}$ . Analogous to (4.11) in the proof of Theorem 4.1, we have

$$\begin{aligned} \mathbb{P}[\widehat{S} \subseteq S] &= 1 - \mathbb{P}\left[\max_{j \in S^c} x(j) > t_p^*\right] = 1 - \mathbb{P}\left[\max_{j \in S^c} \epsilon(j) > t_p^*\right] \\ &\geq 1 - \mathbb{P}\left[\max_{j \in \{1, \dots, p\}} \epsilon(j) > t_p^*\right] \geq 1 - \mathbb{P}\left[\max_{j \in \{1, \dots, \lfloor p^{1-\beta} \rfloor\}} \widetilde{\epsilon}(j) > t_p^*\right] \end{aligned}$$

where  $(\widetilde{\epsilon})_{j=1}^{\lfloor p^{1-\beta} \rfloor}$ 's are independent Gaussian errors; in the last inequality we used the assumption that there are at most  $\lfloor p^{1-\beta} \rfloor$  independently distributed Gaussian errors in  $(\epsilon_p(j))_{j=1}^p$ . By Example 4.1 (with  $\lfloor p^{1-\beta} \rfloor$  taking the role of  $p$ ), we know that the FWER goes to 0 at a rate of  $(2 \log \lfloor p^{1-\beta} \rfloor)^{-1/2}$ . Therefore, the probability of no false inclusion converges to 1.

On the other hand, since the signal sizes are no smaller than  $(\nu \underline{r} \log p)^{1/\nu}$ , similar to (4.13), we obtain

$$\begin{aligned} \mathbb{P}[\widehat{S} \supseteq S] &\geq \mathbb{P}\left[\min_{j \in S} \epsilon(j) > \sqrt{2(1 - \beta) \log p} - \sqrt{2\underline{r} \log p}\right] \\ &= \mathbb{P}\left[\max_{j \in S} (-\epsilon(j)) < \sqrt{2 \log p} \left(\sqrt{\underline{r}} - \sqrt{1 - \beta}\right)\right] \\ &= \mathbb{P}\left[\frac{\max_{j \in S} (-\epsilon(j))}{u_{|S|}} < \frac{\sqrt{\underline{r}} - \sqrt{1 - \beta}}{\sqrt{1 - \beta}} (1 + o(1))\right], \end{aligned} \quad (6.45)$$

where in the last line we used the quantiles (2.33). Since the minimum signal size is bounded below by  $\underline{r} > 4(1 - \beta)$ , the right-hand-side of the inequality in (6.45) converges to a constant strictly larger than 1. While the left-hand-side, by Slepian's Lemma Slepian (1962), is stochastically smaller than a r.v. going to 1, i.e.,

$$\frac{1}{u_{|S|}} \max_{j \in S} (-\epsilon(j)) \stackrel{d}{\leq} \frac{1}{u_{|S|}} \max_{j \in S} \epsilon^*(j) \xrightarrow{\mathbb{P}} 1, \quad (6.46)$$

where  $(\epsilon^*)_{j=1}^{\lfloor D^{1-\beta} \rfloor}$ 's are independent Gaussian errors. Therefore the probability in (6.45) must also converge to 1.  $\square$



## Chapter 7

# Fundamental Statistical Limits in Genome-wide Association Studies

We investigate the fundamental limits of multiple testing problems in high-dimensional chi-square models, and in genome-wide association studies as introduced in Section 1.2.

In Section 7.1, we shall establish the phase transitions of the sparse chi-square model (1.3). Recall that in large-scale screening studies where a large number of association tests are conducted, resulting statistics may be approximated by

$$x(i) \sim \chi_\nu^2(\lambda(i)), \quad i = 1, \dots, p, \quad (7.1)$$

where  $\chi_\nu^2(\lambda(i))$  is a chi-square distributed random variable with  $\nu$  degrees of freedom and non-centrality parameter  $\lambda(i)$ . In parallel to results in Chapter 3, we show that several commonly used family-wise error rate-control procedures — including Bonferroni’s procedure — are asymptotically optimal for the exact, and exact-approximate support recovery problems (as defined in Definition 2.5) in idealized chi-square models (7.1) with independent components. We further show that the BH procedure is asymptotically optimal for the approximate, and approximate-exact support recovery problems. Under appropriate parametrizations of the signal sizes and sparsity, they establish the phase transitions of support recovery problems in the chi-square model. Remarkably, the degree-of-freedom parameter does not affect the asymptotic boundaries in any of the four support recovery problems.

All phase transition boundaries coincide with those in the additive error models obtained in Chapter 3 under suitable parametrizations. This indicates vanishing differences between the difficulties of the one-sided and two-sided alternatives in the Gaussian additive error model (1.1).

We then return to association screenings of categorical variables in Section 7.2, and present the consequences of the phase transition in the exact-approximate problem in large-scale genetic association studies. We do so by characterizing the relationship between the signal size  $\lambda$  and the marginal frequencies, odds ratio, and sample sizes for association tests on 2-by-2 contingency tables. This result, establishing the relationship between sample sizes and signal sizes, is made precise in Section 7.2.

We elaborate on the implications of this relationship on optimal study designs for association studies in Section 7.3. Perhaps surprisingly, our analysis reveals that balanced designs with equal number of cases and controls are often statistically inefficient. Practical consequences of these results in power analysis will be illustrated with data examples in Section 7.4.

The phase transitions in the chi-square models are demonstrated with numerical simulations in Section 7.5. Proofs of results in this Chapter are collected in Section 7.6.

## 7.1 Support recovery problems in chi-squared models

Similar to the analysis of additive error models in Chapter 3, we will work with triangular arrays of chi-square models (1.3) indexed by  $p$ . We adopt the same parametrization for the sparsity of the non-centrality parameter vectors  $\lambda = \lambda_p$ ,

$$|S_p| = \lfloor p^{1-\beta} \rfloor, \quad \beta \in (0, 1] \quad (7.2)$$

where  $\beta$  parametrizes the problem sparsity. The closer  $\beta$  is to 1, the sparser the support  $S_p$ ; conversely, when  $\beta$  is close to 0, the support is dense with many non-null signals.

We parametrize the range of the non-zero and perhaps unequal signals in the chi-square model with

$$\underline{\Delta} = 2\underline{r} \log p \leq \lambda(i) \leq \bar{\Delta} = 2\bar{r} \log p, \quad \text{for all } i \in S_p, \quad (7.3)$$

for some constants  $0 < \underline{r} \leq \bar{r} \leq +\infty$ .

### 7.1.1 The exact support recovery problem

The first main result characterizes the phase transition phenomenon in the exact support recovery problem under the chi-square model.

**Theorem 7.1.** *Consider the high-dimensional chi-squared model (1.3) with signal sparsity and size as described in (7.2) and (7.3). The function*

$$g(\beta) = \left(1 + \sqrt{1 - \beta}\right)^2 \quad (7.4)$$

*characterizes the phase transition of exact support recovery problem. Specifically, if  $\underline{r} > g(\beta)$ , then Bonferroni's, Sidák's, Holm's, and Hochberg's procedures with slowly vanishing (see Definition 3.1) nominal FWER levels all achieve asymptotically exact support recovery in the sense of (2.25).*

Conversely, if  $\bar{r} < g(\beta)$ , then for any thresholding procedure  $\hat{S}_p$ , we have  $\mathbb{P}[\hat{S}_p = S_p] \rightarrow 0$ . Therefore, in view of Lemma 2.1, exact support recovery asymptotically fails for all thresholding procedures in the sense of (2.26).

The procedures listed in Theorem 7.1 were reviewed in Section 2.2. Proof of the theorem can be found in Section 7.6.2.

It is evident that the exact support recovery boundary (7.4) coincides with that in parallel results for the Gaussian additive error models (1.1) in Chapter 3. Implications of these results will be discussed in Section 7.1.5 below.

*Remark 7.1.* Theorem 7.1 predicts that the asymptotic boundaries are the same for all values of the parameter  $\nu$ . In simulations (Section 7.5), we find this asymptotic prediction to be quite accurate for  $\nu \leq 3$  even in moderate dimensions ( $p = 100$ ). For  $\nu > 3$ , the phase transitions take place somewhat above the boundary  $g$ . The behavior is qualitatively similar for the other three phase transitions (see Theorems 7.2, 7.3, and 7.4 below).

### 7.1.2 The exact-approximate support recovery problem

The next theorem describes the phase transition in the exact-approximate support recovery problem.

**Theorem 7.2.** *In the context of Theorem 7.1, the function*

$$\tilde{g}(\beta) = 1 \tag{7.5}$$

*characterizes the phase transition of exact-approximate support recovery problem. Specifically, if  $\underline{r} > \tilde{g}(\beta)$ , then the procedures listed in Theorem 7.1 with slowly vanishing nominal FWER levels achieve asymptotically exact-approximate support recovery in the sense of (2.25).*

*Conversely, if  $\bar{r} < \tilde{g}(\beta)$ , then for any thresholding procedure  $\hat{S}_p$ , the exact-approximate support recovery fails in the sense of (2.26).*

Theorem 7.2 is proved in Section 7.6.4.

### 7.1.3 The approximate support recovery problem

Our third main result characterizes the phase transition phenomenon in the approximate support recovery problem in the chi-square model.

**Theorem 7.3.** *Consider the high-dimensional chi-squared model (1.3) with signal sparsity and size as described in (7.2) and (7.3). The function*

$$h(\beta) = \beta \tag{7.6}$$

characterizes the phase transition of approximate support recovery problem. Specifically, if  $\underline{r} > h(\beta)$ , then the BH procedure  $\hat{S}_p$  (defined in Section 2.2) with slowly vanishing (see Definition 3.1) nominal FDR levels achieves asymptotically approximate support recovery in the sense of (2.25).

Conversely, if  $\bar{r} < h(\beta)$ , then approximate support recovery asymptotically fails in the sense of (2.26) for all thresholding procedures.

Theorem 7.3 is proved in Section 7.6.4 below.

### 7.1.4 The approximate-exact support recovery problem

A counterpart of Theorem 3.5 also holds in the chi-square models.

**Theorem 7.4.** *In the context of Theorem 7.3, the function*

$$\tilde{h}(\beta) = \left( \sqrt{\beta} + \sqrt{1 - \beta} \right)^2 \quad (7.7)$$

characterizes the phase transition of approximate-exact support recovery problem. Specifically, if  $\underline{r} > \tilde{h}(\beta)$ , then the Benjamini-Hochberg procedure with slowly vanishing nominal FDR levels achieves asymptotically approximate-exact support recovery in the sense of (2.25).

Conversely, if  $\bar{r} < \tilde{h}(\beta)$ , then for any thresholding procedure  $\hat{S}_p$ , the approximate-exact support recovery fails in the sense of (2.26).

Theorem 7.4 is proved in Section 7.6.2.

Notice that all phase transitions boundaries are identical to those in the Gaussian additive error model (1.1) under one-side alternative. We refer readers to Figure 3.2 in Section 3.2 for a visualization of the results in Theorems 7.1 through 7.4.

The all four Theorems so far focus only on the idealized models (1.3) where statistics are *independent*. Support recovery problems under dependent observations remain to be explored. Recall in Chapter 3 we showed that the boundary for the exact support recovery problem in the additive error model (1.1) continues to hold even under severe dependence and general distributional assumptions. We conjecture that similar results would also hold, under classes of dependence structures that are “not too different from independence”, in the chi-square models. As an example, in the GWAS application, dependence among the genetic markers at different locations (known as linkage disequilibrium) decay as a function of their physical distances on the genome (Bush and Moore, 2012), resulting in locally dependent test statistics. It would be of great interest to extend the current theory to cover important dependence structures that arise in such applications.



### 7.1.5 Comparison of one- versus two-sided alternatives in additive error models

As alluded to in Section 1.2 in the introduction, we draw explicit comparisons between the one-sided and two-sided alternatives in Gaussian additive error models (1.1).

The exact support recovery problem in the dependent Gaussian additive error model (1.1) was studied in Chapter 3, with parametrization of sparsity identical to that in (7.2), whereas the range of the non-zero (and perhaps unequal) mean shifts  $\mu(i)$  was parametrized as

$$\underline{\Delta} = \sqrt{2\underline{r} \log p} \leq \mu(i) \leq \bar{\Delta} = \sqrt{2\bar{r} \log p}, \quad \text{for all } i \in S_p,$$

for some constants  $0 < \underline{r} \leq \bar{r} \leq +\infty$ . Under this one-sided alternative, a phase transition in the  $r$ - $\beta$  plane was described, where the boundary was found to be identical to (7.4) in Theorem 7.1 for the chi-square models (7.1).

As discussed in Section 1.2, support recovery problems in the chi-square model with  $\nu = 1$  correspond to the support recovery problems in the additive model under two-sided alternatives. This implies that the asymptotic signal size requirements are identical between the two-sided alternative and its one-sided counterpart, in order to achieve exact support recovery. As we shall see in numerical experiments (in Section 7.5 below), the difference is not very pronounced even in moderate dimensions, and vanishes as  $p \rightarrow \infty$ , in accordance with Theorem 7.1.

Comparisons can also be drawn in the approximate, approximate-exact, and exact approximate support recovery problems between the two types of alternatives.

Specifically, the approximate support recovery problem in the Gaussian additive error model (1.1) under one-sided alternatives exhibits a phase transition phenomenon characterized by a boundary that coincides with (7.6) in Theorem 7.3. Similar to the exact support recovery problem, this indicates vanishing difference in the difficulties of the two types alternatives in approximate support recovery problems.

Comparing Theorems 7.2 to 3.4 and Theorems 7.4 to 3.5, we see that the phase transition boundaries under the two types of alternatives are also identical in the exact-approximate and approximate-exact support recovery problems. The additional uncertainty in the two-sided alternatives do not call for larger signal sizes asymptotically in these problems.

To complete the comparisons, we point out that the phase transition boundaries for the sparse signal detection problem in the two types of alternatives are both identical to (3.4). This was analyzed in Donoho and Jin (2004).

## 7.2 Odds ratios and statistical power

We return to the application of association screenings for categorical variables, and put the results in the previous section to use. In particular, we focus on the exact-approximate support recovery problem, and demonstrate the consequences of its phase transition (Theorem 7.2) in genetic association studies.

In order to do so, we must first connect the concept of “statistical signal size”  $\lambda$  with some key quantities in association tests. While “signal size” likely sounds foreign to most practitioners, it is intimately linked with the concept of “effect sizes” — or odds ratios — in association studies, which are frequently estimated and reported in GWAS catalogs. We characterize the relationship between the two quantities in the special, but fairly common case of association tests on 2-by-2 contingency tables in Section 7.2.

Consider a 2-by-2 multinomial distribution with marginal probabilities of phenotypes  $(\phi_1, \phi_2)$  and genotypes  $(\theta_1, \theta_2)$ . The *probability* table (as opposed to the table of multinomial *counts* in the introduction) is as follows.

Probabilities	Genotype		Total by phenotype
	Variant 1	Variant 2	
Cases	$\mu_{11}$	$\mu_{12}$	$\phi_1$
Controls	$\mu_{21}$	$\mu_{22}$	$\phi_2$
Total by genotype	$\theta_1$	$\theta_2$	1

The odds ratio (i.e., “effect size”) is defined as the ratio of the phenotype frequencies between the two genotype variants,

$$R := \frac{\mu_{11}}{\mu_{21}} \bigg/ \frac{\mu_{12}}{\mu_{22}} = \frac{\mu_{11}\mu_{22}}{\mu_{12}\mu_{21}}. \quad (7.8)$$

The multinomial distribution is fully parametrized by the trio  $(\theta_1, \phi_1, R)$ . Odds ratios further away from 1 indicate greater contrasts between the probability of outcomes. Independence between the genotypes and phenotypes would imply an odds ratio of one, and hence  $\mu_{jk} = \phi_j\theta_k$ , for all  $j, k \in \{1, 2\}$ .

For a sequence of local alternatives  $\mu^{(1)}, \mu^{(2)}, \dots$ , such that  $\sqrt{n}(\mu_{jk}^{(n)} - \phi_j\theta_k)$  converges to a constant table  $\delta = (\delta_{jk})$ , the chi-square test statistics converge in distribution to the non-central chi-squared distribution with non-centrality parameter

$$\lambda = \sum_{j=1}^2 \sum_{k=1}^2 \delta_{jk}^2 / (\phi_j\theta_k).$$

See, e.g., Ferguson (2017). Hence, for large samples from a fixed distribution  $(\mu_{ij})$ , the statistic is well approximated by a  $\chi_1^2(\lambda)$  distribution, where

$$\lambda = n \sum_{j=1}^2 \sum_{k=1}^2 \frac{(\mu_{jk} - \phi_j\theta_k)^2}{\phi_j\theta_k}. \quad (7.9)$$

Power calculations therefore only depend on the  $\mu_{jk}$ 's through  $\lambda = nw^2$ , where we define

$$w^2 := \lambda/n \quad (7.10)$$

to be the *signal size per sample*. Statistical power would be increasing in  $w^2$  for fixed sample sizes.

The next proposition states that the statistical signal size per sample can be parametrized by the odds ratio and the marginals in the probability table.

**Proposition 7.1.** *Consider a 2-by-2 multinomial distribution with marginal distributions  $(\phi_1, \phi_2 = 1 - \phi_2)$  and  $(\theta_1, \theta_2 = 1 - \theta_1)$ . Let signal size  $w^2$  be defined as in (7.10), and odds ratio  $R$  be defined as in (7.8). If  $R = 1$ , we have  $w^2 = 0$ ; if  $R \in (0, 1) \cup (1, +\infty)$ , then we have*

$$w^2(R) = \frac{1}{4A(R-1)^2} \left( B + CR - \sqrt{(B + CR)^2 - 4A(R-1)^2} \right)^2, \quad (7.11)$$

where  $A = \phi_1\theta_1\phi_2\theta_2$ ,  $B = \phi_1\theta_1 + \phi_2\theta_2$ , and  $C = \phi_1\theta_2 + \phi_2\theta_1$ .

Proposition 7.1 is derived in Section 7.6.5.

To understand Proposition 7.1, we illustrate Relation (7.11) for selected values of marginals  $\theta_1$  and  $\phi_1$  in Figure 7.1. Observe in the figure that an odds ratio further away from one corresponds to stronger statistical signal per sample, ceteris paribus. However, this “valley” pattern is in general not symmetric around 1, except for balanced marginal distributions ( $\phi_1 = 1/2$  or  $\theta_1 = 1/2$ ). While the odds ratio  $R$  can be arbitrarily close to 0 or diverge to  $+\infty$  for any marginal distribution, the signal sizes  $w^2$  are bounded from above by constants that depend only on the marginals.

**Corollary 7.1.** *The signal size as a function of the odds ratio  $w^2(R)$  is decreasing on  $(0, 1)$  and increasing on  $(1, \infty)$ , with limits*

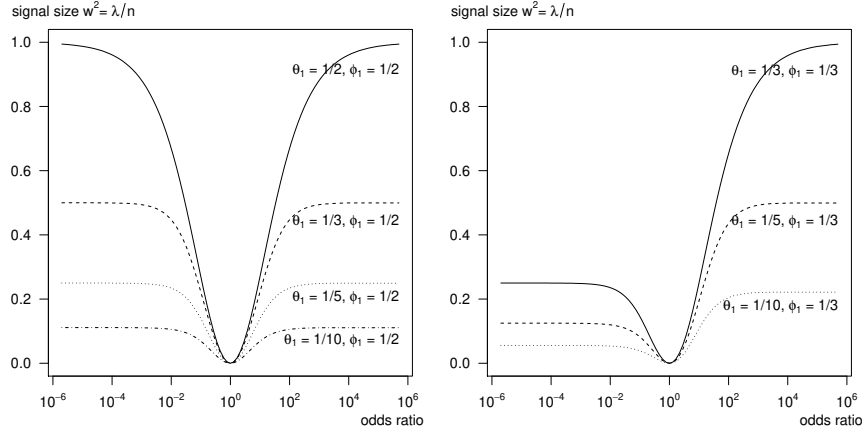
$$\lim_{R \rightarrow 0_+} w^2(R) = \min \left\{ \frac{\phi_1\theta_1}{\phi_2\theta_2}, \frac{\phi_2\theta_2}{\phi_1\theta_1} \right\}, \quad (7.12)$$

and

$$\lim_{R \rightarrow +\infty} w^2(R) = \min \left\{ \frac{\phi_1\theta_2}{\phi_2\theta_1}, \frac{\phi_2\theta_1}{\phi_1\theta_2} \right\}. \quad (7.13)$$

Corollary 7.1 immediately implies that balanced designs with roughly equal number of cases and controls are not necessarily the most informative.

For example, in a study where a third of the recruited subjects carry the genetic variant positively correlated with the trait (i.e.,  $\theta_1 = 1/3$ ), an unbalanced design with  $\phi_1 = 1/3$  would maximize  $w^2$  at large odds ratios. This unbalanced design is much more efficient compared to, say, a balanced design with  $\phi_1 = 1/2$ . In the first case, we have  $w^2 \rightarrow 1$  as  $R \rightarrow \infty$ ; whereas in the second design,  $w^2 < 1/2$  no matter how large  $R$  is. This difference can also be read by comparing the dashed curve ( $\theta_1 = 1/3$ ,  $\phi_1 = 1/2$ ) in the left panel of Figure 7.1, with the solid curve ( $\theta_1 = 1/3$ ,  $\phi_1 = 1/3$ ) in the right panel of Figure 7.1.



**Fig. 7.1** Signal sizes per sample  $w^2$  as functions of odds ratios in 2-by-2 multinomial distributions for selected genotype marginals in balanced (left) and unbalanced (right) designs; see Relation (7.11) in Proposition 7.1. For given marginal distributions, extreme odds ratios imply stronger statistical signals at a given sample size. However, the signal sizes are bounded above by constants that depend on the marginal distributions; see Relations (7.12) and (7.13).

### 7.3 Optimal study designs and rare variants

For a study with a fixed budget, i.e., a fixed total number of subjects  $n$ , the researcher is free to choose the fraction of cases  $\phi_1$  to be included in the study. A natural question is how this budget should be allocated to maximize the statistical power of discovery, or equivalently, the signal sizes  $\lambda = nw^2$ .

In principal, Relation (7.11) can be optimized with respect to the fraction of cases  $\phi_1$  in order to find optimal designs, if  $\theta_1$  is known and held constant. In practice, this is not the case. While the fraction of cases can be controlled, the distributions of genotypes *in the study* are often unknown prior to data collection, and can change with the case-to-control ratio.

Fortunately, the conditional distributions of genotypes in the healthy control groups are often estimated by existing studies, and are made available by consortia such as the NHGRI-EBI GWAS catalog (MacArthur et al., 2016). We denote the conditional frequency of the first genetic variant in the control group as  $(f, 1 - f)$ , where

$$f := \mu_{21}/\phi_2. \quad (7.14)$$

The multinomial probability is fully parametrized by the new trio:  $(f, \phi_1, R)$ .

Probabilities	Genotype		Total by phenotype
	Variant 1	Variant 2	
Cases	$\frac{\phi_1 f R}{f R + 1 - f}$	$\frac{\phi_1 (1 - f)}{f R + 1 - f}$	$\phi_1$
Controls	$f(1 - \phi_1)$	$(1 - f)(1 - \phi_1)$	$1 - \phi_1$

Proposition 7.1 may also be re-stated in terms of the new parametrization.

**Corollary 7.2.** *In the 2-by-2 multinomial distribution with marginals  $(\phi_1, \phi_2 = 1 - \phi_1)$ , and conditional distribution of the variants in the control group  $(f, 1 - f)$ , Relation (7.11) holds with  $\theta_1 = \phi_1 f R / (f R + 1 - f) + f(1 - \phi_1)$  and  $\theta_2 = 1 - \theta_1$ .*

The choice of  $\phi_1$  now has a practical solution.

**Corollary 7.3.** *In the context of Corollary 7.2, the optimal design  $(\phi_1^*, \phi_2^*)$  that maximizes the signal size per sample  $w^2$  is prescribed by*

$$\phi_1^* = \frac{fR + 1 - f}{fR + 1 - f + \sqrt{R}}, \quad \text{and} \quad \phi_2^* = 1 - \phi_1^*. \quad (7.15)$$

Corollary 7.3 is proved in Section 7.6.5.

Of particular interest in the genetics literature are genetic variants with very low allele frequencies in the control group (i.e.,  $f \approx 0$ ), known as rare variants. In such cases, Equation (7.15) can be approximated using the Taylor expansion,

$$\phi_1^* = \frac{1}{1 + \sqrt{R}} + \frac{(R - \sqrt{R})f}{1 + \sqrt{R}} + O(f^2). \quad (7.16)$$

To illustrate, for rare and adversarial factors ( $f \approx 0$  and  $R > 1$ ), the optimal  $\phi_1^*$  is less than  $1/2$ . Therefore, for studies under a fixed budget, controls should constitute the majority of the subjects, in order to maximize power. On the other hand, for rare and protective factors ( $f \approx 0$  and  $R < 1$ ), the optimal  $\phi_1^*$  is greater than  $1/2$ , and cases should be the majority.

## 7.4 Phase transitions in large-scale association screening studies

Returning to the problem of *high-dimensional* marginal screenings for categorical covariates, we explore the manifestation of the phase transition in the exact-approximate support recovery problem in the genetic context.

Recall Theorem 7.2 predicts that FWER and FNR can be simultaneously controlled in large dimensions if and only if

$$r = \frac{\lambda}{2 \log p} = \frac{w^2 n}{2 \log p} > 1. \quad (7.17)$$

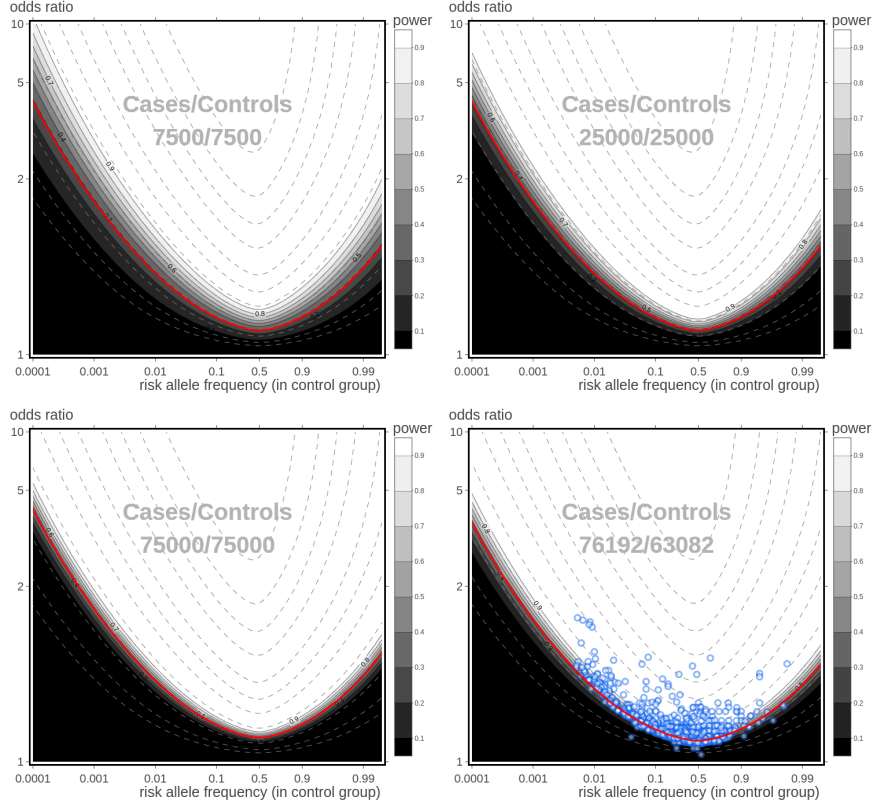
Therefore, if we were to apply FWER-controlling procedures at low nominal levels (say, 5%), then the FNR would experience a phase transition in the sense that, if

$$r > 1 \iff w^2 > \frac{2 \log p}{n}, \quad (7.18)$$

then the FNR can be close to 0; otherwise, FNR must be close to 1.

Using the parametric relationship described in Corollary 7.2 (and Proposition 7.1), the inequalities in (7.18) implicitly define regions of  $(f, R)$  where associations

are discoverable with high power, for a given  $\phi_1$ . Further, the boundary of such discoverable regions sharpens as dimensionality diverges. We illustrate this phase transition through a numerical example next.



**Fig. 7.2** The OR-RAF diagram visualizing the marginal power of discovery in genetic association studies, after applying Bonferroni's procedure with nominal FWER at 5% level. Sample sizes are marked in each panel, and the problem dimensions are, respectively,  $p = 4$  (upper-left),  $p = 10^2$  (upper-right), and  $p = 10^6$  (lower-left), so that  $n/\log p$  are roughly constant. Red curves mark the boundaries ( $r = 1$ ) of the phase transition for the exact-approximate support recovery problem; dashed curves are the equi-signal (equi-power) curves. The phase transition in signal sizes  $\lambda$  translates into the phase transition in terms of  $(f, R)$ , and sharpens as  $p \rightarrow \infty$ ; see Example 7.1. In the lower-right panel, we visualize discovered associations (blue circles) in a recent GWA study (Michailidou et al. (2017)); the estimated odds ratios and risk allele frequencies in a recent GWA study are subject to survival bias and should not be taken at their face values; see Remark 7.2.

*Example 7.1.* Consider association tests on  $2 \times 2$  contingency tables at  $p$  locations as introduced in Section 1.2, where the counts follow a multinomial distribution parametrized by  $(f, R, \phi_1)$  as in Section 7.3. Assume that the phenotype marginals are fixed at  $\phi_1 = \phi_2 = 1/2$ . Applying Bonferroni's procedure with nominal FWER

at  $\alpha = 5\%$  level, we can approximate the marginal power of association tests by

$$\mathbb{P}[\chi_1^2(\lambda) > \chi_{1,\alpha/p}^2], \quad (7.19)$$

where  $\chi_{1,\alpha/p}^2$  is the upper  $(\alpha/p)$ -quantile of a central chi-squared distribution with 1 degree of freedom. We calculate this marginal power as a function of the parameters  $(f, R)$  in three scenarios:

- $p = 4, n = 3 \times 10^4$
- $p = 10^2, n = 1 \times 10^5$
- $p = 10^6, n = 3 \times 10^6$

and visualize the results as heatmaps<sup>1</sup> (referred to as OR-RAF diagrams) in Figure 7.2. These parameter values are chosen so that  $\log p/n$  are roughly constant (around  $4.6 \times 10^{-5}$ ).

We also overlay “equi-signal” curves, i.e., functions implicitly defined by the equations  $r = c$  for a range of  $c$  (dashed curves), and highlight the predicted boundary of phase transition for the exact-approximate support recovery problem  $r = 1$  (red curves). The change in marginal power clearly sharpens around the predicted boundary  $r = 1$  as dimensionality diverges.

*Remark 7.2.* In an attempt to find empirical evidence of our theoretical predictions, we chart the genetic variants associated with breast cancer, discovered in a 2017 study by Michailidou et al. (2017) in an OR-RAF diagram. The estimated risk allele frequencies ( $f$ ) and odds ratios ( $R$ ) are taken from the NHGRI-EBI GWAS catalog MacArthur et al. (2016), and plotted against a power heatmap calculated according to the reported sample sizes. See lower-right panel of Figure 7.2.

It is tempting to believe, on careless inspection, that roughly *all* discovered associations fall inside the high power region of the diagram, therefore demonstrating the phase transition in statistical power. Unfortunately, the estimates here are subject to survival bias — the study in fact uses the same dataset for *both* support estimation and parameter estimation, without adjusting the latter for the selection process. The seemingly striking agreement between the power calculations and the estimated effects of reported associations *should not* be taken as evidence for the validity of our theory. We conjecture, as the theory predicts, that accurate and unbiased parameter estimates from an independent replication will still place the associations in the high power region of the diagram.

Finally, we demonstrate with an example how results in Sections 7.1 and 7.2 may be used for planning prospective association studies.

*Example 7.2.* In a GWAS with  $p = 10^6$  genomic marker locations, researchers wish to locate genetic associations with the trait of interest. Specifically, they wish to

<sup>1</sup> Since genetic variants can always be relabelled such that Variant 1 is positively associated with Cases, we only produce part of the diagram where  $R > 1$ . Sample sizes marked in the figure are adjusted by a factor of 1/2, to reflect the genetic context where a pair of alleles are measured for every individual at every genomic location.

maximize power in the region where genetic variants have risk allele frequencies of 0.01 and odds ratios of 1.2. By Corollary 7.3, the optimal design has a fraction of cases  $\phi^* = 0.478$ , yielding the statistical signal size per sample  $w^2 \approx 9.00 \times 10^{-5}$  according to Corollary 7.2.

If we wish to achieve exact-approximate support recovery in the sense of (2.25), Theorem 7.2 predicts that the signal size parameter  $r$  has to be at least  $\tilde{g}(\beta) = 1$ . This signal size calls for a sample size of  $n = \lambda/w^2 = 2r \log(p)/w^2 \approx 307,011$ . In a typically GWAS, a pair of alleles are sequenced for every marker location, bringing the required number of subjects in the study to  $n/2 \approx 153,509$ .

In comparison, a more accurate power calculation directly using (7.19) predicts that  $n/2 = 165,035$  subjects are needed, under the set of parameters ( $p = 10^6$ ,  $f = 0.01$ ,  $R = 1.2$ ) and  $\text{FWER} = 0.05$ ,  $\text{FNR} = 0.5$ ; this is 7% higher than our crude asymptotic approximation. In general, we recommend using the more precise calculations over the back-of-the-envelope asymptotics for planning prospective studies and performing systematic reviews; a user-friendly web application implementing the more precise approximations is provided in Gao et al. (2019). Nevertheless, the theoretical results on phase transitions generate simple, accurate, and powerful insights that cannot be easily derived from numerical calculations.

## 7.5 Numerical illustrations of the phase transitions in chi-square models

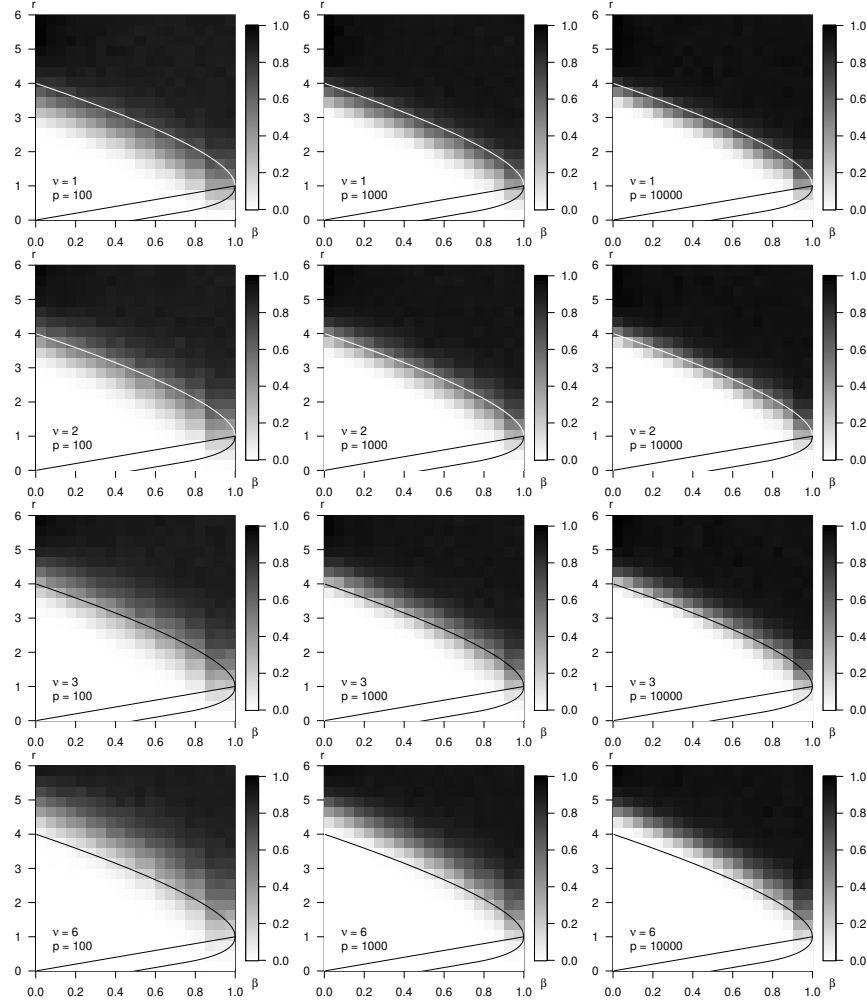
We illustrate with simulations the phase transition phenomena in the chi-square model, and compare numerically the required signal sizes in support recovery problems between the two types of alternatives in the additive error model.

### 7.5.1 Exact support recovery

The sparsity of the signal vectors in the experiments are parametrized as in (7.2). Signal sizes are assumed equal with magnitude  $\lambda(i) = 2r \log p$  for  $i \in S$ . We estimate the support set  $S$  using Bonferroni's procedure with nominal FWER level set at  $1/(5 \log p)$ . The nominal FWER levels vanishes slowly, in line with the assumptions in Theorem 7.1. Experiments were repeated 1000 times at each of the 400 sparsity-signal-size combinations, for dimensions  $p = 10^2, 10^3$ , and  $10^4$ .

The empirical probabilities of exact support recovery under Bonferroni's procedure are shown in Figure 7.3. The numerical results suggest not only good accuracy of the predicted boundaries in high-dimensions ( $p = 10^4$ , right panels of Figure 7.3), but also practical relevance of the theoretical predictions in moderate dimensions ( $p = 100$ , left panels of Figure 7.3).

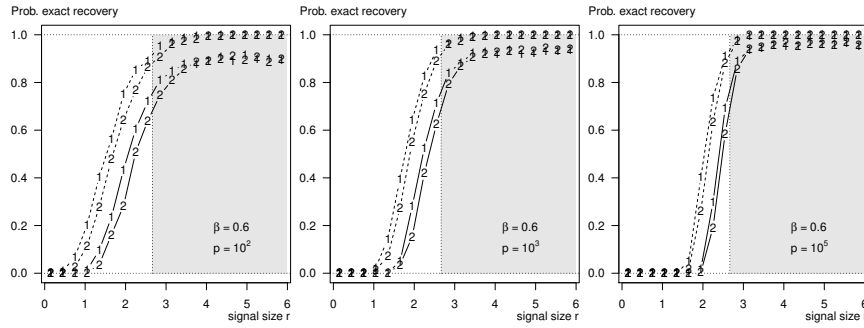




**Fig. 7.3** The empirical probability of exact support recovery of Bonferroni's procedure in the chi-squared model (1.3). We simulate  $\nu = 1, 2, 3, 6$  (first to last row), at dimensions  $p = 10^2, 10^3, 10^4$  (left to right column), for a grid of sparsity levels  $\beta$  and signal sizes  $r$ . The experiments were repeated 1000 times for each sparsity-signal size combination; darker color indicates higher probability of exact support recovery. Numerical results are in general agreement with the boundaries described in Theorem 7.1; for large  $\nu$ 's, the phase transitions take place somewhat above the predicted boundaries. The boundary for the approximate support recovery (Theorem 7.3) and the detection boundary (see Donoho and Jin (2004)) are plotted for comparison.

We conduct further experiments to examine the optimality claims in Theorem 7.1 by comparing with the oracle procedure with thresholds  $t_p = \min_{i \in S} x(i)$ . We also examine the claims in Section 7.1.5, and compare the one-sided alternatives in Gaussian additive models with the two-sided alternatives (or equivalently, the chi-square model with  $\nu = 1$ ). We apply Bonferroni's procedure and the oracle thresholding procedure in both settings.

Experiments were repeated 1000 times for a grid of signal size values ranging from  $r = 0$  to 6, and for dimensions  $10^2, 10^3$ , and  $10^5$ . Results of the experiments, shown in Figure 7.4, suggest vanishing difference between difficulties of two-sided vs one-sided alternatives in the additive error models, as well as vanishing difference between the powers of Bonferroni's procedures and the oracle procedures as  $p \rightarrow \infty$ .



**Fig. 7.4** The empirical probability of exact support recovery of Bonferroni's procedure (solid curves) and the oracle procedure (dashed curves) in the chi-squared model with one degree of freedom (marked '2') in the additive Gaussian error model and under one-sided alternatives (marked '1'). We simulate at dimensions  $p = 10^2, 10^3, 10^5$  (left to right) for a grid of signal sizes  $r$  and sparsity level  $\beta = 0.6$ . The experiments were repeated 1000 times for each method-model-signal-size combination. Numerical results show evidence of convergence to the 0-1 law as predicted by Theorem 7.1; regions where asymptotically exact support recovery can be achieved are shaded in grey. The difference in power between Bonferroni's procedure and the oracle procedure, as well as in the two types of alternatives both decrease as dimensionality increases.

### 7.5.2 Approximate, and approximate-exact support recovery

Similar experiments are conducted to examine the optimality claims in Theorem 7.3, and in Section 7.1.5. We define an oracle thresholding procedure for approximate support recovery, where the threshold is chosen to minimize the empirical risk. That is,

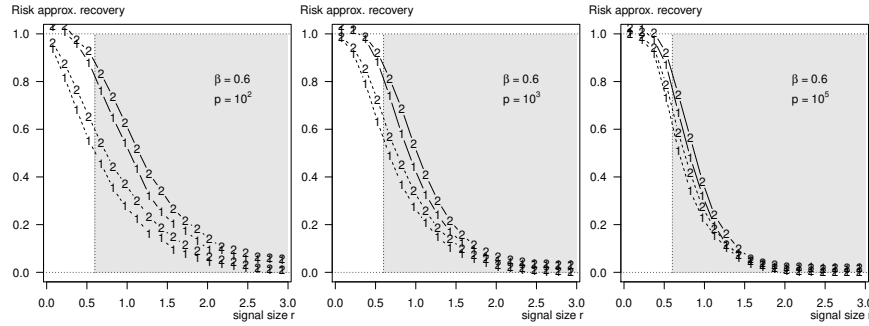
$$t_p(x, S) \in \arg \min_{t \in \mathbb{R}} \frac{|\widehat{S}(t) \setminus S|}{\max\{|\widehat{S}(t)|, 1\}} + \frac{|S \setminus \widehat{S}(t)|}{\max\{|S|, 1\}},$$

where  $\widehat{S}(t) = \{i \mid x(i) \geq t\}$ ; in implementation, we only need to scan the values of observations  $t \in \{x(1), \dots, x(p)\}$ . The nominal FDR level for the BH procedure is set at  $1/(5 \log p)$ , therefore slowly vanishing, in line with the assumptions in Theorem 7.3; all other parameters are identical to that in the experiments for exact support recovery. Results of the experiments are shown in Figure 7.5 and Figure 7.6.

We also examine the boundary described in Theorem 7.2. Experimental settings are identical to that in the experiments for approximate support recovery. We compare the performance of the BH procedure with an oracle procedure with threshold

$$t_p(x, S) \in \min_{i \in S} x(i),$$

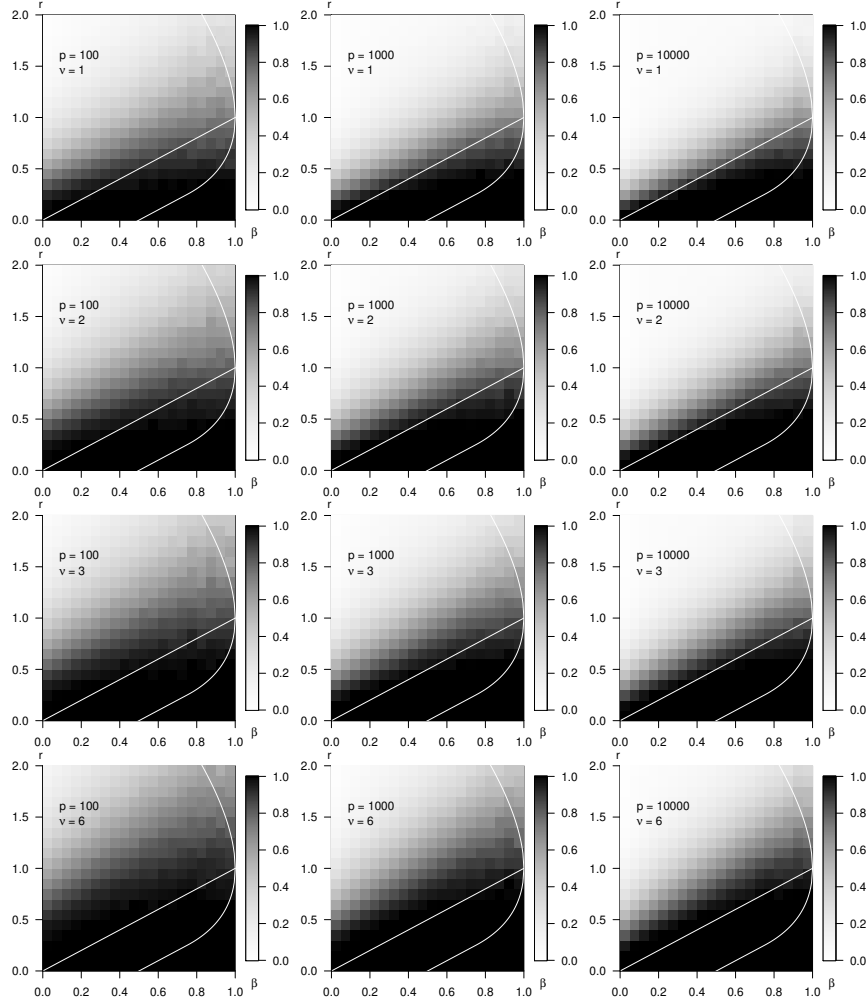
and visualize results of the experiments in Figure 7.7. Notice that the BH procedure sets its threshold somewhat higher than the oracle, especially for small  $\beta$ 's. The empirical risk of the oracle procedure (not shown here in the interest of space) follows much more closely the predicted boundary (7.7).



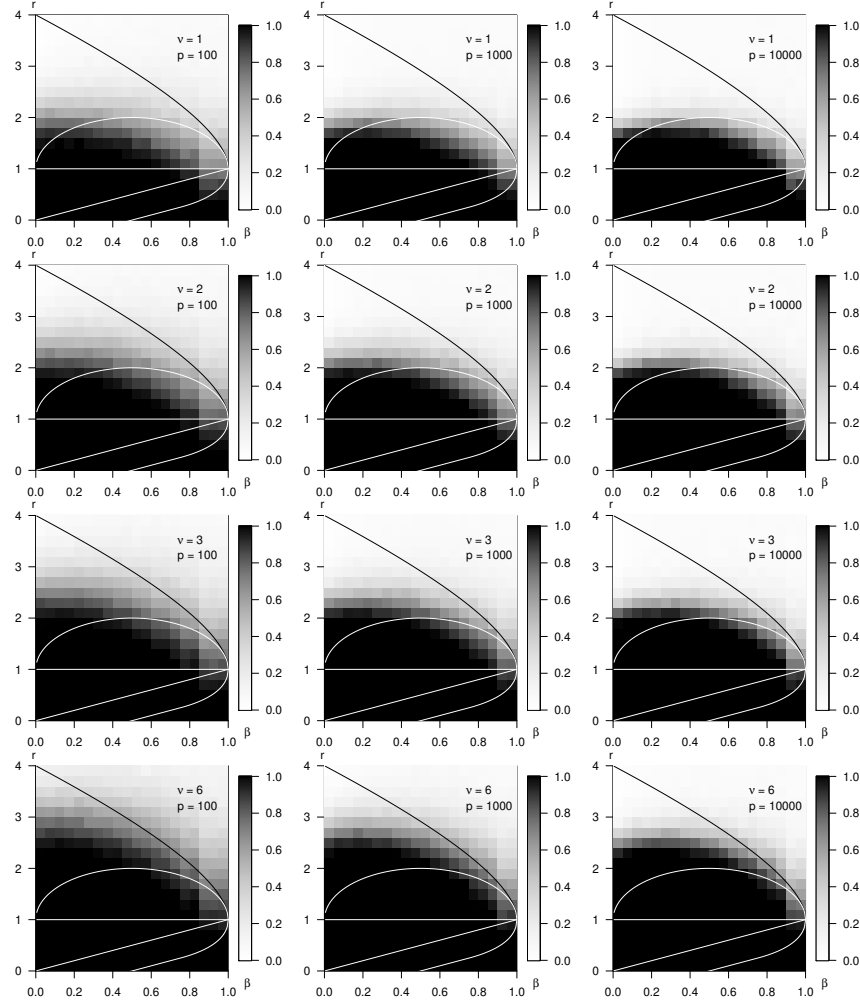
**Fig. 7.5** The empirical risk of approximate support recovery of Benjamini-Hochberg's procedure (solid curves) and the oracle procedure (dashed curves) in the chi-squared model with one degree of freedom (marked '2') and in the additive Gaussian error model under one-sided alternatives (marked '1'). We simulate at dimensions  $p = 10^2, 10^3, 10^5$  (left to right) for a grid of signal sizes  $r$  and sparsity level  $\beta = 0.6$ . The experiments were repeated 1000 times for each method-model-signal-size combination. Numerical results show evidence of convergence to the 0-1 law as predicted by Theorem 7.3; regions where asymptotically approximate support recovery can be achieved are shaded in grey. The difference in risks between Benjamini-Hochberg's procedure and the oracle procedure, as well as in the two types of alternatives, both decrease as dimensionality increases.

## 7.6 Proofs

We review some properties of the chi-square distributions in Section 7.6.1, before presenting the proofs of the main theorems on phase transitions in Sections 7.6.2,



**Fig. 7.6** The estimated risk of approximate support recovery risk<sup>A</sup> (see (2.7)) of the Benjamini-Hochberg procedure in the chi-squared model (1.3). We simulate  $\nu = 1, 2, 3, 6$  (first to last row), at dimensions  $p = 10^2, 10^3, 10^4$  (left to right column), for a grid of sparsity levels  $\beta$  and signal sizes  $r$ . The experiments were repeated 1000 times for each sparsity-signal size combination; darker color indicates higher larger risk<sup>A</sup>. Numerical results are generally in agreement with the boundaries described in Theorem 7.3; for large  $\nu$ 's, the phase transitions take place somewhat above the predicted boundaries. The boundary for the exact support recovery problem (Theorem 7.1) and the detection boundary (see Donoho and Jin (2004)) are plotted for comparison.



**Fig. 7.7** The estimated risk of approximate-exact support recovery risk<sup>EA</sup> (see (2.12)) of the Benjamini-Hochberg procedure in the chi-squared model (1.3). We simulate  $\nu = 1, 2, 3, 6$  (first to last row), at dimensions  $p = 10^2, 10^3, 10^4$  (left to right column), for a grid of sparsity levels  $\beta$  and signal sizes  $r$ . The experiments were repeated 1000 times for each sparsity-signal size combination; darker color indicates higher larger risk<sup>EA</sup>. Numerical results are generally in agreement with the boundaries described in Theorem 7.4; for small  $\beta$ 's and large  $\nu$ 's, the phase transitions take place somewhat above the predicted boundaries. Other boundaries in the support recovery and the detection problems are plotted for comparison.

7.6.3, and 7.6.4. Results relating signal sizes and effect sizes in association tests will be justified in Section 7.6.5.

### 7.6.1 Auxiliary facts of chi-square distributions

We shall recall, and establish, some auxiliary facts about chi-square distributions. These facts will be used in the proofs of Theorem 7.1 and Theorem 7.3.

**Lemma 7.1 (Rapid variation of chi-square distribution tails).** *The central chi-square distribution with  $\nu$  degrees of freedom has rapidly varying tails. That is,*

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[\chi_\nu^2(0) > tx]}{\mathbb{P}[\chi_\nu^2(0) > x]} = \begin{cases} 0, & t > 1 \\ 1, & t = 1 \\ \infty, & 0 < t < 1 \end{cases}, \quad (7.20)$$

where we overloaded the notation  $\chi_\nu^2(0)$  to represent a random variable with the chi-square distribution.

*Proof (Proof of Lemma 7.1).* When  $\nu = 1$ , the chi-square distribution reduces to a squared Normal, and (7.20) follows from the rapid variation of the standard Normal distribution. For  $\nu \geq 2$ , we recall the following bound on tail probabilities (see, e.g., (Ingolot, 2010)),

$$\frac{1}{2}\mathcal{E}_\nu(x) \leq \mathbb{P}[\chi_\nu^2(0) > x] \leq \frac{x}{(x - \nu + 2)\sqrt{\pi}}\mathcal{E}_\nu(x), \quad \nu \geq 2, \ x > \nu - 2,$$

where  $\mathcal{E}_\nu(x) = \exp\{-\frac{1}{2}[(x - \nu - (\nu - 2)\log(x/\nu) + \log \nu)]\}$ . Therefore, we have

$$\frac{(x - \nu + 2)\sqrt{\pi}}{2x} \frac{\mathcal{E}_\nu(tx)}{\mathcal{E}_\nu(x)} \leq \frac{\mathbb{P}[\chi_\nu^2(0) > tx]}{\mathbb{P}[\chi_\nu^2(0) > x]} \leq \frac{2tx}{(tx - \nu + 2)\sqrt{\pi}} \frac{\mathcal{E}_\nu(tx)}{\mathcal{E}_\nu(x)},$$

where  $\mathcal{E}_\nu(tx)/\mathcal{E}_\nu(x) = \exp\{-\frac{1}{2}[(t-1)x - (\nu-2)\log t]\}$  converges to 0 or  $\infty$  depending on whether  $t > 1$  or  $0 < t < 1$ . The case where  $t = 1$  is trivial.  $\square$

Lemma 7.1 and Proposition 2.2 yield the following Corollary.

**Corollary 7.4.** *Maxima of independent observations from central chi-square distributions with  $\nu$  degrees of freedom are relatively stable. Specifically, let  $\epsilon_p = (\epsilon_p(i))_{i=1}^p$  be independently and identically distributed (iid)  $\chi_\nu^2(0)$  random variables. Then the triangular array  $\mathcal{E} = \{\epsilon_p, p \in \mathbb{N}\}$  has relatively stable (RS) maxima in the sense of (2.38).*

**Lemma 7.2 (Stochastic monotonicity).** *The non-central chi-square distribution is stochastically monotone in its non-centrality parameter. Specifically, for two non-central chi-square distributions both with  $\nu$  degrees of freedom, and non-centrality*

parameters  $\lambda_1 \leq \lambda_2$ , we have  $\chi_\nu^2(\lambda_1) \stackrel{d}{\leq} \chi_\nu^2(\lambda_2)$ . That is,

$$\mathbb{P}[\chi_\nu^2(\lambda_1) \leq t] \geq \mathbb{P}[\chi_\nu^2(\lambda_2) \leq t], \quad \text{for any } t \geq 0. \quad (7.21)$$

where we overloaded the notation  $\chi_\nu^2(\lambda)$  to represent a random variable with the chi-square distribution with non-centrality parameter  $\lambda$  and degree-of-freedom parameter  $\nu$ .

*Proof (Proof of Lemma 7.2).* Recall that non-central chi-square distributions can be written as sums of  $\nu - 1$  standard normal random variables and a non-central normal random variable with mean  $\sqrt{\lambda}$  and variance 1,

$$\chi_\nu^2(\lambda) \stackrel{d}{=} Z_1^2 + \dots + Z_{\nu-1}^2 + (Z_\nu + \sqrt{\lambda})^2.$$

Therefore, it suffices to show that  $\mathbb{P}[(Z + \sqrt{\lambda})^2 \leq t]$  is non-increasing in  $\lambda$  for any  $t \geq 0$ , where  $Z$  is a standard normal random variable. We rewrite this expression in terms of standard normal probability function  $\Phi$ ,

$$\begin{aligned} \mathbb{P}[(Z + \sqrt{\lambda})^2 \leq t] &= \mathbb{P}[-\sqrt{\lambda} - \sqrt{t} \leq Z \leq -\sqrt{\lambda} + \sqrt{t}] \\ &= \Phi(-\sqrt{\lambda} + \sqrt{t}) - \Phi(-\sqrt{\lambda} - \sqrt{t}). \end{aligned} \quad (7.22)$$

The derivative of the last expression (with respect to  $\lambda$ ) is

$$\frac{1}{2\sqrt{\lambda}} \left( \phi(\sqrt{\lambda} + \sqrt{t}) - \phi(\sqrt{\lambda} - \sqrt{t}) \right) = \frac{1}{2\sqrt{\lambda}} \left( \phi(\sqrt{\lambda} + \sqrt{t}) - \phi(\sqrt{t} - \sqrt{\lambda}) \right), \quad (7.23)$$

where  $\phi$  is the density of the standard normal distribution. Notice that we have used the symmetry of  $\phi$  around 0 in the last expression.

Since  $0 \leq \max\{\sqrt{\lambda} - \sqrt{t}, \sqrt{t} - \sqrt{\lambda}\} < \sqrt{t} + \sqrt{\lambda}$  when  $t > 0$ , by monotonicity of the normal density on  $(0, \infty)$ , we conclude that the derivative (7.23) is indeed negative. Therefore, (7.22) is decreasing in  $\lambda$ , and (7.21) follows for  $t > 0$ . For  $t = 0$ , equality holds in (7.21) with both probabilities being 0.  $\square$

Finally, we derive asymptotic expressions for chi-square quantiles.

**Lemma 7.3 (Chi-square quantiles).** *Let  $F$  be the central chi-square distributions with  $\nu$  degrees of freedom, and let  $u(y)$  be the  $(1 - y)$ -th generalized quantile of  $F$ , i.e.,*

$$u(y) = F^{\leftarrow}(1 - y). \quad (7.24)$$

*Then*

$$u(y) \sim 2 \log(1/y), \quad \text{as } y \rightarrow 0. \quad (7.25)$$

*Proof (Proof of Lemma 7.3).* The case where  $\nu = 1$  follows from the well-known formula for Normal quantiles (see, e.g., Proposition 1.1 in Gao and Stoev (2020))

$$F^{\leftarrow}(1 - y) = \Phi^{\leftarrow}(1 - y/2) \sim \sqrt{2 \log(2/y)} \sim \sqrt{2 \log(1/y)}.$$

The case where  $\nu \geq 2$  follows from the following estimates of high quantiles of chi-square distributions (see, e.g., (Ingolot, 2010)),

$$\nu + 2 \log(1/y) - 5/2 \leq u(y) \leq \nu + 2 \log(1/y) + 2\sqrt{\nu \log(1/y)}, \quad \text{for all } y \leq 0.17,$$

where both the lower and upper bound are asymptotic to  $2 \log(1/y)$ .  $\square$

### 7.6.2 Proof of Theorem 7.1

*Proof (Proof of Theorem 7.1).* We first prove the sufficient condition. The Bonferroni procedure sets the threshold at  $t_p = F^{\leftarrow}(1-\alpha/p)$ , which, by Lemma 7.3, is asymptotic to  $2 \log p - 2 \log \alpha$ . By the assumption on  $\alpha$  in (3.8), for any  $\delta > 0$ , we have  $p^{-\delta} = o(\alpha)$ . Therefore, we have  $-\log \alpha \leq \delta \log p$  for large  $p$ , and

$$1 \leq \limsup_{p \rightarrow \infty} \frac{2 \log p - 2 \log \alpha}{2 \log p} \leq 1 + \delta,$$

for any  $\delta > 0$ . Hence,  $t_p \sim 2 \log p$ .

The condition  $\underline{r} > g(\beta)$  implies, after some algebraic manipulation,  $\sqrt{\underline{r}} - \sqrt{1-\beta} > 1$ . Therefore, we can pick  $q > 1$  such that

$$\sqrt{\underline{r}} - \sqrt{1-\beta} > \sqrt{q} > 1. \quad (7.26)$$

Setting the  $t^* = t_p^* = 2q \log p$ , we have  $t_p < t_p^*$  for large  $p$ .

On the one hand,  $\text{FWER} = 1 - \mathbb{P}[\widehat{S}_p \subseteq S_p]$  vanishes under the Bonferroni procedure with  $\alpha \rightarrow 0$ . On the other hand, for large  $p$ , the probability of no missed detection is bounded from below by

$$\mathbb{P}[\widehat{S}_p \supseteq S_p] = \mathbb{P}[\min_{i \in S} x(i) \geq t_p] \geq \mathbb{P}[\min_{i \in S} x(i) \geq t^*] \geq 1 - p^{1-\beta} \mathbb{P}[\chi_\nu^2(\underline{\Delta}) < t^*], \quad (7.27)$$

where we have used the fact that signal sizes are bounded below by  $\underline{\Delta}$ , and the stochastic monotonicity of chi-square distributions (Lemma 7.2) in the last inequality. Writing

$$\chi_\nu^2(\underline{\Delta}) \stackrel{d}{=} Z_1^2 + \dots + Z_{\nu-1}^2 + (Z_\nu + \sqrt{\underline{\Delta}})^2$$

where  $Z_i$ 's are iid standard normal variables, we have

$$\begin{aligned} \mathbb{P}[\chi_\nu^2(\underline{\Delta}) < t^*] &\leq \mathbb{P}[(Z_\nu + \sqrt{\underline{\Delta}})^2 < t^*] = \mathbb{P}[|Z_\nu + \sqrt{\underline{\Delta}}| < \sqrt{t^*}] \\ &\leq \mathbb{P}\left[Z_\nu < -\sqrt{\underline{\Delta}} + \sqrt{t^*}\right] \\ &= \mathbb{P}\left[Z_\nu < \sqrt{2 \log p} (\sqrt{q} - \sqrt{\underline{r}})\right]. \end{aligned} \quad (7.28)$$



By our choice of  $q$  in (7.26), the last probability in (7.28) can be bounded from above by

$$\begin{aligned} \mathbb{P}\left[Z_\nu < -\sqrt{2(1-\beta)\log p}\right] &\sim \frac{\phi\left(-\sqrt{2(1-\beta)\log p}\right)}{\sqrt{2(1-\beta)\log p}} \\ &= \frac{1}{\sqrt{2(1-\beta)\log p}} p^{-(1-\beta)}, \end{aligned}$$

where the first line uses Mill's ratio for Gaussian distributions (see Section A.1 and Relation (A.3)). This, combined with (7.27), completes the proof of the sufficient condition for the Bonferroni's procedure.

Under the assumption of independence, Sidák's, Holm's, and Hochberg's procedures are strictly more powerful than Bonferroni's procedure, while controlling FWER at the nominal levels. Therefore, the risks of exact support recovery for these procedures also vanishes. This completes the proof for the first part of Theorem 7.1.

We now show the necessary condition. We first normalize the maxima by the chi-square quantiles  $u_p = F^{\leftarrow}(1 - 1/p)$ , where  $F$  is the distribution of a (central) chi-square random variable,

$$\mathbb{P}[\widehat{S}_p = S_p] \leq \mathbb{P}[M_{S^c} < t_p \leq m_S] \leq \mathbb{P}\left[\frac{M_{S^c}}{u_p} < \frac{m_S}{u_p}\right], \quad (7.29)$$

where  $M_{S^c} = \max_{i \in S^c} x(i)$  and  $m_S = \min_{i \in S} x(i)$ . By the relative stability of chi-square random variables (Corollary 7.4), we know that  $M_{S^c}/u_{|S^c|} \rightarrow 1$  in probability. Further, using the expression for  $u_p$  (Lemma 7.3), we obtain

$$\frac{u_{p-p^{1-\beta}}}{u_p} \sim \frac{2 \log(p - p^{1-\beta})}{2 \log p} = \frac{\log p + \log(1 - p^{-\beta})}{\log p} \sim 1.$$

Therefore, the left-hand-side of the last probability in (7.29) converges to 1,

$$\frac{M_{S^c}}{u_p} = \frac{M_{S^c}}{u_{p-p^{1-\beta}}} \frac{u_{p-p^{1-\beta}}}{u_p} \xrightarrow{\mathbb{P}} 1. \quad (7.30)$$

Meanwhile, for any  $i \in S$ , by Lemma 7.2 and the fact that signal sizes are bounded above by  $\bar{\Delta}$ , we have,

$$\chi_\nu^2(\lambda(i)) \stackrel{d}{\leq} \chi_\nu^2(\bar{\Delta}) \stackrel{d}{=} Z_1^2 + \dots + Z_{\nu-1}^2 + \left(Z_\nu + \sqrt{\bar{\Delta}}\right)^2.$$

Dividing through by  $u_p$ , and taking minimum over  $S$ , we obtain

$$\frac{m_S}{u_p} = \min_{i \in S} \frac{\chi_\nu^2(\lambda(i))}{u_p} \stackrel{d}{\leq} \min_{i \in S} \left\{ \frac{Z_1^2(i) + \dots + Z_{\nu-1}^2(i)}{u_p} + \frac{(Z_\nu(i) + \sqrt{\bar{\Delta}})^2}{u_p} \right\}. \quad (7.31)$$

Let  $i^\dagger = i_p^\dagger$  be the index minimizing the second term in (7.31), i.e.,

$$i^\dagger := \arg \min_{i \in S} \frac{(Z_\nu(i) + \sqrt{\Delta})^2}{u_p} = \arg \min_{i \in S} f_p(Z_\nu(i)), \quad (7.32)$$

where  $f_p(x) := (x + \sqrt{\Delta})^2 / (2 \log p)$ . We shall first show that

$$\mathbb{P}[f_p(Z_\nu(i^\dagger)) < 1 - \delta] \rightarrow 1, \quad (7.33)$$

for some small  $\delta > 0$ . On the one hand, we know (by solving a quadratic inequality) that

$$f_p(x) < 1 - \delta \iff \frac{x}{\sqrt{2 \log p}} \in (-\sqrt{\bar{r}} + \sqrt{1 - \delta}, -(\sqrt{\bar{r}} - \sqrt{1 - \delta})). \quad (7.34)$$

On the other hand, we know (by relative stability of iid Gaussians Gao and Stoev (2020)) that

$$\frac{\min_{i \in S} Z_\nu(i)}{\sqrt{2 \log p}} \rightarrow -\sqrt{1 - \beta} \quad \text{in probability.} \quad (7.35)$$

Further, by the assumption on the signal sizes  $\bar{r} < (1 + \sqrt{1 - \beta})^2$ , we have,

$$-(\sqrt{\bar{r}} + 1) < -1 < -\sqrt{1 - \beta} < -(\sqrt{\bar{r}} - 1).$$

Therefore we can pick a small  $\delta > 0$  such that

$$-(\sqrt{\bar{r}} + 1) < -(\sqrt{\bar{r}} + \sqrt{1 - \delta}) < -\sqrt{1 - \beta} < -(\sqrt{\bar{r}} - \sqrt{1 - \delta}) < -(\sqrt{\bar{r}} - 1). \quad (7.36)$$

Combining (7.34), (7.35), and (7.36), we obtain

$$\begin{aligned} \mathbb{P} \left[ \min_{i \in S} f_p(Z_\nu(i)) < 1 - \delta \right] &= \mathbb{P} \left[ f_p(Z_\nu(i^\dagger)) < 1 - \delta \right] \\ &\geq \mathbb{P} \left[ f_p \left( \min_{i \in S} Z_\nu(i) \right) < 1 - \delta \right] \rightarrow 1, \end{aligned}$$

and we arrive at (7.33). As a corollary, since  $u_p \sim 2 \log p$ , it follows that

$$\mathbb{P} \left[ \min_{i \in S} \frac{(Z_\nu(i) + \sqrt{\Delta})^2}{u_p} < 1 - \delta \right] \rightarrow 1. \quad (7.37)$$

Finally, by independence between  $Z_1^2(i) + \dots + Z_{\nu-1}^2(i)$  and  $(Z_\nu^2(i) + \sqrt{\Delta})^2$ , and the fact that  $i^\dagger$  is a function of only the latter, we have

$$Z_1^2(i^\dagger) + \dots + Z_{\nu-1}^2(i^\dagger) \stackrel{d}{=} Z_1^2(i) + \dots + Z_{\nu-1}^2(i) \quad \text{for all } i \in S.$$

Therefore,  $Z_1^2(i^\dagger) + \dots + Z_{v-1}^2(i^\dagger) = O_{\mathbb{P}}(1)$ , and

$$\frac{Z_1^2(i^\dagger) + \dots + Z_{v-1}^2(i^\dagger)}{u_p} \rightarrow 0 \quad \text{in probability.} \quad (7.38)$$

Together, (7.37) and (7.38) imply that

$$\begin{aligned} \mathbb{P} \left[ \frac{m_S}{u_p} < 1 - \delta \right] &\geq \mathbb{P} \left[ \min_{i \in S} \left\{ \frac{Z_1^2(i) + \dots + Z_{v-1}^2(i)}{u_p} + \frac{(Z_v(i) + \sqrt{\Delta})^2}{u_p} \right\} < 1 - \delta \right] \\ &\geq \mathbb{P} \left[ \frac{Z_1^2(i^\dagger) + \dots + Z_{v-1}^2(i^\dagger)}{u_p} + \frac{(Z_v(i^\dagger) + \sqrt{\Delta})^2}{u_p} < 1 - \delta \right] \rightarrow 1. \end{aligned} \quad (7.39)$$

In view of (7.29), (7.30), and (7.39), we conclude that exact recovery cannot succeed with any positive probability. The proof of the necessary condition is complete.  $\square$

### 7.6.3 Proof of Theorem 7.3

We first show the necessary condition. That is, when  $\bar{r} < \beta$ , no thresholding procedure is able to achieve approximate support recovery.

The proof follows the ideas in Arias-Castro and Chen (2017), and is very similar to the proof of Theorem 3.3. One could in principle obtain the proofs in this section by referencing arguments that have appeared in Section ?? . We choose to present the proof here in full to make this section self-contained.

*Proof (Proof of necessary condition in Theorem 7.3).* Denote the distributions of  $\chi_v^2(0)$ ,  $\chi_v^2(\underline{\Delta})$  and  $\chi_v^2(\bar{\Delta})$  as  $F_0$ ,  $F_{\underline{\Delta}}$ , and  $F_{\bar{\Delta}}$  respectively.

Recall that thresholding procedures are of the form

$$\widehat{S}_p = \{i \mid x(i) > t_p(x)\}.$$

Denote  $\widehat{S} := \{i \mid x(i) > t_p(x)\}$ , and  $\widehat{S}(u) := \{i \mid x(i) > u\}$ . For any threshold  $u \geq t_p$  we must have  $\widehat{S}(u) \subseteq \widehat{S}$ , and hence

$$\text{FDP} := \frac{|\widehat{S} \setminus S|}{|\widehat{S}|} \geq \frac{|\widehat{S} \setminus S|}{|\widehat{S} \cup S|} = \frac{|\widehat{S} \setminus S|}{|\widehat{S} \setminus S| + |S|} \geq \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|}. \quad (7.40)$$

On the other hand, for any threshold  $u \leq t_p$  we must have  $\widehat{S}(u) \supseteq \widehat{S}$ , and hence

$$\text{NDP} := \frac{|S \setminus \widehat{S}|}{|S|} \geq \frac{|S \setminus \widehat{S}(u)|}{|S|}. \quad (7.41)$$

Since either  $u \geq t_p$  or  $u \leq t_p$  must take place, putting (7.40) and (7.41) together, we have

$$\text{FDP} + \text{NDP} \geq \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|} \wedge \frac{|S \setminus \widehat{S}(u)|}{|S|}, \quad (7.42)$$

for any  $u$ . Therefore it suffices to show that for a suitable choice of  $u$ , the RHS of (7.42) converges to 1 in probability; the desired conclusion on FDR and FNR follows by the dominated convergence theorem.

Let  $t^* = 2q \log p$  for some fixed  $q$ , we obtain an estimate of the tail probability

$$\begin{aligned} \overline{F}_0(t^*) &= \mathbb{P}[\chi_\nu^2(0) > t^*] = \frac{2^{-\nu/2}}{\Gamma(\nu/2)} \int_{2q \log p}^{\infty} x^{\nu/2-1} e^{-x/2} dx \\ &\sim \frac{2^{-\nu/2}}{\Gamma(\nu/2)} 2(2q \log p)^{\nu/2-1} p^{-q}. \end{aligned} \quad (7.43)$$

where  $a_p \sim b_p$  is taken to mean  $a_p/b_p \rightarrow 1$ ; this tail estimate was also obtained in Donoho and Jin (2004). Observe that  $|\widehat{S}(t^*) \setminus S|$  has distribution  $\text{Binom}(p-s, \overline{F}_0(t^*))$  where  $s = |S|$ , denote  $X = X_p := |\widehat{S}(t^*) \setminus S|/|S|$ , and we have

$$\mu := \mathbb{E}[X] = \frac{(p-s)\overline{F}_0(t^*)}{s}, \quad \text{and} \quad \text{Var}(X) = \frac{(p-s)\overline{F}_0(t^*)F_0(t^*)}{s^2} \leq \mu/s.$$

Therefore for any  $M > 0$ , we have, by Chebyshev's inequality,

$$\mathbb{P}[X < M] \leq \mathbb{P}[|X - \mu| > \mu - M] \leq \frac{\mu/s}{(\mu - M)^2} = \frac{1/(\mu s)}{(1 - M/\mu)^2}. \quad (7.44)$$

Now, from the expression of  $\overline{F}_0(t^*)$  in (7.43), we obtain

$$\mu = (p^\beta - 1)\overline{F}_0(t^*) \sim \frac{2^{1-\nu/2}}{\Gamma(\nu/2)} (2q \log p)^{\nu/2-1} p^{\beta-q}.$$

Since  $\bar{r} < \beta$ , we can pick  $q$  such that  $\bar{r} < q < \beta$ . In turn, we have  $\mu \rightarrow \infty$ , as  $p \rightarrow \infty$ . Therefore the last expression in (7.44) converges to 0, and we conclude that  $X \rightarrow 0$  in probability, and hence

$$\frac{|\widehat{S}(t^*) \setminus S|}{|\widehat{S}(t^*) \setminus S| + |S|} = \frac{X}{X+1} \rightarrow 1 \quad \text{in probability.} \quad (7.45)$$

On the other hand, we show that with the same choice of  $u = t^*$ ,

$$\frac{|S \setminus \widehat{S}(t^*)|}{|S|} \rightarrow 1 \quad \text{in probability.} \quad (7.46)$$

By the stochastic monotonicity of chi-square distributions (Lemma 7.2), the probability of missed detection for each signal is lower bounded by  $\mathbb{P}[\chi_\nu^2(\lambda_i) \leq t^*] \geq F_{\bar{a}}(t^*)$ .

Therefore,  $|S \setminus \widehat{S}(t^*)| \stackrel{d}{\geq} \text{Binom}(s, F_{\bar{a}}(t^*))$ , and it suffices to show that  $F_{\bar{a}}(t^*)$  converges to 1. This is indeed the case, since

$$\begin{aligned} F_{\bar{a}}(t^*) &= \mathbb{P}[Z_1^2 + \dots + Z_{\nu}^2 + 2\sqrt{2\bar{r} \log p} Z_{\nu} + 2\bar{r} \log p \leq 2q \log p] \\ &\geq \mathbb{P}[Z_1^2 + \dots + Z_{\nu}^2 \leq (q - \bar{r}) \log p, 2\sqrt{2\bar{r} \log p} Z_{\nu} \leq (q - \bar{r}) \log p], \end{aligned}$$

and both events in the last line have probability going to 1 as  $p \rightarrow \infty$ . The necessary condition is shown.  $\square$

We now turn to the sufficient condition. That is, when  $\underline{r} > \beta$ , the Benjamini-Hochberg procedure with slowly vanishing FDR levels achieves asymptotic approximate support recovery. The structure for the proof of sufficient condition follows that of Theorem 2 in Arias-Castro and Chen (2017).

*Proof (Proof of sufficient condition in Theorem 7.3).* The FDR vanishes by our choice of  $\alpha$  and the FDR-controlling property of the BH procedure. It only remains to show that FNR also vanishes.

To do so we compare the FNR under the alternative specified in Theorem 7.3 to one with all of the signal sizes equal to  $\underline{\Delta}$ . Let  $x(i)$  be vectors of independent observations with  $p - s$  nulls having  $\chi_{\nu}^2(0)$  distributions, and  $s$  signals having  $\chi_{\nu}^2(\underline{\Delta})$  distributions. By Lemma A.1, it suffices to show that the FNR under the BH procedure in this setting vanishes.

Let  $\widehat{G}$  denote the empirical survival function as in (A.5). Define the empirical survival functions for the null part and signal part

$$\widehat{W}_{\text{null}}(t) = \frac{1}{p-s} \sum_{i \notin S} \mathbb{1}\{x(i) \geq t\}, \quad \widehat{W}_{\text{signal}}(t) = \frac{1}{s} \sum_{i \in S} \mathbb{1}\{x(i) \geq t\}, \quad (7.47)$$

where  $s = |S|$ , so that

$$\widehat{G}(t) = \frac{p-s}{p} \widehat{W}_{\text{null}}(t) + \frac{s}{p} \widehat{W}_{\text{signal}}(t).$$

Apply Lemma A.2 to the two summands in  $\widehat{G}$ , we obtain  $\widehat{G}(t) = G(t) + \widehat{R}(t)$ , where

$$G(t) = \frac{p-s}{p} \overline{F}_0(t) + \frac{s}{p} \overline{F}_a(t), \quad (7.48)$$

where  $\overline{F}_0$  and  $\overline{F}_a$  are the survival functions of  $\chi_{\nu}^2(0)$  and  $\chi_{\nu}^2(\underline{\Delta})$  respectively, and

$$\widehat{R}(t) = O_{\mathbb{P}} \left( \xi_p \sqrt{\overline{F}_0(t) F_0(t)} + \frac{s}{p} \xi_s \sqrt{\overline{F}_a(t) F_a(t)} \right), \quad (7.49)$$

uniformly in  $t$ .

Recall (see proof of Lemma A.1) that the BH procedure is the thresholding procedure with threshold set at  $\tau$  (defined in (A.6)). The NDP may also be re-written as

$$\text{NDP} = \frac{|S \setminus \widehat{S}|}{|S|} = \frac{1}{s} \sum_{i \in S} \mathbb{I}\{x(i) < \tau\} = 1 - \widehat{W}_{\text{signal}}(\tau),$$

so that it suffices to show that

$$\widehat{W}_{\text{signal}}(\tau) \rightarrow 1 \quad (7.50)$$

in probability. Applying Lemma A.2 to  $\widehat{W}_{\text{signal}}$ , we know that

$$\widehat{W}_{\text{signal}}(\tau) = \overline{F}_a(\tau) + O_{\mathbb{P}} \left( \xi_s \sqrt{\overline{F}_a(\tau) F_a(\tau)} \right) = \overline{F}_a(\tau) + o_{\mathbb{P}}(1).$$

So it suffices to show that  $F_a(\tau) \rightarrow 0$  in probability. Now let  $t^* = 2q \log(p)$  for some  $q$  such that  $\beta < q < \underline{r}$ . We have

$$\begin{aligned} F_a(t^*) &= \mathbb{P}[\chi_{\nu}^2(\underline{\Delta}) \leq t^*] \leq \mathbb{P} \left[ 2\sqrt{\underline{\Delta}} Z_{\nu} \leq t^* - \underline{\Delta} \right] \\ &= \mathbb{P} \left[ Z_{\nu} \leq \frac{t^*}{2\sqrt{\underline{\Delta}}} - \frac{\sqrt{\underline{\Delta}}}{2} \right] = \mathbb{P} \left[ Z_{\nu} \leq \frac{q-r}{2\sqrt{\underline{r}}} \sqrt{2 \log p} \right] \rightarrow 0. \end{aligned} \quad (7.51)$$

Hence in order to show (7.50), it suffices to show

$$\mathbb{P}[\tau \leq t^*] \rightarrow 1. \quad (7.52)$$

By (7.48), the mean of the empirical process  $\widehat{G}$  evaluated at  $t^*$  is

$$G(t^*) = \frac{p-s}{p} \overline{F}_0(t^*) + \frac{s}{p} \overline{F}_a(t^*). \quad (7.53)$$

The first term, using Relation (7.43), is asymptotic to  $p^{-q} L(p)$ , where  $L(p)$  is the logarithmic term in  $p$ . The second term, since  $\overline{F}_a(t^*) \rightarrow 1$  by Relation (7.51), is asymptotic to  $p^{-\beta}$ . Therefore,  $G(t^*) \sim p^{-q} L(p) + p^{-\beta} \sim p^{-\beta}$ , since  $p^{\beta-q} L(p) \rightarrow 0$  where  $q > \beta$ .

The fluctuation of the empirical process at  $t^*$ , by Relation (7.49), is

$$\begin{aligned} \widehat{R}(t^*) &= O_{\mathbb{P}} \left( \xi_p \sqrt{\overline{F}_0(t^*) F_0(t^*)} + \frac{s}{p} \xi_s \sqrt{\overline{F}_a(t^*) F_a(t^*)} \right) \\ &= O_{\mathbb{P}} \left( \xi_p \sqrt{\overline{F}_0(t^*)} \right) + o_{\mathbb{P}} \left( p^{-\beta} \right). \end{aligned}$$

By (7.43) and the expression for  $\xi_p$ , the first term is  $O_{\mathbb{P}}(p^{-(q+1)/2} L(p))$  where  $L(p)$  is a poly-logarithmic term in  $p$ . Since  $\beta < \min\{q, 1\}$ , we have  $\beta < (q+1)/2$ , and hence  $\widehat{R}(t^*) = o_{\mathbb{P}}(p^{-\beta})$ .

Putting the mean and the fluctuation of  $\widehat{G}(t^*)$  together, we obtain

$$\widehat{G}(t^*) = G(t^*) + \widehat{R}(t^*) \sim_{\mathbb{P}} G(t^*) \sim p^{-\beta},$$

and therefore, together with (7.43), we have

$$\overline{F}_0(t^*)/\widehat{G}(t^*) = p^{\beta-q} L(p)(1 + o_{\mathbb{P}}(1)),$$

which is eventually smaller than the FDR level  $\alpha$  by the assumption (3.8) and the fact that  $\beta < q$ . That is,

$$\mathbb{P} \left[ \overline{F}_0(t^*)/\widehat{G}(t^*) < \alpha \right] \rightarrow 1.$$

By definition of  $\tau$  (recall (A.6)), this implies that  $\tau \leq t^*$  with probability tending to 1, and (7.52) is shown. The proof for the sufficient condition is complete.  $\square$

#### 7.6.4 Proof of Theorems 7.2 and 7.4

As with the proof of Theorem 7.3, one could shorten the presentations in this section by referencing arguments Section ?? . Again, we choose to present the proof in full to make this section self-contained.

*Proof (Proof of Theorem 7.2).* We first show the sufficient condition. Similar to the proof of Theorem 7.3, it suffices to show that

$$\text{NDP} = 1 - \widehat{W}_{\text{signal}}(t_p) \rightarrow 0, \quad (7.54)$$

where  $t_p$  is the threshold of Bonferroni's procedure.

Since  $\underline{r} > \widetilde{g}(\beta) = 1$ , we can pick  $q$  such that  $1 < q < \underline{r}$ . Let  $t^* = 2q \log p$ , we have  $t_p < t_p^*$  for large  $p$  as in the proof of Theorem 7.1. Therefore for large  $p$ , we have

$$\widehat{W}_{\text{signal}}(t_p) \geq \widehat{W}_{\text{signal}}(t^*) \geq \overline{F}_a(t^*) + o_{\mathbb{P}}(1),$$

where the last inequality follows from the stochastic monotonicity of the chi-square family (Lemma 7.2), and Lemma A.2. Indeed,  $F_a(t^*) \rightarrow 0$  by (7.51) and our choice of  $q < \underline{r}$ . The proof of the sufficient condition is complete.

Proof of the necessary condition follows a similar structure to that of Theorem 7.3. That is, we show that FWER + FNR has liminf at least 1 by working with the lower bound

$$\text{FWER}(\mathcal{R}) + \text{FNR}(\mathcal{R}) \geq \mathbb{P} \left[ \max_{i \in S^c} x(i) > u \right] \wedge \mathbb{E} \left[ \frac{|S \setminus \widehat{S}(u)|}{|S|} \right], \quad (7.55)$$

which holds for any thresholding procedure  $\mathcal{R}$  and for arbitrary  $u \in \mathbb{R}$ . By the assumption that  $\bar{r} < \widetilde{g}(\beta) = 1$ , we can pick  $q$  such that  $\bar{r} < q < 1$  and let  $u = t^* = 2q \log p$ . By relative stability of chi-squared random variables (Lemma 7.1), we have

$$\mathbb{P} \left[ \frac{\max_{i \in S^c} x(i)}{2 \log p} > \frac{t^*}{2 \log p} \right] \rightarrow 1. \quad (7.56)$$

where the first fraction in (7.56) converges to 1, while the second converges to  $q < 1$ . On the other hand, by our choice of  $q > \bar{r}$ , the second term in (7.55) also converges to 1 as in (7.46). This completes the proof of the necessary condition.  $\square$

*Proof (Proof of Theorem 7.4).* We first show the sufficient condition. Since FDR control is guaranteed by the BH procedure, we only need to show that the FWNR also vanishes, that is,

$$\mathbb{P} \left[ \min_{i \in S} x(i) \geq \tau \right] \rightarrow 1, \quad (7.57)$$

where  $\tau$  is the threshold for the BH procedure.

By the assumption that  $\underline{r} > \bar{h}(\beta) = (\sqrt{\beta} + \sqrt{1-\beta})^2$ , we have  $\sqrt{\underline{r}} - \sqrt{1-\beta} > \sqrt{\beta}$ , so we can pick  $q > 0$ , such that

$$\sqrt{\underline{r}} - \sqrt{1-\beta} > \sqrt{q} > \sqrt{\beta}. \quad (7.58)$$

Let  $t^* = 2q \log p$ , we claim that

$$\mathbb{P} [\tau \leq t^*] \rightarrow 1. \quad (7.59)$$

Indeed, by our choice of  $q > \beta$ , (7.59) follows in the same way that (7.52) did.

With this  $t^*$ , we have

$$\mathbb{P} \left[ \min_{i \in S} x(i) \geq \tau \right] \geq \mathbb{P} \left[ \min_{i \in S} x(i) \geq t^*, t^* \geq \tau \right]. \quad (7.60)$$

However, by our choice of  $\sqrt{q} < \sqrt{\underline{r}} - \sqrt{1-\beta}$ , the probability of the first event on the right-hand side of (7.60) also goes to 1 according to (7.27) and (7.28). Together with (7.59), this proves (7.57), and completes proof of the sufficient condition.

The necessary condition follows from the lower bound

$$\text{FDR}(\mathcal{R}) + \text{FWNR}(\mathcal{R}) \geq \mathbb{E} \left[ \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|} \right] \wedge \mathbb{P} \left[ \min_{i \in S} x(i) < u \right], \quad (7.61)$$

which holds for any thresholding procedure  $\mathcal{R}$  and for arbitrary  $u \in \mathbb{R}$ .

By the assumption that  $\bar{r} < \bar{h}(\beta) = (\sqrt{\beta} + \sqrt{1-\beta})^2$ , we can pick a constant  $q > 0$ , such that

$$\sqrt{\bar{r}} - \sqrt{1-\beta} < \sqrt{q} < \sqrt{\beta}. \quad (7.62)$$

Let also  $u = t^* = 2q \log p$ . By our choice of  $q < \beta$ , we know from (7.45) that the first term on the right-hand-side of (7.61) converges to 1. It remains to show that the second term in (7.61) also converges to 1.

For the second term in (7.61), dividing through by  $2 \log p$ , we obtain

$$\mathbb{P} \left[ \min_{i \in S} x(i) < t^* \right] = \mathbb{P} \left[ \frac{m_S}{2 \log p} < q \right]. \quad (7.63)$$

Similar to (7.31), we have



$$\frac{m_S}{2 \log p} \stackrel{d}{\leq} \min_{i \in S} \frac{Z_1^2(i) + \dots + Z_{v-1}^2(i)}{2 \log p} + \frac{(Z_v(i) + \sqrt{\Delta})^2}{2 \log p}. \quad (7.64)$$

Define  $i^\dagger = i_p^\dagger$  to be the index minimizing the second term in (7.64), i.e.,

$$i^\dagger := \arg \min_{i \in S} f_p(Z_v(i)), \quad (7.65)$$

where  $f_p(x) := (x + \sqrt{\Delta})^2 / (2 \log p)$ .

Since  $\sqrt{q} > \sqrt{r} - \sqrt{1 - \beta}$  and  $q > 0$ , we have  $\frac{\sqrt{r} - \sqrt{q}}{\sqrt{1 - \beta}} < 1$ . Also, since

$$\frac{\sqrt{r} + \sqrt{q}}{\sqrt{1 - \beta}} > 0, \quad \text{and} \quad \frac{\sqrt{r} - \sqrt{q}}{\sqrt{1 - \beta}} < \frac{\sqrt{r} + \sqrt{q}}{\sqrt{1 - \beta}},$$

we can further pick a constant  $\beta_0 \in (0, 1]$  such that

$$\frac{\sqrt{r} - \sqrt{q}}{\sqrt{1 - \beta}} < \sqrt{\beta_0} < \frac{\sqrt{r} + \sqrt{q}}{\sqrt{1 - \beta}}. \quad (7.66)$$

Let  $Z_{[1]} \leq Z_{[2]} \leq \dots \leq Z_{[s]}$  be the order statistics of  $\{Z_v(i)\}_{i \in S}$  and define  $k = \lfloor s^{1 - \beta_0} \rfloor$ . Applying Lemma 7.4 (stated below), we obtain

$$\frac{Z_{[k]}}{\sqrt{2 \log p}} = \frac{Z_{[k]}}{\sqrt{2 \log s}} \frac{\sqrt{2 \log s}}{\sqrt{2 \log p}} \rightarrow -\sqrt{\beta_0(1 - \beta)} \quad \text{in probability.} \quad (7.67)$$

Since we know (by solving a quadratic inequality) that

$$f_p(x) < q \iff \frac{x}{\sqrt{2 \log p}} \in \left( -(\sqrt{r} + \sqrt{q}), -(\sqrt{r} - \sqrt{q}) \right), \quad (7.68)$$

combining (7.66), (7.67), and (7.68), it follows that

$$\mathbb{P} \left[ f_p \left( Z_v(i^\dagger) \right) < q \right] \geq \mathbb{P} \left[ f_p \left( Z_{[k]} \right) < q \right] \rightarrow 1.$$

Finally, using (7.38), we conclude that

$$\mathbb{P} \left[ \min_{i \in S} x(i) < t^* \right] = \mathbb{P} \left[ \frac{m_S}{2 \log p} < q \right] \geq \mathbb{P} \left[ o_{\mathbb{P}}(1) + f_p \left( Z_v(i^\dagger) \right) < q \right] \rightarrow 1.$$

Therefore, the two terms on the right-hand-side of (7.61) both converge 1. This completes the proof of the necessary condition.  $\square$

It only remains to justify (7.67).

**Lemma 7.4 (Relative stability of order statistics).** *Let  $Z_{[1]} \leq \dots \leq Z_{[s]}$  be the order statistics of  $s$  iid standard Gaussian random variables. Let  $\beta_0 \in (0, 1]$  and*

define  $k = \lfloor s^{1-\beta_0} \rfloor$ , then we have

$$\frac{Z_{[k]}}{\sqrt{2 \log s}} \rightarrow -\sqrt{\beta_0} \quad \text{in probability.} \quad (7.69)$$

*Proof (Proof of Lemma 7.4).* Using the Renyi representation for order statistics, we write

$$Z_{[i]} = \Phi^{\leftarrow}(U_{[i]}), \quad (7.70)$$

where  $U_{[i]}$  is the  $i^{\text{th}}$  (smallest) order statistic of  $s$  independent uniform random variables over  $(0, 1)$ . Since  $U_{[i]}$  has a Beta( $i, s+1-i$ ) distribution, with mean and standard deviation,

$$\mathbb{E}[U_{[k]}] = k/(s+1) \sim s^{-\beta_0}, \quad \text{and} \quad \text{sd}(U_{[k]}) = \frac{1}{s+1} \sqrt{\frac{k(s+1-k)}{s+2}} \sim s^{-\frac{1+\beta_0}{2}},$$

we obtain by Chebyshev's inequality

$$\mathbb{P} \left[ s^{-\beta_0}(1-\epsilon) < U_{[k]} < s^{-\beta_0}(1+\epsilon) \right] \rightarrow 1,$$

where  $\epsilon$  is an arbitrary positive constant. This implies, by representation (7.70),

$$\mathbb{P} \left[ \Phi^{\leftarrow} \left( s^{-\beta_0}(1-\epsilon) \right) < Z_{[k]} < \Phi^{\leftarrow} \left( s^{-\beta_0}(1+\epsilon) \right) \right] \rightarrow 1. \quad (7.71)$$

Using the expression for standard Gaussian quantiles (see, e.g., Proposition 1.1. in Gao and Stoev (2020)), we know that

$$\begin{aligned} \Phi^{\leftarrow} \left( s^{-\beta_0}(1-\epsilon) \right) &\sim -\sqrt{2 \log (s^{\beta_0}/(1-\epsilon))} \\ &= -\sqrt{2(\beta_0 \log s - \log(1-\epsilon))} \sim -\sqrt{2\beta_0 \log s}, \end{aligned}$$

and similarly  $\Phi^{\leftarrow} \left( s^{-\beta_0}(1+\epsilon) \right) \sim -\sqrt{2\beta_0 \log s}$ . Since both ends of the interval in (7.71) are asymptotic to  $-\sqrt{2\beta_0 \log s}$ , the desired conclusion follows.  $\square$

### 7.6.5 Proof of Proposition 7.1 and Corollary 7.1

*Proof (Proof of Proposition 7.1 and Corollary 7.1).* We parametrize the 2-by-2 multinomial distribution with the parameter  $\delta$ ,

$$\mu_{11} = \phi_1 \theta_1 + \delta, \quad \mu_{12} = \phi_1 \theta_2 - \delta, \quad \mu_{21} = \phi_2 \theta_1 - \delta, \quad \mu_{22} = \phi_2 \theta_2 + \delta. \quad (7.72)$$

By relabelling of categories, we may assume  $0 < \theta_1, \phi_1 \leq 1/2$  without loss of generality. Note that  $\delta$  must lie within the range  $[\delta_{\min}, \delta_{\max}]$ , where

$$\delta_{\min} := \max\{-\phi_1 \theta_1, -\phi_2 \theta_2, \phi_1 \theta_2 - 1, \phi_2 \theta_1 - 1\} = -\phi_1 \theta_1,$$

and

$$\delta_{\max} := \min\{1 - \phi_1\theta_1, 1 - \phi_2\theta_2, \phi_1\theta_2, \phi_2\theta_1\} = \min\{\phi_1\theta_2, \phi_2\theta_1\},$$

in order for  $\mu_{ij} \geq 0$  for all  $i, j \in \{1, 2\}$ . Under this parametrization, Relation (7.8) then becomes

$$R = \frac{\mu_{11}\mu_{22}}{\mu_{12}\mu_{21}} = \frac{\phi_1\theta_1\phi_2\theta_2 + \delta(\phi_1\theta_1 + \phi_2\theta_2) + \delta^2}{\phi_1\theta_1\phi_2\theta_2 - \delta(\phi_1\theta_2 + \phi_2\theta_1) + \delta^2}, \quad (7.73)$$

which is one-to-one and increasing in  $\delta$  on  $(\delta_{\min}, \delta_{\max})$ . Equation (7.10) becomes

$$w^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(\mu_{ij} - \phi_i\theta_j)^2}{\phi_i\theta_j} = \delta^2 \sum_i \sum_j \frac{1}{\phi_i\theta_j} = \frac{\delta^2}{\phi_1\theta_1\phi_2\theta_2}, \quad (7.74)$$

Solving for  $\delta$  in (7.73), and plugging into the expression for signal size (7.74) yields Relation (7.11). Corollary 7.1 follows from the fact that  $w^2(\delta)$  is decreasing on  $[\delta_{\min}, 0)$ , increasing on  $(0, \delta_{\max}]$ , with limits

$$\lim_{\delta \rightarrow \delta_{\min}} w^2(\delta) = \frac{\phi_1\theta_1}{\phi_2\theta_2}, \quad \text{and} \quad \lim_{\delta \rightarrow \delta_{\max}} w^2(\delta) = \min \left\{ \frac{\phi_1\theta_2}{\phi_2\theta_1}, \frac{\phi_2\theta_1}{\phi_1\theta_2} \right\}.$$

The other three cases ( $1/2 \leq \theta_1, \phi_1 \leq 1$ ;  $0 < \theta_1 \leq 1/2 \leq \phi_1 \leq 1$ ; and  $0 \leq \phi_1 \leq 1/2 \leq \theta_1 \leq 1$ ) may be obtained similarly, or by appealing to the symmetry of the problem.  $\square$

*Proof (Proof of Corollary 7.3).* Using the parametrization in (7.72) and in Corollary 7.2, we solve for  $\delta$  in (7.73) to obtain

$$\begin{aligned} \delta &= \frac{\phi_1 f R}{f R + 1 - f} - \left( \frac{\phi_1 f R}{f R + 1 - f} + f(1 - \phi_1) \right) \phi_1 \\ &= \frac{f(1 - f)\phi_1(1 - \phi_1)(R - 1)}{f R + 1 - f}. \end{aligned} \quad (7.75)$$

Substituting (7.75) into the expression (7.74), after some simplification, yields

$$w^2 = \frac{f(1 - f)\phi_1(1 - \phi_1)(R - 1)^2}{[\phi_1 R + (1 - \phi_1)D][\phi_1 + (1 - \phi_1)D]}, \quad (7.76)$$

where  $D = fR + 1 - f > 0$ . Therefore, the derivative of (7.76) with respect to  $\phi_1$  is

$$\frac{dw^2}{d\phi_1} = \frac{f(1 - f)(R - 1)^2}{[\phi_1 R + (1 - \phi_1)D]^2 [\phi_1 + (1 - \phi_1)D]^2} [(D^2 - R)\phi_1^2 - 2D^2\phi_1 + D^2]. \quad (7.77)$$

Further, we obtain the second derivative with respect to  $\phi_1$ ,

$$\frac{d^2 w^2}{d\phi_1^2} = h(R, f) [(\phi_1 - 1)D^2 - \phi_1 R], \quad (7.78)$$

where  $h$  is some function of  $(R, f)$  taking on strictly positive values.

Since  $[(\phi_1 - 1)D^2 - \phi_1 R] < 0$ , the second derivative (7.78) must be strictly negative on  $[0, 1]$ . This implies that the first derivative (7.77) is strictly decreasing on  $[0, 1]$ . Since the first derivative (7.77) is strictly positive at  $\phi_1 = 0$ , and strictly negative at  $\phi_1 = 1$ , it must have a unique zero between 0 and 1, and hence, the solution to  $(D^2 - R)\phi_1^2 - 2D^2\phi_1 + D^2 = 0$  in the interval of  $[0, 1]$  must be the maximizer of (7.76) — when  $D^2 - R > 0$ , the smaller of the two roots maximizes (7.76), and when  $D^2 - R < 0$ , it is the larger of the two. They share the same expression  $D/(D + \sqrt{R})$ , which coincides with (7.15). Finally, when  $D^2 = R$ , the only root  $\phi_1^* = 1/2$ , which also coincides with (7.15), is the maximizer of (7.76).  $\square$

# Bibliography

## References

- Adler, R. J. and Taylor, J. E. (2009). *Random fields and geometry*. Springer Science & Business Media.
- Agresti, A. (2018). *An introduction to categorical data analysis*. Wiley.
- Anderson, T. W. and Darling, D. A. (1952). Asymptotic theory of certain "goodness of fit" criteria based on stochastic processes. *The annals of mathematical statistics*, pages 193–212.
- Arias-Castro, E. and Chen, S. (2017). Distribution-free multiple testing. *Electronic Journal of Statistics*, 11(1):1983–2001.
- Arias-Castro, E. and Wang, M. (2017). Distribution-free tests for sparse heterogeneous mixtures. *Test*, 26(1):71–94.
- Arias-Castro, E. and Ying, A. (2019). Detection of sparse mixtures: higher criticism and scan statistic. *Electronic Journal of Statistics*, 13(1):208–230.
- Barber, R. F. and Candès, E. J. (2015). Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085.
- Barndorff-Nielsen, O. (1963). On the limit behaviour of extreme order statistics. *The Annals of Mathematical Statistics*, 34(3):992–1002.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, pages 289–300.
- Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188.
- Berman, S. M. (1964). Limit theorems for the maximum term in stationary sequences. *The Annals of Mathematical Statistics*, pages 502–516.
- Bogdan, M., Chakrabarti, A., Frommlet, F., and Ghosh, J. K. (2011). Asymptotic bayes-optimality under sparsity of some multiple testing procedures. *The Annals of Statistics*, 39(3):1551–1579.
- Bush, W. S. and Moore, J. H. (2012). Genome-wide association studies. *PLoS computational biology*, 8(12):e1002822.

- Butucea, C., Ndaoud, M., Stepanova, N. A., and Tsybakov, A. B. (2018). Variable selection with Hamming loss. *The Annals of Statistics*, 46(5):1837–1875.
- Cai, T. T., Jeng, X. J., and Jin, J. (2011). Optimal detection of heterogeneous and heteroscedastic mixtures. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(5):629–662.
- Cai, T. T., Jin, J., and Low, M. G. (2007). Estimation and confidence sets for sparse normal mixtures. *The Annals of Statistics*, 35(6):2421–2449.
- Cai, T. T. and Wu, Y. (2014). Optimal detection of sparse mixtures against a given null distribution. *IEEE Transactions on Information Theory*, 60(4):2217–2232.
- Candès, E. J. (2018). Lecture 3: Global testing, chi-square test, optimality of chi-square test for distributed mild effects. In *Stats 300C: Theory of Statistics (Spring 2018), Stanford Lecture notes*, <https://statweb.stanford.edu/~candes/teaching/stats300c/index.html>.
- Chatterjee, S. (2014). *Superconcentration and related topics*, volume 15. Springer.
- Commings, L. and Dalalyan, A. S. (2012). Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics*, 40(5):2667–2696.
- Conlon, D., Fox, J., and Sudakov, B. (2015). Recent developments in graph Ramsey theory. In *Surveys in combinatorics 2015*, volume 424 of *London Math. Soc. Lecture Note Ser.*, pages 49–118. Cambridge Univ. Press, Cambridge.
- Cramér, H. (1928). On the composition of elementary errors: First paper: Mathematical deductions. *Scandinavian Actuarial Journal*, 1928(1):13–74.
- De Haan, L. and Ferreira, A. (2007). *Extreme value theory: an introduction*. Springer Science & Business Media.
- Dedecker, J., Doukhan, P., Lang, G., León R., J. R., Louhichi, S., and Prieur, C. (2007). *Weak dependence: with examples and applications*, volume 190 of *Lecture Notes in Statistics*. Springer, New York.
- Domingos, P. and Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. *Machine learning*, 29(2-3):103–130.
- Donoghue, W. F. (2014). *Distributions and Fourier transforms*, volume 32. Academic Press.
- Donoho, D. and Jin, J. (2004). Higher Criticism for detecting sparse heterogeneous mixtures. *The Annals of Statistics*, 32(3):962–994.
- Donoho, D. and Jin, J. (2015). Special invited paper: Higher criticism for large-scale inference, especially for rare and weak effects. *Statistical Science*, pages 1–25.
- Dudoit, S., Shaffer, J. P., and Boldrick, J. C. (2003). Multiple hypothesis testing in microarray experiments. *Statistical Science*, pages 71–103.
- Dunn, O. J. (1961). Multiple comparisons among means. *Journal of the American statistical association*, 56(293):52–64.
- Efron, B. (2004). Large-scale simultaneous hypothesis testing: the choice of a null hypothesis. *Journal of the American Statistical Association*, 99(465):96–104.
- Eicker, F. (1979). The asymptotic distribution of the suprema of the standardized empirical processes. *The Annals of Statistics*, pages 116–138.
- Embrechts, P., Klüppelberg, C., and Mikosch, T. (2013). *Modelling extremal events: for insurance and finance*, volume 33. Springer Science & Business Media.

- Erdős, P. and Szekeres, G. (1935). A combinatorial problem in geometry. *Compositio mathematica*, 2:463–470.
- Fan, J. (1996). Test of significance based on wavelet thresholding and neyman’s truncation. *Journal of the American Statistical Association*, 91(434):674–688.
- Ferguson, T. S. (2017). *A course in large sample theory*. Routledge.
- Fox, J. (2009). Lecture 5: Ramsey Theory. In *MAT 307: Combinatorics (Spring 2009), MIT Lecture notes*, <http://math.mit.edu/~fox/MAT307.html>.
- Gao, Z. (2019). U-PASS: <https://power.stat.lsa.umich.edu/u-pass/>. An R Shiny App accompanying the paper “U-PASS: unified power analysis and forensics for qualitative traits in genetic association studies”.
- Gao, Z. and Stoev, S. (2020). Fundamental limits of exact support recovery in high dimensions. *Bernoulli*, 26(4):2605–2638.
- Gao, Z., Terhorst, J., Van Hout, C. V., and Stoev, S. (2019). U-PASS: unified power analysis and forensics for qualitative traits in genetic association studies. *Bioinformatics*, 36(3):974–975.
- Genovese, C. and Wasserman, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):499–517.
- Genovese, C. R., Jin, J., Wasserman, L., and Yao, Z. (2012). A comparison of the lasso and marginal regression. *Journal of Machine Learning Research*, 13(Jun):2107–2143.
- Gnedenko, B. (1943). Sur la distribution limite du terme maximum d’une serie aleatoire. *Annals of mathematics*, pages 423–453.
- Hall, P. and Jin, J. (2010). Innovated higher criticism for detecting sparse signals in correlated noise. *The Annals of Statistics*, 38(3):1686–1732.
- He, Y., Xu, G., Wu, C., and Pan, W. (2018). Asymptotically independent u-statistics in high-dimensional testing. *arXiv preprint arXiv:1809.00411*.
- Hochberg, Y. (1988). A sharper bonferroni procedure for multiple tests of significance. *Biometrika*, 75(4):800–802.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*, pages 65–70.
- Hsing, T. (1995). A note on the asymptotic independence of the sum and maximum of strongly mixing stationary random variables. *The Annals of Probability*, pages 938–947.
- Inglot, T. (2010). Inequalities for quantiles of the chi-square distribution. *Probability and Mathematical Statistics*, 30(4):339–351.
- Ingster, Y. I. (1998). Minimax detection of a signal for  $l_n^p$ -balls. *Mathematical Methods of Statistics*, 7(4):401–428.
- Ji, P. and Jin, J. (2012). UPS delivers optimal phase diagram in high-dimensional variable selection. *The Annals of Statistics*, 40(1):73–103.
- Jin, J., Zhang, C.-H., and Zhang, Q. (2014). Optimality of graphlet screening in high dimensional variable selection. *The Journal of Machine Learning Research*, 15(1):2723–2772.
- Kallitsis, M., Stoev, S. A., Bhattacharya, S., and Michailidis, G. (2016). AMON: An open source architecture for online monitoring, statistical analysis, and forensics

- of multi-gigabit streams. *IEEE Journal on Selected Areas in Communications*, 34(6):1834–1848.
- Kartsioukas, R., Gao, Z., and Stoev, S. (2019). On the rate of concentration of maxima in gaussian arrays. *arXiv preprint arXiv:1910.04259*.
- Klass, M. J. (1984). The minimal growth rate of partial maxima. *Ann. Probab.*, 12(2):380–389.
- Leadbetter, M. R., Lindgren, G., and Rootzén, H. (1983). *Extremes and related properties of random sequences and processes*. Springer Series in Statistics. Springer-Verlag, New York-Berlin.
- MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., and Morales, J. (2016). The new nhgri-ebi catalog of published genome-wide association studies (gwas catalog). *Nucleic acids research*, 45(D1):D896–D901.
- McCormick, W. and Mittal, Y. (1976). *On weak convergence of the maximum*. Stanford University. Department of Statistics.
- Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, 34(3):1436–1462.
- Michailidou, K., Lindström, S., Dennis, J., Beesley, J., Hui, S., Kar, S., Lemaçon, A., Soucy, P., Glubb, D., and Rostamianfar, A. (2017). Association analysis identifies 65 new breast cancer risk loci. *Nature*, 551(7678):92.
- Naveau, P. (2003). Almost sure relative stability of the maximum of a stationary sequence. *Advances in Applied Probability*, 35(3):721–736.
- Neuvial, P. and Roquain, E. (2012). On false discovery rate thresholding for classification under sparsity. *The Annals of Statistics*, 40(5):2572–2600.
- Nichols, T. and Hayasaka, S. (2003). Controlling the familywise error rate in functional neuroimaging: a comparative review. *Statistical methods in medical research*, 12(5):419–446.
- Pipiras, V. and Taqqu, M. S. (2017). *Long-range dependence and self-similarity*. Cambridge Series in Statistical and Probabilistic Mathematics, [45]. Cambridge University Press, Cambridge.
- Ramsey, F. P. (2009). On a problem of formal logic. In *Classic Papers in Combinatorics*, pages 1–24. Springer.
- Resnick, S. I. (2013). *Extreme values, regular variation and point processes*. Springer.
- Resnick, S. I. and Tomkins, R. (1973). Almost sure stability of maxima. *Journal of Applied Probability*, 10(2):387–401.
- Robbins, H. (1955). A remark on Stirling’s formula. *The American mathematical monthly*, 62(1):26–29.
- Šidák, Z. (1967). Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American Statistical Association*, 62(318):626–633.
- Skorokhod, A. V. (1956). Limit theorems for stochastic processes. *Theory of Probability & Its Applications*, 1(3):261–290.
- Slepian, D. (1962). The one-sided barrier problem for gaussian noise. *Bell Labs Technical Journal*, 41(2):463–501.



- Smirnov, N. (1948). Table for estimating the goodness of fit of empirical distributions. *The annals of mathematical statistics*, 19(2):279–281.
- Storey, J. D. (2007). The optimal discovery procedure: a new approach to simultaneous significance testing. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(3):347–368.
- Sun, W. and Cai, T. T. (2007). Oracle and adaptive compound decision rules for false discovery rate control. *Journal of the American Statistical Association*, 102(479):901–912.
- Tanguy, K. (2015a). Some superconcentration inequalities for extrema of stationary Gaussian processes. *Statist. Probab. Lett.*, 106:239–246.
- Tanguy, K. (2015b). Some superconcentration inequalities for extrema of stationary gaussian processes. *Statistics & Probability Letters*, 106:239–246.
- Taqqu, M. S. (2003). Fractional Brownian motion and long-range dependence. In Doukhan, P., Oppenheim, G., and Taqqu, M. S., editors, *Theory and Applications of Long-range Dependence*, pages 5–38. Birkhäuser.
- Tsagris, M., Beneki, C., and Hassani, H. (2014). On the folded normal distribution. *Mathematics*, 2(1):12–28.
- Tukey, J. W. (1976). T13N: The higher criticism. In *Statistics 411, Princeton University Lecture notes*.
- Wainwright, M. J. (2009a). Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting. *IEEE Transactions on Information Theory*, 55(12):5728–5741.
- Wainwright, M. J. (2009b). Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, UK.
- Wasserman, L. and Roeder, K. (2009). High dimensional variable selection. *Annals of statistics*, 37(5A):2178.
- Wu, C., Xu, G., and Pan, W. (2019). An adaptive test on high-dimensional parameters in generalized linear models. *Statistica Sinica*, 29:2163–2186.
- Xu, G., Lin, L., Wei, P., and Pan, W. (2016). An adaptive two-sample test for high-dimensional means. *Biometrika*, 103(3):609–624.
- Zhang, J. (2002). Powerful goodness-of-fit tests based on the likelihood ratio. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(2):281–294.
- Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine learning research*, 7(Nov):2541–2563.
- Zhong, P.-S., Chen, S. X., and Xu, M. (2013). Tests alternative to higher criticism for high-dimensional means under sparsity and column-wise dependence. *The Annals of Statistics*, 41(6):2820–2851.



## Appendix A

### Proofs for Chapter 3

We first recall some basic properties of the Gaussian distribution in Section A.1. Section A.2 states and proves an interesting property of the BH procedure which may be of independent interest. Results on the signal detection problem (Theorem 3.1) are proved in Section A.3, and the phase transition results on the support recovery problems (Theorems 3.2 through 3.5) are shown in Sections A.4 and A.5.

#### A.1 Auxiliary facts of Gaussian distributions

We recall three facts of Gaussian distributions that will be used in the proofs later.

We first state the relative stability of iid standard Gaussian random variables. Since the standard Gaussian distribution falls in the class of asymptotically generalized Gaussians (AGG; see Definition 2.6), by Example 2.1, we know that the triangular array  $\mathcal{E} = \{(\epsilon_p(i))_{i=1}^p, p \in \mathbb{N}\}$  has relatively stable (RS) maxima in the sense of (2.38), i.e.,

$$\frac{1}{u_p} \max_{i=1, \dots, p} \epsilon_p(i) \xrightarrow{\mathbb{P}} 1, \quad \text{as } p \rightarrow \infty, \quad (\text{A.1})$$

where  $u_p$  is the  $(1/p)$ -th upper quantile as defined in (2.33). Similarly, since the array  $\mathcal{E}$  has distributions symmetric around 0, it also has relatively stable minima

$$\frac{1}{u_p} \min_{i=1, \dots, p} \epsilon_p(i) \xrightarrow{\mathbb{P}} -1, \quad \text{as } p \rightarrow \infty. \quad (\text{A.2})$$

The second fact is on the well-known bounds for the Mill's ratio of Gaussian tails. Let  $\Phi$  denote the CDF of the standard Gaussian distribution and  $\phi$  its density. One can show that for all  $x > 0$  we have

$$\frac{x}{1+x^2} \phi(x) \leq \bar{\Phi}(x) = 1 - \Phi(x) \leq \frac{1}{x} \phi(x), \quad (\text{A.3})$$

using e.g., integration by parts.

The third fact is the stochastic monotonicity of the Gaussian location family. In fact, for all location families  $\{F_\delta(x)\}_\delta$  where  $F_\delta(x) = F(x - \delta)$ , we have,

$$F_{\delta_1}(t) \geq F_{\delta_2}(t), \quad \text{for all } t \in \mathbb{R} \quad \text{and all } \delta_1 \leq \delta_2. \quad (\text{A.4})$$

Relation (A.4) holds, of course, when  $F$  is the standard Gaussian distribution.

## A.2 Monotonicity of the Benjamini-Hochberg procedure

We make a connection between power of the BH procedure and the stochastic ordering of distributions under the alternative. This result, though natural, seems new.

**Lemma A.1 (Monotonicity of the BH procedure).** *Consider  $p$  independent observations  $x(i)$ ,  $i \in \{1, \dots, p\}$ , where the  $(p - s)$  coordinates in the null part have common distribution  $F_0$ , and the remaining  $s$  signals have alternative distributions  $F_j^i$ ,  $i \in S$ , respectively. Compare the two alternatives  $j \in \{1, 2\}$ , where the distributions in Alternative 2 are stochastically larger than those in Alternative 1, i.e.,*

$$F_2^i(t) \leq F_1^i(t), \quad \text{for all } t \in \mathbb{R}, \text{ and for all } i \in S.$$

*If the BH procedure is applied at the same nominal level of FDR, then the FNR of the BH procedure under Alternative 2 is bounded above by the FNR under Alternative 1. Further, the threshold of the BH procedure under Alternative 2 is stochastically smaller than that under Alternative 1.*

Loosely put, the power of the BH procedure is monotone increasing with respect to the stochastic ordering of the alternatives, yet (the distribution of) the BH threshold is monotone decreasing in the distributions of the alternatives.

*Proof (Proof of Lemma A.1).* We first re-express the BH procedure in a different form. Recall that on observing  $x(i)$ ,  $i \in \{1, \dots, p\}$ , the BH procedure is the thresholding procedure with threshold set at  $x_{[i^*]}$ , where  $i^* := \max\{i \mid \overline{F}_0(x_{[i]}) \leq \alpha i/p\}$ , and  $x_{[1]} \geq \dots \geq x_{[p]}$  are the order statistics.

Let  $\widehat{G}$  denote the left-continuous empirical survival function

$$\widehat{G}(t) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}\{x(i) \geq t\}. \quad (\text{A.5})$$

By the definition, we know that  $\widehat{G}(x_{[i]}) = i/p$ . Therefore, by the definition of  $i^*$ , we have

$$\overline{F}_0(x_{[i]}) > \alpha \widehat{G}(x_{[i]}) = \alpha i/p \quad \text{for all } i > i^*.$$

Since  $\widehat{G}$  is constant on  $(x_{[i^*+1]}, x_{[i^*]}]$ , the fact that  $\overline{F}_0(x_{[i^*]}) \leq \alpha \widehat{G}(x_{[i^*]})$  and  $\overline{F}_0(x_{[i^*+1]}) > \alpha \widehat{G}(x_{[i^*+1]})$  implies that  $\alpha \widehat{G}$  and  $\overline{F}_0$  must “intersect” on the interval

by continuity of  $F_0$ . We denote this “intersection” as

$$\tau = \inf\{t \mid \overline{F_0}(t) \leq \alpha \widehat{G}(t)\}. \quad (\text{A.6})$$

Note that  $\tau$  cannot be equal to  $x_{[i^*+1]}$  since  $\overline{F_0}$  is càdlàg. Since there is no observation in  $[\tau, x_{[i^*+1]})$ , we can write the BH procedure as the thresholding procedure with threshold set at  $\tau$ .

Now, denote the observations under Alternatives 1 and 2 as  $x_1(i)$  and  $x_2(i)$ . Since  $x_2(i)$  stochastically dominates  $x_1(i)$  for all  $i \in \{1, \dots, p\}$ , there exists a coupling  $(\tilde{x}_1, \tilde{x}_2)$  of  $x_1$  and  $x_2$  such that  $\tilde{x}_1(i) \leq \tilde{x}_2(i)$  almost surely for all  $i$ . We will replace  $\tilde{x}_1$  and  $\tilde{x}_2$  with  $x_1$  and  $x_2$  in what follows. Since we will compare the FNR's, i.e., expectations with respect to the marginals of  $x$ 's in the last step, this replacement does not affect the conclusions. To simplify notation, we still write  $x_1$  and  $x_2$  in place of  $\tilde{x}_1$  and  $\tilde{x}_2$ .

Let  $\widehat{G}_k$  be the left-continuous empirical survival function under Alternative  $k$ , i.e.,

$$\widehat{G}_k(t) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}\{x_k(i) \geq t\}, \quad k \in \{1, 2\}. \quad (\text{A.7})$$

We define the BH thresholds  $\tau_1$  and  $\tau_2$  by replacing  $\widehat{G}$  in (A.6) with  $\widehat{G}_1$  and  $\widehat{G}_2$ , respectively. Denote the set estimates of signal support  $\widehat{S}_k = \{i \mid x_k(i) \geq \tau_k\}$  by the BH procedure. We claim that

$$\tau_2 \leq \tau_1 \quad \text{with probability 1.} \quad (\text{A.8})$$

Indeed, by definition of the empirical survival function (A.7) and the fact that  $x_1(i) \leq x_2(i)$  almost surely for all  $i$ , we have  $\widehat{G}_1(t) \leq \widehat{G}_2(t)$  for all  $t$ . Hence,  $\overline{F_0}(t) \leq \alpha \widehat{G}_1(t)$  implies  $\overline{F_0}(t) \leq \alpha \widehat{G}_2(t)$ , and Relation (A.8) follows from the definition of  $\tau$  in (A.6). The claim of stochastic ordering of the BH thresholds in Lemma A.1 follows from (A.8).

Finally, when  $\tau_2 \leq \tau_1$ , we have  $\tau_2 \leq x_1(i) \leq x_2(i)$  with probability 1 for all  $i \in \widehat{S}_1$ . Therefore, it follows that  $\widehat{S}_1 \subseteq \widehat{S}_2$  and hence  $|S \setminus \widehat{S}_2| \leq |S \setminus \widehat{S}_1|$  almost surely. The first conclusion in Lemma A.1 follows from the last inequality.  $\square$

### A.3 Proof of Theorem 3.1

*Proof (Proof of Theorem 3.1).* Statements about  $L_1$ ,  $L_2$ , and sum statistics  $S$  in the case of diverging signal sizes (3.3) can be found in Fan (1996) and Candès (2018). We prove here the statements for the case where signals are dense and small, as parametrized in (3.5).

We first show that the sum statistic  $S$ , or equivalently, the simple arithmetic mean attains the sparse signal detection boundary.

Consider the case of vanishing signals as prescribed in (3.5), by normality of the summands, we have,

$$\frac{1}{\sqrt{p}} \sum_{i=1}^p x(i) \sim \begin{cases} N(0, 1), & \text{under } H_0 \\ N(p^{(r-\beta)+1/2}, 1), & \text{under } H_1. \end{cases} \quad (\text{A.9})$$

It immediately follows that the two distributions can be distinguished perfectly if  $p^{r-(\beta-1/2)}$  diverges, i.e.,  $r > \beta - 1/2$ . This can be seen by simply setting the rejection region at  $(p^{(r-\beta)+1/2}/2, +\infty)$  for the scaled statistic  $\sum_{i=1}^p x(i)/\sqrt{p}$ . According to the lower bound on the performance limit in detection problems (see Theorem 8 in Cai et al. (2011)), we have shown that  $S$  attains the optimal detection boundary (3.6).

We now turn to the  $L_2$ -norms. Recall a non-central chi-square random variable  $\chi_k^2(\lambda)$  has mean  $(k + \lambda)$  and variance  $2(k + 2\lambda)$ . Since the observations have distributions  $N(0, 1)$  under the null and  $N(p^r, 1)$  under the alternative, we have  $x^2(i) \sim \chi_1^2(0)$  for  $i \notin S$  and  $x^2(i) \sim \chi_1^2(p^{2r})$  for  $i \in S$ . Therefore, mean and variance of the (centered and scaled)  $L_2$  statistics are

$$\mathbb{E} \left[ \frac{1}{\sqrt{p}} \sum_{i=1}^p (x(i)^2 - 1) \right] = \begin{cases} 0, & \text{under } H_0 \\ p^{1-\beta} p^{2r} p^{-1/2} = p^{1/2-\beta+2r}, & \text{under } H_1, \end{cases} \quad (\text{A.10})$$

and

$$\text{Var} \left( \frac{1}{\sqrt{p}} \sum_{i=1}^p (x(i)^2 - 1) \right) = \begin{cases} \frac{1}{p} 2p = 2, & \text{under } H_0 \\ \frac{1}{p} (2p + 2p^{1-\beta+2r}) = 2(1 + p^{2r-\beta}), & \text{under } H_1, \end{cases} \quad (\text{A.11})$$

respectively. By the (Lyapunov) central limit theorem, we have

$$\frac{1}{2p} \sum_{i=1}^p (x(i)^2 - 1) \implies N(0, 1), \quad (\text{A.12})$$

under the null, and

$$\frac{1}{2p} \left( \sum_{i=1}^p (x(i)^2 - 1) - p^{1/2-\beta+2r} \right) \implies N(0, 1), \quad (\text{A.13})$$

under the alternative since  $p^{2r-\beta} \rightarrow 0$  for all  $r < 0$  and  $\beta > 0$ . Hence, perfect detection with the  $L_2$ -norm is possible if  $p^{1/2-\beta+2r}$  diverges, i.e.,  $r > \beta/2 - 1/4$ . On the other hand, if  $r < \beta/2 - 1/4$ , the distributions of the (scaled) statistics merge under the null and the alternative.

The case of  $L_1$ -norm is treated similarly. Let  $Y = |X|$  where  $X \sim |N(\mu, 1)|$ . Using the expressions for the mean and variance of  $Y$  (see, e.g., Tsagris et al. (2014)),

$$\mu_Y = \mathbb{E}[Y] = \sqrt{\frac{2}{\pi}} e^{-\mu^2/2} + \mu(1 - \Phi(-\mu)), \quad (\text{A.14})$$

$$\sigma_Y^2 = \text{Var}(Y) = \mu^2 + 1 - \mu_Y^2, \quad (\text{A.15})$$

where  $\Phi$  is the CDF of a standard normal random variable, we have, regardless of the value of  $\mu$ ,

$$\sigma_Y^2 = \text{Var}(Y) = \mathbb{E}(Y - \mathbb{E}Y)^2 \leq \mathbb{E}(X - \mathbb{E}X)^2 = 1, \quad (\text{A.16})$$

where the inequality holds because absolute value is a Lipschitz function with Lipschitz constant 1.

By the central limit theorem, we have,

$$\frac{1}{\sqrt{p}} \left( \sum_{i=1}^p |x(i)| - \sqrt{\frac{2}{\pi}} \right) \Rightarrow N(0, 1 - 2/\pi) \quad (\text{A.17})$$

under the null. On the other hand, when the alternative hypothesis holds, we have

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{\sqrt{p}} \left( \sum_{i=1}^p |x(i)| - \sqrt{\frac{2}{\pi}} \right) \right] &= \frac{p^{1-\beta}}{\sqrt{p}} \left[ \left( \sqrt{\frac{2}{\pi}} e^{-\mu^2/2} + \mu (1 - 2\Phi(-\mu)) \right) - \sqrt{\frac{2}{\pi}} \right] \\ &= p^{1/2-\beta} \left[ \sqrt{\frac{2}{\pi}} (e^{-p^{2r}/2} - 1) + p^r (1 - 2\Phi(-\mu)) \right] \\ &= p^{1/2-\beta} \left[ \sqrt{\frac{2}{\pi}} (-p^{2r}/2 - O(p^{4r})) + p^r \left( \sqrt{\frac{2}{\pi}} p^r + O(p^{3r}) \right) \right] \\ &= p^{1/2-\beta} \sqrt{\frac{2}{\pi}} (p^{2r}/2 + O(p^{4r})) \\ &= p^{1/2-\beta+2r} \sqrt{1/2\pi} + O(p^{1/2-\beta+4r}), \end{aligned}$$

and

$$\begin{aligned} \text{Var} \left( \frac{1}{\sqrt{p}} \left( \sum_{i=1}^p |x(i)| - \sqrt{\frac{2}{\pi}} \right) \right) &= (1 - p^{1-\beta})(1 - 2/\pi) + p^{1-\beta} \sigma_Y^2 \\ &\rightarrow 1 - 2/\pi, \end{aligned}$$

by boundedness of  $\sigma_Y^2$  shown in (A.16). Again, by the (Lyapunov) central limit theorem, we conclude asymptotic normality of the centered and scaled  $L_1$ -norms under the alternative. In an entirely analogous argument to the  $L_2$ -norm case, asymptotically perfect detection can be achieved if  $p^{1/2-\beta+2r}$  diverges, i.e.,  $r > \beta/2 - 1/4$ . On the other hand, when  $r < \beta/2 - 1/4$ , the two hypotheses cannot be told apart by the  $L_1$ -norms since the distributions of the (scaled) statistics merge under the two hypotheses.  $\square$

### A.4 Proof of Theorem 3.3

We first show the necessary condition. That is, when  $\bar{r} < \beta$ , no thresholding procedure is able to achieve approximate support recovery. The arguments are similar to that in Theorem 1 of Arias-Castro and Chen (2017), although we allow for unequal signal sizes.

*Proof (Proof of necessary condition in Theorem 3.3).* Denote the distributions of  $N(0, 1)$ ,  $N(\underline{\Delta}, 1)$ , and  $N(\bar{\Delta}, 1)$  as  $F_0$ ,  $F_{\underline{a}}$ , and  $F_{\bar{a}}$  respectively.

Recall that thresholding procedures are of the form

$$\widehat{S}_p = \{i \mid x(i) > t_p(x)\}.$$

Denote  $\widehat{S} := \{i \mid x(i) > t_p(x)\}$ , and  $\widehat{S}(u) := \{i \mid x(i) > u\}$ . For any threshold  $u \geq t_p$  we must have  $\widehat{S}(u) \subseteq \widehat{S}$ , and hence

$$\text{FDP} := \frac{|\widehat{S} \setminus S|}{|\widehat{S}|} \geq \frac{|\widehat{S} \setminus S|}{|\widehat{S} \cup S|} = \frac{|\widehat{S} \setminus S|}{|\widehat{S} \setminus S| + |S|} \geq \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|}. \quad (\text{A.18})$$

On the other hand, for any threshold  $u \leq t_p$  we must have  $\widehat{S}(u) \supseteq \widehat{S}$ , and hence

$$\text{NDP} := \frac{|S \setminus \widehat{S}|}{|S|} \geq \frac{|S \setminus \widehat{S}(u)|}{|S|}. \quad (\text{A.19})$$

Since either  $u \geq t_p$  or  $u \leq t_p$  must take place, putting (A.18) and (A.19) together, we have

$$\text{FDP} + \text{NDP} \geq \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|} \wedge \frac{|S \setminus \widehat{S}(u)|}{|S|}, \quad (\text{A.20})$$

for any  $u$ . Therefore it suffices to show that for a suitable choice of  $u$ , the RHS of (A.20) converges to 1 in probability; the desired conclusion on FDR and FNR follows by the dominated convergence theorem.

Let  $t^* = \sqrt{2q \log p}$  for some fixed  $q$ , we obtain an estimate of the tail probability by Mill's ratio (A.3),

$$\overline{F}_0(t^*) \sim \frac{1}{t^*} \phi(t^*) = \frac{1}{2\sqrt{\pi q \log p}} p^{-q}, \quad (\text{A.21})$$

where  $a_p \sim b_p$  is taken to mean  $a_p/b_p \rightarrow 1$ . Observe that  $|\widehat{S}(t^*) \setminus S|$  has distribution  $\text{Binom}(p - s, \overline{F}_0(t^*))$  where  $s = |S|$ , denote  $X = X_p := |\widehat{S}(t^*) \setminus S|/|S|$ , and we have

$$\mu := \mathbb{E}[X] = \frac{(p - s)\overline{F}_0(t^*)}{s}, \quad \text{and} \quad \text{Var}(X) = \frac{(p - s)\overline{F}_0(t^*)F_0(t^*)}{s^2} \leq \mu/s.$$

Therefore for any  $M > 0$ , we have, by Chebyshev's inequality,



$$\mathbb{P}[X < M] \leq \mathbb{P}[|X - \mu| > \mu - M] \leq \frac{\mu/s}{(\mu - M)^2} = \frac{1/(\mu s)}{(1 - M/\mu)^2}. \quad (\text{A.22})$$

Now, from the expression of  $\overline{F}_0(t^*)$  in (A.21), we obtain

$$\mu = (p^\beta - 1)\overline{F}_0(t^*) \sim \frac{1}{2\sqrt{\pi q \log p}} p^{\beta-q}.$$

Since  $\bar{r} < \beta$ , we can pick  $q$  such that  $\bar{r} < q < \beta$ . In turn, we have  $\mu \rightarrow \infty$ , as  $p \rightarrow \infty$ . Therefore the last expression in (A.22) converges to 0, and we conclude that  $X \rightarrow \infty$  in probability, and hence

$$\frac{|\widehat{S}(t^*) \setminus S|}{|\widehat{S}(t^*) \setminus S| + |S|} = \frac{X}{X + 1} \rightarrow 1 \quad \text{in probability.} \quad (\text{A.23})$$

On the other hand, we show that with the same choice of  $u = t^*$ , we have,

$$\frac{|S \setminus \widehat{S}(t^*)|}{|S|} \rightarrow 1 \quad \text{in probability.} \quad (\text{A.24})$$

By the stochastic monotonicity of Gaussian location family (A.4), we have the following lower bound for the probability of missed detection for each signal  $\mu(i)$ ,  $i \in S$ ,

$$\mathbb{P}[\mathbf{N}(\mu(i), 1) \leq t^*] \geq F_{\bar{a}}(t^*). \quad (\text{A.25})$$

Since  $|S \setminus \widehat{S}(t^*)|$  can be written as the sum of  $s$  independent Bernoulli random variables,

$$|S \setminus \widehat{S}(t^*)| = \sum_{i \in S} \mathbb{1}_{(-\infty, t^*]}(x(i)),$$

using with (A.25), we conclude that  $|S \setminus \widehat{S}(t^*)| \stackrel{d}{\geq} \text{Binom}(s, F_{\bar{a}}(t^*))$ . Finally, we know that  $F_{\bar{a}}(t^*)$  converges to 1 by our choice of diverging  $t^*$ , and the necessary condition is shown.  $\square$

We now turn to the sufficient condition. That is, when  $\underline{r} > \beta$ , the Benjamini-Hochberg procedure with slowly vanishing FDR levels achieves asymptotic approximate support recovery.

*Proof (Proof of the sufficient condition in Theorem 3.3).* The FDR vanishes by our choice of  $\alpha$  and the FDR-controlling property of the BH procedure (Benjamini and Hochberg, 1995). It only remains to show that FNR also vanishes.

To do so we compare the FNR under the alternative specified in Theorem 3.3 to one with all of the signal sizes equal to  $\underline{\Delta}$ . By Lemma A.1, it suffices to show that the FNR under the BH procedure in this setting vanishes. Let  $x(i)$  be vectors of independent observations with  $p - s$  nulls having standard Gaussian distributions, and  $s$  signals having  $\mathbf{N}(\underline{\Delta}, 1)$  distributions.

Denote the null and the alternative distributions as  $F_0$  and  $F_a$  respectively. Let  $\widehat{G}$  denote the empirical survival function as in (A.5). Define the empirical survival

functions for the null part and signal part

$$\widehat{W}_{\text{null}}(t) = \frac{1}{p-s} \sum_{i \notin S} \mathbb{1}\{x(i) \geq t\}, \quad \widehat{W}_{\text{signal}}(t) = \frac{1}{s} \sum_{i \in S} \mathbb{1}\{x(i) \geq t\}, \quad (\text{A.26})$$

where  $s = |S|$ , so that

$$\widehat{G}(t) = \frac{p-s}{p} \widehat{W}_{\text{null}}(t) + \frac{s}{p} \widehat{W}_{\text{signal}}(t).$$

We need the following result to describe the deviations of the empirical distributions.

**Lemma A.2 (Theorem 1 of Eicker (1979)).** *Let  $Z_1, \dots, Z_k$  be iid with continuous survival function  $Q$ . Let  $\widehat{Q}_k$  denote their empirical survival function and define  $\xi_k = \sqrt{2 \log \log(k)/k}$  for  $k \geq 3$ . Then*

$$\frac{1}{\xi_k} \sup_z \frac{|\widehat{Q}_k(z) - Q(z)|}{\sqrt{Q(z)(1-Q(z))}} \rightarrow 1,$$

in probability as  $k \rightarrow \infty$ . In particular,

$$\widehat{Q}_k(z) = Q(z) + O_{\mathbb{P}}\left(\xi_k \sqrt{Q(z)(1-Q(z))}\right),$$

uniformly in  $z$ .

Apply Lemma A.2 to the two summands in  $\widehat{G}$ , we obtain  $\widehat{G}(t) = G(t) + \widehat{R}(t)$ , where

$$G(t) = \frac{p-s}{p} \overline{F}_0(t) + \frac{s}{p} \overline{F}_a(t), \quad (\text{A.27})$$

and

$$\widehat{R}(t) = O_{\mathbb{P}}\left(\xi_p \sqrt{\overline{F}_0(t)F_0(t)} + \frac{s}{p} \xi_s \sqrt{\overline{F}_a(t)F_a(t)}\right), \quad (\text{A.28})$$

uniformly in  $t$ .

Recall (see proof of Lemma A.1) that the BH procedure is the thresholding procedure with threshold set at

$$\tau = \inf\{t \mid \overline{F}_0(t) \leq \alpha \widehat{G}(t)\}. \quad (\text{A.29})$$

The NDP may also be re-written as

$$\text{NDP} = \frac{|S \setminus \widehat{S}|}{|S|} = \frac{1}{s} \sum_{i \in S} \mathbb{1}\{x(i) < \tau\} = 1 - \widehat{W}_{\text{signal}}(\tau),$$

so that it suffices to show that

$$\widehat{W}_{\text{signal}}(\tau) \rightarrow 1 \quad (\text{A.30})$$

in probability. Applying Lemma A.2 to  $\widehat{W}_{\text{signal}}$ , we know that

$$\widehat{W}_{\text{signal}}(\tau) = \overline{F}_a(\tau) + O_{\mathbb{P}} \left( \xi_s \sqrt{\overline{F}_a(\tau) F_a(\tau)} \right) = \overline{F}_a(\tau) + o_{\mathbb{P}}(1).$$

So it suffices to show that  $F_a(\tau) \rightarrow 0$  in probability. Now let  $t^* = \sqrt{2q \log(p)}$  for some  $q$  such that  $\beta < q < \underline{r}$ . We have

$$F_a(t^*) = \Phi(t^* - \underline{\Delta}) = \Phi(\sqrt{2(q - \underline{r}) \log p}) \rightarrow 0. \quad (\text{A.31})$$

Hence in order to show (A.30), it suffices to show

$$\mathbb{P}[\tau \leq t^*] \rightarrow 1. \quad (\text{A.32})$$

By (A.27), the mean of the empirical process  $\widehat{G}$  evaluated at  $t^*$  is

$$G(t^*) = \frac{p-s}{p} \overline{F}_0(t^*) + \frac{s}{p} \overline{F}_a(t^*). \quad (\text{A.33})$$

The first term, using Relation (A.21), is asymptotic to  $p^{-q} L(p)$ , where  $L(p)$  is the logarithmic term in  $p$ . The second term, since  $\overline{F}_a(t^*) \rightarrow 1$  by Relation (A.31), is asymptotic to  $p^{-\beta}$ . Therefore,  $G(t^*) \sim p^{-q} L(p) + p^{-\beta} \sim p^{-\beta}$ , since  $p^{\beta-q} L(p) \rightarrow 0$  where  $q > \beta$ .

The fluctuation of the empirical process at  $t^*$ , by Relation (A.28), is

$$\begin{aligned} \widehat{R}(t^*) &= O_{\mathbb{P}} \left( \xi_p \sqrt{\overline{F}_0(t^*) F_0(t^*)} + \frac{s}{p} \xi_s \sqrt{\overline{F}_a(t^*) F_a(t^*)} \right) \\ &= O_{\mathbb{P}} \left( \xi_p \sqrt{\overline{F}_0(t^*)} \right) + o_{\mathbb{P}}(p^{-\beta}). \end{aligned}$$

By (A.21) and the expression for  $\xi_p$ , the first term is  $O_{\mathbb{P}}(p^{-(q+1)/2} L(p))$  where  $L(p)$  is a poly-logarithmic term in  $p$ . Since  $\beta < \min\{q, 1\}$ , we have  $\beta < (q+1)/2$ , and hence  $\widehat{R}(t^*) = o_{\mathbb{P}}(p^{-\beta})$ .

Putting the mean and the fluctuation of  $\widehat{G}(t^*)$  together, we obtain

$$\widehat{G}(t^*) = G(t^*) + \widehat{R}(t^*) \sim_{\mathbb{P}} G(t^*) \sim p^{-\beta},$$

and therefore, together with (A.21), we have

$$\overline{F}_0(t^*)/\widehat{G}(t^*) = p^{\beta-q} L(p)(1 + o_{\mathbb{P}}(1)),$$

which is eventually smaller than the FDR level  $\alpha$  by the assumption (3.8) and the fact that  $\beta < q$ . That is,

$$\mathbb{P} \left[ \overline{F}_0(t^*)/\widehat{G}(t^*) < \alpha \right] \rightarrow 1.$$

By definition of  $\tau$  (recall (A.29)), this implies that  $\tau \leq t^*$  with probability tending to 1, and (A.32) is shown. The proof for the sufficient condition is complete.  $\square$

## A.5 Proof of Theorems 3.4 and 3.5

Proof of Theorem 3.4 uses ideas from the proof of Theorem 3.3 and is substantially shorter.

*Proof (Proof of Theorem 3.4).* We first show the sufficient condition. Vanishing FWER is guaranteed by the properties of the procedures, and we only need to show that FNR also goes to zero. Similar to the proof of Theorem 3.3, it suffices to show that

$$\text{NDP} = 1 - \widehat{W}_{\text{signal}}(t_p) \rightarrow 0, \quad (\text{A.34})$$

where  $t_p$  is the threshold of Bonferroni's procedure.

Since  $\alpha$  vanishes slowly (see Definition 3.8), for any  $\delta > 0$ , we have  $p^{-\delta} = o(\alpha)$ . Therefore, we have  $-\log \alpha \leq \delta \log p$  for large  $p$ , and

$$1 \leq \limsup_{p \rightarrow \infty} \frac{2 \log p - 2 \log \alpha}{2 \log p} \leq 1 + \delta,$$

for any  $\delta > 0$ . Therefore, by the expression for normal quantiles, we know that

$$t_p = F^{\leftarrow}(1 - \alpha/p) \sim (2 \log p - 2 \log \alpha)^{1/2} \sim (2 \log p)^{1/2}.$$

Since  $\underline{r} > \widetilde{g}(\beta) = 1$ , we can pick  $q$  such that  $1 < q < \underline{r}$ . Let  $t^* = \sqrt{2q \log p}$ , we know that  $t_p < t_p^*$  for large  $p$ . Therefore for large  $p$ , we have

$$\widehat{W}_{\text{signal}}(t_p) \geq \widehat{W}_{\text{signal}}(t^*) \geq \overline{F}_a(t^*) + o_{\mathbb{P}}(1),$$

where  $\overline{F}_a$  is the survival function of  $N(\sqrt{2\underline{r} \log p}, 1)$ ; the last inequality follows from the stochastic monotonicity of the Gaussian location family (A.4), and Lemma A.2. Indeed, by our choice of  $q < \underline{r}$ , we obtain

$$F_a(t^*) = \Phi\left(\sqrt{2(q - \underline{r}) \log p}\right) \rightarrow 0,$$

and (A.34) is shown. This completes the proof of the sufficient condition.

The proof of the necessary condition follows similar structure as in the proof of Theorem 3.3, and uses the lower bound

$$\text{FWER}(\mathcal{R}) + \text{FNR}(\mathcal{R}) \geq \mathbb{P}\left[\max_{i \in S^c} x(i) > u\right] \wedge \mathbb{E}\left[\frac{|S \setminus \widehat{S}(u)|}{|S|}\right], \quad (\text{A.35})$$

which holds for any arbitrary thresholding procedure  $\mathcal{R}$  and arbitrary real  $u \in \mathbb{R}$ .

By the assumption that  $\bar{r} < \tilde{g}(\beta) = 1$ , we can pick  $q$  such that  $\bar{r} < q < 1$  and let  $u = t^* = \sqrt{2q \log p}$  in (A.35). By relative stability of iid Gaussian random variables (A.1), we have

$$\mathbb{P} \left[ \frac{\max_{i \in S^c} x(i)}{\sqrt{2 \log p}} > \frac{t^*}{\sqrt{2 \log p}} \right] \rightarrow 1. \quad (\text{A.36})$$

since the first fraction in (A.36) converges to 1, while the second converges to  $q < 1$ . Therefore, the first term on the right-hand side of (A.35) converges to 1.

On the other hand, by the stochastic monotonicity of Gaussian location family (A.4), the probability of missed detection for each signal is lower bounded by  $\mathbb{P}[Z + \mu(i) \leq t^*] \geq F_{\bar{a}}(t^*)$ , where  $Z$  is a standard Gaussian r.v., and  $F_{\bar{a}}$  is the cdf of  $N(\sqrt{2\bar{r} \log p}, 1)$ . Therefore,  $|S \setminus \widehat{S}(t^*)| \stackrel{d}{\geq} \text{Binom}(s, F_{\bar{a}}(t^*))$ , and it suffices to show that  $F_{\bar{a}}(t^*)$  converges to 1. Indeed,

$$F_{\bar{a}}(t^*) = \Phi(\sqrt{2(q - \bar{r}) \log p}) \rightarrow 1,$$

by our choice of  $q > \bar{r}$ . Hence both quantities in the minimum on the right-hand side of (A.35) converge to 1 in the limit, and the necessary condition is shown.  $\square$

*Proof (Proof of Theorem 3.5).* We first show the sufficient condition. Since FDR control is guaranteed by the BH procedure, we only need to show that the FWNH also vanishes, that is,

$$\mathbb{P} \left[ \min_{i \in S} x(i) \geq \tau \right] \rightarrow 1, \quad (\text{A.37})$$

where  $\tau$  is the threshold for the BH procedure.

By the assumption that  $\underline{r} > \tilde{h}(\beta) = (\sqrt{\beta} + \sqrt{1 - \beta})^2$ , we have  $\sqrt{\underline{r}} - \sqrt{1 - \beta} > \sqrt{\beta}$ , so we can pick  $q > 0$ , such that

$$\sqrt{\underline{r}} - \sqrt{1 - \beta} > \sqrt{q} > \sqrt{\beta}. \quad (\text{A.38})$$

We only need to show that with a specific choice of  $t^* = \sqrt{2q \log p}$  where

$$\sqrt{\underline{r}} - \sqrt{1 - \beta} > \sqrt{q} > \sqrt{\beta}, \quad (\text{A.39})$$

we have both

$$\mathbb{P}[\tau \leq t^*] \rightarrow 1, \quad (\text{A.40})$$

and

$$\mathbb{P} \left[ \min_{i \in S} x(i) \geq t^* \right] \rightarrow 1, \quad (\text{A.41})$$

so that

$$\mathbb{P} \left[ \min_{i \in S} x(i) \geq \tau \right] \geq \mathbb{P} \left[ \min_{i \in S} x(i) \geq t^*, t^* \geq \tau \right] \rightarrow 1.$$

Relation (A.40) follows in exactly the same way (A.32) did on page 127.

Dividing the left-hand-side in Relation (A.41) by  $\sqrt{2 \log p}$ , we have,

$$\begin{aligned} \frac{\min_{i \in S} x(i)}{\sqrt{2 \log p}} &= \frac{\min_{i \in S} \mu(i) + \epsilon(i)}{\sqrt{2 \log p}} \stackrel{d}{\geq} \frac{\sqrt{2r} \log p + \min_{i \in S} \epsilon(i)}{\sqrt{2 \log p}} \\ &\rightarrow -\sqrt{1 - \beta} + \sqrt{r}, \end{aligned}$$

where the last convergence follows from the relative stability of iid Gaussians minima (A.2). On the other hand,  $t^*/\sqrt{2 \log p} = \sqrt{q} < \sqrt{r} - \sqrt{1 - \beta}$  by our choice of  $q$ , and Relation (A.41) follows.

The necessary condition follows from the lower bound

$$\text{FDR}(\mathcal{R}) + \text{FWNR}(\mathcal{R}) \geq \mathbb{E} \left[ \frac{|\widehat{S}(u) \setminus S|}{|\widehat{S}(u) \setminus S| + |S|} \right] \wedge \mathbb{P} \left[ \min_{i \in S} x(i) < u \right], \quad (\text{A.42})$$

which holds for any thresholding procedure  $\mathcal{R}$  and for arbitrary  $u \in \mathbb{R}$ . In particular, we show that both terms in the minimum in (A.42) converge to 1 when we set  $u = t^* = \sqrt{2q \log p}$  where

$$\sqrt{r} - \sqrt{1 - \beta} < \sqrt{q} < \sqrt{\beta}. \quad (\text{A.43})$$

On the one hand, we have,

$$\frac{\min_{i \in S} x(i)}{\sqrt{2 \log p}} \stackrel{d}{\leq} \frac{\min_{i \in S} \epsilon(i) + \sqrt{2r} \log p}{\sqrt{2 \log p}} \rightarrow \sqrt{r} - \sqrt{1 - \beta},$$

by relative stability of iid Gaussians (A.2). On the other hand,  $t^*/\sqrt{2 \log p} = \sqrt{q} > \sqrt{r} - \sqrt{1 - \beta}$  by our choice of  $q$ ; this shows that the second term on the right-hand side of (A.42) converges to 1.

Observe that  $|\widehat{S}(t^*) \setminus S|$  has distribution  $\text{Binom}(p - s, \overline{\Phi}(t^*))$ , and define  $X = X_p := |\widehat{S}(t^*) \setminus S|/|S|$ , we obtain,

$$\begin{aligned} \mu &:= \mathbb{E}[X] = (p^\beta - 1) \overline{\Phi}(t^*) \sim (p^\beta - 1) \frac{\phi(t^*)}{t^*} \\ &\sim \frac{1}{\sqrt{2\pi}} (2q \log p)^{-1/2} p^{\beta-q} \rightarrow \infty, \end{aligned}$$

where the divergence follows from our choice of  $q < \beta$ . Using again Relations (A.22) and (A.23), we conclude that the first term on the right-hand side of (A.42) also converges to 1. This completes the proof of the necessary condition.  $\square$

## Appendix B

### Exact support recovery in non AGG models

#### B.1 Strong classification boundaries in other light-tailed error models

The strong classification boundaries extend beyond the AGG models. As our analysis in Section ?? suggests, all additive error models where the errors have URS maxima demonstrate this phase transition phenomenon under appropriate parametrization of the sparsity and signal sizes. We derive explicit boundaries for two additional classes of models under the general form of the additive noise models (1.1), with heavier and lighter tails than the AGG models, respectively.

We would like to point out that the sparsity and signal sizes can be re-parametrized for the boundaries to have different shapes. For example in the case of Gaussian errors, if we re-parametrize sparsity  $s$  with  $\tilde{\beta} = 2 - \left(1 + \sqrt{1 - \beta}\right)^2$  where  $\tilde{\beta} \in (0, 1)$ , then the signal sparsity would have a slightly more complicated form:

$$|S_p| = \lfloor p^{1-\beta} \rfloor = \left\lfloor p^{\left(\sqrt{2-\tilde{\beta}}-1\right)^2} \right\rfloor,$$

while the strong classification boundary would take on the simpler form:

$$g(\beta) = \tilde{g}(\tilde{\beta}) = 2 - \tilde{\beta}. \quad (\text{B.1})$$

In the next two classes of models we will adopt parametrizations such that the boundaries are of the form  $\tilde{g}$  in (B.1).

##### *B.1.1 Additive error models with heavier-than-AGG tails*

Distributions such as the log-normal have heavier tails than the AGG model, yet the tails are nevertheless rapidly-varying. Therefore, Proposition 2.2 applies, and

we expect to see phase-transition-type results when the additive errors have these heavier-than-AGG tails.

*Example B.1 (Heavier than AGG).* Let  $\gamma > 1$ ,  $c > 0$ , and suppose that

$$\log \bar{F}(x) = -(\log x)^\gamma (c + M(x)), \quad (\text{B.2})$$

where  $\lim_{x \rightarrow \infty} M(x) \log^\gamma x = 0$ . Then, Relation (2.39) holds under model (B.2). Further, if the entries in the array are independent, the maxima are relatively stable.

The behavior of the quantiles  $u_p$  in this model is as follows. As  $p \rightarrow \infty$ ,

$$u_p \sim \exp \left\{ \left( c^{-1} \log p \right)^{1/\gamma} \right\} \iff c (\log u_p)^\gamma + o(1) = \log(p) = -\log \bar{F}(u_p).$$

since  $u_p$  diverges, and  $M(u_p)$  is  $o((\log^\gamma u_p)^{-1})$ .

Following Example B.1, assume that the errors in Model (1.1) have rapidly varying right tails

$$\log \bar{F}(x) = -(\log x)^\gamma (c + M(x)), \quad (\text{B.3})$$

as  $x \rightarrow \infty$ , and left tails

$$\log F(x) = -(\log(-x))^\gamma (c + M(-x)), \quad (\text{B.4})$$

as  $x \rightarrow -\infty$ .

**Theorem B.1.** *Suppose the marginals  $F$  follows (B.3) and (B.4). Let*

$$k(\beta) = \log p - \left( (\log p)^{1/\gamma} + \log(1 - \beta) \right)^\gamma,$$

and let the signal  $\mu$  have

$$|S_p| = \left\lfloor p e^{-k(\beta)} \right\rfloor$$

non-zero entries. Assume the magnitudes of non-zero signal entries are in the range between

$$\underline{\Delta} = \exp \left\{ (\log p)^{1/\gamma} \right\} \underline{r} \quad \text{and} \quad \bar{\Delta} = \exp \left\{ (\log p)^{1/\gamma} \right\} \bar{r}.$$

If  $\underline{r} > \tilde{g}(\beta) = 2 - \beta$ , then Bonferroni's procedure  $\hat{S}_p$  (defined in (2.21)) with appropriately calibrated FWER  $\alpha \rightarrow 0$  achieves asymptotic perfect support recovery, under arbitrary dependence of the errors.

On the other hand, when the errors are uniformly relatively stable, if  $\bar{r} < \tilde{g}(\beta) = 2 - \beta$ , then no thresholding procedure can achieve asymptotic perfect support recovery with positive probability.



### B.1.2 Additive error models with lighter-than-AGG tails

Similar to how Proposition 2.2 applies to models with heavier-than-AGG tails, it also to error models with lighter tails than the AGG class.

*Example B.2 (Lighter than AGG).* With  $\nu > 0$ , and  $L(x)$  a slowly varying function, the class of distributions

$$\log \bar{F}(x) = -\exp\{x^\nu L(x)\}, \quad (\text{B.5})$$

is rapidly varying. The quantiles can be derived explicitly in a subclass of (B.5) where  $L(x) \rightarrow 1$ , or equivalently, when  $\log |\log \bar{F}(x)| \sim x^\nu$ ,

$$u_p \sim (\log \log p)^{1/\nu} \iff \exp\{u_p^\nu (1 + o(1))\} = \log(p) = -\log \bar{F}(u_p).$$

Following Example B.2, assume that errors in Model (1.1) has rapidly varying right tails

$$\log \bar{F}(x) = -\exp\{x^\nu L(x)\}, \quad (\text{B.6})$$

where  $L(x)$  is a slowly varying function, as  $x \rightarrow \infty$ , and left tails

$$\log \bar{F}(x) = -\exp\{-x^\nu L(-x)\}, \quad (\text{B.7})$$

as  $x \rightarrow -\infty$ .

The phase transition results in multiple testing problems under such tail assumptions is characterizes as follows.

**Theorem B.2.** *Suppose marginals  $F$  follow (B.6) and (B.7). Let*

$$k(\beta) = \log p - (\log(p))^{(1-\beta)^\nu},$$

*and let the signal  $\mu$  have*

$$|S_p| = \left\lfloor p e^{-k(\beta)} \right\rfloor$$

*non-zero entries. Assume the magnitudes of non-zero signal entries are in the range between*

$$\underline{\Delta} = \log \log p^{1/\nu} \underline{r} \quad \text{and} \quad \bar{\Delta} = \log \log p^{1/\nu} \bar{r}.$$

*If  $\underline{r} > \bar{g}(\beta) = 2 - \beta$ , then Bonferroni's procedure  $\hat{S}_p$  (defined in (2.21)) with appropriately calibrated FWER  $\alpha \rightarrow 0$  achieves asymptotic perfect support recovery, under arbitrary dependence of the errors.*

*On the other hand, when the errors are uniformly relatively stable, if  $\bar{r} < \bar{g}(\beta) = 2 - \beta$ , then no thresholding procedure can achieve asymptotic perfect support recovery with positive probability.*

## B.2 Thresholding procedures under heavy-tailed errors

We analyze the performance of thresholding estimators under heavy-tailed models in this section, and illustrate its lack of phase transition. Suppose we have iid errors with Pareto tails in Model (1.1), that is,  $\epsilon(i)$ 's have common marginal distribution  $F$  where

$$\overline{F}(x) \sim x^{-\alpha} \quad \text{and} \quad F(-x) \sim x^{-\alpha}, \quad (\text{B.8})$$

as  $x \rightarrow \infty$ . It is well-known (see, e.g., Theorem 1.6.2 of (Leadbetter et al., 1983)) that the maxima of iid Pareto random variables have Frechet-type limits. Specifically, we have

$$\frac{\max_{i \in \{1, \dots, p\}} \epsilon(i)}{u_p} \implies Y, \quad (\text{B.9})$$

in distribution, where  $u_p = F^{\leftarrow}(1 - 1/p) \sim p^{1/\alpha}$ , and  $Y$  is a standard  $\alpha$ -Frechet random variable, i.e.,

$$\mathbb{P}[Y \leq t] = \exp\{-t^{-\alpha}\}, \quad t > 0.$$

By symmetry in our assumptions, the same argument applies to the minima as well.

**Theorem B.3.** *Let errors in Model (1.1) be as described in Relation (B.8). Let the signal have  $s = |S| = fp$  non-zero entries, with magnitude  $\Delta = rp^{1/\alpha}$ , where both  $f \in (0, 1)$  and  $r \in (0, +\infty)$  may depend on  $p$ , so that no generality is lost.*

*Under these assumptions, the necessary condition for thresholding procedures  $\widehat{S}$  to achieve exact support recovery ( $\mathbb{P}[\widehat{S} = S] \rightarrow 1$ ) is*

$$\liminf_{p \rightarrow \infty} r = \infty. \quad (\text{B.10})$$

*Condition (B.10) is also sufficient for the oracle thresholding procedure to succeed in the exact support recovery problem.*

*On the other hand, the necessary and sufficient condition for all thresholding procedures to fail exact support recovery ( $\mathbb{P}[\widehat{S} = S] \rightarrow 0$ ) is*

$$\limsup_{p \rightarrow \infty} r = 0.$$

In other words, Theorem B.3 states that there does not exist a non-trivial phase transition for thresholding procedures when errors have (two-sided)  $\alpha$ -Pareto tails.

*Proof (Proof of Theorem B.3).* Recall the oracle thresholding procedure  $\widehat{S}^* = \{i : x(i) \geq x_{[s]}\}$ , and the set of all thresholding procedures, denoted  $\mathcal{S}$  (see Definition 2.20). The probability of exact support recovery by any thresholding procedure  $\widehat{S} \in \mathcal{S}$  is bounded above by that of  $\widehat{S}^*$ , that is,

$$\begin{aligned}
\max_{\widehat{S} \in S} \mathbb{P}[\widehat{S} = S] &= \mathbb{P}[\widehat{S}^* = S] = \mathbb{P}\left[\max_{i \in S^c} x(i) \leq \min_{i \in S} x(i)\right] \\
&= \mathbb{P}\left[\frac{\max_{i \in S^c} x(i)}{u_p} \leq \frac{\min_{i \in S} x(i)}{u_p}\right] \\
&= \mathbb{P}\left[\frac{M_{S^c}}{u_p} \leq \frac{m_S}{u_p} + r_p\right], \tag{B.11}
\end{aligned}$$

where  $M_{S^c} = \max_{i \in S^c} \epsilon(i)$  and  $m_S = \min_{i \in S} \epsilon(i)$ . For any  $\alpha > 0$ , the following elementary relations hold,

$$0 < L \leq (1-f)^{1/\alpha} + f^{1/\alpha} \leq U < \infty, \quad \text{for all } f \in (0, 1),$$

where  $L = \min\{1, 2(1/2)^{1/\alpha}\}$  and  $U = \max\{1, 2(1/2)^{1/\alpha}\}$ . Therefore we have,

$$U \max\left\{\frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p}\right\} < r_p \implies (1-f)^{1/\alpha} \frac{M_{S^c}}{u_p} - f^{1/\alpha} \frac{m_S}{u_p} < r_p, \tag{B.12}$$

and

$$L \min\left\{\frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p}\right\} < r_p \iff (1-f)^{1/\alpha} \frac{M_{S^c}}{u_p} - f^{1/\alpha} \frac{m_S}{u_p} < r_p. \tag{B.13}$$

Putting together (B.11), (B.12), and (B.13), we have

$$\mathbb{P}\left[\max\left\{\frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p}\right\} < r_p/U\right] \leq \mathbb{P}[\widehat{S}^* = S] \leq \mathbb{P}\left[\min\left\{\frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p}\right\} < r_p/L\right]. \tag{B.14}$$

We know from the weak convergence result (B.9) that for any  $\epsilon > 0$  there is a constant  $N$  such that for all  $p > N$  we have

$$\mathbb{P}\left[\max\left\{\frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p}\right\} < r_p/U\right] \geq \mathbb{P}\left[\max\{Y^{(1)}, Y^{(2)}\} < r_p/U\right] - \epsilon, \tag{B.15}$$

where  $Y^{(1)}$  and  $Y^{(2)}$  are independent  $\alpha$ -Frechet random variables with scale coefficients  $(1-f)^{1/\alpha}$  and  $f^{1/\alpha}$  respectively. That is,

$$\mathbb{P}[Y^{(1)} \leq t] = \exp\{-(1-f)/t^\alpha\}, \quad \text{and} \quad \mathbb{P}[Y^{(2)} \leq t] = \exp\{-f/t^\alpha\}.$$

Since the distributional limit in (B.15) has a density (with respect to the Lebesgue measure), we know that density is bounded above by a finite constant, say,  $K$ . For the same choice of  $\epsilon$  as before, we can find a further constant  $N'$  such that for all  $p > \max\{N, N'\}$  we have

$$\liminf r_p < \epsilon/K + r_p,$$

so that the right hand side of (B.15) is bounded by

$$\mathbb{P}\left[\max\{Y^{(1)}, Y^{(2)}\} < r_p/U\right] - \epsilon \geq \mathbb{P}\left[\max\{Y^{(1)}, Y^{(2)}\} < \frac{\liminf r_p}{U}\right] - 2\epsilon. \tag{B.16}$$

By the arbitrariness in the choice of  $\epsilon$ , we conclude from (B.15) and (B.16) that

$$\liminf \mathbb{P} \left[ \max \left\{ \frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p} \right\} < r_p/U \right] \geq \mathbb{P} \left[ \max \{Y^{(1)}, Y^{(2)}\} < \frac{\liminf r_p}{U} \right]. \quad (\text{B.17})$$

Combining Relations (B.14) and (B.17), we know that if  $\liminf r_p = \infty$ , we must have

$$\liminf \mathbb{P} \left[ \widehat{S}^* = S \right] \geq \mathbb{P} \left[ \max \{Y^{(1)}, Y^{(2)}\} < \frac{\liminf r_p}{U} \right] = 1.$$

Conversely, if  $\liminf \mathbb{P} \left[ \widehat{S}^* = S \right] < 1$ , we must have  $\liminf r_p < \infty$ .

Similarly, we can obtain the upper bound of exact support recovery probability for the optimal thresholding procedure,

$$\limsup \mathbb{P} \left[ \min \left\{ \frac{M_{S^c}}{u_p}, -\frac{m_S}{u_p} \right\} < r_p/L \right] \leq \mathbb{P} \left[ \min \{Y^{(1)}, Y^{(2)}\} < \frac{\limsup r_p}{L} \right]. \quad (\text{B.18})$$

The conclusions of the second part of Theorem B.3 follow from (B.14) and (B.18).

□

The probability of exact recovery can be approximated if the parameters  $r$  and  $f$  converge. The next result follows from a small modification of the arguments in the proof of Theorem B.3.

**Corollary B.1.** *Under the assumptions in Theorem B.3, if  $\lim r = r^*$ , and  $\lim f = f^*$ , for some constant  $r^* \geq 0$  and  $f^* \in [0, 1]$ , then*

$$\lim \mathbb{P}[\widehat{S}^* = S] = \mathbb{P} \left[ (1 - f^*)^{1/\alpha} Z_1 + (f^*)^{1/\alpha} Z_2 < r^* \right].$$

where  $Z_1$  and  $Z_2$  are independent standard  $\alpha$ -Frechet random variables, i.e.,  $\mathbb{P}[Z_i \leq x] = \exp \{-x^{-\alpha}\}$ ,  $x > 0$ .

*Remark B.1.* Of course one might wonder if it would be meaningful to derive a “phase transition” under a different parametrization of the signal sizes, say

$$\Delta = p^{r/\alpha}. \quad (\text{B.19})$$

In this case, Theorem B.3 suggests that a “phase transition” takes place at  $r = 1$ . However, this non-multiplicative parametrization of the signal sizes would make power analysis (like in Example 3.1) dimension-dependent.

To illustrate, in the case of Gaussian errors with variance 1, if we were interested in small signals of size  $\sqrt{2r \log p}$ , where  $r < 1$  is below the boundary (4.5), then we only need  $n > 2/r$  samples to guarantee discovery of their support. In the Pareto case with parametrization (B.19), however, if we were interested in small signals of size  $p^{r/\alpha}$ , where  $r < 1$ , then the “boundary” says that we will need  $n > p^{2(1-r)/\alpha}$  samples, which is exponential in the dimension  $p$  and quickly diverges. Recall that the “boundary” is really an asymptotic result in  $p$ . Such an approximation in finite dimensions becomes invalid.