# The Price of Competition: Effect Size Heterogeneity Matters in High Dimensions

Hua Wang<sup>®</sup>, Yachong Yang<sup>®</sup>, and Weijie J. Su<sup>®</sup>, Member, IEEE

Abstract—In high-dimensional sparse regression, would increasing the signal-to-noise ratio while fixing the sparsity level always lead to better model selection? For high-dimensional sparse regression problems, surprisingly, in this paper we answer this question in the negative in the regime of linear sparsity for the Lasso method, relying on a new concept we term effect size heterogeneity. Roughly speaking, a regression coefficient vector has high effect size heterogeneity if its nonzero entries have significantly different magnitudes. From the viewpoint of this new measure, we prove that the false and true positive rates achieve the optimal trade-off uniformly along the Lasso path when this measure is maximal in a certain sense, and the worst trade-off is achieved when it is minimal in the sense that all nonzero effect sizes are roughly equal. Moreover, we demonstrate that the first false selection occurs much earlier when effect size heterogeneity is minimal than when it is maximal. The underlying cause of these two phenomena is, metaphorically speaking, the "competition" among variables with effect sizes of the same magnitude in entering the model. Taken together, our findings suggest that effect size heterogeneity shall serve as an important complementary measure to the sparsity of regression coefficients in the analysis of high-dimensional regression problems. Our proofs use techniques from approximate message passing theory as well as a novel technique for estimating the rank of the first false variable.

Index Terms—Approximate message passing, false discovery rate, high-dimensional sparse regression, model selection, signal-to-noise ratio.

#### I. INTRODUCTION

E CONSIDER high-dimensional sparse regression problems in which we observe an n-dimensional response vector y that is generated by a linear model

$$y = X\beta + z, \tag{I.1}$$

where X is an  $n \times p$  design matrix of features,  $\beta \in \mathbb{R}^p$  denotes an unknown vector of regression coefficients, and

Manuscript received 6 March 2021; revised 20 November 2021; accepted 5 March 2022. Date of publication 13 April 2022; date of current version 13 July 2022. This work was supported in part by NSF through CAREER under Grant DMS-1847415, Grant CCF-1763314, and Grant CCF-1934876; in part by the Wharton Dean's Research Fund; and in part by the Facebook Faculty Research Award. (Hua Wang and Yachong Yang are co-first authors.) (Corresponding author: Weijie J. Su.)

The authors are with the Department of Statistics and Data Science, University of Pennsylvania, Philadelphia, PA 19104 USA (e-mail: suw@wharton.upenn.edu).

Communicated by A. Maleki, Associate Editor for Signal Processing and Source Coding.

Color versions of one or more figures in this article are available at https://doi.org/10.1109/TIT.2022.3166720.

Digital Object Identifier 10.1109/TIT.2022.3166720

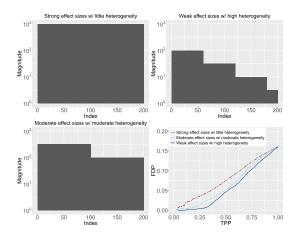


Fig. 1. The TPP-FDP trade-off along the entire Lasso path, with three different sets of regression coefficients. Note that the TPP-FDP trade-off is equivalent to the receiver operating characteristic curve. The sparsity of  $\boldsymbol{\beta}$  is fixed to k=200 (throughout this paper, we use k to denote the sparsity level) and the 200 true effects are plotted in the logarithmic scale in the three panels. For example, in the "Strong effect sizes w/ little heterogeneity" setting,  $\beta_1=\dots=\beta_{200}=10^3$ , and  $\beta_{201}=\dots=\beta_{1000}=0$ . The design matrix  $\boldsymbol{X}\in\mathbb{R}^{n\times p}$  has independent  $\mathcal{N}(0,1/n)$  entries, where n=p=1000, and the noise term  $\boldsymbol{z}$  has independent  $\mathcal{N}(0,\sigma^2)$  entries with  $\sigma=0.01$ . The bottom-right panel shows the plot of FDP as a function of TPP, averaged over 100 independent runs. For completeness, we remark that effect size heterogeneity influences model selection in a more complex manner at a higher noise level (see Figure 10).

 $z \in \mathbb{R}^n$  is a noise term. In the big data era, this model has been increasingly applied to high-dimensional settings where the number of variables p is comparable to or even much larger than the number of observational units n. While this reality poses challenges to the regression problem, in many scientific problems there are good reasons to suspect that truly relevant variables account for a small fraction of all the observed variables or, equivalently,  $\beta$  is sparse in the sense that many of its components are zero or nearly zero. Indeed, a very impressive body of *theoretical* work shows that the difficulty of variable selection in the high-dimensional setting relies crucially on how sparse the regression coefficients are [1], [2].

This paper, however, asks whether there are other measures concerning the regression coefficients that have a *practical* impact on variable selection for the linear model (I.1). To address this question, we present a simulation study in Figure 1. Notably, the sparsity—or, put differently, the number of nonzero components—of the regression coefficients  $\boldsymbol{\beta}$  is *fixed* to 200 across three experimental settings, but with vary-

0018-9448 © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

ing magnitudes of the 200 true effect sizes. The method we use for variable selection is the Lasso [3], which is perhaps the most popular model selector in the high-dimensional setting. Given a penalty parameter  $\lambda>0$ , this method finds the solution to the convex optimization program

$$\widehat{\boldsymbol{\beta}}(\lambda) = \underset{\boldsymbol{b} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{b}\|^2 + \lambda \|\boldsymbol{b}\|_1, \quad (I.2)$$

where  $\|\cdot\|$  and  $\|\cdot\|_1$  denote the  $\ell_2$  and the  $\ell_1$  norms, respectively. A variable j is selected by this method at  $\lambda$  if  $\widehat{\beta}_j(\lambda) \neq 0$ , and a false selection occurs if it is a noise variable in the sense that  $\beta_j = 0$ . Formally, we use the false discovery proportion (FDP) and true positive proportion (TPP) as measures of the type I error and power, respectively, to assess the quality of the selected model  $\{1 \leq j \leq p : \widehat{\beta}_j(\lambda) \neq 0\}$ :

$$FDP_{\lambda} = \frac{\#\{j : \beta_j = 0 \text{ and } \widehat{\beta}_j(\lambda) \neq 0\}}{\#\{j : \widehat{\beta}_j(\lambda) \neq 0\}}, \quad (I.3)$$

$$TPP_{\lambda} = \frac{\#\{j : \beta_j \neq 0 \text{ and } \widehat{\beta}_j(\lambda) \neq 0\}}{\#\{j : \beta_j \neq 0\}}.$$
 (I.4)

As is clear, we wish to select a model with a small FDP and large TPP.

Despite weaker effect sizes, strikingly, Figure 1 shows that the Lasso can achieve *better* model selection in terms of the TPP–FDP trade-off and, in particular, this counterintuitive behavior holds uniformly along the entire Lasso path or, equivalently, over all values of  $\lambda$ . Existing theory often analyzes how the worst-case performance of the Lasso and other related procedures depends on the regression coefficients through the sparsity of the regression coefficients (see, for example, [4]). However, the sparsity level is fixed across the experimental settings of Figure 1. In light of this, therefore, one would expect that the strong signals and weak signals would yield the best and worst model selection results, respectively. Figure 1 shows that this is not necessarily the case.

Thus, a *finer-grained* structural analysis of the effect sizes is needed to better understand the Lasso in some settings. In this paper, we address this important question by proposing a concept that we term effect size heterogeneity concerning the regression coefficients in high dimensions. Roughly speaking, a regression coefficient vector has higher effect size heterogeneity than another vector (of the same sparsity) if the nonzero entries of the former are more heterogeneous than those of the latter in terms of magnitude. As a complement to sparsity, effect size heterogeneity will be shown to have a significant impact on how the Lasso performs in terms of the false and true positive rates trade-off: while the sparsity level of the regression coefficients is fixed, the higher the effect size heterogeneity is, the better the Lasso performs. Turning back to Figure 1, we note that the strong effect sizes are the least heterogeneous in magnitude, and the weak effect sizes are the most heterogeneous. Therefore, the comparisons made in Figure 1 match well the implication of effect size heterogeneity.

Concretely, the main thrust of this paper lies in the development of two complementary perspectives to precisely quantify

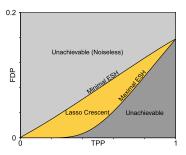


Fig. 2. The Lasso Crescent diagram specified by the parameters n/p=1 and k/p=0.2, following the setting in Figure 1. The lower/upper smooth curve is asymptotically achieved with maximal/minimal effect size heterogeneity (ESH) in the regime of infinite signal-to-noise ratio. The explicit expressions of the curves are given in Section II. Our Theorem 1 implies that nowhere on the Lasso path we can find any (TPP, FDP) pairs in the region below the Lasso Crescent (also see [5]). In the noiseless setting, moreover, this impossibility result continues to hold in the region above the Lasso Crescent (Theorem 2), which is labeled "Unachievable (Noiseless)."

the impact of effect size heterogeneity. First, following the setup of Figure 1, we consider the full possible range of the asymptotic trade-off between the TPP and FDP along the Lasso path, while varying the level of effect size heterogeneity. Assuming a random design with independent Gaussian entries and working in the regime of linear sparsity—meaning that the fraction of true effect sizes tends to a constant—we formally show that the full possible range is enclosed by two smooth curves in the (TPP, FDP) plane, which we referred to as the Lasso Crescent. Figure 2 presents an instance of the Lasso Crescent. More precisely, having excluded the impact of noise by taking z = 0 in the linear model (I.1), the lower curve is asymptotically achieved when effect size heterogeneity is maximal in the sense that all true effect sizes are widely different from each other, while the upper curve is asymptotically achieved when the heterogeneity is minimal in the sense that all true effects are of the same size. In general, the (TPP, FDP) pairs computed from the entire Lasso path must be asymptotically sandwiched between the two curves in the noiseless setting or, equivalently, in the regime of the infinite signal-to-noise ratio. The gap between the two curves is fundamental in the sense that it persists no matter how strong the effects are.

While the TPP-FDP trade-off essentially examines the "bulk" of the Lasso solution path, the second perspective we take extends to the "edge": when does the first noise variable enter the model along the Lasso path? More precisely, we decrease  $\lambda$  from  $\infty$  to 0 and find the first "time" a false selection occurs. To indicate the difficulty of consistent model selection, formally, we consider the rank of the first noise variable or, put concretely, one plus the number of the true variables before the Lasso selects the first noise variable. Intuitively, a large rank is desirable. As with the first perspective, assuming a Gaussian random design and regression coefficients with linear sparsity, we prove that the rank of the first false selection is bounded above by  $(1+o(1))n/(2\log p)$ . This upper bound, which approximately equals 72 in the setting of Figure 1, holds no matter how strong the effect sizes are. Interestingly, this upper bound is *exactly* achieved when effect

size heterogeneity is maximal and the noise level tends to zero. On the other hand, [6] has obtained a sharp prediction of the rank of the first false variable in the case of minimal effect size heterogeneity, which, together with our new result, shows that the first noise variable occurs much earlier with minimal effect size heterogeneity than with maximal effect size heterogeneity. Although not entirely related, the two perspectives consistently demonstrate that effect size heterogeneity is an important and useful concept for understanding the performance of the Lasso.

The fact that effect size heterogeneity matters, as shown above, is due to the bias introduced by the shrinkage nature of the Lasso. This bias in turn makes the residuals absorb many of the true effects that act as what we may want to call "shrinkage noise". Metaphorically, variables yet to be selected tend to "compete" with each other in entering the Lasso path and contribute to the shrinkage noise. The "competition" is particularly intensive among variables having about the same effect sizes, which is the case when effect size heterogeneity is low. As a price, the shrinkage noise gets inflated and some noise variables may be selected early due to their high correlations with the shrinkage noise. This is why false selections occur with a good chance and early. In contrast, when the heterogeneity is high, the largest true effect yet to be selected tends to have a significant correlation with the residuals, thereby having a better chance to be selected sooner. To appreciate this heuristic explanation, it is instructive to note that the least-squares estimator, if available, does not exhibit this price-of-competition phenomenon, as it is unbiased for the regression coefficients. An alternative but less direct way to appreciate effect size heterogeneity is to relate it to the restricted eigenvalue condition [2]. Roughly speaking, this condition is concerned with a vector of regression coefficients such that its  $\ell_1$  norm is largely contributed by a few components, and such approximate sparse regression coefficients can be well estimated by the Lasso and the Dantzig selector under certain designs [2], [7]. From the viewpoint of this condition, therefore, regression coefficients with high effect size heterogeneity can be thought of as having a smaller effective sparsity level, which is favored by the Lasso.

As a final remark, the price-of-competition phenomenon does not appear if the sparsity is sub-linear in the ambient dimension p, which is often assumed in the copious body of literature on high-dimensional regression. In this regime of sparsity, effect size heterogeneity has a vanishing impact on the performance of the Lasso if the signal-to-noise ratio is sufficiently strong or the beta-min condition is satisfied. Our paper also departs from this line of literature from a technical standpoint. Indeed, the proofs of the results in this paper make heavy use of approximate message passing (AMP) theory [8]-[10], with nontrivial extensions.

# A. Organization

The remainder of this paper is organized as follows. In Section II, we formalize the Lasso Crescent diagram by presenting our theoretical results that predict the TPP-FDP

 $^{1}$ If the effect sizes are sufficiently strong, variable selection using the t values of the least-squares estimator can lead to full power without any type  $L_{\rm eff}$ 

trade-off with respect to effect size heterogeneity. Section III extends the investigation of effect size heterogeneity to the problem of the first false variable along the Lasso path. Section IV is devoted to proving the results in Section II, whereas technical details of the proofs are deferred to the appendix. In Section V, we provide numerical studies to demonstrate the impact of effect size heterogeneity in general settings. We conclude the paper in Section VI with a few directions for future research.

#### II. THE LASSO CRESCENT

In this section, we derive the full possible range of the asymptotic trade-off between the TPP and FDP along the Lasso path, with a focus on its dependence on effect size heterogeneity. Specifically, our results can be pictorially presented by the Lasso Crescent as in Figure 2, hence the title of this section. The proofs are deferred to Section IV.

Throughout this section, and indeed the entire paper, we assume the following *working hypotheses* to specify the linear model (I.1). For ease of reading, we use boldface letters to denote vectors and matrices.

Gaussian Design Matrix: We consider a sequence of designs  $X \in \mathbb{R}^{n_l \times p_l}$  consisting of i.i.d.  $\mathcal{N}(0, 1/n_l)$  entries so that each column has an approximate unit  $\ell_2$  norm. As the index  $l \to \infty$ , we assume  $p_l, n_l \to \infty$  with  $n_l/p_l \to \delta$  for a constant  $\delta > 0$ . The index l is often omitted for the sake of simplicity.

Regression Coefficients: Let the regression coefficients  $\beta_1,\ldots,\beta_p$  be i.i.d. copies of a random variable  $\Pi$  that satisfies  $\mathbb{E}\,\Pi^2<\infty$ . Of particular interest to this paper is an  $\epsilon$ -sparse prior  $\Pi$  in the sense that  $\mathbb{P}(\Pi\neq 0)=\epsilon$  for a constant  $0<\epsilon<1$ . Thus, the realized  $\beta_1,\ldots,\beta_p$  are in the linear sparsity regime since the sparsity is approximately equal to  $\epsilon p$ .

*Noise*: The noise term z consists of i.i.d. elements drawn from  $\mathcal{N}(0, \sigma^2)$ , where the noise level  $\sigma \geq 0$  is fixed.

For completeness,  $X, \beta$ , and z are jointly independent. These assumptions are used in the literature on AMP theory and its applications (see, for example, [9]–[13]) and, more recently, have been commonly made in the high-dimensional regression literature [14]–[17]. On top of that, we adopt some adjustments made by [5] that slightly simplify the assumptions on  $\beta$  and z. Regarding the assumption on the noise, it is worth noting that we do not exclude the case  $\sigma=0$ , which corresponds to noiseless observations. For some of the results in this section, the price-of-competition phenomenon manifests itself most clearly in the noiseless setting.

# A. Most Heterogeneous Effect Sizes

Our first main theorem considers regression coefficients that are drawn from the following prior distribution:

Definition 2.1: For M > 0 and an integer m > 0, we call

$$\Pi^{\Delta} = \begin{cases}
0 & \text{w.p. } 1 - \epsilon \\
M & \text{w.p. } \frac{\epsilon}{m} \\
M^2 & \text{w.p. } \frac{\epsilon}{m} \\
\cdots & \cdots \\
M^m & \text{w.p. } \frac{\epsilon}{m}
\end{cases}$$
(II.1)

the  $(\epsilon, m, M)$ -heterogeneous prior.

For notational convenience, we suppress the dependence of  $\Pi^{\Delta}$  on  $\epsilon, m, M$ . This prior is  $\epsilon$ -sparse in the sense of the working hypotheses. As is clear, larger values of m, M would render the prior more heterogeneous. Indeed, this paper is primarily concerned with the case where both  $M, m \to \infty$ . This corresponds to the regime where the signal-to-noise ratio tends to infinity and, in addition, the true effect sizes are increasingly different. To be complete, the  $(\epsilon, m, M)$ -heterogeneous prior is only a specific example that attains increasing heterogeneity. See Remark 2.2 for more examples.

Following (I.3),  $\text{FDP}_{\lambda}(\Pi)$  and  $\text{TPP}_{\lambda}(\Pi)$  denote the (random) false discovery proportion and true positive proportion, respectively, of the Lasso estimate at  $\lambda$  when the regression coefficients in (I.1) are i.i.d. draws from a prior  $\Pi$ . For ease of reading, we say a pair (TPP, FDP) *outperforms* another pair (TPP', FDP') if TPP > TPP' and FDP < FDP'. As noted earlier, all theoretical results in this paper are obtained under the working hypotheses. For conciseness, the statements of our theorems shall not mention this fact anymore.

Theorem 1: Let C>c>0 be fixed. For any  $\epsilon$ -sparse prior  $\Pi$ , if both m and M are sufficiently large in the  $(\epsilon,m,M)$ -heterogeneous prior  $\Pi^{\Delta}$ , then the following conclusions are true:

(a) The event

$$\bigcup_{c<\lambda,\lambda'< C} \bigg\{ \big( \mathrm{TPP}_{\lambda'}\big(\Pi\big), \mathrm{FDP}_{\lambda'}\big(\Pi\big) \big) \text{ outperforms} \\ \\ \big( \mathrm{TPP}_{\lambda}\big(\Pi^{\Delta}\big), \mathrm{FDP}_{\lambda}\big(\Pi^{\Delta}\big) \big) \bigg\}$$

happens with probability tending to **zero** as  $n, p \to \infty$ . (b) For any constant  $\nu > 0$ , no matter how we choose  $\widehat{\lambda}'(\boldsymbol{y}, \boldsymbol{X}) \geq c$  adaptively as long as it always satisfies  $\operatorname{TPP}_{\widehat{\lambda}'}(\Pi) > \nu$ , with probability approaching **one** there exists  $\widehat{\lambda} > 0$  such that  $(\operatorname{TPP}_{\widehat{\lambda}}(\Pi^{\Delta}), \operatorname{FDP}_{\widehat{\lambda}}(\Pi^{\Delta}))$  outperforms  $(\operatorname{TPP}_{\widehat{\lambda}'}(\Pi), \operatorname{FDP}_{\widehat{\lambda}'}(\Pi))$ .

Remark 2.2: The priors for which the theorem holds can be extended in the following way. Consider a sequence of priors  $\Pi^{\Delta}$  satisfying  $\Pi^{\Delta}=0$  with probability  $1-\epsilon$  and  $\Pi^{\Delta}=M_i\neq 0$  with probability  $\epsilon\gamma_i$  for  $i=1,\ldots,m$  such that  $\gamma_1+\cdots+\gamma_m=1,\max_i\gamma_i\to 0$ , and  $\min_{1\leq i\leq m}|M_i/M_{i-1}|\to\infty$  (set  $M_0=1$ ). Alternatively, the nonzero component of the prior can be drawn from a continuous random variable with cumulative distribution function of form  $\frac{\log_M x}{m}$  for  $1\leq x\leq M^m$ . While the theorem statement is restricted to  $(\epsilon,m,M)$ -heterogeneous priors for brevity, its proof considers the general case.

This theorem demonstrates the optimality of heterogeneous and strong effects in terms of the trade-off between the TPP and FDP. Importantly, this optimality is *uniform* in the sense that it holds along the entire Lasso path, no matter how strong the true effects coming from  $\Pi$  are. To be sure, the event as a union in (a) is taken over any (TPP, FDP) pair from the prior  $\Pi$  and any pair from the prior  $\Pi^{\Delta}$ . Although each conclusion alone is not a consequence of the other, as we will see from the proof in Section IV, the two conclusions are built on top of the fact that the pairs  $(\text{TPP}_{\lambda}, \text{FDP}_{\lambda})$  with varying  $\lambda$  converge

uniformly to a deterministic smooth curve for both  $\Pi$  and  $\Pi^{\Delta}$ . This fact allows us to obtain the following byproduct:

Proposition 2.3: Under the assumptions of Theorem 1, for any sufficiently small constant  $\nu>0$ , the following statement holds with probability tending to one: for any  $\lambda,\lambda'>c$  such that  $\left|\mathrm{TPP}_{\lambda}(\Pi^{\Delta})-\mathrm{TPP}_{\lambda'}(\Pi)\right|<\nu$  and  $\mathrm{TPP}_{\lambda'}(\Pi)>0.001$ , we have

$$FDP_{\lambda}(\Pi^{\Delta}) < FDP_{\lambda'}(\Pi).$$

This result makes it self-evident why the prior  $\Pi^{\Delta}$  is a most favorable for the entire Lasso path, though literally, we should interpret this favorability in the limit  $m, M \to \infty$ . More precisely, this result implies that given a required power level, the smallest possible FDP is achieved when the effects are increasingly heterogeneous and strong. Of note, the number 0.001 above can be replaced by any small positive constant, and it does not impede the interpretability of the theorem since we are generally not interested in a model that includes only a tiny fraction of true variables.

An interesting yet unaddressed question is to find an expression of the asymptotic minimum of FDP given  $\mathrm{TPP}_{\lambda}(\Pi^{\Delta}) = u$  in the limit  $m, M \to \infty$ . Call this function  $q^{\Delta}(u; \delta, \epsilon)$ . From our results, one can easily see that  $q^{\Delta}$  is nothing but the lower *envelope* of instance-specific TPP–FDP trade-off curves over all  $\epsilon$ -sparse priors. To see this, first note that one can prove that as  $n, p \to \infty$ , the pairs  $(\mathrm{TPP}_{\lambda}(\Pi), \mathrm{FDP}_{\lambda}(\Pi))$  over all  $\lambda$  converge to a smooth curve, which is denoted by  $q^{\Pi}(u)$  (see Section IV). Recognizing that  $\Pi^{\Delta}$  is also  $\epsilon$ -sparse and assuming  $\lim_{m,M\to\infty}q^{\Pi^{\Delta}}$  exists, we must have

$$q^{\Delta}(u) := \lim_{m, M \to \infty} q^{\Pi^{\Delta}}(u) \ge \inf_{\Pi : \epsilon\text{-sparse}} q^{\Pi}(u). \tag{II.2}$$

On the other hand, it follows from Theorem 1 and in particular Proposition 2.3 that

$$q^{\Pi}(u) \ge \lim_{m, M \to \infty} \left( q^{\Pi^{\Delta}}(u) + o(1) \right) = q^{\Delta}(u)$$

for any  $\epsilon$ -sparse prior  $\Pi$ . This display, together with (II.2), gives

$$q^{\Delta}(u) = \inf_{\Pi: \epsilon \text{-sparse}} q^{\Pi}(u). \tag{II.3}$$

Interestingly, the right-hand side of (II.3) has been tackled in [5], leading to a precise expression. To describe this expression, let  $t^{\Delta}(u)$  be the largest positive root of the following equation in t,

$$\begin{split} \frac{2(1-\epsilon)\left[(1+t^2)\Phi(-t)-t\phi(t)\right]+\epsilon(1+t^2)-\delta}{\epsilon\left[(1+t^2)(1-2\Phi(-t))+2t\phi(t)\right]} \\ &=\frac{1-u}{1-2\Phi(-t)}, \quad \text{(II.4)} \end{split}$$

where  $\Phi(\cdot)$  and  $\phi(\cdot)$  denote the cumulative distribution function and probability density function of the standard normal distribution, respectively. Theorem 2.1 in [5] shows that

$$\inf_{\Pi:\epsilon\text{-sparse}} q^\Pi(u) = \frac{2(1-\epsilon)\Phi(-t^\Delta(u))}{2(1-\epsilon)\Phi(-t^\Delta(u)) + \epsilon u}. \tag{II.5}$$

Taken together, (II.3) and (II.5) yield

$$q^{\Delta}(u) = \frac{2(1-\epsilon)\Phi(-t^{\Delta}(u))}{2(1-\epsilon)\Phi(-t^{\Delta}(u)) + \epsilon u}.$$
 (II.6)

Remark 2.4: If u=0, treat  $\infty$  as a root of the equation and set 0/0=0 in (II.5). As such,  $q^{\Delta}$  satisfies  $q^{\Delta}(0)=0$ . If  $\delta<1$  and  $\epsilon$  is larger than a threshold determined by  $\delta$ , the function  $q^{\Delta}$  is defined only for u between 0 and a certain number strictly smaller than 1. This is where the celebrated Donoho–Tanner phase transition occurs [18] (also see Section B.2). Throughout this paper, however, we focus on the regime that is below the Donoho–Tanner phase transition—that is, the case where  $\delta \geq 1$ , or  $\delta < 1$  and  $\epsilon$  is small so that the range of u is the unit interval [0,1]. In contrast, above the phase transition, the mapping from the TPP to FDP might not be unique (see Figure 3.1 in [19] and [20]).

In summary, we have the following corollary, which addresses the aforementioned question.

Corollary 2.5: Under the assumptions of Theorem 1, we have

$$\lim_{m,M\to\infty}\lim_{n,p\to\infty}\sup_{\lambda>c}\left|\mathrm{FDP}_{\lambda}(\Pi^{\Delta})-q^{\Delta}\left(\mathrm{TPP}_{\lambda}(\Pi^{\Delta})\right)\right|=0,$$

where  $\lim_{n,p\to\infty}$  is taken in probability. Moreover, for any  $\epsilon$ -sparse prior  $\Pi$ , we have

$$\text{FDP}_{\lambda}(\Pi) \geq q^{\Delta} \left( \text{TPP}_{\lambda}(\Pi) \right) - 0.001$$

for all  $\lambda > c$  with probability tending to one.

Remark 2.6: As  $\lambda \to \infty$ , both  $\text{TPP}_{\lambda}(\Pi^{\Delta})$  and  $\text{FDP}_{\lambda}(\Pi^{\Delta})$  tend to 0. Hence, there is no need to impose an upper bound on  $\lambda$  when taking the supremum  $\sup_{\lambda>c}$ . The second conclusion of Corollary 2.5 follows from Proposition 2.3 in conjunction with the continuity of  $q^{\Delta}$ . As earlier, 0.001 can be replaced by any positive constant.

The second conclusion of Corollary 2.5 is part of Theorem 2.1 in [5] and demonstrates that true variables and irrelevant variables are always interspersed along the Lasso path. In particular, this is true when the regularization parameter  $\lambda$  tends to 0. In this case, indeed, the Lasso would select a significant fraction of false variables with vanishing but nonzero estimated coefficients. This fact necessitates a form of calibration of the Lasso estimates for variable selection [19], [20].

The significance of Theorem 1 and Corollary 2.5, however, extends beyond earlier results. Precisely, [5] derived the expression (II.5) by constructing a different signal prior II for each power level u. Indeed, the nonzero component of the prior constructed in [5] has two different magnitudes with weights depending on u, as opposed to an increasing number of different magnitudes as in the  $(\epsilon, m, M)$ -heterogeneous prior. The increasing level of heterogeneity allows us to give a one-shot construction of most heterogeneous priors at all power levels.

#### B. Least Heterogeneous Effect Sizes

We now turn to the opposite question: which effect sizes lead to the worst trade-off between the TPP and FDP along the Lasso path? Inspired by the interpretation of effect size heterogeneity, it is natural to consider the following signal prior as a candidate:

Definition 2.7: For M > 0, we call

$$\Pi^{\nabla} = \begin{cases}
0 & \text{w.p. } 1 - \epsilon \\
M & \text{w.p. } \epsilon
\end{cases}$$
(II.7)

the  $(\epsilon, M)$ -homogeneous prior.

This prior would render all true effect sizes equal, thereby being least heterogeneous or most homogeneous among all  $\epsilon$ -sparse priors. The following theorem confirms our intuition that this homogeneous prior is least favorable for the Lasso as the resulting effect sizes give the least optimal trade-off between false positives and power.

Theorem 2: Let C>c>0 be fixed. In the noiseless setting — that is, z=0 — for any  $\epsilon$ -sparse prior  $\Pi$  that is non-constant conditional on  $\Pi \neq 0$ , the following conclusions are true for the  $(\epsilon,M)$ -homogeneous prior  $\Pi^{\nabla}$ :

(a) The event

$$\bigcup_{c<\lambda,\lambda'< C} \bigg\{ (\mathsf{TPP}_{\lambda}(\Pi^{\nabla}), \mathsf{FDP}_{\lambda}(\Pi^{\nabla})) \text{ outperforms} \\ \\ (\mathsf{TPP}_{\lambda'}(\Pi), \mathsf{FDP}_{\lambda'}(\Pi)) \bigg\}$$

happens with probability tending to **zero** as  $n, p \to \infty$ .

(b) For any constant  $\nu>0$ , no matter how we choose  $\widehat{\lambda}'(\boldsymbol{y},\boldsymbol{X})\geq c$  adaptively as long as it always satisfies  $\mathrm{TPP}_{\widehat{\lambda}'}(\Pi)>\nu$ , with probability tending to **one** there exists  $\widehat{\lambda}>0$  such that  $\left(\mathrm{TPP}_{\widehat{\lambda}'}(\Pi),\mathrm{FDP}_{\widehat{\lambda}'}(\Pi)\right)$  outperforms  $\left(\mathrm{TPP}_{\widehat{\lambda}}(\Pi^{\nabla}),\mathrm{FDP}_{\widehat{\lambda}}(\Pi^{\nabla})\right)$ .

This theorem is similar, but in the opposite sense, to Theorem 1. One distinction between the two theorems is that Theorem 2 assumes the noiseless setting, as opposed to the noisy setting considered in Theorem 1. The noiseless setting is equivalent to an infinite value of the signal-to-noise ratio, which allows us to better isolate the impact of effect size heterogeneity from that of the noise term. That said, this theorem remains true in the presence of noise by setting a sufficiently large magnitude M for the true effect sizes.

Just as Proposition 2.3 does, the following result follows from the proof of Theorem 2 presented in Section IV.

Proposition 2.8: Under the assumptions of Theorem 2, for any sufficiently small constant  $\nu > 0$ , the following statement holds with probability tending to one: if  $\lambda, \lambda' > c$  satisfy  $\text{TPP}_{\lambda'}(\Pi) > 0.001$  and  $|\text{TPP}_{\lambda}(\Pi^{\nabla}) - \text{TPP}_{\lambda'}(\Pi)| < \nu$ , then

$$FDP_{\lambda}(\Pi^{\nabla}) > FDP_{\lambda'}(\Pi).$$

As is clear, this result demonstrates that the prior  $\Pi^{\nabla}$  is least favorable for the entire Lasso path in the noiseless case. Roughly speaking, this proposition shows that if the two Lasso problems agree on the value of the TPP along their paths, then the  $(\epsilon, M)$ -homogeneous prior  $\Pi^{\nabla}$  must yield a higher FDP. As with Proposition 2.3, 0.001 can be replaced by any positive constant. On a related note, the prior (II.7) is known to be least favorable for certain estimation problems both in the noiseless and noisy settings [21] (see also Lemma 4.4.1 and Corollary 4.4.3 in [22]). An interesting direction for future research is to study the relationship between estimation and variable selection with regard to the least favorability of the prior distribution.

The sharp distinction between Propositions 2.3 and 2.8 must be attributed to the priors  $\Pi^{\Delta}$  and  $\Pi^{\nabla}$ . The cause is, loosely speaking, due to the "competition" among variables with about the same effect sizes in entering the Lasso model. However, we find it easier to elucidate the underlying cause when studying the rank of the first false variable and thus defer the detailed discussion to Section III.

We now proceed to specify the curve on which  $(\text{TPP}_{\lambda}(\Pi^{\nabla}), \text{FDP}_{\lambda}(\Pi^{\nabla}))$  lies in the limit. For a fixed  $\alpha$ , let  $\varsigma = \varsigma(\alpha)$  denote the largest root of the equation

$$\delta = 2(1 - \epsilon)[(1 + \alpha^2)\Phi(-\alpha) - \alpha\phi(\alpha)] - \epsilon(2\alpha + \varsigma)\phi(\varsigma) + \epsilon\varsigma\phi(2\alpha + \varsigma) + \epsilon(1 + \alpha^2)[\Phi(\varsigma) + \Phi(-2\alpha - \varsigma)] + \epsilon(\varsigma + \alpha)^2[\Phi(-\varsigma) + \Phi(-2\alpha - \varsigma)],$$

and let  $t^{\nabla} = t^{\nabla}(u; \delta, \epsilon)$  be the largest root of the following equation in  $\alpha$ :

$$\Phi(\varsigma(\alpha)) + \Phi(-2\alpha - \varsigma(\alpha)) = u.$$

With all of these in place, define

$$q^{\nabla}(u; \delta, \epsilon) = \frac{2(1 - \epsilon)\Phi(-t^{\nabla}(u))}{2(1 - \epsilon)\Phi(-t^{\nabla}(u)) + \epsilon u}.$$
 (II.8)

The derivation of the expression is given in Lemma 1.16 in Section B.2. The following result shows that this function describes the limiting trade-off between the TPP and FDP in the case of minimal effect size heterogeneity:

Corollary 2.9: Under the assumptions of Theorem 2, we have

$$\lim_{n,p\to\infty} \sup_{\lambda>c} \left| \mathsf{FDP}_{\lambda}(\Pi^{\nabla}) - q^{\nabla} \left( \mathsf{TPP}_{\lambda}(\Pi^{\nabla}) \right) \right| = 0.$$

Moreover, for any  $\epsilon$ -sparse prior  $\Pi$ , we have

$$\mathrm{FDP}_{\lambda}(\Pi) \leq q^{\nabla} \left( \mathrm{TPP}_{\lambda}(\Pi) \right) + 0.001$$

for all  $\lambda > c$  with probability tending to one.

As with Theorem 2, Corollary 2.9 holds for any M > 0because of the noiseless setting.

Taken together, Corollaries 2.5 and 2.9 give the following

Theorem 3: Let c > 0 be any small constant. In the noiseless setting, for any  $\epsilon$ -sparse prior  $\Pi$ , we have

$$q^{\Delta} \left( \text{TPP}_{\lambda}(\Pi) \right) - 0.001 \le \text{FDP}_{\lambda}(\Pi) \le q^{\nabla} \left( \text{TPP}_{\lambda}(\Pi) \right) + 0.001$$

for all  $\lambda>c$  with probability tending to one. The two curves  $q^\Delta$  and  $q^\nabla$  enclose a crescent-shaped region, which we call the Lasso Crescent. This theorem shows any (TPP, FDP) pairs along the entire Lasso path would essentially lie in the corresponding Lasso Crescent that is specified by the shape n/p of the design and the sparsity ratio k/p of the effect sizes, and this region is tight. Figure 3 presents two instances of the Lasso Crescent, with simulations showing good agreement between the predicted and observed behaviors.2

 $^2$ R and Matlab code to calculate  $q^{\Delta}$  and  $q^{
abla}$  is available at https://github.com/huawang-wharton/effectsizeheterogeneity.

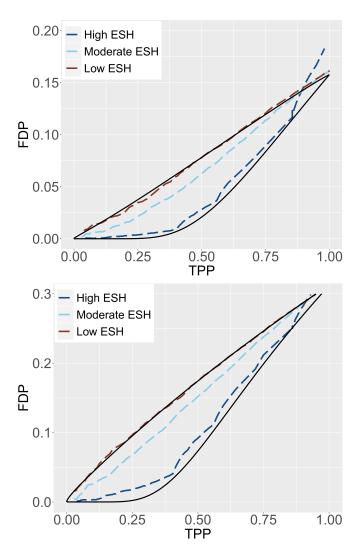


Fig. 3. Illustration of the interpretation of the Lasso Crescent via Theorem 3. The design of size  $n \times p$  has i.i.d.  $\mathcal{N}(0, 1/n)$  entries and the noise level is set to 0. Specifically, we use n = p = 1000, and sparsity k = 200 in the left panel, and n = 800, p = 1200, k = 200 in the right panel. The "high effect size heterogeneity (ESH)" setting: the 200 coefficients take 4 different values; The "moderate ESH" setting: the first 100 coefficients are set to 100 and the second 100 coefficients are set to 50; The "low ESH" setting: the 200 coefficients are set to 100. The dashed lines are averaged over 200 independent runs of the Lasso path. The two boundaries  $q^{\Delta}$  and  $q^{\nabla}$  are in solid black lines.

# III. THE FIRST FALSE SELECTION

In this section, we examine the impact of effect size heterogeneity on model selection by the Lasso from a different perspective: when is the first false variable selected? Intuitively, the later the first false variable occurs, the better the method performs. Using a mix of new and old results, this section will show that the first false variable occurs much earlier when effect size heterogeneity is minimal than when it is maximal.

Denote the rank of the first falsely selected variable by

$$T := \#\{j : \widehat{\beta}_i(\lambda^* - 0) \neq 0\} = \#\{j : \widehat{\beta}_i(\lambda^*) \neq 0\} + 1.$$

Above,  $\lambda^*$  is the first time along the Lasso path that a false variable is about to be selected:

$$\lambda^* = \sup\{\lambda : \text{there exists } 1 \le i \le p, \widehat{\beta}_i(\lambda) \ne 0, \beta_i = 0\},\$$

and  $\lambda^* - 0$  informally represents a value that is infinitesimally smaller than  $\lambda^*$ . In words, T is equal to one plus the number of true variables before the first false variable.

The problem of the rank of the first false selection has been considered by [6] in the case where all nonzero regression coefficients are equal. This corresponds to minimal effect size heterogeneity. While we continue employing the working hypotheses as earlier, in this section the regression coefficients  $\beta$  are assumed to be deterministic.

Proposition 3.1: [6, Theorem 2] Under the working hypotheses, let  $\beta_j=M$  for  $1\leq j\leq k$  and  $\beta_j=0$  for  $k+1\leq j\leq p$ , where  $k/p\to\epsilon$  and  $M\to\infty$  as  $n,p\to\infty$ . Then, the rank T of the first false variable selected by the Lasso satisfies

$$\log T = (1 + o_{\mathbb{P}}(1)) \sqrt{\frac{2\delta \log p}{\epsilon}},$$

where  $o_{\mathbb{P}}(1)$  tends to 0 in probability.

This result also applies to forward stepwise regression and least angle regression [23]. Note that this proposition considers the regime where the signal-to-noise ratio  $M/\sigma \to \infty$  as  $\sigma$  is fixed. If  $M/\sigma$  is bounded, one has  $\log T \le (1+o_{\mathbb{P}}(1))\sqrt{\frac{2\delta\log p}{\epsilon}}$  [6, Theorem 1]. Indeed, the original theorem predicts that

$$\log T \le (1 + o_{\mathbb{P}}(1)) \left( \sqrt{2n(\log p)/k} - n/(2k) + \log(n/(2p\log p)) \right),$$

which is reduced to the upper bound above since  $n/p \to \delta$  and  $k/p \to \epsilon$  under our working hypotheses.

Turning to most heterogeneous effect sizes, we have the result below.

Proposition 3.2: Under the working hypotheses, let  $\beta_j = M^{k+1-j}$  for  $1 \le j \le k$  and  $\beta_j = 0$  for  $k+1 \le j \le p$ , where  $k/p \to \epsilon$ . If M is sufficiently large, then there exists  $\lambda$  depending on n such that

$$\begin{split} \#\left\{j: \widehat{\beta}_j(\lambda) \neq 0, \beta_j \neq 0\right\} &= (1+o_{\mathbb{P}}(1))\frac{n}{2\log p}, \\ \text{and } \#\left\{j: \widehat{\beta}_j(\lambda) \neq 0, \beta_j = 0\right\} &= 0 \end{split}$$

as  $n, p \to \infty$ .

The proof of this proposition is given in the appendix. Regarding how large M should be, precisely, this proposition holds if M satisfies  $M \geq n^a$  as  $n \to \infty$  for any constant  $a > \frac{1}{2}$ . It is also worth mentioning that the proof is adapted from the proof of Theorem 1 in [24]. The effect sizes in Proposition 3.2 are essentially the same as an  $(\epsilon, m, M)$ -heterogeneous prior (II.1) with  $m \to \infty$ .

Proposition 3.2 asserts that all  $(1 + o_{\mathbb{P}}(1)) \frac{n}{2 \log p}$  selected variables are true at some point along the Lasso path. If the Lasso does not kick out any selected variables before that point,<sup>3</sup> this result implies that  $T \geq (1 + o_{\mathbb{P}}(1)) \frac{n}{2 \log p}$ .

Recognizing the fact that

$$e^{(1+o_{\mathbb{P}}(1))\sqrt{\frac{2\delta\log p}{\epsilon}}} \ll (1+o_{\mathbb{P}}(1))\frac{n}{2\log p}$$

most heterogeneous effect sizes are more *favorable* than least heterogeneous effect sizes for the Lasso not only in terms of the TPP–FDP trade-off as shown in the previous section, but also in terms of the rank of the first false variable.

Unlike Theorem 1 and Theorem 2, the two propositions here are silent on whether their bounds can be extended to general  $\epsilon p$ -sparse effect sizes. The following theorem gives a partial *affirmative* answer to this question, which broadly applies to all regression coefficients with sparsity no more than  $\epsilon p$ , as opposed to the exact sparsity level  $\epsilon p$ .

Theorem 4: Under the working hypotheses, for arbitrary regression coefficients  $\beta$  with sparsity satisfying  $k \leq \epsilon p$ , the rank T of the first false variable selected by the Lasso satisfies

$$T \le (1 + o_{\mathbb{P}}(1)) \frac{n}{2 \log p}$$

as  $n, p \to \infty$ .

Together with Proposition 3.2, this theorem indicates that maximal effect size heterogeneity is most favorable for the Lasso in terms of the rank of the first false variable. Importantly, the sharp bound  $(1+o_{\mathbb{P}}(1))\frac{n}{2\log p}$  is the maximum number of true variables before a false selection for essentially all sparsity levels, no matter how strong and how heterogeneous the effect sizes are. This novel result is a contribution of independent interest to high-dimensional statistics. The proof is given in Section III-A and does not involve any elements from AMP theory.

In regard to Proposition 3.1, however, it is tempting to ask whether minimal effect size heterogeneity is least favorable from the same standpoint; that is, whether or not

$$\log T \ge (1 + o_{\mathbb{P}}(1)) \sqrt{\frac{2\delta \log p}{\epsilon}}$$

for any  $\epsilon p$ -sparse regression coefficients in the noiseless case. We leave this question for future work.

In passing, we briefly explain how and why effect size heterogeneity has a significant impact on model selection by the Lasso, shedding light on the price-of-competition phenomenon. To ease the elaboration, we assume the noiseless setting (z=0) and denote by S the set of all true variables. Consider the Lasso solution  $\widehat{\beta}(\lambda)$  at some  $\lambda$  where no false selection occurs (that is, the support  $\widehat{S}$  of  $\widehat{\beta}$  is a subset of S). Our explanation relies crucially on the fact that a variable  $X_j$   $(j \notin \widehat{S})$  is likely to be the next selected variable if its inner product with the residual,  $X_j^{\top}(y-X\widehat{\beta})$ , is the largest in magnitude. Note that (denote by  $X_Q$  the matrix that is formed by the columns corresponding to Q for a subset Q of  $\{1,\ldots,p\}$ )

$$\begin{split} \boldsymbol{X}_{j}^{\top}(\boldsymbol{y} - \boldsymbol{X}\widehat{\boldsymbol{\beta}}) = & \boldsymbol{X}_{j}^{\top}(\boldsymbol{y} - \boldsymbol{X}_{\widehat{S}}\widehat{\boldsymbol{\beta}}_{\widehat{S}}) \\ = & \boldsymbol{X}_{j}^{\top}\boldsymbol{X}_{S\backslash\widehat{S}}\boldsymbol{\beta}_{S\backslash\widehat{S}} + \boldsymbol{X}_{j}^{\top}\boldsymbol{X}_{\widehat{S}}(\boldsymbol{\beta}_{\widehat{S}} - \widehat{\boldsymbol{\beta}}_{\widehat{S}}). \end{split}$$

Now, we argue that the largest  $\boldsymbol{X}_j^{\top}\boldsymbol{X}_{S\backslash\widehat{S}}\boldsymbol{\beta}_{S\backslash\widehat{S}}$  in absolute value in the case of high effect size heterogeneity is likely to be from a true variable  $\boldsymbol{X}_j$   $(j\in S\setminus\widehat{S})$  and, conversely,

<sup>&</sup>lt;sup>3</sup>It is well-known that along the Lasso path, a selected variable can be dropped out as  $\lambda$  decreases [23]. However, we did not observe this phenomenon before the first  $(1+o_{\mathbb{P}}(1))\frac{n}{2\log p}$  variables are selected in all of our simulations.

it is likely to be from an irrelevant variable  $\boldsymbol{X}_j$   $(j \notin S)$  if effect size heterogeneity is low. Informally, regarding  $\widehat{S}$  as deterministic, then  $\boldsymbol{X}_j^{\top} \boldsymbol{X}_{S \setminus \widehat{S}} \boldsymbol{\beta}_{S \setminus \widehat{S}}$  is approximately normally distributed with variance  $\|\boldsymbol{X}_{S \setminus \widehat{S}} \boldsymbol{\beta}_{S \setminus \widehat{S}}\|^2/n \approx \|\boldsymbol{\beta}_{S \setminus \widehat{S}}\|^2/n$  and mean

$$\begin{cases} 0 & \text{if } j \notin S \\ \beta_j & \text{if } j \in S \setminus \widehat{S}. \end{cases}$$

In the setting of Proposition 3.2 where true effect sizes are widely different from each other, the standard deviation  $\|\boldsymbol{\beta}_{S\setminus\widehat{S}}\|/\sqrt{n}$  is much smaller than  $\sup_{j\in S\setminus\widehat{S}}\beta_j$ . Consequently, the unselected variable with the largest effect size  $\sup_{j \in S \setminus \widehat{S}} \beta_j$ tends to stand out, with essentially no "competition" among all unselected variables, thereby being the next selected variable. In the setting of Proposition 3.1, however, the standard deviation  $\|\boldsymbol{\beta}_{S\setminus\widehat{S}}\|/\sqrt{n}$  is comparable to the largest unselected effect sizes, which are in fact of the same size. Another way to put this is that the overall effect is evenly distributed across true variables, and the resulted competition renders any variable dwarfed by the noise. Accordingly, some noise variable  $X_i$ is very likely to have a larger inner product  $X_i^{\top}(y - X\widehat{\beta})$ in magnitude than any unselected true variable does. As such, a false selection is likely to occur very early when effect size heterogeneity is low.

# A. Proof of Theorem 4

Let  $\nu > 0$  be any small constant and denote by  $\mathcal{A}_{\nu}$  the event that the rank of the first false variable

$$T \ge (1+\nu) \frac{n}{2\log p}.$$

The proof follows if one can show  $\mathbb{P}(\mathcal{A}_{\nu}) \to 0$  for all  $\nu > 0$  as  $n, p \to \infty$ . Recall that S denotes the support  $\operatorname{supp}(\beta)$ . If the sparsity  $|S| = k < (1 + \nu) \frac{n}{2 \log p} - 1 = (1 + \nu + o(1)) \frac{n}{2 \log p}$ , the event  $\mathcal{A}_{\nu}$  is an empty set because T is always no greater than  $|S| + 1 < (1 + \nu) \frac{n}{2 \log n}$ , leading to  $\mathbb{P}(\mathcal{A}_{\nu}) = 0$ .

than  $|S|+1<(1+\nu)\frac{n}{2\log p}$ , leading to  $\mathbb{P}(\mathcal{A}_{\nu})=0$ . Now, we turn to the more challenging case where  $k\geq (1+\nu)\frac{n}{2\log p}-1$  and the remainder of the proof aims to show  $\mathbb{P}(\mathcal{A}_{\nu})\to 0$ . Consider the solution  $\widetilde{\boldsymbol{\beta}}(\lambda)$  to the restricted Lasso problem

$$\widetilde{\boldsymbol{\beta}}(\lambda) := \underset{\boldsymbol{b} \in \mathbb{R}^k}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{X}_S \boldsymbol{b}\|^2 + \lambda \|\boldsymbol{b}\|_1.$$
 (III.1)

Let

$$\overline{\lambda} = \sup \left\{ \lambda : \|\widetilde{\beta}(\lambda)\|_0 \ge (1+\nu) \frac{n}{2\log p} - 1 \right\}$$

be the first time that the restricted Lasso selects  $(1+\nu)\frac{n}{2\log p}-1$  variables and denote by  $\widehat{S}$  the support of  $\widetilde{\beta}(\overline{\lambda}-0)$  (here  $\overline{\lambda}-0$  is infinitesimally smaller than  $\overline{\lambda}$ ). In particular, this set must satisfy

$$(1+\nu)\frac{n}{2\log p} - 1 \le |\widehat{S}| \le (1+\nu)\frac{n}{2\log p}.$$
 (III.2)

On the event  $A_{\nu}$ , the support of the full Lasso solution is a subset of S. Therefore,  $\beta(\overline{\lambda})$  defined in (III.1) is also the solution to the full Lasso problem at  $\overline{\lambda}$ :

$$\widetilde{\boldsymbol{\beta}}(\overline{\lambda}) = \operatorname*{argmin}_{\boldsymbol{b} \in \mathbb{R}^p} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{b}\|^2 + \overline{\lambda} \|\boldsymbol{b}\|_1.$$

Note that  $\widetilde{\boldsymbol{\beta}}(\overline{\lambda})$  may be k-dimensional as in (III.1) or p-dimensional by setting the remaining p-k entries to zero, depending on the context. Writing  $\widetilde{\boldsymbol{\beta}}$  as a shorthand for  $\widetilde{\boldsymbol{\beta}}(\overline{\lambda})$ , as a consequence, we have  $\left|\boldsymbol{X}_j^\top(\boldsymbol{y}-\boldsymbol{X}_S\widetilde{\boldsymbol{\beta}})\right| \leq \overline{\lambda}$  for all  $j \notin S$  on the event  $\mathcal{A}_{\nu}$ , thereby certifying

$$\mathbb{P}(\mathcal{A}_{\nu}) \leq \mathbb{P}\left(\left|\boldsymbol{X}_{j}^{\top}(\boldsymbol{y} - \boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})\right| \leq \overline{\lambda} \text{ for all } j \notin S\right).$$

To prove  $\mathbb{P}(A_{\nu}) \to 0$ , therefore, it suffices to show that

$$\max_{j \notin S} \left| \boldsymbol{X}_{j}^{\top} (\boldsymbol{y} - \boldsymbol{X}_{S} \widetilde{\boldsymbol{\beta}}) \right| > \overline{\lambda}$$
 (III.3)

with probability tending to one. Making use of the independence between  $X_j$  and  $y - X_S \widetilde{\boldsymbol{\beta}}$ ,  $X_j^{\top} (y - X_S \widetilde{\boldsymbol{\beta}})$ 's are p - k i.i.d. normal random variables with mean 0 and variance  $\|\boldsymbol{y} - \boldsymbol{X}_S \widetilde{\boldsymbol{\beta}}\|^2 / n$ , conditional on  $y - \boldsymbol{X}_S \widetilde{\boldsymbol{\beta}}$ . This gives

$$\max_{j \notin S} \left| \boldsymbol{X}_{j}^{\top} (\boldsymbol{y} - \boldsymbol{X}_{S} \widetilde{\boldsymbol{\beta}}) \right| = (1 + o_{\mathbb{P}}(1)) \frac{\|\boldsymbol{y} - \boldsymbol{X}_{S} \widetilde{\boldsymbol{\beta}}\|}{\sqrt{n}} \sqrt{2 \log(p - k)} \\
\geq (1 + o_{\mathbb{P}}(1)) \frac{\|\boldsymbol{X}_{\widehat{S}} (\boldsymbol{X}_{\widehat{S}}^{\top} \boldsymbol{X}_{\widehat{S}})^{-1} \boldsymbol{X}_{\widehat{S}}^{\top} (\boldsymbol{y} - \boldsymbol{X}_{S} \widetilde{\boldsymbol{\beta}})\|}{\sqrt{n}} \sqrt{2 \log(p - k)}, \tag{III.4}$$

where the inequality follows since  $X_{\widehat{S}}(X_{\widehat{S}}^{\top}X_{\widehat{S}})^{-1}X_{\widehat{S}}^{\top}$  is a projection. For the moment, take the inequality

$$\|\boldsymbol{X}_{\widehat{S}}(\boldsymbol{X}_{\widehat{S}}^{\top}\boldsymbol{X}_{\widehat{S}})^{-1}\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y}-\boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})\| \geq (1+c)\overline{\lambda}\sqrt{\frac{n}{2\log p}}$$
(III.5)

as given for some constant c>0 possibly depending on  $\nu$ , with probability tending to one. Combining (III.4) and (III.5) yields

$$\begin{aligned} \max_{j \notin S} \left| \boldsymbol{X}_{j}^{\top} (\boldsymbol{y} - \boldsymbol{X}_{S} \widetilde{\boldsymbol{\beta}}) \right| &\geq (1 + o_{\mathbb{P}}(1)) \sqrt{2 \log(p - k)} \frac{(1 + c) \overline{\lambda}}{\sqrt{2 \log p}} \\ &= (1 + c + o_{\mathbb{P}}(1)) \overline{\lambda} \sqrt{\frac{\log(p - k)}{\log p}} \\ &\geq (1 + c + o_{\mathbb{P}}(1)) \overline{\lambda} \sqrt{\frac{\log(p - \epsilon p)}{\log p}} \\ &= (1 + c + o_{\mathbb{P}}(1)) \overline{\lambda} \end{aligned}$$

with probability tending to one, which ensures (III.3).

We proceed to complete the proof of this theorem by verifying (III.5). The Karush-Kuhn-Tucker condition for the Lasso asserts that

$$\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y}-\boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})=\overline{\lambda}\operatorname{sgn}(\widetilde{\boldsymbol{\beta}}_{\widehat{S}})\in\overline{\lambda}\{1,-1\}^{|\widehat{S}|},$$

from which we get

$$\|\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y} - \boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})\| = \overline{\lambda}\sqrt{|\widehat{S}|}.$$

A classical result in random matrix theory (see Lemma 1.1 in the appendix) shows that the singular values of

 $X_{\widehat{S}}(X_{\widehat{S}}^{\top}X_{\widehat{S}})^{-1}$  are all bounded below by  $\frac{1}{\sqrt{1+\theta}}$  with probability  $1-1/p^2$ , where

$$\theta = C\sqrt{\frac{(1+\nu)\frac{n}{2\log p}\log(p/((1+\nu)\frac{n}{2\log p}))}{n}} \asymp \sqrt{\frac{\log\log p}{\log p}}$$
 (III.6)

for an absolute constant C. This allows us to get

$$\|\boldsymbol{X}_{\widehat{S}}(\boldsymbol{X}_{\widehat{S}}^{\top}\boldsymbol{X}_{\widehat{S}})^{-1}\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y} - \boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})\|$$

$$\geq \frac{1}{\sqrt{1+\theta}}\|\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y} - \boldsymbol{X}_{S}\widetilde{\boldsymbol{\beta}})\|$$

$$= \sqrt{\frac{|\widehat{S}|}{1+\theta}} \cdot \overline{\lambda}$$
(III.7)

with probability tending to one. Recognizing that  $\theta < \nu/2$  for sufficiently large p and plugging (III.6) and (III.2) into (III.7), we obtain

$$\|\boldsymbol{X}_{\widehat{S}}(\boldsymbol{X}_{\widehat{S}}^{\top}\boldsymbol{X}_{\widehat{S}})^{-1}\boldsymbol{X}_{\widehat{S}}^{\top}(\boldsymbol{y}-\boldsymbol{X}_{\widehat{S}}\widetilde{\boldsymbol{\beta}})\| \geq \sqrt{\frac{(1+\nu)\frac{n}{2\log p}-1}{1+\nu/2}} \cdot \overline{\lambda}$$
$$= (1+c)\overline{\lambda}\sqrt{\frac{n}{2\log p}}$$

with probability approaching one, where  $c=\sqrt{\frac{1+\nu}{1+\nu/2}}-1>0$ . This proves (III.5), thereby completing the proof of Theorem 4.

# IV. PROOFS FOR THE LASSO CRESCENT

To prove Theorems 1 and 2, we start by introducing AMP theory at a *minimal* level. In the case of the Lasso, loosely speaking, tools from AMP theory enable the characterization of the asymptotic joint distribution of the Lasso estimate  $\hat{\beta}(\lambda)$  and the true regression coefficients  $\beta$  under the working hypotheses [8]–[10]. The distribution is determined by several parameters that can be solved from two equations (see (IV.1) below). It is important to note, however, that this body of literature only allows for the analysis of the Lasso at a *fixed* value of  $\lambda$ . As such, these tools are not directly applicable to the full Lasso path that this paper deals with.

To overcome this difficulty, we leverage a recent extension on AMP theory that allows us to work on the Lasso problem uniformly over its penalty parameter [5]. Under the working hypotheses, let  $\tau>0$  and  $\alpha>\alpha_0$  be the unique solution to

$$\tau^{2} = \sigma^{2} + \frac{1}{\delta} \mathbb{E}(\eta_{\alpha\tau}(\Pi + \tau W) - \Pi)^{2}$$
$$\lambda = \left(1 - \frac{1}{\delta} \mathbb{P}(|\Pi + \tau W| > \alpha\tau)\right) \alpha\tau, \quad (IV.1)$$

where  $\eta_c(x) := \mathrm{sgn}(x) \cdot \max\{|x| - c, 0\}$  is the soft-thresholding function, W is a standard normal random variable that is independent of  $\Pi$ , and  $\alpha_0 = 0$  if  $\delta > 1$  and otherwise is the unique root of  $(1+t^2)\Phi(-t) - t\phi(t) = \frac{\delta}{2}$  in  $t \geq 0$ . Let  $\Pi^\star$  be

distributed the same as  $\Pi$  conditional on  $\Pi \neq 0$ , and define the two deterministic functions

$$\begin{aligned} \operatorname{tpp}_{\lambda}^{\infty}(\Pi) &= \mathbb{P}(|\Pi^{\star} + \tau W| > \alpha \tau) \\ \operatorname{fdp}_{\lambda}^{\infty}(\Pi) &= \frac{2(1 - \epsilon)\Phi(-\alpha)}{2(1 - \epsilon)\Phi(-\alpha) + \epsilon \mathbb{P}(|\Pi^{\star} + \tau W| > \alpha \tau)}. \end{aligned} \tag{IV.2}$$

Above,  $\Pi^*$  remains independent of W. For convenience, we use  $\stackrel{\mathbb{P}}{\longrightarrow}$  to denote convergence in probability. With the notation in place, now we state a lemma that our proofs rely on.

Lemma 4.1 ([25, Lemma A.2]): Fix 0 < c < C. Under the working hypotheses, we have

$$\begin{split} \sup_{c<\lambda< C} |\mathrm{TPP}_{\lambda}(\Pi) - \mathrm{tpp}_{\lambda}^{\infty}(\Pi)| &\stackrel{\mathbb{P}}{\longrightarrow} 0, \\ \text{and} \quad \sup_{\Omega \subset C} |\mathrm{FDP}_{\lambda}(\Pi) - \mathrm{fdp}_{\lambda}^{\infty}(\Pi)| &\stackrel{\mathbb{P}}{\longrightarrow} 0. \end{split} \tag{IV.3}$$

Lemma 4.1 offers all the elements the present paper needs from AMP theory. In addition to the use of this lemma, notably, our proofs of Theorems 1 and 2 involve several technical novelties that we shall highlight in Sections IV-A and IV-B. Relating to the literature, the convergence of  $\text{TPP}_{\lambda}(\Pi)$  and  $\text{FDP}_{\lambda}(\Pi)$  for a single  $\lambda$  has been established earlier in [10], [26].

We use  $q^{\Pi}(\cdot)$  to represent the  $\lambda$ -parameterized curve  $(\operatorname{tpp}_{\lambda}^{\infty}, \operatorname{fdp}_{\lambda}^{\infty})$  in the sense that

$$\mathrm{fdp}_{\lambda}^{\infty}(\Pi) = q^{\Pi}(\mathrm{tpp}_{\lambda}^{\infty}(\Pi)).$$

Formally, Lemma 1.11 in Section B.1 demonstrates that the instance-specific trade-off curve  $q^\Pi$  is continuously differentiable and strictly increasing. Relating to Section II, Corollary 2.5 implies that, taking the  $(\epsilon,m,M)$ -heterogeneous prior  $\Pi^\Delta, \ q^{\Pi^\Delta}$  converges to  $q^\Delta$  as  $m,M\to\infty$ . Likewise, from Corollary 2.9 it follows that  $q^{\Pi^\nabla}(\cdot)$  is identical to  $q^\nabla(\cdot)$  in the noiseless setting.

# A. The Upper Boundary

Our first aim is to prove Theorem 2 along with Proposition 2.8. The proof is built on top of the following important lemma, which considers a non-constant  $\Pi^*$ .

Lemma 4.2: Let  $\Pi$  be any  $\epsilon$ -sparse prior that is non-constant conditional on  $\Pi \neq 0$ . In the noiseless setting  $\sigma = 0$ , we have

$$q^{\Pi}(u) < q^{\nabla}(u)$$

for all 0 < u < 1.

Taking this lemma as given for the moment, we prove part (a) of Theorem 2.

*Proof of Theorem 2(a):* We start by pointing out the following fact: there exists a constant v > 0 such that for all  $c < \lambda, \lambda' < C$ , the two inequalities

$$\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\nabla}) > \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) - \upsilon \text{ and } \operatorname{fdp}_{\lambda}^{\infty}(\Pi^{\nabla}) < \operatorname{fdp}_{\lambda'}^{\infty}(\Pi) + \upsilon \tag{IV.4}$$

cannot hold simultaneously.

Assuming this fact for the moment, it is a stone's throw away to prove Theorem 2. Lemma 4.1 ensures that,

with probability tending to one as  $n,p \to \infty$ , the four terms  $|\mathrm{TPP}_{\lambda}(\Pi^{\nabla}) - \mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\nabla})|$ ,  $|\mathrm{FDP}_{\lambda}(\Pi^{\nabla}) - \mathrm{fdp}_{\lambda}^{\infty}(\Pi^{\nabla})|$ ,  $|\mathrm{TPP}_{\lambda'}(\Pi) - \mathrm{tpp}_{\lambda'}^{\infty}(\Pi)|$ , and  $|\mathrm{FDP}_{\lambda'}(\Pi) - \mathrm{fdp}_{\lambda'}^{\infty}(\Pi)|$  are all smaller than v/2 for all  $c < \lambda, \lambda' < C$ . In this event,

$$TPP_{\lambda}(\Pi^{\nabla}) > TPP_{\lambda'}(\Pi)$$

implies

$$\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\nabla}) > \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) - \upsilon$$

and, likewise,  $FDP_{\lambda}(\Pi^{\nabla}) < FDP_{\lambda'}(\Pi)$  implies  $fdp_{\lambda}^{\infty}(\Pi^{\nabla}) < fdp_{\lambda'}^{\infty}(\Pi) + \upsilon$ . Recognizing that the two inequalities in (IV.4) cannot both hold, therefore, in this event the following inequalities

$$\operatorname{TPP}_{\lambda}(\Pi^{\nabla}) > \operatorname{TPP}_{\lambda'}(\Pi) \quad \text{and} \quad \operatorname{FDP}_{\lambda}(\Pi^{\nabla}) < \operatorname{FDP}_{\lambda'}(\Pi)$$

cannot hold simultaneously for all  $c<\lambda,\lambda'< C$ . In words,  $\left(\operatorname{TPP}_{\lambda}(\Pi^{\nabla}),\operatorname{FDP}_{\lambda}(\Pi^{\nabla})\right)$  does not outperform  $\left(\operatorname{TPP}_{\lambda'}(\Pi),\operatorname{FDP}_{\lambda'}(\Pi)\right)$ , and this applies to all  $c<\lambda,\lambda'< C$  with probability tending to one.

We conclude the proof by verifying (IV.4). To this end, first find  $0 < u_1 < u_2 < 1$  such that the asymptotic powers  $\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\nabla}), \operatorname{tpp}_{\lambda'}^{\infty}(\Pi)$  are always between  $u_1$  and  $u_2$  for  $c < \lambda$ ,  $\lambda' < C$ . Next, set

$$v' := \inf_{u_1 \le u_< u_2} (q^{\nabla}(u) - q^{\Pi}(u)).$$
 (IV.5)

From Lemma 4.2, we must have v' > 0. Since  $q^{\nabla}$  is a continuous function on the closed interval [0,1], its uniform continuity gives

$$\left| q^{\nabla}(u) - q^{\nabla}(u') \right| < \frac{v'}{2}$$
 (IV.6)

as long as  $|u - u'| \le v''$  for some v'' > 0.

As the final step, we show that (IV.4) cannot hold simultaneously by taking  $v=\min\{v'/2,v''\}$ . To see this, suppose we already have  $\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\nabla})>\operatorname{tpp}_{\lambda'}^{\infty}(\Pi)-v$ , from which we get

$$\begin{split} \mathrm{fdp}_{\lambda}^{\infty}(\Pi^{\nabla}) &= q^{\nabla}(\mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\nabla})) \\ &\geq q^{\nabla}\left(\mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\nabla}) + \upsilon\right) - \frac{\upsilon'}{2} \\ &> q^{\nabla}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi)\right) - \frac{\upsilon'}{2}. \end{split}$$

Above, the first inequality follows from (IV.6). We proceed by leveraging (IV.5) and obtain

$$\begin{split} \operatorname{fdp}_{\lambda}^{\infty}(\Pi^{\nabla}) &> q^{\nabla} \left( \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) \right) - \frac{\upsilon'}{2} \\ &\geq q^{\Pi} \left( \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) \right) + \upsilon' - \frac{\upsilon'}{2} \\ &= q^{\Pi} \left( \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) \right) + \frac{\upsilon'}{2}. \end{split}$$

Finally, note that

$$q^\Pi\left(\operatorname{tpp}_{\lambda'}^\infty(\Pi)\right) + \frac{\upsilon'}{2} \geq q^\Pi\left(\operatorname{tpp}_{\lambda'}^\infty(\Pi)\right) + \upsilon = \operatorname{fdp}_{\lambda'}^\infty(\Pi) + \upsilon.$$

Taken together, these calculations reveal that the condition  $\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\nabla}) > \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) - v$  implies  $\operatorname{fdp}_{\lambda}^{\infty}(\Pi^{\nabla}) > \operatorname{fdp}_{\lambda'}^{\infty}(\Pi) + v$ .

As such, the two inequalities in (IV.4) cannot hold at the same time. This completes the proof.  $\Box$ 

The same reasoning in the proof above can be used to prove part (b) of Theorem 2 and Proposition 2.8. More precisely, the first step is to establish the desired result for the deterministic functions  $\operatorname{tpp}_{\lambda}^{\infty}$  and  $\operatorname{fdp}_{\lambda}^{\infty}$  using Lemma 4.2, followed by the second step that shows the uniform convergence using Lemma 4.1. In particular, part (b) of Theorem 2 relies on the strictly increasing property of  $q^{\nabla}$ . Moreover, note that a lower bound on  $\operatorname{TPP}_{\lambda'}(\Pi)$  can be translated into an upper bound on  $\lambda'$  [25, Lemma D.1].

Before turning to the proof of Lemma 4.2, we propose the following preparatory lemma.

Lemma 4.3 ([5, Lemma C.1]): For any fixed  $\alpha > 0$ , define a function y = f(x) in the parametric form

$$x(t) = \mathbb{P}(|t + W > \alpha|)$$
  
$$y(t) = \mathbb{E}(\eta_{\alpha}(t + W) - t)^{2}$$

for  $t \geq 0$ , where W is a standard normal random variable. Then f is strictly concave.

*Proof of Lemma 4.2:* We parameterize the curve  $(tpp_{\lambda}^{\infty}, fdp_{\lambda}^{\infty})$  using  $\alpha > \alpha_0$ . Explicitly, treating  $\alpha$  as the free parameter instead of  $\lambda$ , we can solve  $\tau$  from the AMP equation (IV.1). Define

$$\begin{split} \mathrm{fd}_\alpha^\infty(\Pi) &= 2(1-\epsilon)\Phi(-\alpha) \\ \mathrm{td}_\alpha^\infty(\Pi) &= \epsilon\,\mathbb{P}(|\Pi^\star + \tau W| > \alpha\tau). \end{split}$$

This allows us to express the asymptotic power and FDP as functions of  $\alpha$ :

$$\begin{split} & \operatorname{tpp}_{\alpha}^{\infty}(\Pi) = \frac{\operatorname{td}_{\alpha}^{\infty}(\Pi)}{\epsilon} \\ & \operatorname{fdp}_{\alpha}^{\infty}(\Pi) = \frac{\operatorname{fd}_{\alpha}^{\infty}(\Pi)}{\operatorname{fd}_{\alpha}^{\infty}(\Pi) + \operatorname{td}_{\alpha}^{\infty}(\Pi)}. \end{split}$$

To prove Lemma 4.2, for each  $\alpha > \alpha_0$ , it is sufficient to find a certain value of M such that

$$\mathrm{fd}_\alpha^\infty(\Pi)=\mathrm{fd}_\alpha^\infty(\Pi^\nabla) \quad \mathrm{and} \quad \mathrm{td}_\alpha^\infty(\Pi)>\mathrm{td}_\alpha^\infty(\Pi^\nabla), \qquad (\mathrm{IV}.7)$$

where  $\Pi^{\nabla}$  is the  $(\epsilon, M)$ -homogeneous prior (II.7). To see this fact, suppose on the contrary that

$$q^{\Pi}(u) \ge q^{\nabla}(u)$$
 (IV.8)

for some 0 < u < 1. Let  $\alpha$  satisfy  $\operatorname{tpp}_{\alpha}^{\infty}(\Pi) = u$ . From (IV.7) we obtain

$$u=\operatorname{tpp}_{\alpha}^{\infty}(\Pi)=\frac{\operatorname{td}_{\alpha}^{\infty}(\Pi)}{\epsilon}>\frac{\operatorname{td}_{\alpha}^{\infty}(\Pi^{\nabla})}{\epsilon}=\operatorname{tpp}_{\alpha}^{\infty}(\Pi^{\nabla}):=u^{\nabla}$$
 (IV.9)

and

$$\begin{split} \mathrm{fdp}_{\alpha}^{\infty}(\Pi) &= \frac{\mathrm{fd}_{\alpha}^{\infty}(\Pi)}{\mathrm{fd}_{\alpha}^{\infty}(\Pi) + \mathrm{td}_{\alpha}^{\infty}(\Pi)} < \frac{\mathrm{fd}_{\alpha}^{\infty}(\Pi^{\nabla})}{\mathrm{fd}_{\alpha}^{\infty}(\Pi^{\nabla}) + \mathrm{td}_{\alpha}^{\infty}(\Pi^{\nabla})} \\ &= \mathrm{fdp}_{\alpha}^{\infty}(\Pi^{\nabla}), \end{split}$$

which gives

$$q^\Pi(u) = \mathrm{fdp}_\alpha^\infty(\Pi) < \mathrm{fdp}_\alpha^\infty(\Pi^\nabla) = q^\nabla(u^\nabla).$$

This inequality combined with (IV.8) gives

$$q^{\nabla}(u) < q^{\nabla}(u^{\nabla}),$$

which, together with the fact that  $q^{\nabla}$  is an increasing function, leads to  $u < u^{\nabla}$ . This is a contradiction to (IV.9). Therefore, (IV.8) cannot hold for any 0 < u < 1.

The remainder of the proof aims to establish (IV.7) by constructing a certain prior  $\Pi^{\nabla}$ . Explicitly, it suffices to show

$$\operatorname{td}_{\alpha}^{\infty}(\Pi) > \operatorname{td}_{\alpha}^{\infty}(\Pi^{\nabla})$$
 (IV.10)

because the equality in (IV.7) holds regardless of the choice of  $\Pi^{\nabla}$ . To construct  $\Pi^{\nabla}$ , we first write  $\mathrm{td}_{\alpha}^{\infty}(\Pi)$  as

$$\operatorname{td}_{\alpha}^{\infty}(\Pi) = \epsilon \int \mathbb{P}(|t+W| > \alpha) d\pi(t),$$
 (IV.11)

where  $\mathrm{d}\pi(t)$  denotes the measure of  $\Pi^\star/\tau$ . Since  $\mathbb{P}(|t+W|>\alpha)$  is a strictly increasing function of t, there must exist t'>0 such that

$$\operatorname{td}_{\alpha}^{\infty}(\Pi) = \epsilon \, \mathbb{P}(|t' + W| > \alpha). \tag{IV.12}$$

Following (II.7), we let  $\Pi^{\nabla} = t'\tau$  with probability  $\epsilon$  and  $\Pi^{\nabla} = 0$  otherwise.

Now, let  $\tau^{\nabla}$  denote the solution to (IV.1) given  $\alpha$  and  $\Pi^{\nabla}$ . That is (note that  $\sigma = 0$ ),

$$(1 - \epsilon) \mathbb{E} \eta_{\alpha}(W)^{2} + \epsilon \mathbb{E} \left( \eta_{\alpha} \left( \frac{t'\tau}{\tau^{\nabla}} + W \right) - \frac{t'\tau}{\tau^{\nabla}} \right)^{2} = \delta.$$

Our next step is to show

$$\tau^{\nabla} > \tau$$
.

To this end, we invoke Lemma 4.3 and the strict concavity of f gives

$$\mathbb{E}\left(\eta_{\alpha}\left(\frac{t'\tau}{\tau} + W\right) - \frac{t'\tau}{\tau}\right)^{2} \equiv f\left(\mathbb{P}(|t' + W| > \alpha)\right)$$

$$= f\left(\int \mathbb{P}(|t + W| > \alpha) d\pi(t)\right)$$

$$> \int f\left(\mathbb{P}(|t + W| > \alpha)\right) d\pi(t)$$

$$= \int \mathbb{E}\left(\eta_{\alpha}\left(t + W\right) - t\right)^{2} d\pi(t)$$

$$= \mathbb{E}\left(\eta_{\alpha}\left(\frac{\Pi^{\star}}{\tau} + W\right) - \frac{\Pi^{\star}}{\tau}\right)^{2}, \quad \text{(IV.13)}$$

where the second equality follows from the definition of t' in (IV.11) and (IV.12), and the inequality is strict because  $\Pi^{\star}$  is not constant. Together with the AMP equation for  $\Pi$ 

$$(1 - \epsilon) \mathbb{E} \eta_{\alpha}(W)^{2} + \epsilon \mathbb{E} \left( \eta_{\alpha} \left( \frac{\Pi^{\star}}{\tau} + W \right) - \frac{\Pi^{\star}}{\tau} \right)^{2} = \delta,$$

(IV.13) implies

$$(1 - \epsilon) \mathbb{E} \eta_{\alpha}(W)^{2} + \epsilon \mathbb{E} \left( \eta_{\alpha} \left( \frac{t'\tau}{\tau} + W \right) - \frac{t'\tau}{\tau} \right)^{2} > \delta$$

or, equivalently,

$$(1 - \epsilon) \mathbb{E} \eta_{\alpha}(W)^{2} + \mathbb{E} \left[ \left( \eta_{\alpha} \left( \frac{\Pi^{\nabla}}{\tau} + W \right) - \frac{\Pi^{\nabla}}{\tau} \right)^{2}; \Pi^{\nabla} \neq 0 \right] > \delta. \quad \text{(IV.14)}$$

By definition, however,  $\tau^{\nabla}$  must satisfy

$$(1 - \epsilon) \mathbb{E} \eta_{\alpha}(W)^{2} + \mathbb{E} \left[ \left( \eta_{\alpha} \left( \frac{\Pi^{\nabla}}{\tau^{\nabla}} + W \right) - \frac{\Pi^{\nabla}}{\tau^{\nabla}} \right)^{2}; \Pi^{\nabla} \neq 0 \right] = \delta.$$
(IV.15)

A comparison between (IV.14) and (IV.15) immediately gives  $\tau^{\nabla} > \tau$ .

Having shown  $\tau^{\nabla} > \tau$ , we complete the proof by noting

$$\begin{split} \operatorname{td}_{\alpha}^{\infty}(\Pi) &= \epsilon \, \mathbb{P}(|t'+W| > \alpha) \\ &> \epsilon \, \mathbb{P}\left( \left| \frac{t'\tau}{\tau^{\nabla}} + W \right| > \alpha \right) \\ &= \epsilon \, \mathbb{P}\left( \left| \Pi^{\nabla} + \tau^{\nabla} W \right| > \alpha \tau^{\nabla} \middle| \Pi^{\nabla} \neq 0 \right) \\ &= \operatorname{td}_{\alpha}^{\infty}(\Pi^{\nabla}). \end{split}$$

This verifies (IV.10).

# B. The Lower Boundary

Now we turn to the proof of Theorem 1. As with the architecture of the proof of Theorem 2, our strategy is to first prove the theorem for the deterministic functions  $tpp_{\lambda}^{\infty}$  and  $fdp_{\lambda}^{\infty}$ , and then apply Lemma 4.1 to carry over the results to the random functions  $TPP_{\lambda}$  and  $FDP_{\lambda}$ . Having said this, it is important to note that the proof presents a novel element to the literature. Below, we shall highlight the novel part of the proof of Theorem 1 and leave the rest to the appendix.

As shown in [5], the trade-off curve  $q^{\Pi}$  of any  $\epsilon$ -sparse prior  $\Pi$  obeys

$$q^{\Pi}(u) > q^{\Delta}(u)$$

for 0 < u < 1 in both the noiseless and noisy settings, where the curve  $q^{\Delta}$  is defined in (II.6). If the (TPP, FDP) pairs from the  $(\epsilon, m, M)$ -heterogeneous prior  $\Pi^{\Delta}$  form the curve  $q^{\Delta}$  asymptotically as  $n, p \to \infty$ , the proof of Theorem 1 would follow immediate, just as Theorem 2. For any values of m and M, however, the  $\lambda$ -parameterized curve  $(\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\Delta}), \operatorname{fdp}_{\lambda}^{\infty}(\Pi^{\Delta}))$  does not agree with  $q^{\Delta}$ . This is in contrast to the proof of Theorem 2, where the (TPP, FDP) pairs from the  $(\epsilon, M)$ -homogeneous prior (II.7) converge to the curve  $q^{\nabla}$  for any value of  $M \neq 0$ , thanks to the assumed noiseless setting.

To tackle this challenge, our strategy is to uniformly approximate  $q^\Delta$  using a more general prior for effect sizes that takes the form

$$\Pi^{\Delta}(\boldsymbol{M}, \boldsymbol{\gamma}) = \begin{cases}
0 & \text{w.p. } 1 - \epsilon \\
M_1 & \text{w.p. } \epsilon \gamma_1 \\
M_2 & \text{w.p. } \epsilon \gamma_2 \\
M_3 & \text{w.p. } \epsilon \gamma_3 \\
\dots & \dots \\
M_m & \text{w.p. } \epsilon \gamma_m,
\end{cases} (IV.16)$$

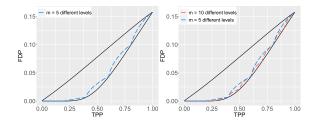


Fig. 4. Illustration of Lemma 4.4, showing the convergence to the lower curve  $q^{\Delta}$ . Left: m=5 different levels in the prior (IV.16) with  $\gamma_1=\cdots=\gamma_5=0.2$ , and the associated trade-off curve touches the lower boundary at 4 points; Right: the case m=10 and  $\gamma_1=\cdots=\gamma_{10}=0.1$  is added as a comparison to the left case.

where  $0 < M_1 < M_2 < \cdots < M_m$  and  $\gamma_1 + \cdots + \gamma_m = 1$  with  $\gamma_i > 0$ . Fixing  $\gamma = (\gamma_1, \dots, \gamma_m)$  while letting  $M_1 \to \infty$  and  $M_{i+1}/M_i \to \infty$  for all i, we have the following lemma:

 $M_{i+1}/M_i \to \infty$  for all i, we have the following lemma: Lemma 4.4: The curve  $q^{\Pi^{\Delta}(M,\gamma)}$  converges to a function that agrees with  $q^{\Delta}$  at m-1 points on (0,1).

For convenience, denote by  $q^{\Delta(\gamma)}$  the limiting curve of  $q^{\Pi^{\Delta}(M,\gamma)}$  as  $M_1 \to \infty$  and  $M_{i+1}/M_i \to \infty$ . Figure 4 provides an illustration of this limiting curve. To see why  $q^{\Delta(\gamma)}$  is close to  $q^{\Delta}$ , note that Lemma 4.4 ensures that there exist  $0 < u_1 < u_2 < \cdots < u_{m-1} < 1$  such that

$$q^{\Delta(\gamma)}(u_i) = q^{\Delta}(u_i)$$

for  $i=1,\ldots,m-1$ . In fact, the two functions also agree at  $u_0:=0$  and  $u_m:=1$ . Recognizing that both functions are increasing, for any  $u_i \leq u \leq u_{i+1}$  we get

$$0 \le q^{\Delta(\gamma)}(u) - q^{\Delta}(u) \le q^{\Delta(\gamma)}(u_{i+1}) - q^{\Delta}(u_i) = q^{\Delta}(u_{i+1}) - q^{\Delta}(u_i).$$

Making use of the uniform continuity of  $q^{\Delta}$ , the desired conclusion follows if we show that the gaps  $u_{i+1}-u_i$  are small for all  $i=0,\ldots,m-1$ . The proof of Lemma 4.4, indeed, reveals that this is true if  $\max \gamma_i$  is sufficiently small. See the proof of this lemma and the remaining details in Section B.

In passing, we remark that (IV.16) in the special case m=2 has been considered in [5]. Explicitly, the lower boundary  $q^{\Delta}$  is formed as the lower envelope of the instance-specific trade-off curves induced by the  $\epsilon$ -sparse priors. See the discussion following (II.3) in Section II.

#### V. ILLUSTRATIONS

In this section, we present simulation studies to illustrate the impact of effect size heterogeneity beyond the working hypotheses, with a focus on how the impact depends on the design matrix and the noise level.

# A. Design Matrix

We perform four simulation studies to examine the impact of effect size heterogeneity on the Lasso method under various synthetic design matrices. Overall, the simulation results show that effect size heterogeneity remains an influential factor in determining the performance of the Lasso far beyond Gaussian designs.

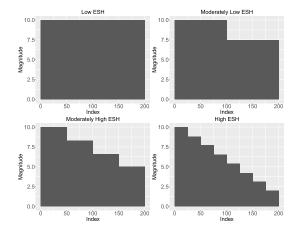


Fig. 5. Four sets of effect sizes ranked in increasing order of their effect size heterogeneity. The corresponding regression coefficients in  $\mathbb{R}^{1000}$  with sparsity 200 are used in the experiments of Figures 6, 7, and 10.

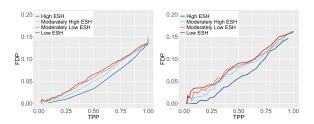


Fig. 6. The TPP-FDP trade-off along the Lasso path under a correlated Gaussian design and a Bernoulli design (Study 1). We set n=p=1000, k=200 and  $\sigma=0$  in both simulations. Left: Gaussian design matrix, each row having covariance  $\Sigma$  taking the form  $\Sigma_{ij}=0.5^{\lfloor i-j\rfloor}$ . Right: design matrix with i.i.d. Bernoulli entries taking the value  $1/\sqrt{1000}$  or  $-1/\sqrt{1000}$  with equal probability. The four sets of regression coefficients are described in Figure 5. The mean FDP is obtained by averaging over 200 replicates.

Study 1. We consider a design matrix of size  $1000 \times 1000$  that has each row independently drawn from  $\mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ , where  $\Sigma_{ij} = 0.5^{|i-j|}/1000$  and another design matrix of the same size that has independent Bernoulli entries, which take the value  $1/\sqrt{1000}$  with probability half and otherwise  $-1/\sqrt{1000}$ . The sparsity is fixed to k=200 while we consider four scenarios of the 200 true effects corresponding to low, moderately low, moderately high, and high effect size heterogeneity (see Figure 5). The results on the TPP-FDP trade-off are presented in Figure 6

Study 2. In this study, we use a dataset of size  $1000 \times 892$  that is simulated from the admixture of the African-American and European populations, based on the HapMap genotype data [27] (see more details in [28], [29]). The variables can only take 0,1, or 2 according to the genotype of a genetic marker. To improve the conditioning of the design matrix, we add i.i.d.  $\mathcal{N}(0,1/1000)$  perturbations to all the entries. Each column is further standardized to have mean 0 and unit norm. We use the effect sizes described in Figure 5 to generate a synthetic response y following the linear model (I.1), with noise z=0. The results are plotted in Figure 7.

Study 3. Working under Gaussian and Bernoulli designs, we now empirically examine the rank of the first false variable. This study considers a varying sparsity level k and sets the

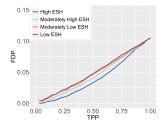


Fig. 7. The TPP–FDP trade-off for the genotype dataset (Study 2). The four curves correspond to the four sets of effect sizes described in Figure 5. The noise term is set to be **0**. The results are obtained by averaging over 200 replicates.

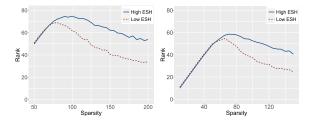


Fig. 8. The rank of the first spurious variable with varying sparsity (Study 3). Left: design matrix of size  $1000 \times 1000$  consists of i.i.d.  $\mathcal{N}(0, \frac{1}{1000})$  entries. Right: design matrix of size  $800 \times 1200$ , with i.i.d. Bernoulli entries that take the value  $1/\sqrt{500}$  with probability 1/2 and value  $-1/\sqrt{500}$  otherwise. Each curve is averaged over 200 independent replicates.

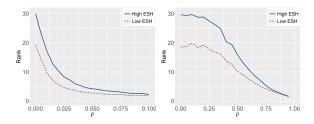


Fig. 9. The rank of the first spurious variable (Study 4). Left: Gaussian design with an equi-correlation covariance matrix, with the non-diagonal correlation  $\rho$  varying from 0 to 0.1. Right: Gaussian design with covariance  $\Sigma$  taking the form  $\Sigma_{ij}=\rho^{|i-j|}/1000$ . Each curve is averaged over 200 independent replicates.

effect sizes to  $\beta_j=100$  for  $j=1,\ldots,k$  (low effect size heterogeneity) or  $\beta_j=j$  for  $j=1,\ldots,k$  (high effect size heterogeneity). Each noise component  $z_i$  follows  $\mathcal{N}(0,1)$  independently. Figure 8 shows the results under an independent Gaussian random design and an independent Bernoulli design.

Study 4. This scenario uses  $500 \times 1000$  design matrices that have each row drawn independently from  $\mathcal{N}(0, \Sigma)$ . In the left panel of Figure 9, the  $1000 \times 1000$  covariance matrix  $\Sigma$  is set to  $\Sigma_{ij} = \rho/1000$  if  $i \neq j$  and  $\Sigma_{jj} = 1/1000$ . In the right panel, the covariance satisfies  $\Sigma_{ij} = \rho^{|i-j|}/1000$  for all i, j, with  $\rho$  varying from 0 to 0.95. The effect sizes are set to  $\beta_j = 100\sqrt{2\log p}$  for  $j \leqslant k$  (low effect size heterogeneity) or, in the low effect size heterogeneity case, the true effect sizes are set to a decreasing sequence from  $100\sqrt{2\log p}$  to 0. The noise z consists of independent standard normal variables. As is clear, both Figure 8 and Figure 9 show that the rank of

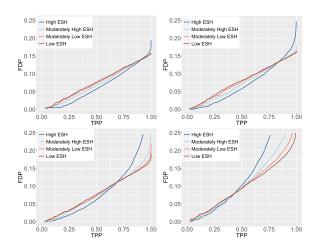


Fig. 10. The TPP-FDP trade-off plot with varying noise levels. The design matrix is specified by n=p=1000, with i.i.d. Gaussian entries. The regression coefficients are from Figure 5, and the noise vector has i.i.d.  $\mathcal{N}(0,\sigma^2)$  entries, where  $\sigma$  is set to 0.1,0.2,0.5 and 1 in the top-left, top-right, bottom-left, and bottom-right panels, respectively. The mean FDP is obtained by averaging over 100 replicates.

the first false variable is larger when effect size heterogeneity is high, aligning with our analysis in Section III.

#### B. Noise Level

While Theorem 2 concerning the regime of low effect size heterogeneity only applies to the noiseless case, we make an attempt to show the impact of effect size heterogeneity in the noisy setting via simulations. Under an independent Gaussian random design of size  $1000 \times 1000$ , we set the nonzero regression coefficients to the four sets of effect sizes as depicted in Figure 5. The noise term z consists of independent  $\mathcal{N}(0,\sigma^2)$  entries with  $\sigma=0.1,0.2,0.5,1.0$ . The results are displayed in Figure 10.

As with our previous simulation results, higher effect size heterogeneity tends to give rise to a better trade-off between the TPP and FDP from the beginning of the Lasso path. Interestingly, we observe a crossing point in each of the four panels of Figure 10 where higher heterogeneity undergoes a transition from giving a better trade-off down to a worse trade-off. In particular, the crossing point occurs earlier as the noise level  $\sigma$  goes up. While it requires further research to understand this transition in a concrete manner, our observation is that the unselected effect sizes in the late stage of the Lasso path tend to be relatively small compared to the noise level, especially the effect sizes depicted in the bottomright panel of Figure 10, which have relatively high effect size heterogeneity. Intuitively, this crossing point is where signalto-noise ratio becomes the dominant factor in place of effect size heterogeneity.

# VI. DISCUSSION

In this paper, we have proposed a concept termed effect size heterogeneity for measuring how diverse the nonzero regression coefficients are. Working under Gaussian random designs, we demonstrate that effect size heterogeneity has a significant impact on model selection consistency of the Lasso when the sparsity is linear in the ambient dimension. In short, we prove that the Lasso attains the optimal trade-off between true and false positives uniformly along its path when the effect sizes are strong and heterogeneous, and attains the worst trade-off when the effects are about the same size in the noiseless case. We also identify similar dependence of the rank of the first noise variable on effect size heterogeneity. While the two theoretical results are proved under certain assumptions, our simulations show that effect size heterogeneity has a significant impact on the Lasso estimate in a much wider range of settings.

Moving forward, this paper opens up several directions for future research. First, it is important to develop methods that incorporate the level of effect size heterogeneity for solving high-dimensional regression problems. In particular, one would be tempted to improve on the Lasso when the level of effect size heterogeneity is low. Interestingly, the SLOPE method has inadvertently addressed this question as its sorted  $\ell_1$  penalty generally increases as the heterogeneity gets higher [30]-[32]. Another related method developed from a Bayesian angle is the spike-and-slab Lasso procedure [33], which enables the adaptation to a mixture of large and small effects. Nevertheless, it is highly desirable to have methods that leverage effect size heterogeneity more directly. Moreover, a pressing question is to give a quantitative and formal definition of effect size heterogeneity. From a practical standpoint, regression coefficients are seldom exactly zero and thus it might be more appropriate to consider the Type S error, which occurs when a nonzero effect is selected but with the incorrect sign [34], [35]. This reality should prompt one to investigate how effect size heterogeneity interacts with the trade-off between the resulted directional FDP and power. Another question of practical importance is to examine carefully how the impact of effect size heterogeneity depends on the noise level. As an aside, given that Proposition 3.1 remains true for forward stepwise regression and least angle regression, we conjecture that Proposition 3.2 and Theorem 4 also hold for the two model selection procedures. More broadly, it is of interest to investigate whether effect size heterogeneity retains its impact on other  $\ell_1$  regularized methods such as the two-stage Lasso [19] and the Dantzig selector.

#### **APPENDIX**

A. Technical Proofs for Section III

1) Proof of Proposition 3.2:

*Proof of Proposition 3.2:* We use the "primal-dual witness" argument in the Lasso literature (for example, see Theorem 2 in [24]). As a reminder, here we consider the standard form of Lasso as in (I.2).

$$\widehat{\boldsymbol{\beta}} = \underset{\boldsymbol{b} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{b}\|^2 + \lambda \|\boldsymbol{b}\|_1$$

with the model specified by (I.1).

$$y = X\beta + z$$
.

We define a pair  $(\widehat{\beta}, \widehat{w}) \in \mathbb{R}^p \times \mathbb{R}^p$  to be *primal-dual optimal* if  $\widehat{\beta}$  is a minimizer of (I.2), and  $\widehat{w} \in \partial \|\widehat{\beta}\|_1$ , satisfying the

zero-subgradient condition

$$\boldsymbol{X}^{\top}(\boldsymbol{X}\widehat{\boldsymbol{\beta}} - \boldsymbol{y}) + \lambda \widehat{\boldsymbol{w}} = 0.$$

For the convenience of analysis, we denote  $\lambda_n = \frac{\lambda}{n}$ . Thus the condition above is equivalent to

$$\frac{1}{n} \mathbf{X}^{\top} (\mathbf{X} \widehat{\boldsymbol{\beta}} - \mathbf{y}) + \lambda_n \widehat{\mathbf{w}} = 0.$$
 (A.1)

By the sufficiency of KKT condition, we know that if there exists some  $\hat{w}$  such that the pair  $(\hat{\beta}, \hat{w}) \in \mathbb{R}^p \times \mathbb{R}^p$  satisfies (A.1), then  $\hat{\beta}$  is the solution to the Lasso. So  $\hat{w}$  can be seen as a "dual witness" showing  $\hat{\beta}$  is indeed a solution. We are therefore going to construct a "dual witness" vector  $\hat{w}$  to prove a certain  $\hat{\beta}$  is the solution to the Lasso.

To concretely give our construction of  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\boldsymbol{w}}$ , we fix an arbitrary small  $\xi>0$ , and then let  $s=[(1-\xi)\frac{n-1}{1-\xi+2\log p}].^4$  Denote  $S_0\equiv\{1,2,\ldots,s\},\,S_1\equiv\{s+1,s+2,\ldots,k\},\,$  and  $S=S_0\cup S_1$ , thus we have  $S^C=\{k+1,\ldots,p\}.$  Let  $M(n)=n^a$  for some  $a>\frac12$ , and let  $\lambda_n=n^b$  for some b that satisfies  $(k-s)a-1< b<(k-s+1)a-\frac32.$  We omit the dependence of M on n in the following proof. For clarity, for any subset of  $T\subset\{1,2,\ldots,p\},$  we always use the notation  $\boldsymbol{w}_T$  to denote the restricted vector  $(w_i)_{i\in T}$  of a vector  $\boldsymbol{w}$ , and the notation  $\boldsymbol{X}_T$  to denote the restricted column matrix  $(x_{i,j})_{j\in T}$  of a matrix  $\boldsymbol{X}$ . We consider the following procedure to construct the pair  $(\widehat{\boldsymbol{\beta}},\widehat{\boldsymbol{w}}),$ 

- 1) Let  $\widehat{\boldsymbol{\beta}}_{S_{0,\widehat{\cdot}}^{C}} = 0;$
- 2) Solve  $(\widehat{\beta}_{S_0}, \widehat{w}_{S_0}) \in \mathbb{R}^s \times \mathbb{R}^s$  from the following oracle sub-problem

$$\widehat{oldsymbol{eta}}_{S_0} \in \operatorname*{argmin}_{oldsymbol{b} \in \mathbb{R}^s} \left\{ rac{1}{2} \|oldsymbol{y} - oldsymbol{X}_{S_0} oldsymbol{b}\|_2^2 + \lambda \|oldsymbol{b}\|_1 
ight\}, \quad (A.2)$$

and choose  $\widehat{\boldsymbol{w}}_{S_0} \in \partial \|\widehat{\boldsymbol{\beta}}_{S_0}\|_1$  such that

$$\frac{1}{n} \boldsymbol{X}_{S_0}^{\top} (\boldsymbol{X}_{S_0} \widehat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{y}) + \lambda_n \widehat{\boldsymbol{w}}_{S_0} = 0;$$

3) Given  $\widehat{\boldsymbol{\beta}}_{S_0}$ ,  $\widehat{\boldsymbol{w}}_{S_0}$ , and  $\widehat{\boldsymbol{\beta}}_{S_0^C} = 0$ , compute  $\widehat{\boldsymbol{w}}_{S_0^C} \in \mathbb{R}^{p-s}$  by equation (A.1), and check whether the *strict dual feasibility* condition  $\|\widehat{\boldsymbol{w}}_{S_0^C}\|_{\infty} < 1$  holds.

The primal-dual witness construction guarantees that if a pair  $(\widehat{\boldsymbol{\beta}},\widehat{\boldsymbol{w}})$  satisfies all the three conditions above, then  $\widehat{\boldsymbol{\beta}}$  is the unique solution of the Lasso [24]. Once we prove our construction satisfies the conditions above, the second claim of Proposition 3.2 is an easy corollary as we explicitly require  $\widehat{\boldsymbol{\beta}}_j = 0$  for all  $j \in S_0^C$ , and this gives

$$\#\left\{j:\widehat{\beta}_j(\lambda)\neq 0, \beta_j=0\right\}=0.$$

And from this construction, it is also not hard to prove the first claim of the Proposition 3.2. With this protocol in mind, we proceed to prove that we can construct such a pair of  $(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{w}})$ . Now, we solve  $\widehat{\boldsymbol{\beta}}_{S_0}, \widehat{\boldsymbol{w}}_{S_0}$  from the subproblem in condition 2. Then, we set  $\widehat{\boldsymbol{\beta}}_{S_0^C} = 0$  as in condition 1, and solve  $\widehat{\boldsymbol{w}}_{S^C} \in \mathbb{R}^{p-s}$  from (A.1). To prove Proposition 3.2, we only need to prove that with this construction, the strict dual

<sup>4</sup>This is only for technical convenience. One can easily verify that this condition is equivalent to  $s=(1-o_{\mathbb{P}}(1))\frac{2n}{\log p}$ .

feasibility condition holds with high probability as  $n, p \to \infty$ . To prove this, we first simplify condition (A.1) by substituting  $\hat{\boldsymbol{\beta}}_{S_{c}^{C}} = 0$ , and write it in block matrix form as follows,

$$\frac{1}{n} \begin{bmatrix} \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} & \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_1} & \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S^c} \\ \boldsymbol{X}_{S_1}^{\top} \boldsymbol{X}_{S_0} & \boldsymbol{X}_{S_1}^{\top} \boldsymbol{X}_{S_1} & \boldsymbol{X}_{S_1}^{\top} \boldsymbol{X}_{S^c} \\ \boldsymbol{X}_{S^c}^{\top} \boldsymbol{X}_{S_0} & \boldsymbol{X}_{S^c}^{\top} \boldsymbol{X}_{S_1} & \boldsymbol{X}_{S^c}^{\top} \boldsymbol{X}_{S^c} \end{bmatrix} \cdot \begin{bmatrix} \hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0} \\ -\boldsymbol{\beta}_{S_1} \\ 0 \end{bmatrix} \\
-\frac{1}{n} \begin{bmatrix} \boldsymbol{X}_{S_0}^{\top} \boldsymbol{z} \\ \boldsymbol{X}_{S_1}^{\top} \boldsymbol{z} \\ \boldsymbol{X}_{S^c}^{\top} \boldsymbol{z} \end{bmatrix} + \lambda_n \begin{bmatrix} \hat{\boldsymbol{w}}_{S_0} \\ \hat{\boldsymbol{w}}_{S_1} \\ \hat{\boldsymbol{w}}_{S^c} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix},$$

or equivalently,

$$\frac{1}{n} \boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{X}_{S_0} (\widehat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}) - \frac{1}{n} \boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1} - \frac{1}{n} \boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{z} + \lambda_n \widehat{\boldsymbol{w}}_{S_0} = 0,$$

$$(A.3)$$

$$\frac{1}{n} \boldsymbol{X}_{S_1}^{\mathsf{T}} \boldsymbol{X}_{S_0} (\widehat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}) - \frac{1}{n} \boldsymbol{X}_{S_1}^{\mathsf{T}} \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1} - \frac{1}{n} \boldsymbol{X}_{S_1}^{\mathsf{T}} \boldsymbol{z} + \lambda_n \widehat{\boldsymbol{w}}_{S_1} = 0,$$

$$\frac{1}{n} \boldsymbol{X}_{SC}^{\top} \boldsymbol{X}_{S_0} (\widehat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}) - \frac{1}{n} \boldsymbol{X}_{SC}^{\top} \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1} - \frac{1}{n} \boldsymbol{X}_{SC}^{\top} \boldsymbol{z} + \lambda_n \widehat{\boldsymbol{w}}_{SC} = 0$$

(A.5)

By (A.3), we have

$$\widehat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0} = \left(\boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{X}_{S_0}\right)^{-1} \left[\boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1} + \boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{z}\right] - \lambda_n n \left(\boldsymbol{X}_{S_0}^{\mathsf{T}} \boldsymbol{X}_{S_0}\right)^{-1} \widehat{\boldsymbol{w}}_{S_0}.$$
(A.6)

By substituting (A.6) into (A.4) and (A.5), we can solve  $\widehat{w}_j$  for any  $j \in S_0^C$  as

$$\widehat{w}_{j} = -\frac{1}{\lambda_{n}n} X_{j}^{\top} \boldsymbol{X}_{S_{0}} (\widehat{\boldsymbol{\beta}}_{S_{0}} - \boldsymbol{\beta}_{S_{0}}) + \frac{1}{\lambda_{n}n} \left[ X_{j}^{\top} \boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + X_{j}^{\top} \boldsymbol{z} \right] \quad \text{A si}$$

$$= X_{j}^{\top} \boldsymbol{X}_{S_{0}} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}} \qquad \qquad \mathbb{P}$$

$$-\frac{1}{\lambda_{n}n} X_{j}^{\top} \boldsymbol{X}_{S_{0}} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \left[ \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{z} \right] \qquad \leq (p + \frac{1}{\lambda_{n}n} \left[ X_{j}^{\top} \boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + X_{j}^{\top} \boldsymbol{z} \right] \qquad \text{Con}$$

$$= \underbrace{X_{j}^{\top} \boldsymbol{X}_{S_{0}} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}}}_{v_{j}} + \underbrace{X_{j}^{\top} \mathcal{P}_{S_{0}^{\perp}} \left[ \frac{\boldsymbol{z}}{\lambda_{n}n} + \boldsymbol{X}_{S_{1}} \frac{\boldsymbol{\beta}_{S_{1}}}{\lambda_{n}n} \right]}_{u_{j}}, \qquad (A.7)$$

where  $\mathcal{P}_{S_0^{\perp}} = \mathbf{I} - \mathbf{X}_{S_0} \left( \mathbf{X}_{S_0}^{\top} \mathbf{X}_{S_0} \right)^{-1} \mathbf{X}_{S_0}^{\top}$ . As mentioned previously, our goal is to show the strict dual feasibility condition  $\max_{j \in S_0^C} |\widehat{w}_j| < 1$  holds with high probability. We will prove it by analyzing the two terms  $u_j$  and  $v_j$  separately. Specifically, we prove that  $v_j < 1 - \frac{\xi}{16}$  and  $u_j \to 0$  with high probability.

Denote  $M_n = \frac{1}{n} \widehat{\boldsymbol{w}}_{S_0}^{\top} \left( \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} \right)^{-1} \widehat{\boldsymbol{w}}_{S_0}$ . Conditioning on the event  $E = \{ \widehat{\boldsymbol{w}}_{S_0} = \operatorname{sgn}(\boldsymbol{\beta}_{S_0}) \}$  and its complement gives us

$$\mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge c\right) \le \mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge c \mid E\right) + \mathbb{P}(E^C).$$

It can be seen through Lemma 1.10 that the second term of the last display tends to zero; For the first term, we let  $T(\vartheta)$  denote the event  $\{|M_n - \mathbb{E} M_n| \ge \vartheta \mathbb{E} M_n\}$ . Similar as before, conditioning on the event  $T(\vartheta)$  and its complement gives for any  $c \in (0,1)$ ,

$$\mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge c \middle| E\right) \le \mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge c \middle| T(\vartheta)^C \cap E\right) + \mathbb{P}(T(\vartheta) \cap E).$$

By Lemma 1.2 and Lemma 1.4, the second term in the last display goes to 0 as  $n\to\infty$  faster than  $\frac{2}{\vartheta^2(n-s-3)}$ . And for the first term, we tackle it by considering  $\max_{j\in S_0^C} v_j$  and  $\min_{j\in S_0^C} v_j$  separately. Denote the event  $T=T(\vartheta)^C\cap E$  for convenience, we have

$$\mathbb{P}\left(\max_{j \in S_0^C} v_j \ge c \middle| T\right) \le \mathbb{P}\left(\max_{j \in S_0^C} \widetilde{v}_j \ge c\right),$$

where  $\widetilde{v}_j$  are i.i.d. from  $\mathcal{N}(0,(1+\vartheta)\mathbb{E}[M_n|E])=\mathcal{N}(0,(1+\vartheta)\frac{s}{n-s-1})$ . This inequality follows from Lemma 1.9, which states that the probability of the event  $\{\max_{i\in S_0^C}v_i\geq c\}$  increases as the mean and variance of  $v_i$  increase for Gaussian variables  $v_i$ . Given the event T, the maximum variance of  $v_j$  is simply  $(1+\vartheta)\mathbb{E}[M_n|E]$ , and thus we have the inequality above. Set  $c=\varpi+\mathbb{E}\max_{j\in S_0^C}\widetilde{v}_j$ . From Lemma 1.5, we have

$$\mathbb{P}\left(\max_{j \in S_0^C} v_j \ge c \middle| T\right) \le \mathbb{P}\left(\max_{j \in S_0^C} \widetilde{v}_j > \varpi + \mathbb{E} \max_{j \in S_0^C} \widetilde{v}_j\right)$$
$$\le (p - s) \exp\left(-\frac{\varpi^2}{2(1 + \vartheta)\mathbb{E}[M_n | E]}\right).$$

A similar argument for  $\min_{j \in S_0^C} \widetilde{v}_j$  gives us

$$\mathbb{P}\left(\min_{j \in S_0^C} v_j < -c \middle| T\right) = \mathbb{P}\left(-\min_{j \in S_0^C} v_j > \varpi + \mathbb{E} \max_{j \in S_0^C} \widetilde{v}_j \middle| T\right) \\
\leq (p-s) \exp\left(-\frac{\varpi^2}{2(1+\vartheta)\mathbb{E}[M_n|E]}\right).$$

Combining the two inequalities above yields

$$\mathbb{P}\left(\max_{j \in S_0^C} |v_j| > \varpi + \mathbb{E} \max_{j \in S_0^C} \widetilde{v}_j \middle| T\right) \\
\leq 2(p-s) \exp\left(-\frac{\varpi^2}{2(1+\vartheta)\mathbb{E}[M_n|E]}\right) \to 0.$$
(A.8)

By Lemma 1.4 and (A.8), we get

$$\mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge \varpi + \mathbb{E} \max_{j \in S_0^C} \widetilde{v}_j\right) \le \frac{2}{\vartheta^2(n-s-3)} + 2(p-s) \exp\left(-\frac{\varpi^2(n-s-1)}{2s(1+\vartheta)}\right). \tag{A.9}$$

With the relation of  $s=[(1-\xi)\frac{n-1}{1-\xi+2\log p}]$ , we can set  $\varpi=\frac{\xi}{8}, \vartheta=\frac{(1-\frac{\xi}{4})^2}{1-\xi}-1>0$ , and obtain

$$\varpi + \sqrt{(1+\vartheta)\frac{s}{n-s-1}2\log(p-s)}$$

$$\leq \varpi + \sqrt{(1+\vartheta)\frac{(1-\xi)\frac{n-1}{1-\xi+2\log p}+1}{n-(1-\xi)\frac{n-1}{1-\xi+2\log p}-1}2\log p}$$

$$= \varpi + \sqrt{(1+\vartheta)\frac{(1-\xi)(n-1)+(1-\xi+\log p)}{2(n-1)\log p}2\log p}$$

$$= \varpi + \sqrt{(1+\vartheta)\left((1-\xi)+\frac{(1-\xi+\log p)}{(n-1)}\right)}$$

$$\leq \frac{\xi}{8} + (1-\frac{\xi}{4}) + \frac{C}{n-1}$$

$$< 1 - \frac{\xi}{16} \quad \text{for some large } n. \tag{A.10}$$

Combining this with the well-known fact of the expectation of the maximum of i.i.d. Gaussian variables that  $\mathbb{E}\max_{j\in S_0^C}|\widetilde{v}_j|\leq \sqrt{(1+\vartheta)\frac{s}{n-s-1}}2\log(p-s)$  and (A.9), we know

$$\mathbb{P}(\max_{j \in S_0^C} |v_j| \ge 1 - \frac{\xi}{16}) \le \frac{2}{\vartheta^2(n - s - 3)} + 2(p - s) \exp\left(-\frac{\varpi^2(n - s - 1)}{2s(1 + \vartheta)}\right) + \mathbb{P}(E^C).$$
(A.11)

This bound is good enough for our purpose. We now proceed to obtain a similar bound of  $u_i$ .

By the Cauchy-Schwarz inequality, we have

$$|u_j| \le \left\| X_j^{\top} \mathcal{P}_{S_0^{\perp}} \right\| \cdot \left\| \frac{z}{\lambda_n n} + X_{S_1} \frac{\beta_{S_1}}{\lambda_n n} \right\|.$$
 (A.12)

Therefore, we can bound  $|u_j|$  if we can control the two norms in (A.12) separately. Because all eigenvalues of  $\mathcal{P}_{S_0^{\perp}}$  are less than 1, we have

$$\|X_j^{\top} \mathcal{P}_{S_0^{\perp}}\|^2 \le \|X_j\|^2 = \sum_{i=1}^n \mathbf{X}_{ij}^2 \stackrel{\mathcal{D}}{=} \frac{1}{n} \sum_{i=1}^n W_i^2,$$

where  $W_i \overset{i.i.d.}{\sim} \mathcal{N}(0,1)$ . The summation  $\sum_{i=1}^n W_i^2$  is thus a  $\chi^2$ -distribution with degree of freedom n. By Lemma 1.6, for any  $t_1 > 0$ , we have

$$\mathbb{P}\left(\left|\frac{\|X_j\|^2}{n} - 1\right| \ge t_1\right) \le 2\exp(-nt_1^2/8).$$

Combining this with  $\left\| {X_j}^{ op} \mathcal{P}_{S_0^{\perp}} \right\| \leq \|X_j\|$  gives us

$$\|X_j^{\top} \mathcal{P}_{S_0^{\perp}}\|_2 \le \sqrt{1+t_1}, \quad \text{w.p.} \ge 1-2\exp(-nt_1^2/8).$$

Finally, let  $\widetilde{u}_i$  denote

$$\widetilde{u}_j = \left\| \frac{z}{\lambda_n n} + X_{S_1} \frac{\beta_{S_1}}{\lambda_n n} \right\|.$$

It is easy to realize that

$$e_i \cdot (\boldsymbol{z} + \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1}) \sim \mathcal{N}(0, \sigma'^2), \text{ for all } 1 \leq i \leq s,$$

where  $e_i \in \mathbb{R}^n$  is the *i*-th standard unit vector and  $\sigma'^2 = \sigma^2 + \frac{M^{2(k-s+1)}-1}{n(M^2-1)}$ . By easy calculation, we obtain

$$\mathbb{E}(\|\widetilde{u}_j\|^2) = n \cdot \frac{\sigma'^2}{\lambda_n^2 n^2}.$$

Applying Lemma 1.6 with  $Z_i = e_i \cdot (\boldsymbol{z} + \boldsymbol{X}_{S_1} \boldsymbol{\beta}_{S_1})$  again, we know for any  $t_2 \geq 0$ 

$$\mathbb{P}\left(\left|\frac{\|\widetilde{u}_j\|^2}{\mathbb{E}(\|\widetilde{u}_j\|^2)} - 1\right| \ge t_2\right) \le 2\exp(-nt_2^2/8),$$

which is equivalent to

$$\|\widetilde{u}_j\| \le \sqrt{1+t_2} \cdot \frac{\sigma'}{\lambda_n \sqrt{n}}, \quad \text{w.p.} \ge 1 - 2\exp(-nt_2^2/8).$$
(A.14)

Using (A.12) by combining the two bounds (A.13) and (A.14), we get

$$\mathbb{P}(\max_{j \in S_0^C} |u_j| \ge \sqrt{1 + t_1} \sqrt{1 + t_2} \cdot \frac{\sigma'}{\lambda_n \sqrt{n}})$$

$$\le 2(p - s)(\exp(-nt_1^2/8) + 2\exp(-nt_1^2/8)). \tag{A.15}$$

Now, we can set  $M=n^a$  for some  $a>\frac{1}{2},$   $\lambda_n=n^b$  for some b that satisfies  $(k-s)a-1< b<(k-s+1)a-\frac{3}{2},$  and  $t_1=t_2=1$  to obtain

$$\sqrt{1+t_1}\sqrt{1+t_2} \cdot \frac{\sigma'}{\lambda_n \sqrt{n}} = 2 \frac{\sqrt{\sigma^2 + \frac{M^{2(k-s+1)} - 1}{n(M^2 - 1)}}}{\lambda_n \sqrt{n}}$$

$$\lesssim 2 n^{(k-s)a - 1 - b} \to 0. \quad (A.16)$$

Particularly, when n is large enough,  $\sqrt{1+t_1}\sqrt{1+t_2}\cdot\frac{\sigma'}{\lambda_n\sqrt{n}}$  is less than  $\frac{\xi}{32}$ , which in turn gives

$$\mathbb{P}(\max_{j \in S_0^C} |u_j| \ge \frac{\xi}{32}) \lesssim 2(p-s)(\exp(-nt_1^2/8) + 2\exp(-nt_1^2/8)). \tag{A.17}$$

And thus by a union bound and then (A.7, A.11, A.17), we have

$$\mathbb{P}\left(\max_{j \in S_0^C} |\widehat{w}_j| \ge 1 - \frac{\xi}{32}\right)$$

$$\le \mathbb{P}\left(\max_{j \in S_0^C} |v_j| \ge 1 - \frac{\xi}{16}\right) + \mathbb{P}\left(\max_{j \in S_0^C} |u_j| \ge \frac{\xi}{32}\right) + \mathbb{P}(E^C)$$

$$\le \frac{2}{\vartheta^2(n-s-3)} + 2(p-s)\left(\exp\left(-\frac{\varpi^2(n-s-1)}{2s(1+\vartheta)}\right) + \exp(-nt_2^2/8) + \exp(-nt_1^2/8)\right). \tag{A.18}$$

We observe that as  $n \to \infty$ ,

$$\frac{2}{\vartheta^{2}(n-s-3)} + 2(p-s) \left( \exp\left(-\frac{\varpi^{2}(n-s-1)}{2s(1+\vartheta)}\right) + \exp(-nt_{2}^{2}/8) + \exp(-nt_{1}^{2}/8) \right) \to 0,$$
(A.19)

which simply implies that as  $n \to \infty$ ,

$$\mathbb{P}(\max_{j \in S_0^C} |\widehat{w}_j| \ge 1) \to 0$$

Thus, we have proven under our construction, strict dual feasibility holds. And as we pointed out at the beginning of the proof, the second part of the proposition 3.2 holds, since we set  $\beta_j = 0$  for any  $j \in S_0^C$ . Therefore, we obtain

$$\#\left\{j:\widehat{\beta}_{j}(\lambda)\neq0,\beta_{j}=0\right\}=0.$$

Now, we proceed to prove the first part of the proposition, that

$$\#\left\{j: \widehat{\beta}_j(\lambda) \neq 0, \beta_j \neq 0\right\} = s = (1 - o_{\mathbb{P}}(1)) \frac{n}{2\log p}.$$

Observe that the second equality is due to our assumption on s. So we only need to prove the first equality, that is, for all  $j \in S_0$ ,  $\widehat{\beta}_j$  are non-zero, and thus the total number of non-zero  $\beta$ 's is exactly s. To show this, observe that if  $\beta_j - \beta_j < \beta_j$ , it is clear that  $\beta_i > 0$ . Therefore, it suffices to show the following inequality

$$\max_{j \in S_0} (\beta_j - \widehat{\beta}_j) < \min_{j \in S_0} \beta_j = M^{k-s+1} \equiv \rho$$

holds with probability tending to 1. We denote

$$Y_{i} = -e_{i}^{\mathsf{T}} \cdot \left( \boldsymbol{X}_{S_{0}}^{\mathsf{T}} \boldsymbol{X}_{S_{0}} \right)^{-1} \boldsymbol{X}_{S_{0}}^{\mathsf{T}} \left[ \boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + \boldsymbol{z} \right]$$

$$+ e_{i} \cdot \lambda_{n} n \left( \frac{\boldsymbol{X}_{S_{0}}^{\mathsf{T}} \boldsymbol{X}_{S_{0}}}{n} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}}, \quad (A.20)$$

where  $e_i \in \mathbb{R}^n$  is the *i*-th standard unit vector. By (A.6), we know

$$\max_{j \in S_0} (\beta_j - \widehat{\beta}_j) = \max_{1 \le i \le n} Y_i.$$

So it is equivalent to show that  $\max_{1 \le i \le n} Y_i \ge \rho$  holds with probability tending to zero. By Lemma 1.7, we know for  $E_i$  $\mathbb{E}(Y_i|\boldsymbol{X}_{S_0})$ , and  $V_i = \text{Var}(Y_i|\boldsymbol{X}_{S_0})$ , the event

$$A \equiv \bigcup_{i=1}^{s} \left\{ |E_i| \ge (1 + \sqrt{n}) |\mathbb{E}(E_i)|, \text{ or } |V_i| \ge 2 \mathbb{E}(V_i) \right\}$$

has probability

$$\mathbb{P}(A) \leq \frac{sK}{n-s} \to 0$$
, as  $n \to \infty$ .

By conditioning on the event A and its complement, we have

$$\mathbb{P}(\max_{i \in S_0} Y_i \ge \rho) \le \mathbb{P}(\max_{i \in S_0} Y_i \ge \rho | A^C) + \mathbb{P}(A)$$
$$\le \mathbb{P}(\max_{i \in S_0} \widetilde{Y}_i \ge \rho) + \frac{K}{\frac{n}{s} - 1},$$

where  $\widetilde{Y}_i \overset{i.i.d.}{\sim} \mathcal{N}((1+\sqrt{n})\,\mathbb{E}(E_i), 2\,\mathbb{E}(V_i))$  and the second inequality used the fact in Lemma 1.9 that the probability of the event  $\{\max_{i \in S_0} Y_i \ge \rho\}$  increases as the mean and variance increase as long as the mean is less than  $\rho$ , which can be directly verified by

$$(1+\sqrt{n})\mathbb{E}(E_i) = (1+\sqrt{n})\frac{\lambda_n n^2}{n-s-1} < \rho.$$
 (A.21)

Markov's inequality then gives us

$$\mathbb{P}(\max_{i \in S_0} \widetilde{Y}_i \ge \rho) \le \frac{1}{\rho} \mathbb{E}\left(\max_{i \in S_0} |\widetilde{Y}_i|\right) 
\le \frac{1}{\rho} \left(\mathbb{E}(\widetilde{Y}_i) + \mathbb{E}\left(\max_{i \in S_0} |\widetilde{Y}_i - \mathbb{E}(\widetilde{Y}_i)|\right)\right) 
\le \frac{1}{\rho} \left(\frac{(1+\sqrt{n})\lambda_n n^2}{n-s-1} + 3\sqrt{\frac{2\sigma'^2 \log s}{n-s-1}}\right),$$
(A.22)

where the last inequality uses the bound on Gaussian maxima in Lemma 1.8. By the relation in (A.21), we can easily verify that the probability in (A.22) converges to zero under our conditions of  $M=n^a$ ,  $\rho=M^{k-s+1}$  and  $\lambda_n=n^b$ , where a and b satisfy  $a>\frac{1}{2}$  and  $(k-s)a-1< b<(k-s+1)a-\frac{3}{2}$ ,

2) Miscellaneous Lemmas for Section III: We first state a well-known result in the random matrix theory (see, for example, Theorem 5.2 in [36]) that we use in the proof of Theorem 4. Then we list all the necessary lemmas for proving Proposition 3.2.

Lemma 1.1: Under the working assumptions, for any deterministic  $1 \le m \le p/2$ , the matrix spectrum norm  $\|\cdot\|_2$  satisfies

$$\max_{|S| \le m} \left\| \boldsymbol{X}_S^{\top} \boldsymbol{X}_S - \boldsymbol{I} \right\|_2 \le C \sqrt{\frac{m \log(p/m)}{n}}$$

with probability  $1-1/p^2$ , where C is a universal constant and T is any set of column indices.

Lemma 1.2: For  $v_j$  defined in (A.7), and any  $i, j \in S_0^C$ , we have the following facts,

- 1)  $\mathbb{E}(v_j|X_{S_0})=0;$
- 2)  $\operatorname{Var}(v_{j}|\boldsymbol{X}_{S_{0}}) = \frac{1}{n}\widehat{\boldsymbol{w}}_{S_{0}}^{\top}(\boldsymbol{X}_{S_{0}}^{\top}\boldsymbol{X}_{S_{0}})^{-1}\widehat{\boldsymbol{w}}_{S_{0}};$ 3)  $\operatorname{Cov}(v_{j},v_{i}|\boldsymbol{X}_{S_{0}}) = 0$ , if  $i \neq j$ .

Proof of Lemma 1.2:

Because  $j \in S_0^C$ ,  $X_j \perp X_{S_0}$  and  $v_j = X_j^{\top} X_{S_0} t (X_{S_0}^{\top} X_{S_0} t)^{-1} \hat{\boldsymbol{w}}_{S_0}$ , fact 1 follows from  $X_j$  being a centered Gaussian variable.

For fact 2 and fact 3, we observe

$$\begin{aligned} &\operatorname{Cov}(V_{j}, V_{i} | \boldsymbol{X}_{S_{0}}) \\ = & \mathbb{E} \left( \widehat{\boldsymbol{w}}_{S_{0}}^{\top} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{j} \boldsymbol{X}_{i}^{\top} \boldsymbol{X}_{S_{0}} \right. \\ & \left. \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}} \middle| \boldsymbol{X}_{S_{0}} \right) \\ = & \widehat{\boldsymbol{w}}_{S_{0}}^{\top} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \boldsymbol{X}_{S_{0}}^{\top} \mathbb{E} \left( \boldsymbol{X}_{j} \boldsymbol{X}_{i}^{\top} \middle| \boldsymbol{X}_{S_{0}} \right) \\ & \boldsymbol{X}_{S_{0}} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}} \\ = & \begin{cases} \frac{1}{n} \widehat{\boldsymbol{w}}_{S_{0}}^{\top} \left( \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}} & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases} \end{aligned}$$

Lemma 1.3: Consider  $X_{S_0} \in \mathbb{R}^{n \times s}$ , and suppose each of its column  $\mathbf{X}_j \overset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{\Sigma})$  where  $\mathbf{\Sigma} \in \mathbb{R}^{s \times s}$  is positive definite. Then  ${\boldsymbol{X}_{S_0}}^{\top}{\boldsymbol{X}_{S_0}}$  is a Wishart distribution of degree of freedom n, and  $({\boldsymbol{X}_{S_0}}^{\top}{\boldsymbol{X}_{S_0}})^{-1}$  is the inverse Wishart distribution, with expectation and variance

Authorized licensed use limited to: University of Pennsylvania. Downloaded on June 09,2023 at 01:55:46 UTC from IEEE Xplore. Restrictions apply.

1) 
$$\mathbb{E}(X_{S_0}^{\top}X_{S_0})^{-1} = \frac{\Sigma^{-1}}{n-s-1};$$

2) 
$$\operatorname{Var}[(X_{S_0}^{\top}X_{S_0})_{i,j}^{-1}] = \frac{(n-s+1)(\Sigma_{i,j}^{-1})^2 + (n-s-1)\Sigma_{i,i}^{-1}\Sigma_{j,j}^{-1}}{(n-s)(n-s-1)^2(n-s-3)}$$
. *Proof of Lemma 1.3:* See for example Lemma 7.7.1 of [37]

and the formula for the second moment of the inverse Wishart matrices in [38].

Lemma 1.4: Let  $M_n = \frac{1}{n} \widehat{\boldsymbol{w}}_{S_0}^{\top} \left( \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} \right)^{-1} \widehat{\boldsymbol{w}}_{S_0}$ . Conditioned on the event E, that is,  $\hat{\boldsymbol{w}}_{S_0} = \operatorname{sgn}(\boldsymbol{\beta}_{S_0})$ , we have the following facts:

1) 
$$\mathbb{E}(M_n|E) = \frac{s}{n-s-1}$$
;

2) 
$$Var(M_n|E) = \frac{2s^2}{(n-s-1)^2(n-s-3)}$$

1) 
$$\mathbb{E}(M_n|E) = \frac{s}{n-s-1};$$
  
2)  $\operatorname{Var}(M_n|E) = \frac{2s^2}{(n-s-1)^2(n-s-3)};$   
3)  $\forall \vartheta > 0, \mathbb{P}[|M_n - \mathbb{E}(M_n)| \ge \vartheta \mathbb{E}(M_n)|E] \le \frac{2}{\vartheta^2(n-s-3)}.$ 

*Proof of Lemma 1.4:* Observe that  $\boldsymbol{X}_{S_0}^{\top}\boldsymbol{X}_{S_0}$  follows the Wishart distribution with variance  $\frac{1}{n}\boldsymbol{I}_{S_0}$ , and thus by Lemma 1.3, the matrix  $(\boldsymbol{X}_{S_0}^{\top}\boldsymbol{X}_{S_0})^{-1}$  is the inverse Wishart distribution with mean

$$\mathbb{E}(X_{S_0}^{\mathsf{T}} X_{S_0})^{-1} = \frac{n}{n - s - 1} I_{S_0}.$$
 (A.23)

Notice that  $\widehat{w}_i = \pm 1$  for all  $i \in S_0$ , and when conditioned on E, it is equal to  $sgn(\beta_{S_0})$ , which is independent of  $X_{S_0}$ . Therefore, we have

$$\mathbb{E}(M_n|E) = \frac{1}{n} \frac{n}{n-s-1} \widehat{\boldsymbol{w}}_{S_0}^\top \boldsymbol{I}_{S_0} \widehat{\boldsymbol{w}}_{S_0} = \frac{s}{n-s-1}.$$

To calculate the second moment of the inverse Wishart matrices ([38]), we have that for n-s-3>0,

$$\mathbb{E}(M_n^2|E) = \frac{1}{n^2} \frac{1}{(n-s)(n-s-3)} (n \cdot \widehat{\boldsymbol{w}}_{S_0}^{\top} \boldsymbol{I}_{S_0} \widehat{\boldsymbol{w}}_{S_0})^2 \frac{n-s}{n-s-1}$$
$$= \frac{s^2}{(n-s-1)(n-s-3)}.$$

Therefore, combining the two equations above, we obtain

$$Var(M_n|E) = \frac{s^2}{(n-s-1)(n-s-3)} - \frac{s^2}{(n-s-1)^2}$$
$$= \frac{2s^2}{(n-s-1)^2(n-s-3)}.$$

For the third statement, Markov's inequality gives us

$$\mathbb{P}(|M_n - \mathbb{E}(M_n)| \ge \vartheta \mathbb{E}(M_n|E) \le \frac{\operatorname{Var}(M_n|E)}{\vartheta^2 (\mathbb{E}(M_n|E))^2}$$
$$= \frac{\frac{2s^2}{(n-s-1)^2(n-s-3)}}{\vartheta^2 \frac{s^2}{\vartheta^2}} = \frac{2}{\vartheta^2(n-s-3)}.$$

Lemma 1.5: Consider i.i.d. Gaussian random variables  $z_i \sim$  $\mathcal{N}(0,\sigma^2)$ , where  $j=1,\ldots,l$  for some  $l\geq 2$ . We have for any  $\varpi > 0$ ,

$$\mathbb{P}\left(\max_{1\leq j\leq l} z_j > \varpi + \mathbb{E}\left(\max_{1\leq j\leq l} z_j\right)\right) \leq \mathrm{e}^{-\frac{\varpi^2}{2\sigma^2}}.$$
From of Lemma 1.5: By the Gaussian tail bound

$$\mathbb{P}(z_j > \varpi) \le \frac{\sigma}{\sqrt{2\pi}\varpi} e^{-\frac{\varpi^2}{2\sigma^2}},$$

and the well-known fact for the expectation of maximum of i.i.d. Gaussian variables

$$\mathbb{E}\left(\max_{1\leq j\leq l} z_j\right) \leq \sigma\sqrt{2\log l},$$

we have the following union bound

$$\mathbb{P}\left(\max_{1 \leq j \leq l} z_{j} > \varpi + \mathbb{E}[\max_{1 \leq j \leq l} z_{j}]\right)$$

$$\leq l \frac{1}{\sqrt{2\pi}(\frac{\varpi}{\sigma} + \sqrt{2\log l})} \cdot \exp\left(\frac{(\varpi + \sqrt{2\log(l)}\sigma)^{2}}{2\sigma^{2}}\right)$$

$$= \frac{1}{\sqrt{2\pi}(\frac{\varpi}{\sigma} + \sqrt{2\log l})} \exp\left(-\frac{\varpi^{2}}{2\sigma^{2}}\right) \cdot \exp\left(\sqrt{2\log l}\frac{\varpi}{\sigma}\right)$$

$$\leq \exp\left(-\frac{\varpi^{2}}{2\sigma^{2}}\right)$$

holds as long as  $l \geq 2$ .

Lemma 1.6: Consider  $Z_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0,\theta^2)$ , and denote Z= $\sum_{i=1}^{n} Z_i^2$ . For t > 0, we have the inequality

$$\mathbb{P}\left(\left|\frac{Z}{\mathbb{E}(Z)} - 1\right| \ge t\right) \le 2\exp(-nt^2/8).$$

*Proof of Lemma 1.6:* Let  $\widetilde{Z}_i$  be defined as

$$\widetilde{Z}_i = \frac{Z_i}{\theta} \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1).$$

We have

$$\mathbb{E}(Z) = \sum_{i=1}^{n} \theta^2 = n\theta^2,$$

and therefore

$$\frac{Z}{\mathbb{E}(Z)} = \frac{\sum_{i=1}^{n} \widetilde{Z}_{i}}{n} \stackrel{\mathcal{D}}{=} \frac{\chi^{2}}{n}.$$

By easy calculation, we obtain

$$\mathbb{E}\left(e^{\lambda(\widetilde{Z}_{i}^{2}-1)}\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{\lambda(z^{2}-1)} e^{-z^{2}/2} dz$$
$$= \frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \le e^{2\lambda^{2}}.$$

This means Z is sub-exponential with parameter (2,4)(definition of a sub-exponential variable is standard, so we refer the reader to, for example, the Definition 2.2 and Example 2.11 in [4]), and thus Z is a sub-exponential variable with parameter  $(2\sqrt{n}, 4)$ . By the Bernstein inequality, we obtain

$$\mathbb{P}\left(\left|\frac{Z}{\mathbb{E}(Z)} - 1\right| \ge t\right) \le 2\exp(-nt^2/8).$$

Lemma 1.7: For  $Y_i$  defined in (A.20), we have the following facts:

- 1) Denote  $E_i = \mathbb{E}(Y_i \big| \boldsymbol{X}_{S_0})$ , then we have  $|\mathbb{E}(Y_i)| = |\mathbb{E}(E_i)| = \frac{\lambda_n n^2}{n-s-1}$ ; 2) Denote  $V_i = \operatorname{Var}(Y_i \big| \boldsymbol{X}_{S_0})$ , then we have  $\mathbb{E}(V_i) = 1$ ;
- 3) For n sufficiently large, the inequality  $\mathbb{P}(|E_i| \geq (1 +$  $\sqrt{n}$   $|\mathbb{E}(E_i)|$ , or  $|V_i| \geq 2\mathbb{E}(V_i)$   $\leq \frac{K}{n-s}$  holds for some constant K independent of n and s.

*Proof of Lemma 1.7:* The idea of the following proof is adapted from Lemma 6 of [24].

Authorized licensed use limited to: University of Pennsylvania. Downloaded on June 09,2023 at 01:55:46 UTC from IEEE Xplore. Restrictions apply.

For part (a), since  $X_{S_1} \perp X_{S_0}$  and  $z \perp X_{S_0}$ , we get

$$E_i = \mathbb{E}(Y_i | \boldsymbol{X}_{S_0}) = -\lambda_n n e_i^{\top} \left( \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} \right)^{-1} \widehat{\boldsymbol{w}}_{S_0}.$$

Thus, we have

$$|\mathbb{E}(Y_i)| = \left| \mathbb{E}\left( -\lambda_n n e_i \left( \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} \right)^{-1} \widehat{\boldsymbol{w}}_{S_0} \right) \right|$$

$$= \left| -\lambda_n n e_i^{\top} \frac{n}{n-s-1} \boldsymbol{I}_{S_0}^{-1} \widehat{\boldsymbol{w}}_{S_0} \right|$$

$$= \frac{\lambda_n n^2}{n-s-1},$$

where the second equality is by (A.23) for the mean of the inverse Wishart distribution.

Next, we turn to prove part (b). We observe that each entry of vector  $(\boldsymbol{X}_{S_1}\boldsymbol{\beta}_{S_1}+\boldsymbol{z})$  is i.i.d. distributed as  $\mathcal{N}(0,\sigma'^2)$ , and is independent of  $\boldsymbol{X}_{S_0}$ , where we denote  $\sigma'^2=\sigma^2+\frac{M^{2(k-s+1)}-1}{n(M^2-1)}$ . So we have

$$\begin{aligned} & \operatorname{Var}(Y_{i} \big| \boldsymbol{X}_{S_{0}}) \\ &= \mathbb{E}\left[\left(\boldsymbol{e}_{i}^{\top} \cdot \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{X}_{S_{0}}^{\top} \left(\boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + \boldsymbol{z}\right)\right)^{2} \big| \boldsymbol{X}_{S_{0}}\right] \\ &= \boldsymbol{e}_{i}^{\top} \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{X}_{S_{0}}^{\top} \mathbb{E}\left[\left(\boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + \boldsymbol{z}\right)\right] \\ & \cdot \left(\boldsymbol{X}_{S_{1}} \boldsymbol{\beta}_{S_{1}} + \boldsymbol{z}\right)^{\top} \big| \boldsymbol{X}_{S_{0}}\right] \boldsymbol{X}_{S_{0}} \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{e}_{i} \\ &= \boldsymbol{\sigma}'^{2} \boldsymbol{e}_{i}^{T} \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}} \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{e}_{i} \\ &= \boldsymbol{\sigma}'^{2} \boldsymbol{e}_{i}^{\top} \left(\boldsymbol{X}_{S_{0}}^{\top} \boldsymbol{X}_{S_{0}}\right)^{-1} \boldsymbol{e}_{i}. \end{aligned}$$

Thus by (A.23) again, we obtain

$$\mathbb{E}(V_i) = \mathbb{E}(\sigma'^2 e_i^{\top} \left( \boldsymbol{X}_{S_0}^{\top} \boldsymbol{X}_{S_0} \right)^{-1} e_i)$$

$$= \sigma'^2 e_i^{\top} \frac{n}{n - s - 1} \boldsymbol{I}_{S_0}^{-1} e_i$$

$$= \frac{n\sigma'^2}{n - s - 1}.$$

To prove part (c), we use the formula for the second moment of the inverse Wishart distribution in the part (2) of Lemma 1.3. With  $E_i = \mathbb{E}(Y_i | \boldsymbol{X}_{S_0})$ , we get

$$\mathbb{E}(E_{i}^{2}) = \mathbb{E}(\mathbb{E}(Y_{i}|\boldsymbol{X}_{S_{0}}))^{2}$$

$$= \frac{\lambda_{n}^{2}n^{2}}{(n-s)(n-s-3)} \left[ \left( e_{i}^{\top} \left( \frac{1}{n} \boldsymbol{I}_{S_{0}} \right)^{-1} \widehat{\boldsymbol{w}}_{S_{0}} \right)^{2} + \frac{n^{2}}{n-s-1} \left( \widehat{\boldsymbol{w}}_{S_{0}}^{\top} \widehat{\boldsymbol{w}}_{S_{0}} \right) \left( e_{i}^{\top} e_{i} \right) \right]$$

$$= \frac{\lambda_{n}^{2}n^{2}}{(n-s)(n-s-3)} \left[ n^{2} + \frac{1}{n-s-1} \cdot ns \cdot n \right]$$

$$= \frac{\lambda_{n}^{2}n^{4}(ns+n-s^{2}-2s+1)}{(n-s)(n-s-3)(n-s-1)}.$$

Thus, we have

$$\operatorname{Var}(E_i) = \frac{\lambda_n^2 n^4 (ns + n - s^2 - 2s + 1)}{(n - s)(n - s - 3)(n - s - 1)} - \frac{\lambda_n^2 n^4}{(n - s - 1)^2}$$
$$= \frac{\lambda_n^2 n^4 (n - 1)}{(n - s)(n - s - 3)(n - s - 1)}. \tag{A.24}$$

By Chebyshev's inequality, we can get

$$\mathbb{P}(|E_{i}| \geq (1+\sqrt{n})\,\mathbb{E}(E_{i})) \leq \mathbb{P}(|E_{i}-\mathbb{E}(E_{i})| \geq \sqrt{n}\,\mathbb{E}(E_{i}))$$

$$\leq \frac{\operatorname{Var}(E_{i})}{n(\mathbb{E}(E_{i}))^{2}}$$

$$= \frac{ns+n-s^{2}-2s+1}{4n(n-s)(n-s-3)}$$

$$\leq \frac{K_{1}}{n-s}, \tag{A.25}$$

for some constant  $K_1$  when n is large enough.

Similarly, by Lemma 1.3 (2) again for i = j, and  $\Sigma = \frac{1}{n}I$ , we have

$$\operatorname{Var}(V_i^2) = \sigma'^4 \operatorname{Var} \left[ (e_i^\top \left( \boldsymbol{X}_{S_0}^\top \boldsymbol{X}_{S_0} \right)^{-1} e_i)^2 \right]$$
$$= \sigma'^4 \frac{(n-s+1+n-s-1)n^2}{(n-s)(n-s-1)^2(n-s-3)}$$
$$= \frac{2\sigma'^4 n^2}{(n-s)(n-s-1)(n-s-3)},$$

and thus

$$\mathbb{P}(V_i \ge 2E(V_i)) = \mathbb{P}(V_i - \mathbb{E}(V_i) \ge E(V_i))$$

$$\le \frac{\operatorname{Var}(V_i)}{(\mathbb{E}(V_i))^2} = \frac{\frac{2\sigma'^4 n^2}{(n-s)(n-s-1)(n-s-3)}}{\left(\frac{n\sigma'^2}{n-s-1}\right)^2}$$

$$\le \frac{K_2}{n-s}, \tag{A.26}$$

for some constant  $K_2$  for large n. Therefore combining (A.25,A.26) with union bound, the statement in part 3 holds with  $K = K_1 + K_2$ .

Lemma 1.8: Let  $(X_1, \ldots, X_n)$  be independent and normally distributed. We have

$$\begin{split} \mathbb{E}[\max_{1 \leq i \leq n} |X_i|] \leq 3\sqrt{\log n} \max_{1 \leq i \leq n} \sqrt{\mathbb{E}\,X_i^2}. \\ \textit{Proof of Lemma 1.8:} \text{ This is a well-known result of} \end{split}$$

Proof of  $\overline{Lemma}$  1.8: This is  $\overline{a}$  well-known result of Gaussian maxima. We omit its proof and refer the reader to, for example, [24] Lemma 9.

Lemma 1.9: Let  $Y \sim \mathcal{N}(\mu, \sigma^2)$ . Suppose  $\mu \leq \mu_0$ ,  $\sigma \leq \sigma_0$ , and  $\rho \geq \mu_0$ , then the probability  $\mathbb{P}(Y \geq \rho) \leq \mathbb{P}(Z \geq \rho)$ , where  $Z \sim \mathcal{N}(\mu_0.\sigma_0)$ .

Proof of Lemma 1.9: By definition, we have

$$\mathbb{P}(Y \ge \rho) = 1 - \Phi\left(\frac{\rho - \mu}{\sigma}\right) = \Phi\left(\frac{\mu - \rho}{\sigma}\right).$$

And similarly, we have

$$\mathbb{P}(Z \ge \rho) = \Phi\left(\frac{\mu_0 - \rho}{\sigma_0}\right).$$

By the assumption, we know that

$$\frac{\mu-\rho}{\sigma} \leq \frac{\mu_0-\rho}{\sigma_0}$$

and thus the lemma follows from the fact that  $\Phi$  is increasing.  $\hfill\Box$ 

Lemma 1.10: Under the working assumptions, the event  $E = \{\widehat{\boldsymbol{w}}_{S_0} = \operatorname{sgn}(\boldsymbol{\beta}_{S_0})\}$  satisfies

$$\mathbb{P}(E^C) = o_n(1).$$

Proof of Lemma 1.10: This is the classic result of Theorem 3 in [24]. For the Lasso problem considered in Equation (A.2), with the identity covariance matrix, all conditions (26a), (26b), (26c) therein are easily satisfied with  $C_{\min} = C_{\max} = 1$ . As long as we have the condition

$$\min_{j=1,\dots,s} \beta_j > g(\lambda) = c_3 \lambda + 20 \sqrt{\frac{\sigma^2 \log s}{n}}, \tag{A.27}$$

for some constant  $c_3 > 0$ , we can guarantee that  $\operatorname{sgn}(\boldsymbol{\beta}_{S_0}) = \operatorname{sgn}(\hat{\boldsymbol{\beta}}_{S_0}) \equiv \widehat{\boldsymbol{w}}_{S_0}$  with probability  $1 - c_1 \exp(-c_2 \min\{s, \log(k-s)\})$  for some constant  $c_1, c_2 > 0$ .

To verify condition (A.27), recall that  $\lambda = n^b$ , s = $O(n/\log p)$ , and  $\min_{i=1,\dots,s} \beta_i = M(n)^{k-s+1}$ . Because  $M(n) = n^a$  and  $b < (k-s+1)a - \frac{3}{2}$ , we have

$$\min_{j=1,\dots,s} \beta_j = M(n)^{k-s+1} = n^{(k-s+1)a}$$

$$\geq n^{b+\frac{3}{2}}$$

$$= \lambda n^{\frac{3}{2}}$$

$$\geq c_3 \lambda + 20 \sqrt{\frac{\sigma^2 \log s}{n}} = g(\lambda).$$

for sufficiently large n. This implies  $P(E) \to 1$ , or  $P(E^C) =$  $o_n(1)$ .

# B. Technical Proofs for Section IV

1) A Property of FDP and TPP: Any Trade-off Curve Is Strictly Increasing: A natural belief on the pair of (TPP, FDP) is that FDP (type-I error) should increase with TPP (power), which may be strengthened by our simulation plots. However, along a single Lasso path, this is in general not necessarily true. It is well-known that Lasso is not monotone [18], so it is possible that with more and more true variables entering the Lasso path, fewer and fewer noise variables retain in the Lasso path. In such a case, FDP is no longer a monotone function of TPP. This possibility complicates our analysis, yet the following lemma asserts that this possibility is impossible. We prove that the asymptotic FDP is strictly increasing with the asymptotic TPP. Formally speaking, as  $\lambda$  varies,  $fdp_{\lambda}^{\infty}$ can be seen as a function of  $tpp_{\lambda}^{\infty}$ , and  $fdp_{\lambda}^{\infty}$  is a strictly increasing function of  $tpp_{\lambda}^{\infty}$ . To be rigorous, in the following lemma—indeed throughout the paper—we consider the regime below the Donoho-Tanner phase transition. We refer interested readers to [19] for results above this phase transition.

Lemma 1.11: Fix  $\epsilon, \delta, \sigma$ , and  $\Pi \neq 0$ . We have that  $\operatorname{fdp}_{\lambda}^{\infty}(\Pi)$  is a strictly increasing function of  $\operatorname{tpp}_{\lambda}^{\infty}(\Pi)$ . That is,  $\left. \mathrm{fdp}_{\lambda}^{\infty}(\mathrm{tpp}_{\lambda}^{\infty}) \text{ is a well-defined function, and } \mathrm{fdp}_{\lambda}^{\infty\prime}(\cdot) \right|_{\mathrm{tpp}_{\lambda}^{\infty}} >$ 0 for any valid value of  $tpp_{\lambda}^{\infty}$ .

To prove this lemma, we need the following characterizations among  $\alpha$ ,  $\lambda$ ,  $fdp_{\lambda}^{\infty}$  and  $tpp_{\lambda}^{\infty}$ .

<sup>5</sup>As we will see in Lemma 1.15, the valid range of  $tpp_{\lambda}^{\infty}$  is the range  $(0, u^*)$ . In this paper, we only focus on the case where  $u^* = 1$ .

Lemma 1.12: Fix  $\epsilon, \delta, \sigma$ , and  $\Pi \neq 0$ . Consider any  $\alpha, \tau, \lambda$ that solve equations (IV.1). We have the following facts

$$\frac{\mathrm{d}\alpha}{\mathrm{d}\lambda} > 0 \tag{A.28}$$

$$\frac{\mathrm{dtpp}^{\infty}}{\mathrm{d}\alpha} < 0 \tag{A.29}$$

$$\frac{\mathrm{d}f\mathrm{d}p^{\infty}}{\mathrm{d}\alpha} < 0 \tag{A.30}$$

$$\frac{\mathrm{dtpp}^{\infty}}{\mathrm{dfdp}^{\infty}} > 0 \tag{A.31}$$

 $\frac{\mathrm{d}\lambda}{\mathrm{d}\mathrm{pp}^{\infty}} < 0 \qquad (A.29)$   $\frac{\mathrm{d}\mathrm{fdp}^{\infty}}{\mathrm{d}\alpha} < 0 \qquad (A.30)$   $\frac{\mathrm{d}\mathrm{tpp}^{\infty}}{\mathrm{d}\alpha} > 0 \qquad (A.31)$ We note that the denotations of  $\mathrm{fdp}^{\infty}$  and  $\mathrm{tpp}^{\infty}$  stand for  $fdp^{\infty}_{\alpha}$  and  $tpp^{\infty}_{\alpha}$ , where we treat  $\alpha$  as the free parameter. We often suppress this dependence on  $\alpha$  in the following proof when it is clear from the context, and use the denotations of  $fdp^{\infty}$  and  $tpp^{\infty}$  for simplicity.

Proof of Lemma 1.12: The (A.28) is a well-known result, and one can refer to, for example, Lemma 4.11 of [15] for a

To prove (A.29), we note that  $tpp^{\infty} = \mathbb{P}(|\Pi^{*} + \tau W| > \alpha \tau)$ , where  $\Pi^*$  is the distribution of an entry of  $\beta$  given it's not zero. For any  $\Pi$  that is a proper distribution and satisfies (IV.1), proving  $\frac{\mathrm{d}}{\mathrm{d}\alpha}\mathbb{P}(|\Pi+\tau W|>\alpha\tau)<0$  will suffice. And this result follows from Lemma 4.10 of [15]. Now we left to prove (A.30) and (A.31). We note that, however, (A.31) follows directly from (A.30), and therefore we only need to prove (A.30). To see this fact, we note that  $tpp^{\infty}$  is a strictly increasing function of  $\alpha$ , so  $\alpha$  is also a function of tpp $^{\infty}$ . By the chain rule, we have

$$\frac{\mathrm{d}f\mathrm{d}p^{\infty}}{\mathrm{d}tpp^{\infty}} = \frac{\mathrm{d}f\mathrm{d}p^{\infty}}{\mathrm{d}\alpha} \cdot \frac{\mathrm{d}\alpha}{\mathrm{d}tpp^{\infty}}.$$

Note that we have already proven  $\frac{\mathrm{dtpp}^{\infty}}{\mathrm{d}\alpha} < 0$ , which implies  $\frac{\mathrm{d}\alpha}{\mathrm{dtpp}^{\infty}} < 0$ . Therefore, we only need to show (A.30) is true to prove (A.31).

Now, to prove (A.30), we observe that

$$\mathrm{fdp}^{\infty}(\alpha) = \frac{1}{1 + \frac{\epsilon \mathbb{P}\left(\left|\frac{\Pi^*}{\tau} + W\right| > \alpha\right)}{2(1 - \epsilon)\Phi(-\alpha)}},$$

and thus

$$\frac{\mathrm{d} \mathrm{f} \mathrm{d} \mathrm{p}^{\infty}}{\mathrm{d} \alpha} = \frac{\epsilon}{2(1-\epsilon)} \cdot \frac{\frac{\mathrm{d} \left(\frac{\mathbb{P}\left(\left|\frac{\Pi^{+}}{\tau}+W\right|>\alpha\right)}{\Phi\left(-\alpha\right)}\right)}{\mathrm{d} \alpha}}{\left(1 + \frac{\epsilon \mathbb{P}\left(\left|\frac{\Pi^{+}}{\tau}+W\right|>\alpha\right)}{2(1-\epsilon)\Phi\left(-\alpha\right)}\right)^{2}}.$$

Since the denominator is positive, we only need to show the numerator is negative. For simplicity, we will abuse the notation a little bit by using  $\Pi$  for  $\Pi^*$  in the rest of the proof. We need to show for all  $\Pi \neq 0$ ,

$$\frac{\mathrm{d}\left(\frac{\mathbb{P}\left(\left|\frac{\Pi}{\tau}+W\right|>\alpha\right)}{\Phi(-\alpha)}\right)}{\mathrm{d}\alpha}<0,$$

or equivalently,

$$\frac{\mathrm{d}\mathbb{P}\left(\left|\frac{\Pi}{\tau} + W\right| > \alpha\right)}{\mathrm{d}\alpha} \cdot \Phi(-\alpha) + \mathbb{P}\left(\left|\frac{\Pi}{\tau} + W\right| > \alpha\right) \cdot \phi(\alpha) > 0.$$
(A.32)

We observe that

$$\mathbb{P}\left(\left|\frac{\Pi}{\tau} + W\right| > \alpha\right) \\
= \mathbb{E}\left[\mathbb{P}_{W}\left(W > \alpha - \frac{\Pi}{\tau}\middle|\Pi\right) + \mathbb{P}_{W}\left(W < -\alpha - \frac{\Pi}{\tau}\middle|\Pi\right)\right] \\
= \mathbb{E}\left[\Phi\left(\frac{\Pi}{\tau} - \alpha\right) + \Phi\left(-\frac{\Pi}{\tau} - \alpha\right)\right].$$
(A.33)

Substituting the expressions of (A.33) into (A.32), we obtain, (A.34), as shown at the bottom of the next page.

Note the denominator is positive, so we only need to prove that the numerator is positive. Let  $g(u) = (\Phi(u-\alpha) + \Phi(-u-\alpha))\phi(\alpha) - (\phi(u-\alpha) + \phi(u+\alpha))\Phi(-\alpha)$ . By Lemma 1.14, g(u) > 0 for  $u \neq 0$ , and therefore we have

$$\Omega_1 + \Omega_3 = \frac{\sigma^2 \delta}{\tau^3} \mathbb{E}_{\Pi} \left[ g \left( \frac{\Pi}{\tau} \right) \right] > 0.$$
 (A.35)

For  $\Omega_5$ , we observe that if  $\Pi>0$ , then  $-\phi\left(\alpha-\frac{\Pi}{T}\right)+\phi\left(\alpha+\frac{\Pi}{T}\right)<0$ ; and if  $\Pi<0$ , then  $-\phi\left(\alpha-\frac{\Pi}{T}\right)+\phi\left(\alpha+\frac{\Pi}{T}\right)>0$ . Therefore, we have

$$\mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau} \left( -\phi \left( \alpha - \frac{\Pi}{\tau} \right) + \phi \left( \alpha + \frac{\Pi}{\tau} \right) \right) \right] \leq 0$$

So, by the fact that  $\Phi(-\alpha) \leq \frac{\phi(\alpha)}{\alpha}$ , the definition of  $\Phi(x)$ , we have

$$\Omega_{5} = \mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau^{2}} \left( -\phi \left( \alpha - \frac{\Pi}{\tau} \right) + \phi \left( \alpha + \frac{\Pi}{\tau} \right) \right) \right] 
\cdot \mathbb{E}_{\Pi} \left[ \left( \int_{\alpha - \frac{\Pi}{\tau}}^{\infty} \phi(w) dw + \int_{-\infty}^{-\alpha - \frac{\Pi}{\tau}} \phi(w) dw \right) \right] \cdot \alpha \Phi(-\alpha) 
\geq \mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau^{2}} \left( -\phi \left( \alpha - \frac{\Pi}{\tau} \right) + \phi \left( \alpha + \frac{\Pi}{\tau} \right) \right) \right] 
\cdot \mathbb{E}_{\Pi} \left[ \left( \int_{\alpha - \frac{\Pi}{\tau}}^{\infty} \phi(w) dw + \int_{-\infty}^{-\alpha - \frac{\Pi}{\tau}} \phi(w) dw \right) \right] \cdot \phi(\alpha) 
= \mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau^{2}} \left( -\phi \left( \alpha - \frac{\Pi}{\tau} \right) + \phi \left( \alpha + \frac{\Pi}{\tau} \right) \right) \right] 
\cdot \mathbb{E}_{\Pi} \left[ \Phi \left( \frac{\Pi}{\tau} - \alpha \right) + \Phi \left( -\frac{\Pi}{\tau} - \alpha \right) \right] \cdot \phi(\alpha) 
= \mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau^{2}} \int_{-\alpha - \frac{\Pi}{\tau}}^{\alpha - \frac{\Pi}{\tau}} w \phi(w) dw \right] 
\cdot \mathbb{E}_{\Pi} \left[ \Phi \left( \frac{\Pi}{\tau} - \alpha \right) + \Phi \left( -\frac{\Pi}{\tau} - \alpha \right) \right] \cdot \phi(\alpha).$$
(A.36)

Similarly, by (A.34) and (A.36), and then by the definition of f(u) in Lemma 1.13, we obtain

$$\Omega_{2} + \Omega_{5} \geq \mathbb{E}_{\Pi} \left[ \Phi \left( \frac{\Pi}{\tau} - \alpha \right) + \Phi \left( -\frac{\Pi}{\tau} - \alpha \right) \right] 
\cdot \mathbb{E}_{\Pi} \left[ \frac{\Pi}{\tau^{2}} \int_{-\alpha - \frac{\Pi}{\tau}}^{\alpha - \frac{\Pi}{\tau}} w \phi(w) dw + \frac{\Pi^{2}}{\tau^{3}} \int_{-\alpha - \frac{\Pi}{\tau}}^{\alpha - \frac{\Pi}{\tau}} \phi(w) dw \right] \cdot \phi(\alpha) 
= \mathbb{E}_{\Pi} \left[ \Phi \left( \frac{\Pi}{\tau} - \alpha \right) + \Phi \left( -\frac{\Pi}{\tau} - \alpha \right) \right] \mathbb{E}_{\Pi} \left[ \frac{1}{\tau} f \left( \frac{\Pi}{\tau} \right) \right] \phi(\alpha).$$

Similarly, we have

$$\Omega_4 + \Omega_6 = -\mathbb{E}_{\Pi} \left[ \left( \alpha - \frac{\Pi}{\tau} \right) + \phi \left( \alpha + \frac{\Pi}{\tau} \right) \right] \mathbb{E}_{\Pi} \left[ \frac{1}{\tau} f \left( \frac{\Pi}{\tau} \right) \right].$$
(A.38)

Combining the last display, (A.37), (A.38), Lemma 1.13 and the well-known result that  $\Phi(-\alpha) \leq \frac{\phi(\alpha)}{\alpha}$ , we obtain

$$\Omega_{2} + \Omega_{4} + \Omega_{5} + \Omega_{6} \ge \mathbb{E}_{\Pi} \left[ g \left( \frac{\Pi}{\tau} \right) \right] \cdot \mathbb{E}_{\Pi} \left[ \frac{1}{\tau} f \left( \frac{\Pi}{\tau} \right) \right] > 0,$$
(A.39)

Put together (A.35), (A.34) and (A.39), we have for all  $\Pi \neq 0$ 

$$\frac{\mathrm{d}\mathbb{P}\left(\left|\frac{\Pi}{\tau}+W\right|>\alpha\right)}{\mathrm{d}\alpha}\cdot\Phi(-\alpha)+\mathbb{P}\left(\left|\frac{\Pi}{\tau}+W\right|>\alpha\right)\cdot\phi(\alpha)>0,$$

which, by (A.32), amounts to (A.30), or

$$\frac{\mathrm{d}f\mathrm{d}p^{\infty}}{\mathrm{d}\alpha} < 0.$$

Therefore, combining with (A.29), we obtain that

$$\frac{\mathrm{d}tpp^{\infty}}{\mathrm{d}fdp^{\infty}}>0.$$

Summarizing the result we have proven, it is very easy to prove Lemma 1.11.

*Proof of Lemma 1.11:* Observe  $\operatorname{tpp}^{\infty}(\alpha)$  is a strictly increasing function of  $\alpha$ , and thus  $\operatorname{tpp}^{\infty}$  is a one-to-one function of  $\alpha$ . The inverse function therefore exists, so  $\alpha$  is a strictly increasing function of  $\operatorname{tpp}^{\infty}$ . Similarly,  $\operatorname{fdp}^{\infty}$  is also a strictly increasing function of  $\alpha$ . Therefore, we conclude that  $\operatorname{fdp}^{\infty} = \operatorname{fdp}^{\infty}(\alpha) = \operatorname{fdp}^{\infty}(\alpha(\operatorname{tpp}^{\infty})) = \operatorname{fdp}^{\infty}(\operatorname{tpp}^{\infty})$  is a strictly increasing function of  $\operatorname{tpp}^{\infty}$ , and that  $\frac{\operatorname{dfdp}^{\infty}}{\operatorname{dtpp}^{\infty}} > 0$  holds for any valid value of  $\operatorname{tpp}^{\infty}$ .

Now, we prove the lemmas that we have used in the proof of Lemma 1.12.

Lemma 1.13: Let  $f(u) = u \int_{-\alpha-u}^{\alpha-u} (w+u)\phi(w) dw$ . We have f(u) > 0, for all  $u \neq 0 \in \mathbb{R}$ .

Proof of Lemma 1.13: Observe that

$$f(u) = u \int_{-\alpha - u}^{\alpha - u} (w + u)\phi(w) dw$$

$$\stackrel{w' = w + u}{=} u \int_{-\alpha}^{\alpha} w' \phi(w' - u) dw'$$

$$= u \int_{0}^{\alpha} w' [\phi(w' - u) - \phi(-w' - u)] dw'.$$

So, if u>0, then  $\phi(w'-u)-\phi(-w'-u)>0$ , for any  $w'\in(0,\alpha]$ , thus f(u)>0; and if u<0, then  $\phi(w'-u)-\phi(-w'-u)<0$ , for any  $w'\in(0,\alpha]$ , thus f(u)>0.  $\square$  Lemma 1.14: For any fixed  $\alpha>0$ , let  $g(u)=(\Phi(u-\alpha)+\Phi(-u-\alpha))\phi(\alpha)-(\phi(u-\alpha)+\phi(u+\alpha))\Phi(-\alpha)$ . Then we have  $g(u)\geq 0$ , and the equality g(u)=0 holds if and only if

Proof of Lemma 1.14: We observe that

$$g(u) = \phi(\alpha) \left( \int_{\alpha - u}^{\infty} \phi(w) dw + \int_{\alpha + u}^{\infty} \phi(w) dw - \int_{\alpha}^{\infty} \phi(w) dw \left( \frac{\phi(u - \alpha)}{\phi(\alpha)} + \frac{\phi(u + \alpha)}{\phi(\alpha)} \right) \right)$$

$$= \phi(\alpha) \left( \int_{\alpha}^{\infty} \phi(w) e^{wu} \cdot e^{\frac{-u^2}{2}} dw + \int_{\alpha}^{\infty} \phi(w) e^{-wu} \cdot e^{\frac{-u^2}{2}} dw - \int_{\alpha}^{\infty} \phi(w) (e^{\alpha u} + e^{-\alpha u}) \cdot e^{\frac{-u^2}{2}} dw \right)$$

$$= \phi(\alpha) e^{\frac{-u^2}{2}} \int_{\alpha}^{\infty} \phi(w) [(e^{wu} + e^{-wu}) - (e^{\alpha u} + e^{-\alpha u})] dw.$$

Since for any  $w > \alpha > 0$ , we have

$$e^{wu} + e^{-wu} > e^{\alpha u} + e^{-\alpha u}$$
, for any  $u \in \mathbb{R}$ .

We obtain  $g(u) \ge 0$ , and it is clear the equality holds if and only if u = 0.

2) Miscellaneous Proofs for Section IV-A: In this section, we prove all the necessary lemmas for Theorem 2. To start with, we state the following lemma that specifies the range of all valid tpp<sup>∞</sup>'s.

Lemma 1.15: [Lemma C.1 and Lemma C.4 in [5]] Put

$$u^{\star}(\delta, \epsilon) := \begin{cases} 1 - \frac{(1 - \delta)(\epsilon - \epsilon^{\star})}{\epsilon (1 - \epsilon^{\star})}, & \delta < 1 \text{ and } \epsilon > \epsilon^{\star}(\delta), \\ 1, & \text{otherwise.} \end{cases}$$

Then

$$tpp^{\infty} < u^{\star}(\delta, \epsilon).$$

$$\begin{split} &\frac{\mathrm{d}\mathbb{P}\left(\left|\frac{\Pi}{\tau}+W\right|>\alpha\right)\cdot\Phi(-\alpha)+\mathbb{P}\left(\left|\frac{\Pi}{\tau}+W\right|>\alpha\right)\cdot\phi(\alpha)}{\mathrm{d}\alpha} \\ &=\frac{1}{\frac{\sigma^{2}\delta}{\tau^{3}}+\mathbb{E}_{\Pi,W}\left[\frac{\Pi^{2}}{\tau^{3}}\mathbf{1}_{\{-\alpha<\frac{\Pi}{\tau}+W<\alpha\}}\right]}}\\ &\cdot\left[\boxed{\mathbb{E}_{\Pi}\left[\Phi\left(\frac{\Pi}{\tau}-\alpha\right)+\Phi(-\frac{\Pi}{\tau}-\alpha)\right]\cdot\frac{\sigma^{2}\delta}{\tau^{3}}\cdot\phi(\alpha)}\right. \\ &+\mathbb{E}_{\Pi}\left[\Phi\left(\frac{\Pi}{\tau}-\alpha\right)+\Phi(-\frac{\Pi}{\tau}-\alpha)\right]\cdot\mathbb{E}_{\Pi,W}\left[\frac{\Pi^{2}}{\tau^{3}}\mathbf{1}_{\{-\alpha<\frac{\Pi}{\tau}+W<\alpha\}}\right]\cdot\phi(\alpha)}\right. \\ &-\mathbb{E}_{\Pi}\left[\phi\left(\alpha-\frac{\Pi}{\tau}\right)+\phi\left(\alpha+\frac{\Pi}{\tau}\right)\right]\cdot\mathbb{E}_{\Pi,W}\left[\frac{\Pi^{2}}{\tau^{3}}\mathbf{1}_{\{-\alpha<\frac{\Pi}{\tau}+W<\alpha\}}\right]\cdot\Phi(-\alpha)}\right. \\ &+\mathbb{E}_{\Pi}\left[\frac{\Pi}{\tau^{2}}\left(-\phi\left(\alpha-\frac{\Pi}{\tau}\right)+\phi\left(\alpha+\frac{\Pi}{\tau}\right)\right)\right]\cdot\mathbb{E}_{\Pi,W}\left[\frac{\Pi^{2}}{\tau^{3}}\mathbf{1}_{\{-\alpha<\frac{\Pi}{\tau}+W<\alpha\}}\right]\cdot\Phi(-\alpha)}\right. \\ &+\mathbb{E}_{\Pi}\left[\frac{\Pi}{\tau^{2}}\left(-\phi\left(\alpha-\frac{\Pi}{\tau}\right)+\phi\left(\alpha+\frac{\Pi}{\tau}\right)\right)\right] \\ &\cdot\mathbb{E}_{\Pi}\left[\frac{\Pi}{\tau^{2}}\left(-\phi\left(\alpha-\frac{\Pi}{\tau}\right)+\phi\left(\alpha+\frac{\Pi}{\tau}\right)\right)\right]\right. \\ &+\mathbb{E}_{\Pi}\left[\frac{\Pi}{\tau^{2}}\left(-\phi\left(\alpha-\frac{\Pi}{\tau}\right)+\phi\left(\alpha+\frac{\Pi}{\tau}\right)\right)\right] \\ &\cdot\mathbb{E}_{\Pi}\left[\int_{\alpha-\frac{\Pi}{\tau}}^{\infty}-w\phi(w)\mathrm{d}w+\int_{-\infty}^{-\alpha-\frac{\Pi}{\tau}}w\phi(w)\mathrm{d}w\right]\cdot\Phi(-\alpha). \end{split}$$

(A.34)

Moreover, any u between 0 and  $u^*$  can be uniquely realized as  $tpp^{\infty}$ , by setting  $\alpha = t^*(u)$  which is the root to (II.4).

From this lemma, we know for  $\delta < 1$  and large  $\epsilon$ , it is possible that the range of  $tpp^{\infty}$  is no longer (0,1). In such a case, we are "above the Donoho–Tanner phase transition (DTPT)"; and symmetrically, when  $tpp^{\infty}$  has the range (0,1), we are "below the DTPT". The purpose of this lemma is mainly for the completeness of the theory. In the following, however, we will always assume the range of  $tpp^{\infty}$  is (0,1) to avoid extra complicity. This assumption will simplify our argument, but the proofs of the theorems can be extended to the case when this assumption is not true.

Now, we prove that the upper curve can be achieved by any  $(\epsilon, M)$ -homogeneous prior (II.7). This implies that the homogeneous effect sizes are the least desired.

Lemma 1.16: Given  $(\epsilon, \delta)$  and  $\sigma = 0$ . Any  $(\epsilon, M)$ -homogeneous prior gives the same unique trade-off curve  $q^{\nabla}$  on  $(0, u^*)$ . Furthermore, this curve has the expression specified in (II.8).

Proof of Lemma 1.16: We start with the proof to show the curve  $q^{\nabla}$  is unique in the sense that any two different  $(\epsilon, M)$ -homogeneous priors give the same trade-off curve. Consider any two  $(\epsilon, M)$ -homogeneous priors  $\Pi_1$  and  $\Pi_2$ . Let their nonzero conditional priors be  $\Pi_1^* \equiv M_1$  and  $\Pi_2^* \equiv M_2$ . Treating  $\alpha > \alpha_0$  as the free parameter, we denote the solution to  $\tau$  in equation (IV.1) with prior  $\Pi_1$  by  $\tau_1$ . We have

$$\delta = (1 - \epsilon) \mathbb{E}(\eta_{\alpha}(W)^{2}) + \epsilon \mathbb{E}\left(\eta_{\alpha}\left(\frac{M_{1}}{\tau_{1}} + W\right) - \frac{M_{1}}{\tau_{1}}\right)^{2}.$$

It is clear from a simple calculation that  $\tau_2 = \tau_1 \frac{M_1}{M_2}$ ,  $\alpha$  and  $\Pi_2$  also solve the first equation in (IV.1), that is,

$$\delta = (1 - \epsilon) \mathbb{E}(\eta_{\alpha}(W)^{2}) + \epsilon \mathbb{E}\left(\eta_{\alpha}\left(\frac{M_{2}}{\tau_{2}} + W\right) - \frac{M_{2}}{\tau_{2}}\right)^{2},$$
(A.40)

which implies  $\tau_2$  is the solution to (IV.1) given  $\alpha$  and prior  $\Pi_2$ . Observe the relationships  $\tau_2 = \tau_1 \frac{M_1}{M_2}$ ,  $\Pi_1^{\star} \equiv M_1$  and  $\Pi_2^{\star} \equiv M_2$ . We have

$$\mathbb{P}\left(\left|\frac{\Pi_1^\star}{\tau_1} + W\right| > \alpha\right) = \mathbb{P}\left(\left|\frac{\Pi_2^\star}{\tau_2} + W\right| > \alpha\right).$$

Therefore, combining the equality above with (IV.2), we obtain

$$tpp_{\alpha}^{\infty}(\Pi_{1}) = \mathbb{P}\left(\left|\frac{\Pi_{1}^{\star}}{\tau_{1}} + W\right| > \alpha\right) = \mathbb{P}\left(\left|\frac{\Pi_{2}^{\star}}{\tau_{2}} + W\right| > \alpha\right)$$
$$= tpp_{\alpha}^{\infty}(\Pi_{2}),$$

and

$$\begin{split} \mathrm{fdp}_{\alpha}^{\infty}(\Pi_{1}) &= \frac{2(1-\epsilon)\Phi(-\alpha)}{2(1-\epsilon)\Phi(-\alpha) + \epsilon \, \mathbb{P}\left(\left|\frac{\Pi_{1}^{\star}}{\tau_{1}} + W\right| > \alpha\right)} \\ &= \frac{2(1-\epsilon)\Phi(-\alpha)}{2(1-\epsilon)\Phi(-\alpha) + \epsilon \, \mathbb{P}\left(\left|\frac{\Pi_{2}^{\star}}{\tau_{2}} + W\right| > \alpha\right)} \\ &= \mathrm{fdp}_{\alpha}^{\infty}(\Pi_{1}). \end{split}$$

This means that any point on  $q^{\Pi_1}(\cdot)$  is also on  $q^{\Pi_2}(\cdot)$ . Similarly, any point on  $q^{\Pi_2}(\cdot)$  is also on  $q^{\Pi_1}(\cdot)$ . By Lemma 1.11, they are both strictly increasing function, and thus must equal everywhere on the entire domain  $(0, u^*)$ .

Now, we proceed to prove that this unique trade-off curve has the expression given by (II.8). Fix some  $(\epsilon, M)$ -homogeneous prior  $\Pi^{\nabla}$ . Let u be some point between 0 and  $u^* = 1$ , and set  $\alpha$  such that  $\operatorname{tpp}_{\alpha}^{\infty}(\Pi^{\nabla})$  equals to u. We have

$$u\!=\!\operatorname{tpp}_{\alpha}^{\infty}\!=\!\mathbb{P}(|\Pi^{\nabla\star}\!+\!\tau W|\!>\!\alpha\tau)\!=\!\Phi(\!-\alpha\!+\!\widetilde{M})\!+\!\Phi(\!-\alpha-\widetilde{M}),$$

where  $\widetilde{M} = \frac{M}{\tau}$ . Let  $\varsigma = -\alpha + \widetilde{M}$ , then the equation above becomes

$$\Phi(\varsigma) + \Phi(-2\alpha - \varsigma) = u. \tag{A.41}$$

According to (IV.1), we have

$$\delta = (1 - \epsilon) \mathbb{E}[\eta_{\alpha}(W)]^2 + \epsilon \mathbb{E}[\eta_{\alpha}(\widetilde{M} + W) - \widetilde{M}]^2. \quad (A.42)$$

By a simple algebraic fact

$$\mathbb{E}[\eta_{\alpha}(W)]^{2} = 2[(1+\alpha^{2})\Phi(-\alpha) - \alpha\phi(\alpha)],$$

and the fact

$$\begin{split} &\mathbb{E}[\eta_{\alpha}(\widetilde{M}+W)-\widetilde{M}]^2\\ &=-(\alpha+\widetilde{M})\phi(\alpha-\widetilde{M})-(\alpha-\widetilde{M})\phi(\alpha+\widetilde{M})\\ &+(1+\alpha^2)[\Phi(-\alpha+\widetilde{M})+\Phi(-\alpha-\widetilde{M})]\\ &+\widetilde{M}^2[\Phi(\alpha-\widetilde{M})-\Phi(-\alpha-\widetilde{M})], \end{split}$$

we can plug-in the definition of  $\varsigma$  into (A.42) and obtain

$$\delta = 2(1 - \epsilon)[(1 + \alpha^2)\Phi(-\alpha) - \alpha\phi(\alpha)]$$

$$+ \epsilon[-(2\alpha + \varsigma)\phi(\varsigma) + \varsigma\phi(2\alpha + \varsigma) +$$

$$(1 + \alpha^2)[\Phi(\varsigma) + \Phi(-2\alpha - \varsigma)]$$

$$+ (\varsigma + \alpha)^2[\Phi(-\varsigma) + \Phi(-2\alpha - \varsigma)]]. \tag{A.43}$$

So  $\varsigma = \varsigma(\alpha; \epsilon, \delta)$  is the solution of the equation above. And finally combining the last equation with (A.41), we get an equation in  $\alpha$ 

$$\Phi(\varsigma(\alpha)) + \Phi(-2\alpha - \varsigma(\alpha)) = u, \tag{A.44}$$

Denote the solution of  $\alpha$  to the equation above by  $t^{\nabla}=t^{\nabla}(u;\epsilon,\delta)$ . We have

$$q^{\nabla}(u;\epsilon,\delta) = \mathrm{fdp}_{t^{\nabla}}^{\infty}(\Pi^{\nabla}) = \frac{2(1-\epsilon)\Phi(-t^{\nabla}(u))}{2(1-\epsilon)\Phi(-t^{\nabla}(u)) + \epsilon u}.$$
(A.45)

Therefore, the expression for the upper boundary is just defined by (A.45), where  $t^{\nabla}$  is solved from (A.43) and (A.44).  $\square$  We comment about the existence of  $\alpha$  in the proof above. Note that both equations (A.43) and (A.44) are derived from the AMP equations, which for any  $\alpha > \alpha_0$ , have unique solution  $\tau$ . Note that  $\varsigma$  is a function of  $\tau$ , and thus it is also unique. Therefore, the solution to (A.43) also uniquely exists by Lemma 1.15.

3) Miscellaneous Proofs for Section IV-B: In this section, we prove all necessary lemmas needed for Theorem 1, and then prove Theorem 1. We start by giving the proof to Lemma 4.4.

*Proof of Lemma 4.4:* We treat  $\tau$  as the free parameter instead of  $\lambda$ . To explicitly express the limiting process of M, we consider a sequence of priors  $\{\Pi^{\Delta}(M^{(t)}, \gamma)\}_t$ , where

 $M_1^{(t)} \to \infty$  and  $M_{i+1}^{(t)}/M_i^{(t)} \to \infty$  as  $t \to \infty$ . From (IV.2), the asymptotic TPP of  $\Pi^{\Delta}(M^{(t)}, \gamma)$  at  $\tau$  can be written as

$$\operatorname{tpp}_{\tau}^{\infty}(\Pi^{\Delta}(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})) = \mathbb{P}\left(\left|\Pi^{\Delta\star}(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma}) + \tau W\right| > \alpha\tau\right)$$

$$= \left[\gamma_{1} \mathbb{P}\left(\left|W + \frac{M_{1}^{(t)}}{\tau}\right| > \alpha\right) + \gamma_{2} \mathbb{P}\left(\left|W + \frac{M_{2}^{(t)}}{\tau}\right| > \alpha\right) + \dots + \gamma_{m} \mathbb{P}\left(\left|W + \frac{M_{m}^{(t)}}{\tau}\right| > \alpha\right)\right], \tag{A.46}$$

where  $\alpha$  is solved from (IV.1). We denote the last display by  $\operatorname{tpp}_{\tau}^{\infty(t)}$  for convenience. Similarly, we denote

$$\begin{split} \mathrm{fdp}_{\tau}^{\infty(t)} &\equiv \mathrm{fdp}_{\tau}^{\infty}(\Pi^{\Delta}(\boldsymbol{M}^{(t)},\boldsymbol{\gamma})) \\ &= \frac{2(1-\epsilon)\Phi(-\alpha)}{2(1-\epsilon)\Phi(-\alpha) + \epsilon \mathrm{tpp}^{\infty(t)}(\tau)}. \end{split} \tag{A.47}$$

In the following proof, we will choose an m-tuple of  $(\tau_i^{(t)})_{i=1}^m$  for each fixed t, such that as  $t \to \infty$ ,  $(\operatorname{tpp}_{\tau_i^{(t)}}^{\infty(t)}, \operatorname{fdp}_{\tau_i^{(t)}}^{\infty(t)}) \to (u_i, q^\Delta(u_i))$  at m different  $u_i$ 's. This implies exactly that the limit of trade-off curves  $q^{\Pi^\Delta(\boldsymbol{M}^{(t)}, \gamma)}$  agrees with  $q^\Delta$  at (at least) different m points in (0, 1]. A natural way to pick such an m-tuple  $(\tau_i^{(t)})_{i=1}^m$  is

$$\tau_i^{(t)} = \sqrt{M_i^{(t)} M_{i+1}^{(t)}}, 1 \le i \le m-1,$$

and  $\tau_m^{(t)}=m\times M_m^{(t)}$  when  $i=m.^6$  Under the regime of  $M_1^{(t)}\to\infty$  and  $M_{i+1}^{(t)}/M_i^{(t)}\to\infty$  for all i, we know  $\tau_i^{(t)}$  satisfies

$$|M_i^{(t)}| = o(\tau_i^{(t)}), \text{ and } \tau_i^{(t)} = o(|M_{i+1}^{(t)}|), \quad 1 \le i \le m-1 \tag{A.48}$$

and

$$|M_m^{(t)}| = o(\tau_m^{(t)}),$$
 (A.49)

as  $t \to \infty$ . Moreover, for any  $\alpha > \alpha_0$  and any  $1 \le i \le m$ , we have

$$\mathbb{P}\left(\left|W + \frac{M_j^{(t)}}{\tau_i^{(t)}}\right| > \alpha\right)$$

$$= \begin{cases} \mathbb{P}(|W| > \alpha) + o_t(1), & \text{for } j \leq i, \\ 1 + o_t(1), & \text{for } j \geq i + 1, \end{cases}$$
(A.50)

and

$$\eta_{\alpha} \left( W + \frac{M_{j}^{(t)}}{\tau_{i}^{(t)}} \right) - \frac{M_{j}^{(t)}}{\tau_{i}^{(t)}} = \begin{cases} \eta_{\alpha}(W) + o_{t}(1), & \text{for } j \leq i, \\ W - \alpha + o_{\mathbb{P}, t}(1), & \text{for } j \geq i + 1. \end{cases}$$
(A.51)

Let  $\gamma^{(j)} = \sum_{i=j+1}^{m} \gamma_i$  and use  $\alpha_i^{(t)}$  to denote the solution of  $\alpha$  to (IV.1) given  $\tau_i^{(t)}$ . We have

$$\left(1 - \frac{\sigma^2}{\tau_i^{(t)2}}\right) \delta \\
= \epsilon \cdot \mathbb{E}\left(\eta_\alpha \left(\frac{\Pi^\Delta(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})}{\tau_i^{(t)}} + W\right) - \frac{\Pi^\Delta(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})}{\tau_i^{(t)}}\right) \\
\stackrel{(*)}{=} \frac{\Pi^\Delta(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})}{\tau_i^{(t)}} + W \right) = \frac{\Pi^\Delta(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})}{\tau_i^{(t)}} + W$$

<sup>6</sup>In fact, one can pick any  $\tau_i^{(t)}$  such that (A.48) and (A.49) hold.

$$+(1-\epsilon)\mathbb{E}(\eta_{\alpha}(W))^{2}$$
. (A.52)

By (A.51), the (\*) part of the last display is

$$(*) = \sum_{j=1}^{i} \gamma_{j} \mathbb{E} \left( \eta_{\alpha} \left( \frac{M_{j}^{(t)}}{\tau_{i}^{(t)}} + W \right) - \frac{M_{i}^{(t)}}{\tau^{(t)}} \right)^{2} + \sum_{j=i+1}^{m} \gamma_{j} \mathbb{E} \left( \eta_{\alpha} \left( \frac{M_{j}^{(t)}}{\tau_{i}^{(t)}} + W \right) - \frac{M_{j}^{(t)}}{\tau_{i}^{(t)}} \right)^{2} = (1 - \gamma^{(i)}) \mathbb{E}(\eta_{\alpha}(W)^{2}) + \gamma^{(i)} \mathbb{E}((W - \alpha)^{2}) + o_{t}(1).$$

Observe the fact that  $\sigma$  is fixed and thus  $\frac{\sigma^2}{\tau_i^{(t)2}} \to 0$ . With some simple calculation, (A.52) can be written as

$$\epsilon \gamma^{(i)}(1+\alpha^2) + 2(1-\epsilon \gamma^{(i)})[(1+\alpha^2)\Phi(-\alpha) - 2\alpha\phi(\alpha)] = \delta + o_t(1).$$

Therefore the solution  $\alpha_i^{(t)}$  of the equation above has a limit<sup>7</sup>

$$\alpha_i^{(t)} \to \alpha^{(i)}, \quad \text{as } t \to \infty,$$
 (A.53)

which solves the equation

$$\epsilon \gamma^{(i)} (1 + \alpha^2) + 2(1 - \epsilon \gamma^{(i)})[(1 + \alpha^2)\Phi(-\alpha) - 2\alpha\phi(\alpha)] = \delta.$$
(A.54)

Note  $\alpha^{(i)}$  is independent of the choice of  $\{M^{(t)}\}_{t=1}^{\infty}$  and  $(\tau_i^{(t)})_{i=1}^m$ . Direct calculation can verify that each solution  $\alpha^{(i)}$  also satisfies the equation (II.4) with setting

$$u = u_i = 2\Phi(-\alpha^{(i)})(1 - \gamma^{(i)}) + \gamma^{(i)}.$$
 (A.55)

This implies  $\alpha^{(i)}$  is also the unique solution of (II.4), so

$$\alpha^{(i)} = t^{\Delta}(u_i). \tag{A.56}$$

Combining (A.46), (A.47), (A.50) and (A.53), the limits of  $\operatorname{tpp}_{\tau^{(t)}}^{\infty(t)}$  and  $\operatorname{fdp}_{\tau^{(t)}}^{\infty(t)}$  are

$$\begin{cases} \operatorname{tpp}_{\tau_{i}^{(t)}}^{\infty(t)} & \to \ 2\Phi(-\alpha^{(j)})(1-\gamma^{(i)}) + \gamma^{(i)}, \\ \operatorname{fdp}_{\tau_{i}^{(t)}}^{\infty(t)} & \to \frac{2(1-\epsilon)\Phi(-\alpha^{(i)})}{2(1-\epsilon)\Phi(-\alpha^{(i)}) + \epsilon(2\Phi(-\alpha^{(i)})(1-\gamma^{(i)}) + \gamma^{(i)})}. \end{cases}$$
(A.57)

By  $\alpha^{(i)}=t^{\Delta}(u_i)$  and (A.55), we obtain from (II.5) that

$$q^{\Delta}(u_i; \delta, \epsilon) = \frac{2(1 - \epsilon)\Phi(-t^{\Delta}(u_i))}{2(1 - \epsilon)\Phi(-t^{\Delta}(u_i)) + \epsilon u_i}$$

$$= \frac{2(1 - \epsilon)\Phi(-\alpha^{(i)})}{2(1 - \epsilon)\Phi(-\alpha^{(i)}) + \epsilon(2\Phi(-\alpha^{(i)})(1 - \gamma^{(i)}) + \gamma^{(i)})}.$$
(A.58)

Combining (IV.3), (IV.2), (A.55), (A.57) and (A.58), we finally obtain as  $t \to \infty$ ,

$$\begin{cases} \operatorname{tpp}_t^{\infty}(\tau_i^{(t)}) \to u_i, \\ \operatorname{fdp}_t^{\infty}(\tau_i^{(t)}) \to q^{\Delta}(u_i; \delta, \epsilon). \end{cases}$$

Therefore, the limiting function of the trade-off curves of priors  $\Pi^{\Delta}(\boldsymbol{M}^{(t)}, \boldsymbol{\gamma})$  agrees with the lower boundary  $q^{\Delta}(\cdot; \delta, \epsilon)$  at  $(u_i, q^{\Delta}(u_i))$ , for  $i = 1, 2, \cdots, m$ . Since the m different

 $<sup>^{7}</sup>$ By the existence asserted by AMP theory, the equation (A.54) has a unique solution, denote it by  $\alpha^{(i)}$ , we know the solution  $\alpha_{i}^{(t)}$  of (A.53) must converge into it, since the left-hand side of (A.53) is continuous in  $\alpha$ .

points are nonzero, there must be at least m-1 points in (0,1).

An important set of equations is (A.55). Recall that (A.56) asserts  $\alpha^{(i)} = t^{\Delta}(u_i)$  for all i, and therefore the equations (A.55) are, for all i,

$$u_i = 2\Phi(-t^{\Delta}(u_i))(1 - \gamma^{(i)}) + \gamma^{(i)}.$$
 (A.59)

The last display allows us to quantify the exact points  $u_i$  the limit of  $\Pi^{\Delta}(M^{(t)}, \gamma)$  agrees with  $q^{\nabla}$ . This fact allows us to set  $\gamma$  cleverly so that the distance between any two consecutive  $u_i$ 's are small enough. This is formalized in the following lemma.

Lemma 1.17: For any  $\xi > 0$ , there is some  $\gamma =$  $\{\gamma_1,\ldots,\gamma_m\}$ , with  $\sum_i \gamma_i = 1$ , such that the m points specified by (A.59), together with  $u_0 = 0$  and  $u_{m+1} = 1$ <sup>8</sup> satisfy the following

$$\max_{1 \le j \le m+1} |u_j - u_{j-1}| < \frac{\xi}{2}. \tag{A.60}$$

 $\max_{1 \leq j \leq m+1} |u_j - u_{j-1}| < \frac{\xi}{2}. \tag{A.60}$  Proof of Lemma 1.17: We first prove that the difference  $u_{m+1} - u_m = 0$ . Since  $\gamma^{(m)} = 1$ , and by (A.55), we obtain that  $u_m = 1$ , and thus  $u_{m+1} - u_m = 1 - 1 = 0.9$  With this in mind, we only need to prove the following two quantities can be arbitrarily small to ensure (A.60),

$$u_1 - 0 = \gamma_m + 2(1 - \gamma_m)\Phi(-\alpha^{(m)}),$$
 (A.61)

and for all  $m \ge j \ge 2$ ,

$$u_{j} - u_{j-1} = \gamma_{m-j+1} + 2(1 - \gamma^{(j)})\Phi(-\alpha^{(j)})$$
$$-2(1 - \gamma^{(j-1)})\Phi(-\alpha^{(j-1)}), \quad (A.62)$$

where we remind the reader that by definition,  $\gamma^{(m)} = \gamma_m$ and  $\gamma^{(j)} - \gamma^{(j-1)} = \gamma_{m-j+1}$ .

For the expression in (A.62), we observe that

$$u_{j} - u_{j-1} \leq \gamma_{m-j+1} | 1 + 2(\Phi(-\alpha^{(j)}) - \Phi(-\alpha^{(j-1)})) |$$

$$+ 2(1 - \gamma^{(j)}) | \Phi(-\alpha^{(j)}) - \Phi(-\alpha^{(j-1)}) |$$

$$\leq 5\gamma_{m-j+1} + 2|\Phi(-\alpha^{(j)}) - \Phi(-\alpha^{(j)} - \gamma_{m-j+1})) |.$$
(A.63)

We observe that in equation (A.55) or (A.59), the dependence of  $\alpha^{(j)}$  on  $\gamma^{(j)}$  is only through linear functional of  $\gamma^{(j)}$ . Therefore, it is not hard to realize that  $\alpha^{(j)}$  is continuous in  $\gamma^{(j)}$ . When all  $\{\gamma_s\}_{s>m-j+1}$  are fixed, the  $\alpha^{(j)}$  is a continuous function in  $\gamma_{m-j+1}$ , and so is the expression in (A.63). So we can pick  $\gamma_{m-j+1}$  sufficiently small to ensure the (A.63) is less than  $\frac{\xi}{2}$ .

For the expression in (A.61), we pick some M sufficiently large such that  $\Phi(-M) < \frac{\xi}{8}$ . By Lemma 1.18, we can pick  $\gamma_m < \frac{\xi}{4}$  such that the solution to (II.4) with u being (A.61)

<sup>8</sup>Technically speaking, it should be  $u_{m+1} = u^*$ , yet as discussed earlier, we will focus on the case when we are below the Donoho-Tanner phase transition, so always  $u^* = 1$ .

<sup>9</sup>One might want to verify the existence of  $\alpha^{(m)}$ . Since we always assume that we are below the DTPT, then for any  $u < u^* = 1$ , the  $\alpha = t^{\Delta}(u_i)$ exists as the solution to (II.4) by Lemma 1.15. By setting  $\gamma^{(m)}=1$ , one can directly verify this corresponds to set  $u = u_m = 1^-$ , and by the continuity of t in equation (II.4), we know  $\alpha^{(m)} = t^{\Delta}(u_m)$  exists and less than infinite. And since all other  $\alpha^{(i)}$ 's are less than  $\alpha^{(m)}$ , they also exist.

satisfies that  $\alpha^{(m)} > M$ , or  $\Phi(-\alpha^{(m)}) < \Phi(-M) < \frac{\xi}{9}$ .

$$\gamma_m + 2(1 - \gamma^{(m)})\Phi(-\alpha(\gamma^{(m)})) < \frac{\xi}{4} + 2 \cdot 1 \cdot \frac{\xi}{8} = \frac{\xi}{2}.$$

Lemma 1.18: For any fixed  $\delta$ ,  $\epsilon$ ,  $\xi > 0$  and M > 0. There exists  $\gamma < \frac{\xi}{4}$ , such that the solution  $\alpha$  to (II.4) with u = $\gamma + 2(1-\gamma)\Phi(-\alpha)$  satisfies  $\alpha > M$ .

Proof of Lemma 1.18: We will use the following fact: there exists  $\gamma < \frac{\xi}{4}$  and large M'' > M' > M, such that:

$$(1 - \epsilon \gamma) A(M') + \epsilon \gamma (1 + M'^2) < \delta,$$
  

$$(1 - \epsilon \gamma) A(M'') + \epsilon \gamma (1 + M''^2) > \delta.$$
 (A.64)

where  $A(M) = 2[(1 + M^2)\Phi(-M) - M\phi(M)].$ 

Taken this as given for the moment, we set  $u = \gamma + 2(1 - 1)$  $\gamma$ ) $\Phi(-\alpha)$  with  $\gamma$  such that (A.64) holds. Then (II.4) becomes <sup>10</sup>

$$\frac{2(1-\epsilon)\left[(1+t^2)\Phi(-t) - t\phi(t)\right] + \epsilon(1+t^2) - \delta}{\epsilon\left[(1+t^2)(1-2\Phi(-t)) + 2t\phi(t)\right]} = \frac{1-\gamma - 2(1-\gamma)\Phi(-t)}{1-2\Phi(-t)}.$$

Observe the right hand side of the last display is just  $1 - \gamma$ , so it is equivalent to

$$2(1 - \epsilon) \left[ (1 + t^2) \Phi(-t) - t \phi(t) \right] + \epsilon (1 + t^2) - \delta$$
  
=  $(1 - \gamma) \epsilon \left[ (1 + t^2) (1 - 2\Phi(-t)) + 2t \phi(t) \right],$ 

$$2(1 - \epsilon \gamma) \left[ (1 + t^2) \Phi(-t) - t \phi(t) \right] + \epsilon \gamma (1 + t^2) - \delta = 0.$$
(A.65)

By relationship (A.64) and the fact that there exists unique solution  $\alpha = t^{\nabla}(u)$  to (II.4) and thus to (A.65), we know the solution  $\alpha \in (M', M'')$ , and thus especially  $\alpha > M$ .

Now, to prove (A.64), we first note that it is direct to verify  $A(t) = E[\eta_t(W)^2]$ , and thus it is decreasing in t. And as  $t \to \infty$ ,  $A(t) \to 0$ . Therefore for any  $\delta > 0$ , we can pick M' large enough such that  $A(M') < \frac{\delta}{2}$ , and now pick  $\gamma < \frac{\xi}{4}$ small enough such that  $\epsilon \gamma (1 + M'^2) < \frac{\delta}{2}$ . Therefore, the lefthand side of the first equation of (A.64) is bounded by  $\delta$ . For the second equation in (A.64), pick M'' large enough so that the term  $\epsilon \gamma (1 + M''^2) > \delta$ , and since  $(1 - \epsilon \gamma) A(M'') > 0$ , the second line also holds.

The agreeing points asserted by Lemma 4.4 are close to each other in their x-coordinate distances. Therefore, by the uniform continuity of the lower curve  $q^{\nabla}$  and Cantor's diagonalization argument, we can extend the result from Lemma 4.4 to uniform convergence.

Lemma 1.19: There exist a sequence of prior of  $\Pi^{(t)} =$  $\Pi^{\Delta}(M^{(t)}, \gamma^{(t)})$ , such that their trade-off curve  $q^{\Pi^{(t)}}$  converge uniformly to  $q^{\Delta}$  on any compact interval in (0,1).

*Proof of Lemma 1.19:* Fix any compact interval  $\mathcal{I} = [a, b]$ in (0,1). As in Lemma 4.4, we first consider prior  $\Pi^{(t)} =$  $\Pi^{\Delta}(M^{(t)}, \gamma^{(t)})$  with  $\gamma^{(t)} = \gamma$  being some fixed mtuple. By Lemma 1.11, we know that both  $q^{\Pi^{(t)}}(u)$  and  $q^{\Delta}(u)$  are continuous and strictly increasing. Consider any

<sup>10</sup>Since the solution to t is  $\alpha$  here, we can plug in  $u = \gamma + 2(1 - \gamma)\Phi(-t)$ .

two adjacent agreeing points  $u_j, u_{j+1}$  specified in (A.55), such that  $q^{\Delta}(u_j) = \lim_{t \to \infty} q^{\Pi^{(t)}}(u_j)$  and  $q^{\Delta}(u_{j+1}) = \lim_{t \to \infty} q^{\Pi^{(t)}}(u_{j+1})$ . Since in the interval  $(u_j, u_{j+1})$  the difference is controlled by

$$q^{\Pi^{(t)}}(u) - q^{\Delta}(u) \le q^{\Delta}(u_{j+1}) - q^{\Delta}(u_j), \text{ for any } u \in (u_j, u_{j+1})$$

by the monotonicity of  $q^{\Pi^{(t)}}$ . This difference will be small as long as the gap  $q^{\Delta}(u_{j+1}) - q^{\Delta}(u_j)$  is small, so we proceed to prove we can select  $\Pi^{(t)}$  to ensure the gaps  $q^{\Delta}(u_{j+1}) - q^{\Delta}(u_j)$  are small for all i.

Fix any  $\theta > 0$ . Since  $q^{\Delta}$  is uniformly continuous on the compact set  $\mathcal{I}$ , there exists  $\xi > 0$  such that for any  $u, v \in \mathcal{I}$ ,

$$|q^{\Delta}(u) - q^{\Delta}(v)| < \frac{\theta}{2}$$
, as long as  $|u - v| < \xi$  (A.66)

By the proof of Lemma 4.4, we can construct  $\gamma^{(t)} = \gamma_{\theta}$  to be specified later, and  $M_{\gamma_{\theta}}^{(t)}$ , such that the limit of  $q^{\Pi^{(t)}}$  agrees with  $q^{\Delta}$  at m points  $u_1, \cdots, u_m$ . This implies there exists some  $T_{\theta}$  such that for all  $t \geq T_{\theta}$ ,

$$\max_{i} \left| q^{\Pi^{(t)}}(u_j) - q^{\Delta}(u_j) \right| < \frac{\theta}{2}. \tag{A.67}$$

To specify the choice of  $\gamma_{\theta}$ , note that by Lemma 1.17, we can choose  $\gamma_{\theta}$  such that  $u_1,...,u_m$  satisfies  $u_1-a < u_1-0 < \frac{\xi}{2},$   $u_m-b < \frac{\xi}{2}$  and

$$\max_{2 \le j \le m} |u_j - u_{j-1}| < \frac{\xi}{2}.$$

With this choice of  $\gamma_{\theta}$  together with (A.66) and (A.67), we obtain

$$\sup_{t \ge T_{\theta}, u \in \mathcal{I}} \left| q^{\Pi^{(t)}}(u) - q^{\Delta}(u) \right| < \theta.$$

Specifically, we have the equation above holds for  $t = T_{\theta}$ ,

$$\sup_{u \in \mathcal{T}} \left| q^{\Pi^{(T_{\theta})}}(u) - q^{\Delta}(u) \right| < \theta.$$

Since  $\Pi^{(T_{\theta})}=\Pi^{\Delta}\left(\boldsymbol{M}_{\boldsymbol{\gamma}_{\theta}}^{(T_{\theta})},\boldsymbol{\gamma}_{\theta}\right)$ , the inequality above is simply

$$\sup_{u \in \mathcal{I}} \left| q^{\Pi^{\Delta}(M_{\gamma_{\theta}}^{(T_{\theta})}, \gamma_{\theta})}(u) - q^{\Delta}(u) \right| < \theta. \tag{A.68}$$

Now we can apply Cantor's diagonalization trick since the last display is true for any  $\theta > 0$ . We set  $\theta_{\zeta} = \frac{1}{\zeta} \to 0, \zeta \ge 1$ , and choose the priors

$$\Pi^{(\zeta)} = \left\{ \Pi^{\Delta} \left( \boldsymbol{M}_{\gamma_{\boldsymbol{\theta}_{\zeta}}}^{(T_{\boldsymbol{\theta}_{\zeta}})}, \gamma_{\boldsymbol{\theta}_{\zeta}} \right) \right\}_{\zeta}.$$

Then we know from (A.68) that  $q^{\Pi^{(\zeta)}}$  converges to  $q^{\Delta}$  uniformly on  $\mathcal{I}$  as  $\zeta \to \infty$ .  $\square$  Now, we can proceed to prove Theorem 1, whose proof is very similar to that of Theorem 2 in Section IV-A.

Proof of Theorem 1 (a): Consider any non-constant prior  $\Pi$ , we first prove that there exists some  $\Pi^{\Delta}$  and  $\nu > 0$  such that for all  $c < \lambda, \lambda' < C$ ,

$$\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\Delta}) < \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) + \nu \text{ and } \operatorname{fdp}_{\lambda}^{\infty}(\Pi^{\Delta}) > \operatorname{fdp}_{\lambda'}^{\infty}(\Pi) - \nu \tag{A.69}$$

cannot hold simultaneously. To see this fact, first find  $0 < u_1 < u_2 < 1$  such that the asymptotic powers  $\operatorname{tpp}^\infty_\lambda(\Pi^\Delta), \operatorname{tpp}^\infty_{\lambda'}(\Pi)$  are always between  $u_1$  and  $u_2$  for c < 0

 $\lambda, \lambda' < C$ . Next, we know from Lemma 1.20 that  $q^{\Delta}$  is the strictly below any trade-off curve. So, for any prior  $\Pi$ , we have  $q^{\Delta}(u) < q^{\Pi}(u)$  for any  $u \in \mathcal{I} = [u_1, u_2]$ . Note that both  $q^{\Delta}$  and  $q^{\Pi}$  are uniformly continuous on  $\mathcal{I}$  and thus one can set  $\nu' > 0$  to be

$$\nu' := \inf_{u_1 \le u_< u_2} \left( q^{\Pi}(u) - q^{\Delta}(u) \right) > 0.$$
 (A.70)

Since  $q^{\Delta}$  is a continuous function on the closed interval [0,1], we can make use of its uniform continuity, which ensures

$$\left| q^{\Delta}(u) - q^{\Delta}(u') \right| < \frac{\nu'}{4} \tag{A.71}$$

as long as  $|u-u'| \le \nu''$  for some  $\nu'' > 0$ . By the assertion of Lemma 1.19, we can choose a prior  $\Pi^{\Delta}$  such that it is  $\frac{\nu'}{4}$ -close to  $q^{\Delta}$  on  $\mathcal{I}$ ,

$$\sup_{u_1 \le u \le u_2} (q^{\Pi^{\Delta}}(u) - q^{\Delta}(u)) < \frac{\nu'}{4}.$$
 (A.72)

Now we can prove (IV.4) cannot hold simultaneously with our choice of  $\Pi^{\Delta}$  and  $\nu = \min\{\nu'/2, \nu''\}$ . To see this, suppose we already have  $\operatorname{tpp}_{\lambda}^{\infty}(\Pi^{\Delta}) < \operatorname{tpp}_{\lambda'}^{\infty}(\Pi) + \nu$ . Now observe that

$$\begin{split} \mathrm{fdp}_{\lambda}^{\infty}(\Pi^{\Delta}) &= q^{\Pi^{\Delta}}(\mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\Delta})) \\ &\leq q^{\Delta}\left(\mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\Delta})\right) + \frac{\nu'}{4} \\ &< q^{\Delta}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi^{\Delta}) + \nu\right) + \frac{\nu'}{4} \\ &\leq q^{\Delta}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi)\right) + \frac{\nu'}{2} \\ &\leq q^{\Pi}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi)\right) - \nu' + \frac{\nu'}{2} \\ &= q^{\Pi}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi)\right) - \frac{\nu'}{2} \\ &\leq q^{\Pi}\left(\mathrm{tpp}_{\lambda'}^{\infty}(\Pi)\right) - \nu \\ &= \mathrm{fdp}_{\lambda'}^{\infty}(\Pi) - \nu, \end{split}$$

where the first inequality follows from (A.72); the second inequality follows from the fact that  $fdp^{\infty}(tpp^{\infty})$  is strictly increasing, and  $tpp^{\infty}_{\lambda}(\Pi^{\Delta}) < tpp^{\infty}_{\lambda'}(\Pi) + \nu$ ; the third inequality is by (A.71); the fourth inequality is by (A.70); the last inequality is by the definition of  $\nu$ . As such, the first inequality in (A.69) leads to the violation of the second inequality. Having shown (A.69), it is easy to prove Theorem 1. Lemma 4.1 ensures that the following four terms

$$\begin{aligned} & \left| TPP_{\lambda}(\Pi^{\Delta}) - tpp_{\lambda}^{\infty}(\Pi^{\Delta}) \right|, \left| FDP_{\lambda}(\Pi^{\Delta}) - fdp_{\lambda}^{\infty}(\Pi^{\Delta}) \right|, \\ & \left| TPP_{\lambda'}(\Pi) - tpp_{\lambda'}^{\infty}(\Pi) \right|, \left| FDP_{\lambda'}(\Pi) - fdp_{\lambda'}^{\infty}(\Pi) \right| \end{aligned} \tag{A.73}$$

are all smaller than  $\nu/2$  for all  $c < \lambda, \lambda' < C$ , with probability tending to one as  $n, p \to \infty$ . On this event, it is easy to check that  $\mathrm{TPP}_{\lambda}(\Pi^{\Delta}) \leq \mathrm{TPP}_{\lambda'}(\Pi)$  implies  $\mathrm{tpp}_{\lambda}^{\infty}(\Pi^{\Delta}) < \mathrm{tpp}_{\lambda'}^{\infty}(\Pi) + \nu$ , and  $\mathrm{FDP}_{\lambda}(\Pi^{\Delta}) \geq \mathrm{FDP}_{\lambda'}(\Pi)$  implies  $\mathrm{fdp}_{\lambda'}^{\infty}(\Pi^{\Delta}) > \mathrm{fdp}_{\lambda'}^{\infty}(\Pi) - \nu$ . As such, in the event (A.73), the impossibility of (A.69) uniformly for all  $c < \lambda, \lambda' < C$  implies the impossibility of

$$\operatorname{TPP}_{\lambda}(\Pi^{\Delta}) \leq \operatorname{TPP}_{\lambda'}(\Pi) \text{ and } \operatorname{FDP}_{\lambda}(\Pi^{\Delta}) \geq \operatorname{FDP}_{\lambda'}(\Pi)$$
 for all  $c < \lambda, \lambda' < C$ .

It is the same as the comment after the proof of Theorem 2, we can prove part (b) of Theorem 1 similarly, and we omit for simplicity.

Authorized licensed use limited to: University of Pennsylvania. Downloaded on June 09,2023 at 01:55:46 UTC from IEEE Xplore. Restrictions apply.

In closing, we present the following lemma to be self-contained. It shows that the lower boundary is strictly below any trade-off curve, on which the proof of Theorem 1 relies.

Lemma 1.20 (Lemma C.3 in [5]): Consider any  $\epsilon$ -sparse prior  $\Pi$ . The lower boundary  $q^{\Delta}$  is strictly below the trade-off curve  $q^{\Pi}(\cdot)$ , that is,  $q^{\Delta}(u) < q^{\Pi}(u)$  for any u.

Proof of Lemma 1.20: This is just a re-statement of Lemma C.3 in [5]. They have proved that for any  $\operatorname{tpp}^{\infty} = u$ ,  $\operatorname{fdp}^{\infty} > q^{\Delta}(u)$ , which implies  $q^{\Pi}(u) > q^{\Delta}(u)$ .

#### ACKNOWLEDGMENT

The authors would like to thank Małgorzata Bogdan, Emmanuel Candès, Edward George, Pragya Sur, and Nancy Zhang for stimulating discussions. They are grateful to the referees and an associate editor for their constructive comments that helped improve the presentation of this work.

#### REFERENCES

- E. J. Candès and T. Tao, "The Dantzig selector: Statistical estimation when p is much larger than n," Ann. Statist., vol. 35, no. 6, pp. 2313–2351, 2007.
- [2] P. J. Bickel, Y. Ritov, and A. B. Tsybakov, "Simultaneous analysis of Lasso and Dantzig selector," *Ann. Statist.*, vol. 37, no. 4, pp. 1705–1732, Aug. 2009.
- [3] R. Tibshirani, "Regression shrinkage and selection via the Lasso," J. Roy. Statist. Soc. B, Methodol., vol. 58, no. 1, pp. 267–288, Feb. 1994.
- [4] M. J. Wainwright, High-Dimensional Statistics: A Non-Asymptotic Viewpoint, vol. 48. Cambridge, U.K.: Cambridge Univ. Press, 2019.
- [5] W. Su, M. Bogdan, and E. J. Candès, "False discoveries occur early on the Lasso path," Ann. Statist., vol. 45, no. 5, pp. 2133–2150, Oct. 2017.
- [6] W. J. Su, "When is the first spurious variable selected by sequential regression procedures?" *Biometrika*, vol. 105, no. 3, pp. 517–527, Sep. 2018.
- [7] G. Raskutti, M. J. Wainwright, and B. Yu, "Restricted eigenvalue properties for correlated Gaussian designs," *J. Mach. Learn. Res.*, vol. 11, pp. 2241–2259, Aug. 2010.
- [8] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proc. Nat. Acad. Sci. USA*, vol. 106, no. 45, pp. 18914–18919, 2009.
- [9] M. Bayati and A. Montanari, "The dynamics of message passing on dense graphs, with applications to compressed sensing," *IEEE Trans.* Int. Theory, vol. 57, pp. 2, pp. 764–785, Feb. 2011.
- Inf. Theory, vol. 57, no. 2, pp. 764–785, Feb. 2011.
  [10] M. Bayati and A. Montanari, "The LASSO risk for Gaussian matrices," IEEE Trans. Inf. Theory, vol. 58, no. 4, pp. 1997–2017, Apr. 2012.
  [11] A. Maleki, L. Anitori, Z. Yang, and R. G. Baraniuk, "Asymptotic
- [11] A. Maleki, L. Anitori, Z. Yang, and R. G. Baraniuk, "Asymptotic analysis of complex LASSO via complex approximate message passing (CAMP)," *IEEE Trans. Inf. Theory*, vol. 59, no. 7, pp. 4290–4308, Jul. 2013.
- [12] M. Bayati, M. A. Erdogdu, and A. Montanari, "Estimating LASSO risk and noise level," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 944–952.
  [13] Z. Bu, J. Klusowski, C. Rush, and W. Su, "Algorithmic analysis and
- [13] Z. Bu, J. Klusowski, C. Rush, and W. Su, "Algorithmic analysis and statistical estimation of SLOPE via approximate message passing," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 9361–9371.
- Proc. Adv. Neural Inf. Process. Syst., 2019, pp. 9361–9371.
  [14] A. Weinstein, R. Barber, and E. J. Candès, "A power and prediction analysis for knockoffs with Lasso statistics," 2017, arXiv:1712.06465.
- [15] A. Mousavi, A. Maleki, and R. G. Baraniuk, "Consistent parameter estimation for LASSO and approximate message passing," *Ann. Statist.*, vol. 46, no. 1, pp. 119–148, 2018.
- [16] H. Weng, A. Maleki, and L. Zheng, "Overcoming the limitations of phase transition by higher order analysis of regularization techniques," *Ann. Statist.*, vol. 46, no. 6A, pp. 3099–3129, Dec. 2018.
- [17] P. Sur, Y. Chen, and E. J. Candès, "The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square," *Probab. Theory Rel. Fields*, vol. 175, nos. 1–2, pp. 487–558, Oct. 2019.
- [18] D. Donoho and J. Tanner, "Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing," *Philos. Trans. Roy. Soc. London A, Math., Phys. Eng. Sci.*, vol. 367, no. 1906, pp. 4273–4293, 2009.
- [19] S. Wang, H. Weng, and A. Maleki, "Which bridge estimator is the best for variable selection?" *Ann. Statist.*, vol. 48, no. 5, pp. 2791–2823, Oct. 2020.

- [20] H. Wang, Y. Yang, Z. Bu, and W. Su, "The complete Lasso tradeoff diagram," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, 2020, pp. 20051–20060.
- [21] D. L. Donoho, A. Maleki, and A. Montanari, "The noise-sensitivity phase transition in compressed sensing," *IEEE Trans. Inf. Theory*, vol. 57, no. 11, pp. 6920–6941, Oct. 2011.
- [22] A. Maleki, "Approximate message passing algorithms for compressed sensing," Ph.D. dissertation, Dept. Elect. Eng., Stanford Univ., Stanford, CA, USA, 2010.
- [23] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, "Least angle regression," Ann. Statist., vol. 32, no. 2, pp. 407–499, 2004.
- [24] M. J. Wainwright, "Sharp thresholds for high-dimensional and noisy sparsity recovery using \(\ell\_1\)-constrained quadratic programming (Lasso)," IEEE Trans. Inf. Theory, vol. 55, no. 5, pp. 2183–2202, May 2009.
- IEEE Trans. Inf. Theory, vol. 55, no. 5, pp. 2183–2202, May 2009.

  [25] W. Su, M. Bogdan, and E. J. Candès, "Supplement to 'False discoveries occur early on the Lasso path," Ann. Statist., 2017. Accessed: Feb. 12, 2021. [Online]. Available: http://www-stat.wharton.upenn.edu/~suw/paper/LassoFDR\_supp.pdf
- [26] M. Bogdan, E. van den Berg, W. Su, and E. J. Candès, "Supplementary material to 'Statistical estimation and testing via the sorted ℓ<sub>1</sub> norm," 2013. [Online]. Available: https://statweb.stanford.edu/~candes/publications/downloads/SortedL1\_SM.pdf
- [27] The International HapMap Consortium, "A second generation human haplotype map of over 3.1 million SNPs," *Nature*, vol. 449, pp. 851–862, Oct. 2007.
- [28] M. F. F. Bogdan, F. Frommlet, P. Szulc, and H. Tang, "Model selection approach for genome wide association studies in admixed populations," Stanford Univ., Stanford, CA, USA, Tech. Rep., 2013.
- [29] M. Bogdan, E. van den Berg, W. Su, and E. J. Candès, "Statistical estimation and testing via the sorted  $\ell_1$  norm," 2013, arXiv:1310.1969.
- [30] M. Bogdan, E. van den Berg, C. Sabatti, W. Su, and E. J. Candès, "SLOPE—Adaptive variable selection via convex optimization," *Ann. Appl. Statist.*, vol. 9, no. 3, pp. 1103–1140, 2015.
- Appl. Statist., vol. 9, no. 3, pp. 1103–1140, 2015.
  [31] W. Su and E. J. Candès, "SLOPE is adaptive to unknown sparsity and asymptotically minimax," Ann. Statist., vol. 44, no. 3, pp. 1038–1068, Jun. 2016.
- [32] Z. Bu, J. Klusowski, C. Rush, and W. J. Su, "Characterizing the SLOPE trade-off: A variational perspective and the Donoho–Tanner limit," 2021, arXiv:2105.13302.
- [33] V. Ročková and E. I. George, "The spike-and-slab LASSO," J. Amer. Stat. Assoc., vol. 113, no. 521, pp. 431–444, Jan. 2018.
- Stat. Assoc., vol. 113, no. 521, pp. 431–444, Jan. 2018.
  [34] A. Gelman and F. Tuerlinckx, "Type S error rates for classical and Bayesian single and multiple comparison procedures," Comput. Statist., vol. 15, no. 3, pp. 373–390, Sep. 2000.
- [35] R. F. Barber and E. J. Candès, "A knockoff filter for high-dimensional selective inference," Ann. Statist., vol. 47, no. 5, pp. 2504–2537, Oct. 2019.
- [36] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices," *Construct. Approx.*, vol. 28, no. 3, pp. 253–263, Dec. 2008.
- [37] T. W. Anderson, An Introduction to Multivariate Statistical Analysis. New York, NY, USA: Wiley, 1962.
- [38] V. Siskind, "Second moments of inverse Wishart-matrix elements," *Biometrika*, vol. 59, no. 3, pp. 690–691, 1972.

**Hua Wang** received the B.S. degree in mathematics and applied mathematics from Fudan University in 2018. He is currently pursuing the Ph.D. degree with the Wharton Statistics and Data Science Department, University of Pennsylvania. His research interests include span high-dimensional statistics, privacy-presevering machine learning, and deep learning theory.

Yachong (Elsa) Yang received the B.S. degree in statistics from the School of Gifted Young, University of Science and Technology of China, in 2019. She is currently pursuing the Ph.D. degree in statistics and data science with the Wharton School, University of Pennsylvania.

Weijie J. Su (Member, IEEE) received the bachelor's degree from Peking University in 2011 and the Ph.D. degree from Stanford University in 2016. He is currently an Assistant Professor with the Wharton Statistics and Data Science Department and the Department of Computer and Information Science, University of Pennsylvania. He is the Co-Director of Penn Research in Machine Learning (https://priml.upenn.edu/). His research interests span privacy-preserving data analysis, optimization, high-dimensional statistics, and deep learning theory. He was a recipient of the Stanford Theodore W. Anderson Dissertation Award in 2016, NSF CAREER Award in 2019, the Alfred P. Sloan Research Fellowship in 2020, and the Society for Industrial and Applied Mathematics (SIAM) Early Career Prize in Data Science in 2022.