

Algorithmic Analysis and Statistical Estimation of SLOPE via Approximate Message Passing

Zhiqi Bu*

Jason Klusowski[†]Cynthia Rush[‡]Weijie Su[§]

July 18, 2019

Abstract

SLOPE is a relatively new convex optimization procedure for high-dimensional linear regression via the sorted ℓ_1 penalty: the larger the rank of the fitted coefficient, the larger the penalty. This non-separable penalty renders many existing techniques invalid or inconclusive in analyzing the SLOPE solution. In this paper, we develop an asymptotically exact characterization of the SLOPE solution under Gaussian random designs through solving the SLOPE problem using approximate message passing (AMP). This algorithmic approach allows us to approximate the SLOPE solution via the much more amenable AMP iterates. Explicitly, we characterize the asymptotic dynamics of the AMP iterates relying on a recently developed state evolution analysis for non-separable penalties, thereby overcoming the difficulty caused by the sorted ℓ_1 penalty. Moreover, we prove that the AMP iterates converge to the SLOPE solution in an asymptotic sense, and numerical simulations show that the convergence is surprisingly fast. Our proof rests on a novel technique that specifically leverages the SLOPE problem. In contrast to prior literature, our work not only yields an asymptotically sharp analysis but also offers an algorithmic, flexible, and constructive approach to understanding the SLOPE problem.

1 Introduction

Consider observing linear measurements $\mathbf{y} \in \mathbb{R}^n$ that are modeled by the equation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{w}, \quad (1.1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a known measurement matrix, $\boldsymbol{\beta} \in \mathbb{R}^p$ is an unknown signal, and $\mathbf{w} \in \mathbb{R}^n$ is the measurement noise. Among numerous methods that seek to recover the signal $\boldsymbol{\beta}$ from the observed data, especially in the setting where $\boldsymbol{\beta}$ is sparse and p is larger than n , SLOPE has recently emerged as a useful procedure that allows for estimation and model selection [9]. This method reconstructs the signal by solving the minimization problem

$$\hat{\boldsymbol{\beta}} := \arg \min_{\mathbf{b}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \sum_{i=1}^p \lambda_i |\mathbf{b}|_{(i)}, \quad (1.2)$$

*Department of Applied Mathematics and Computational Science, University of Pennsylvania, Philadelphia, PA 19104, USA. Email: zbu@sas.upenn.edu

[†]Department of Statistics, Rutgers University, New Brunswick, NJ 08854, USA. Email: jason.klusowski@rutgers.edu

[‡]Department of Statistics, Columbia University, New York, NY 10027, USA. Email: cynthia.rush@columbia.edu

[§]Department of Statistics, University of Pennsylvania, Philadelphia, PA 19104, USA. Email: suw@wharton.upenn.edu This work was supported in part by NSF #1217023.

$$J_{\lambda}(\mathbf{b}) = \sum \lambda_i |\mathbf{b}|_i$$

where $\|\cdot\|$ denotes the ℓ_2 norm, $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ (with at least one strict inequality) is a sequence of thresholds, and $|\mathbf{b}|_{(1)} \geq \dots \geq |\mathbf{b}|_{(p)}$ are the order statistics of the fitted coefficients in absolute value. The regularizer $\sum \lambda_i |\mathbf{b}|_{(i)}$ is a *sorted ℓ_1 -norm* (denoted as $J_{\lambda}(\mathbf{b})$ henceforth), which is *non-separable* due to the sorting operation involved in its calculation. Notably, SLOPE has two attractive features that are not simultaneously present in other methods for linear regression including the LASSO [38] and knockoffs [2]. Explicitly, on the estimation side, SLOPE achieves minimax estimation properties under certain random designs *without* requiring any knowledge of the sparsity degree of β [37, 7]. On the testing side, SLOPE controls the false discovery rate in the case of independent predictors [9, 11]. For completeness, we remark that [10, 39, 19] proposed similar non-separable regularizers to encourage grouping of correlated predictors.

This work is concerned with the algorithmic aspects of SLOPE through the lens of *approximate message passing* (AMP) [4, 16, 23, 31]. AMP is a class of computationally efficient and easy-to-implement algorithms for a broad range of statistical estimation problems, including compressed sensing and the LASSO [5]. When applied to SLOPE, AMP takes the following form: at initial iteration $t = 0$, assign $\beta^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{y}$, and for $t \geq 0$,

$$\beta^{t+1} = \text{prox}_{J_{\theta_t}}(\mathbf{X}^\top \mathbf{z}^t + \beta^t), \quad (1.3a)$$

$$\mathbf{z}^{t+1} = \mathbf{y} - \mathbf{X}\beta^{t+1} + \frac{\mathbf{z}^t}{n} [\nabla \text{prox}_{J_{\theta_t}}(\mathbf{X}^\top \mathbf{z}^t + \beta^t)]. \quad (1.3b)$$

The non-increasing sequence θ_t is proportional to $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_p)$ and will be given explicitly in Section 2. Here, $\text{prox}_{J_{\theta}}$ is the proximal operator of the sorted ℓ_1 norm, that is,

$$\text{prox}_{J_{\theta}}(\mathbf{x}) := \underset{\mathbf{b}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{b}\|^2 + J_{\theta}(\mathbf{b}), \quad (1.4)$$

and $\nabla \text{prox}_{J_{\theta}}$ denotes the divergence of the proximal operator (see an equivalent, but more explicit form, of this algorithm in Section 2 and further discussion of SLOPE and the prox operator in Section 5.1). Compared to the proximal gradient descent (ISTA) [13, 14, 29], AMP has an extra correction term in its residual step that adjusts the iteration in a non-trivial way and seeks to provide improved convergence performance [16].

The *empirical* performance of AMP in solving SLOPE under i.i.d. Gaussian matrix \mathbf{X} is illustrated in Figure 1 and Table 1, which suggest the superiority of AMP over ISTA and FISTA [6]—perhaps the two most popular proximal gradient descent methods—in terms of speed of convergence in this setting. However, the vast AMP literature thus far remains silent on whether AMP *provably* solves SLOPE and, if so, whether one can leverage AMP to get insights into the statistical properties of SLOPE. This vacuum in the literature is due to the *non-separability* of the SLOPE regularizer, making it a major challenge to apply AMP to SLOPE directly. In stark contrast, AMP theory has been rigorously applied to the LASSO [5], showing both good empirical performance and nice theoretical properties of solving the LASSO using AMP. Moreover, AMP in this setting allows for asymptotically exact statistical characterization of its output, which converges to the LASSO solution, thereby providing a powerful tool in fine-grained analyses of the LASSO [3, 36, 28, 35].

Main contributions. In this work, we prove that the AMP algorithm (1.3) solves the SLOPE problem in an asymptotically *exact* sense under independent Gaussian random designs. Our proof uses the recently extended AMP theory for non-separable denoisers [8] and applies this tool to derive the state evolution that describes the asymptotically exact behaviors of the AMP iterates β^t in (1.3). The next step, which is the core of our proof, is to relate the AMP estimates to the

Maybe
compare
with
FISTA

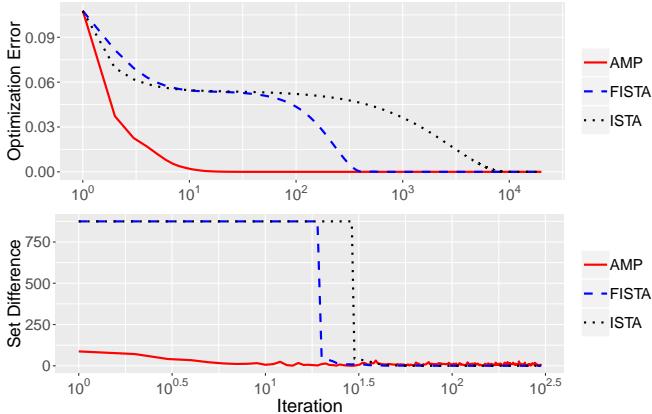


Figure 1: Optimization errors, $\|\beta^t - \hat{\beta}\|^2/p$, and (symmetric) set difference of $\text{supp}(\beta^t)$ and $\text{supp}(\hat{\beta})$.

	Set Diff	Optimization errors				
		10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}
ISTA	60	4048	7326	8569	9007	9161
FISTA	47	275	374	412	593	604
AMP	30	6	13	22	32	40

Table 1: First iteration t for which there is zero set difference or optimization error $\|\beta^t - \hat{\beta}\|^2/p$ falls below a threshold.

Figure 1 and Table 1 Details: Design X is 500×1000 with i.i.d. $\mathcal{N}(0, 1/500)$ entries. True signal β is i.i.d. Gaussian-Bernoulli: $\mathcal{N}(0, 1)$ with probability 0.1 and 0 otherwise. Noise variance $\sigma_w^2 = 0$. A careful calibration between the thresholds θ_t in AMP and λ is SLOPE is used (details in Sec. 2).

SLOPE solution. This presents several challenges that *cannot* be resolved only within the AMP framework. In particular, unlike the LASSO, the number of nonzeros in the SLOPE solution can exceed the number of observations. This fact imposes substantially more difficulties on showing that the distance between the SLOPE solution and the AMP iterates goes to zero than in the LASSO case due to the possible *non-strong convexity* of the SLOPE problem, even restricted to the solution support. To overcome these challenges, we develop novel techniques that are tailored to the characteristics of the SLOPE solution. For example, our proof relies on the crucial property of SLOPE that the *unique* nonzero components of its solution never outnumber the observation units.

As a byproduct, our analysis gives rise to an *exact* asymptotic characterization of the SLOPE solution under independent Gaussian random designs through leveraging the statistical aspect of the AMP theory. In more detail, the probability distribution of the SLOPE solution is completely specified by a few parameters that are the solution to a certain fixed-point equation in an asymptotic sense. This provides a powerful tool for fine-grained statistical analysis of SLOPE as it was for the LASSO problem. We note that a recent paper [20]—which takes an entirely different path—gives an asymptotic characterization of the SLOPE solution that matches our asymptotic analysis deduced from our AMP theory for SLOPE. However, our AMP-based approach is more algorithmic in nature and offers a more concrete connection between the finite-sample behaviors of the SLOPE problem and its asymptotic distribution via the computationally efficient AMP algorithm.

Paper outline. In Section 2 we develop an AMP algorithm for finding the SLOPE estimator in (1.2). Specifically, it is through the threshold values θ_t in the AMP algorithm in (1.3) that one can ensure the AMP estimates converge to the SLOPE estimator with parameter λ , so in Section 2 we provide details for how one should calibrate the thresholds of the AMP iterations in (1.3) in order for the algorithm to solve SLOPE cost in (1.2). Then in Section 3, we state theoretical guarantees showing that the AMP algorithm solves the SLOPE optimization asymptotically and we leverage theoretical guarantees for the AMP algorithm to exactly characterize the mean square error (more generally, any pseudo-Lipschitz error) of the SLOPE estimator in the large system limit. This is done by applying recent theoretical results for AMP algorithms that use a non-separable

non-linearity [8], like the one in (1.3). Finally, Sections 4-7 prove rigorously the theoretical results stated in Section 3 and we end with a discussion in Section 8.

2 Algorithmic Development

To begin with, we state assumptions under which our theoretical results will hold and give some preliminary ideas about SLOPE that will be useful in the development of the AMP algorithm.

Assumptions. Concerning the linear model (1.1) and parameter vector in (1.2), we assume:

- (A1) The measurement matrix \mathbf{X} has independent and identically-distributed (i.i.d.) Gaussian entries that have mean 0 and variance $1/n$.
- (A2) The signal β has elements that are i.i.d. B , with $\mathbb{E}(B^2 \max\{0, \log B\}) < \infty$.
- (A3) The noise \mathbf{w} is elementwise i.i.d. W , with $\sigma_w^2 := \mathbb{E}(W^2) < \infty$.
- (A4) The vector $\lambda(p) = (\lambda_1, \dots, \lambda_p)$ is elementwise i.i.d. Λ , with $\mathbb{E}(\Lambda^2) < \infty$ and $\min\{\lambda(p)\} > 0$.
- (A5) The ratio n/p approaches a constant $\delta \in (0, \infty)$ in the large system limit, as $n, p \rightarrow \infty$.

Remark: (A4) can be relaxed as $\lambda_1, \dots, \lambda_p$ having an empirical distribution that converges weakly to probability measure Λ on \mathbb{R} with $\mathbb{E}(\Lambda^2) < \infty$ and $\|\lambda(p)\|^2/p \rightarrow \mathbb{E}(\Lambda^2)$ and $\min\{\lambda(p)\} > 0$. A similar relaxation can be made for the distributional assumptions (A2) and (A3).

SLOPE preliminaries. For a vector $\mathbf{v} \in \mathbb{R}^p$, the divergence of the proximal operator, $\nabla \text{prox}_f(\mathbf{v})$, is given by the following:

$$\nabla \text{prox}_f(\mathbf{v}) := \sum_{i=1}^p \frac{\partial}{\partial v_i} [\text{prox}_f(\mathbf{v})]_i = \left(\frac{\partial}{\partial v_1}, \frac{\partial}{\partial v_2}, \dots, \frac{\partial}{\partial v_p} \right) \cdot \text{prox}_f(\mathbf{v}), \quad (2.1)$$

where [37, proof of Fact 3.4],

$$\frac{\partial [\text{prox}_{J_\lambda}(\mathbf{v})]_i}{\partial v_j} = \begin{cases} \frac{\text{sign}([\text{prox}_{J_\lambda}(\mathbf{v})]_i) \cdot \text{sign}([\text{prox}_{J_\lambda}(\mathbf{v})]_j)}{\#\{1 \leq k \leq p : |[\text{prox}_{J_\lambda}(\mathbf{v})]_k| = |[\text{prox}_{J_\lambda}(\mathbf{v})]_j|\}}, & \text{if } |[\text{prox}_{J_\lambda}(\mathbf{v})]_j| = |[\text{prox}_{J_\lambda}(\mathbf{v})]_i|, \\ 0, & \text{otherwise.} \end{cases} \quad (2.2)$$

Hence the divergence takes the simplified form

$$\nabla \text{prox}_{J_\lambda}(\mathbf{v}) = \|\text{prox}_{J_\lambda}(\mathbf{v})\|_0^*, \quad (2.3)$$

where $\|\cdot\|_0^*$ counts the unique non-zero magnitudes in a vector, e.g. $\|(0, 1, -2, 0, 2)\|_0^* = 2$. This explicit form of divergence not only waives the need to use approximation in calculation but also speed up the recursion, since it only depends on the proximal operator as a whole instead of on $\theta_{t-1}, \mathbf{X}, \mathbf{z}^{t-1}, \beta^{t-1}$. Therefore, we have

Lemma 2.1. *In AMP, (1.3b) is equivalent to $\mathbf{z}^{t+1} = \mathbf{y} - \mathbf{X}\beta^{t+1} + \frac{\tilde{z}^t}{\delta p} \|\beta^{t+1}\|_0^*$.*

Other details and background on SLOPE and the prox operator are found in Section 5.1. Now we discuss the details of an AMP algorithm that can be used for finding the SLOPE estimator in (1.2).

2.1 AMP Background

An attractive feature of AMP is that its statistical properties can be exactly characterized at each iteration t , at least asymptotically, via a one-dimensional recursion known as state evolution [4, 8, 35, 21]. Specifically, it can be shown that the pseudo-data, meaning the input $\mathbf{X}^\top \mathbf{z}^t + \boldsymbol{\beta}^t$ for the estimate of the unknown signal in (1.3a), is asymptotically equal in distribution to the true signal plus independent, Gaussian noise, i.e. $\boldsymbol{\beta} + \tau_t \mathbf{Z}$, where the noise variance τ_t is defined by the state evolution. For this reason, the function used to update the estimate in (1.3a), in our case, the proximal operator, $\text{prox}_{J_{\theta_t}}(\cdot)$, is usually referred to as a ‘denoiser’ in the AMP literature.

This statistical characterization of the pseudo-data was first rigorously shown to be true in the case of ‘separable’ denoisers by Bayati and Montanari [4], and an analysis of the rate of this convergence was given in [35]. A ‘separable’ denoiser is one that applies the same (possibly non-linear) function to each element of its input. Recent work [8] proves that the pseudo-data has distribution $\boldsymbol{\beta} + \tau_t \mathbf{Z}$ asymptotically, even when the ‘denoisers’ used in the AMP algorithm are non-separable, like the SLOPE prox operator in (1.3a).

As mentioned previously, the dynamics of the AMP iterations are tracked by a recursive sequence referred to as the state evolution, defined as follows. For \mathbf{B} elementwise i.i.d. B independent of $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$, let $\tau_0^2 = \sigma_w^2 + \mathbb{E}[B^2]/\delta$ and for $t \geq 0$,

$$\tau_{t+1}^2 = \sigma_w^2 + \lim_p \frac{1}{\delta p} \mathbb{E} \|\text{prox}_{J_{\theta_t}}(\mathbf{B} + \tau_t \mathbf{Z}) - \mathbf{B}\|^2. \quad (2.4)$$

$\theta_t = \alpha z_t$ is
a stepsize
of the proxy
operator

$\frac{\Delta}{p} \rightarrow \sigma$
 $w \sim \mathcal{N}(0, \sigma_w^2)$

Below we make rigorous the way that the recursion in (2.4) relates to the AMP iteration (1.3).

We note that throughout, we let $\mathcal{N}(\mu, \sigma^2)$ denote the Gaussian density with mean μ and variance σ^2 and we use \mathbb{I}_p to indicate a $p \times p$ identity matrix.

2.2 Analysis of the AMP State Evolution

As the state evolution (2.4) predicts the performance of the AMP algorithm (1.3) (the pseudo-data, $\mathbf{X}^\top \mathbf{z}^t + \boldsymbol{\beta}^t$, is asymptotically equal in distribution $\boldsymbol{\beta} + \tau_t \mathbf{Z}$), it is of interest to study the large t asymptotics of (2.4). Moreover, recall that through the sequence of thresholds $\boldsymbol{\theta}_t$, one can relate the AMP algorithm to the SLOPE estimator in (1.2) for a specific $\boldsymbol{\lambda}$, and the explicit form of this calibration, given in Section 2.3, is motivated by such asymptotic analysis of the state evolution.

It turns out that a finite-size approximation, which we denote $\tau_t^2(p)$, will be easier to analyze than (2.4). The definition of $\tau_{t+1}^2(p)$ is stated explicitly in (2.5) below. Throughout the work, we will define thresholds $\boldsymbol{\theta}_t := \boldsymbol{\alpha} \tau_t(p)$ for every iteration t where the vector $\boldsymbol{\alpha}$ is fixed via a calibration made explicit in Section 2.3. We can interpret this to mean that within the AMP algorithm, $\boldsymbol{\alpha}$ plays the role of the regularizer parameter, $\boldsymbol{\lambda}$. Now we define $\tau_{t+1}^2(p)$, for large p , as a finite-sample approximation to (2.4), namely

$$\tau_{t+1}^2(p) = \sigma_w^2 + \frac{1}{\delta p} \mathbb{E} \|\text{prox}_{J_{\alpha \tau_t(p)}}(\boldsymbol{\beta} + \tau_t(p) \mathbf{Z}) - \boldsymbol{\beta}\|^2, \quad (2.5)$$

where the difference between (2.5) and the state evolution (2.4) is via the large system limit in p . When we refer to the recursion in (2.5), we will always specify the p dependence explicitly as $\tau_t(p)$. An analysis of the limiting properties (in t) of (2.5) is given in Theorem 1 below, after which it is then argued that because interchanging limits and differentiation is justified, the large t analysis of (2.5) holds for (2.4) as well. Before presenting Theorem 1, however, we give the following result which motivates why the AMP iteration should relate at all to the SLOPE estimator.

$$\text{prox}_{\text{ref}}(\mathbf{x}) = \arg\min \left\{ \frac{1}{2} \|\mathbf{y}\|^2 + \frac{1}{2\delta} \|\mathbf{x} - \mathbf{y}\|^2 \right\} = \hat{\mathbf{y}} \Rightarrow \mathbf{o} \in \partial \text{ref}(\hat{\mathbf{y}}) - (\mathbf{x} - \hat{\mathbf{y}}) \Rightarrow \mathbf{x} \in \hat{\mathbf{y}} + \partial \text{ref}(\mathbf{y})$$

$$\text{prox}_h(\beta + v) = \beta \Rightarrow \cancel{\beta + v} \in \cancel{\beta} + \partial h(\beta)$$

Lemma 2.2. Any stationary point $\hat{\boldsymbol{\beta}}$ (with corresponding $\hat{\mathbf{z}}$) in the AMP algorithm (1.3a)-(1.3b) with $\boldsymbol{\theta}_* = \alpha\tau_*$ is a minimizer of the SLOPE cost function in (1.2) with

$$\boldsymbol{\lambda} = \boldsymbol{\theta}_* \left(1 - \frac{1}{\delta p} (\nabla \text{prox}_{J_{\boldsymbol{\theta}_*}}(\hat{\boldsymbol{\beta}} + \mathbf{X}^\top \hat{\mathbf{z}})) \right) = \boldsymbol{\theta}_* \left(1 - \frac{1}{n} \|\text{prox}_{J_{\boldsymbol{\theta}_*}}(\hat{\boldsymbol{\beta}} + \mathbf{X}^\top \hat{\mathbf{z}})\|_0^* \right).$$

Proof of Lemma 2.2. Denote, $\omega := (\nabla \text{prox}_{J_{\boldsymbol{\theta}_*}}(\hat{\boldsymbol{\beta}} + \mathbf{X}^\top \hat{\mathbf{z}})) / (\delta p)$. Now, by stationarity,

$$\hat{\boldsymbol{\beta}} = \text{prox}_{J_{\boldsymbol{\theta}_*}}(\hat{\boldsymbol{\beta}} + \mathbf{X}^\top \hat{\mathbf{z}}), \quad \text{and} \quad \hat{\mathbf{z}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \frac{\hat{\mathbf{z}}}{\delta p} (\nabla \text{prox}_{J_{\boldsymbol{\theta}_*}}(\hat{\boldsymbol{\beta}} + \mathbf{X}^\top \hat{\mathbf{z}})). \quad (2.6)$$

From (2.6), notice that $\hat{\mathbf{z}} = \frac{\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}}{1-\omega}$. By Fact 5.2, $\mathbf{X}^\top \hat{\mathbf{z}} \in \partial J_{\boldsymbol{\theta}_*}(\hat{\boldsymbol{\beta}})$, where $\partial J_{\boldsymbol{\theta}_*}(\hat{\boldsymbol{\beta}})$ is the subgradient of $J_{\boldsymbol{\theta}_*}(\cdot)$ at $\hat{\boldsymbol{\beta}}$ (a precise definition of a subgradient is given in Section 5.1). Then, $\mathbf{X}^\top \hat{\mathbf{z}} = \frac{\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{1-\omega} \in \partial J_{\boldsymbol{\theta}_*}(\hat{\boldsymbol{\beta}})$, and therefore $\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \in \partial J_{\boldsymbol{\theta}_*(1-\omega)}(\hat{\boldsymbol{\beta}})$ which is *exactly* the stationary condition of SLOPE with regularization parameter $\boldsymbol{\lambda} = (1 - \omega)\boldsymbol{\theta}_*$, as desired. \square

Now we present Theorem 1, which provides results about the t asymptotics of the recursion in (2.5) and its proof is given in Appendix A. First, some notation must be introduced: let $\mathbf{A}_{\min}(\delta)$ be the set of solutions to

$$\delta = f(\boldsymbol{\alpha}), \quad \text{where} \quad f(\boldsymbol{\alpha}) := \frac{1}{p} \sum_{i=1}^p \mathbb{E} \left\{ \left(1 - |[\text{prox}_{J_{\boldsymbol{\alpha}}}(\mathbf{Z})]_i| \sum_{j \in I_i} \alpha_j \right) / [\mathbf{D}(\text{prox}_{J_{\boldsymbol{\alpha}}}(\mathbf{Z}))]_i \right\}. \quad (2.7)$$

Here \odot represents elementwise multiplication of vectors and for vector $\mathbf{v} \in \mathbb{R}^p$, \mathbf{D} is defined elementwise as $[\mathbf{D}(\mathbf{v})]_i = \#\{j : |v_j| = |v_i|\}$ if $v_i \neq 0$ and ∞ otherwise. Let $I_i = \{j : 1 \leq j \leq p \text{ and } |[\text{prox}_{J_{\boldsymbol{\alpha}}}(\mathbf{Z})]_j| = |[\text{prox}_{J_{\boldsymbol{\alpha}}}(\mathbf{Z})]_i|\}$. The expectation in (2.7) is taken with respect to \mathbf{Z} , a p -length vector of i.i.d. standard Gaussians. Finally, for $\mathbf{u} \in \mathbb{R}^m$, the notation $\langle \mathbf{u} \rangle := \sum_{i=1}^m u_i/m$ and we say \mathbf{u} is strictly larger than $\mathbf{v} \in \mathbb{R}^m$ if $u_i > v_i$ for all elements $i \in \{1, 2, \dots, m\}$. For the simple case of $p = 2$, we illustrate an example of the set $\mathbf{A}_{\min}(\delta)$ in Figure 2.

Theorem 1. For any $\boldsymbol{\alpha}$ strictly larger than at least one element in the set $\mathbf{A}_{\min}(\delta)$, the recursion in (2.5) has a unique fixed point that we denote as $\tau_*^2(p)$. Then $\tau_t(p) \rightarrow \tau_*(p)$ monotonically for any initial condition. Define a function $F : \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}$ as

$$F(\tau^2(p), \boldsymbol{\alpha}\tau(p)) := \sigma_w^2 + \frac{1}{\delta p} \mathbb{E} \|\text{prox}_{J_{\boldsymbol{\alpha}\tau(p)}}(\mathbf{B} + \tau(p)\mathbf{Z}) - \mathbf{B}\|^2, \quad (2.8)$$

where \mathbf{B} is elementwise i.i.d. B independent of $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$, so that $\tau_{t+1}^2(p) = F(\tau_t^2(p), \boldsymbol{\alpha}\tau_t(p))$. Then $|\frac{\partial F}{\partial \tau^2(p)}(\tau^2(p), \boldsymbol{\alpha}\tau(p))| < 1$ at $\tau(p) = \tau_*(p)$. Moreover, for $f(\boldsymbol{\alpha})$ defined in (2.7), we show that $f(\boldsymbol{\alpha}) = \delta \lim_{\tau(p) \rightarrow \infty} dF/d\tau^2(p)$.

Beyond providing the large t asymptotics of the state evolution sequence, notice that Theorem 1 gives necessary conditions on the calibration vector $\boldsymbol{\alpha}$ under which the recursion in (2.5), and equivalently, the calibration detailed in Section 2.3 below are well-defined.

Recall that it is actually the state evolution in (2.4) (and not that in (2.5)) that predicts the performance of the AMP algorithm, and therefore we would really like a version of Theorem 1 studying the large system limit in p . We argue that because interchanging differentiation and the limit, the proof of Theorem 1 analyzing (2.5), can easily be used to give an analogous result for

(2.4). In particular analyzing (2.4) via the strategy given in the proof of Theorem 1 requires that we study the partial derivative of $\lim_p \mathbb{E}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2/(\delta p)$, with respect to τ^2 . Indeed, to directly make use our proof for the finite- p case given in Theorem 1, it is enough that

$$\frac{\partial}{\partial\tau^2} \lim_p \mathbb{E}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2/(\delta p) = \lim_p \frac{\partial}{\partial\tau^2} \mathbb{E}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2/(\delta p). \quad (2.9)$$

Note that we already have an argument (based on dominated convergence for fixed p , see (A.1) and Lemma A.1) showing that

$$\frac{\partial}{\partial\tau^2} \mathbb{E}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2 = \mathbb{E}\left\{\frac{\partial}{\partial\tau^2}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2\right\}.$$

The next lemma gives us a roadmap for how to proceed (c.f., [34, Theorem 7.17]) to justify the interchange in (2.9).

Lemma 2.3. *Suppose $\{g_m\}$ is a sequence of functions that converge pointwise to g on a compact domain D and whose derivatives $\{g'_m\}$ converge uniformly to a function h on D . Then $h = g'$ on D .*

Therefore, taking $\{g_p\} = \{\mathsf{F}(\tau^2(p), \boldsymbol{\alpha}\tau(p))\}$, it suffices to show that if

$$\frac{\partial\mathsf{F}}{\partial\tau^2}(\tau^2(p), \boldsymbol{\alpha}\tau(p)) = \frac{\partial}{\partial\tau^2(p)} \mathbb{E}\|\text{prox}_{J_{\alpha\tau(p)}}(\mathbf{B} + \tau(p)\mathbf{Z}) - \mathbf{B}\|^2/(\delta p),$$

then the sequence $\{\frac{\partial\mathsf{F}}{\partial\tau^2}(\tau^2, \boldsymbol{\alpha}\tau)\}_p$ converges uniformly as $p \rightarrow \infty$. The main tool for proving such a result is given in the following lemma.

Lemma 2.4. *Suppose $\{g_m\}$ is a sequence of L -Lipschitz functions (where L is independent of m) that converge pointwise to a function g on a compact domain D . Then, the convergence is also uniform on D .*

Using this lemma, the essential idea is to show that there exists a constant $L > 0$, independent of p , such that for all p and all τ_1, τ_2 in a bounded set $D = \{\tau : 0 < r \leq |\tau| \leq R\}$,

$$\left|\frac{\partial\mathsf{F}}{\partial\tau^2}(\tau_1^2, \boldsymbol{\alpha}\tau_1) - \frac{\partial\mathsf{F}}{\partial\tau^2}(\tau_2^2, \boldsymbol{\alpha}\tau_2)\right| \leq L|\tau_1 - \tau_2|.$$

This follows by the mean value theorem and (A.14), with $L = \sup_{p, \tau \in D} \left|\frac{\partial}{\partial\tau^2} \frac{\partial\mathsf{F}}{\partial\tau^2}(\tau^2, \boldsymbol{\alpha}\tau)\right| < +\infty$.

Remark 2.5. *The boundedness of $\{\tau_t(p)\}$ is guaranteed by Proposition 2.6. In particular, since $\boldsymbol{\alpha}$ satisfies the assumption of Theorem 1, Proposition 2.6 guarantees $\boldsymbol{\lambda}$ is bounded and, consequently, so is τ (see the calibration in (2.10) below).*

2.3 Threshold Calibration

Motivated by Lemma 2.2 and the result of Theorem 1, we define a calibration from the regularization parameter $\boldsymbol{\lambda}$, to the corresponding threshold $\boldsymbol{\alpha}$ used to define the AMP algorithm. In practice, we will be given finite-length $\boldsymbol{\lambda}$ and then we want to design the AMP iteration to solve the corresponding SLOPE cost. We do this by choosing $\boldsymbol{\alpha}$ as the vector that solves $\boldsymbol{\lambda} = \boldsymbol{\lambda}(\boldsymbol{\alpha})$ where

$$\boldsymbol{\lambda}(\boldsymbol{\alpha}) := \boldsymbol{\alpha}\tau_*(p)\left(1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{\alpha\tau_*(p)}}(\mathbf{B} + \tau_*(p)\mathbf{Z})\|_0^*\right), \quad (2.10)$$

where \mathbf{B} is elementwise i.i.d. B independent of $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$ and $\tau_*(p)$ is the limiting value defined in Theorem 1. We note the fact that the calibration in (2.10) sets $\boldsymbol{\alpha}$ as a vector *in the same direction* as $\boldsymbol{\lambda}$, but that is scaled by a constant value (for each p), where the scaling constant value is $\tau_*(p)(1 - \mathbb{E} \|\text{prox}_{J_{\boldsymbol{\alpha}\tau_*(p)}}(\mathbf{B} + \tau_*(p)\mathbf{Z})\|_0^*/n)$.

In Proposition 2.6 we show that the calibration (2.10) and its inverse $\boldsymbol{\lambda} \mapsto \boldsymbol{\alpha}(\boldsymbol{\lambda})$ are well-defined and in Algorithm 1 we show that determining the calibration is straightforward in practice.

i.e. we can get $\boldsymbol{\alpha}$ from $\boldsymbol{\lambda}$

Proposition 2.6. *The function $\boldsymbol{\alpha} \mapsto \boldsymbol{\lambda}(\boldsymbol{\alpha})$ defined in (2.10) is continuous on $\{\boldsymbol{\alpha} : f(\boldsymbol{\alpha}) < \delta\}$ for $f(\cdot)$ defined in (2.7) with $\boldsymbol{\lambda}(\mathbf{A}_{\min}) = -\infty$ and $\lim_{\boldsymbol{\alpha} \rightarrow \infty} \boldsymbol{\lambda}(\boldsymbol{\alpha}) = \infty$ (where the limit is taken elementwise). Therefore the function $\boldsymbol{\lambda} \mapsto \boldsymbol{\alpha}(\boldsymbol{\lambda})$ satisfying (2.10) exists. As $p \rightarrow \infty$, the function $\boldsymbol{\alpha} \mapsto \boldsymbol{\lambda}(\boldsymbol{\alpha})$ becomes invertible (given $\boldsymbol{\lambda}, \boldsymbol{\alpha}$ satisfying (2.10) exists uniquely). Furthermore, the inverse function is continuous non-decreasing for any $\boldsymbol{\lambda} > \mathbf{0}$.*

In [5, Proposition 1.4 (first introduced in [17]) and Corollary 1.7] this is proven rigorously for the analogous LASSO calibration and in Appendix A we show how to adapt this proof to SLOPE case. This proposition motivates Algorithm 1 which uses a bisection method to find the unique $\boldsymbol{\alpha}$ for each $\boldsymbol{\lambda}$. It suffices to find two guesses of $\boldsymbol{\alpha}$ parallel to $\boldsymbol{\lambda}$ that, when mapped via (2.10), sandwich the true $\boldsymbol{\lambda}$.

Algorithm 1 Calibration from $\boldsymbol{\lambda} \rightarrow \boldsymbol{\alpha}$

1. Initialize $\alpha_1 = \alpha_{\min}$ such that $\alpha_{\min}\ell \in \mathbf{A}_{\min}$, where $\ell := \boldsymbol{\lambda}/\lambda_1$; Initialize $\alpha_2 = 2\alpha_1$
- while** $L(\alpha_2) < 0$ where $L : \mathbb{R} \rightarrow \mathbb{R}; \alpha \mapsto \text{sign}(\boldsymbol{\lambda}(\alpha\ell) - \boldsymbol{\lambda})$ **do**

 - 2. Set $\alpha_1 = \alpha_2, \alpha_2 = 2\alpha_2$

- end while**
3. **return BISECTION** ($L(\alpha), \alpha_1, \alpha_2$)

Remark: $\text{sign}(\boldsymbol{\lambda}(\cdot) - \boldsymbol{\lambda}) \in \mathbb{R}$ is well-defined since $\boldsymbol{\lambda}(\cdot) \parallel \boldsymbol{\lambda}$ implies all entries share the same sign. The function “**BISECTION**(L, a, b)” finds the root of L in $[a, b]$ via the bisection method.

The calibration in (2.10) is exact when $p \rightarrow \infty$, so we study the mapping between $\boldsymbol{\alpha}$ and $\boldsymbol{\lambda}$ in this limit. Recall from **(A4)**, that the sequence of vectors $\{\boldsymbol{\lambda}(p)\}_{p \geq 0}$ are drawn i.i.d. from distribution Λ . It follows that the sequence $\{\boldsymbol{\alpha}(p)\}_{p \geq 0}$ defined for each p by the finite-sample calibration (2.10) are i.i.d. from a distribution A , where A satisfies $\mathbb{E}(A^2) < \infty$, and is defined via

$$\Lambda = A\tau_*\left(1 - \lim_p \frac{1}{\delta p} \mathbb{E} \|\text{prox}_{J_{\boldsymbol{\alpha}(p)\tau_*}}(\mathbf{B} + \tau_*(p)\mathbf{Z})\|_0^*\right), \quad (2.11)$$

We note, moreover, that the calibrations presented in this section are well-defined:

Fact 2.7. *The limits in (2.4) and (2.11) exist.*

This fact is proven in Appendix C. One idea used in the proof of Fact 2.7 is that the prox operator is *asymptotically separable*, a result shown by [20, Proposition 1]. Specifically, for sequences of input, $\{\mathbf{v}(p)\}$, and thresholds, $\{\boldsymbol{\lambda}(p)\}$, having empirical distributions that weakly converge to distributions

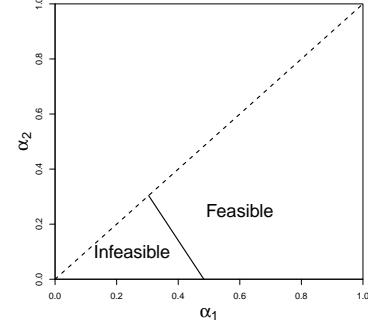


Figure 2: \mathbf{A}_{\min} (black curve) when $p = 2$ and $\delta = 0.6$.

V and Λ , respectively, then there exists a limiting scalar function $h(\cdot) := h(\mathbf{v}(p); V, \Lambda)$ (determined by V and Λ) of the proximal operator $\text{prox}_{J_\lambda}(\mathbf{v}(p))$. Further details are given in Lemma 3.3 in Section 3. Using $h(\cdot) := h(\cdot; B + \tau_* Z, A\tau_*)$, this argument implies that (2.4) can be represented as

$$\tau_*^2 := \sigma_w^2 + \frac{1}{\delta} \mathbb{E}(h(B + \tau_* Z) - B)^2,$$

and if we denote m as the Lebesgue measure, then the limit in (2.11) can be represented as

$$\mathbb{P}\left(B + \tau_* Z \in \left\{x \mid h(x) \neq 0 \quad \text{and} \quad m\{z \mid |h(z)| = |h(x)|\} = 0\right\}\right). \quad (2.12)$$

In other words, the limit in (2.11) is the Lebesgue measure of the domain of the quantile function of h for which the quantile of h assumes unique values (i.e., is not flat).

3 Asymptotic Characterization of SLOPE

3.1 AMP Recovers the SLOPE Estimate

Here we show that the AMP algorithm converges in ℓ_2 to the SLOPE estimator, implying that the AMP iterates can be used as a surrogate for the global optimum of the SLOPE cost function. The schema of the proof is similar to [5, Lemma 3.1], however, major differences lie in the fact that the proximal operator used in the AMP updates (1.3a)-(1.3b) is non-separable. We sketch the proof here, and a forthcoming article will be devoted to giving a complete and detailed argument.

Theorem 2. *Under assumptions (A1) - (A5), for the output of the AMP algorithm in (1.3a) and the SLOPE estimate (1.2),*

$$\plim_{p \rightarrow \infty} \frac{1}{p} \|\hat{\beta} - \beta^t\|^2 = c_t, \quad \text{where} \quad \lim_{t \rightarrow \infty} c_t = 0. \quad (3.1)$$

The proof of Theorem 2 can be found in Section 4. At a high level, the proof requires dealing carefully with the fact that the SLOPE cost function, $\mathcal{C}(\mathbf{b}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + J_\lambda(\mathbf{b})$, given in (1.2) is *not* necessarily strongly convex, meaning that we could encounter the undesirable situation where $\mathcal{C}(\hat{\beta})$ is close to $\mathcal{C}(\beta)$ but $\hat{\beta}$ is not close to β , meaning the statistical recovery of β would be poor.

In the LASSO case, one works around this challenge by showing that the (LASSO) cost function does have nice properties when considering just the elements of the non-zero support of β^t at any (large) iteration t . In the LASSO case, the non-zero support of β has size no larger than $n < p$.

In the SLOPE problem, however, it is possible that the support set has size exceeding n , and therefore the LASSO analysis is not immediately applicable. Our proof develops novel techniques that are tailored to the characteristics of the SLOPE solution. Specifically, when considering the SLOPE problem, one can show nice properties (similar to those in the LASSO case) by considering a support-like set, that being the *unique* non-zeros in the estimate β^t at any (large) iteration t . In other words, if we define an equivalence relation $x \sim y$ when $|x| = |y|$, then entries of AMP estimate at any iteration t are partitioned into equivalence classes. Then we observe from (2.10), and the non-negativity of λ , that the number of equivalence classes is no larger than n . We see an analogy between SLOPE's equivalence class (or 'maximal atom' as described in Appendix 5.1) and LASSO's support set. This approach allows us to deal with the lack of a strongly convex cost.

Theorem 2 ensures that the AMP algorithm solves the SLOPE problem in an asymptotic sense. To better appreciate the convergence guarantee, it calls for elaboration on (3.1). First, it implies that $\|\hat{\beta} - \beta^t\|^2/p$ converges in probability to a constant, say c_t . Next, (3.1) says $c_t \rightarrow 0$ as $t \rightarrow \infty$.

$\lim_{t \rightarrow \infty} c_t = c_p$
 $\lim_{p \rightarrow \infty} c_p = 0$
 would be
 probably
 better ?

3.2 Exact Asymptotic Characterization of the SLOPE Estimate

A consequence of Theorem 4.1, is that the SLOPE estimator $\hat{\beta}$ inherits performance guarantees provided by the AMP state evolution, in the sense of Theorem 3 below. Theorem 3 provides an asymptotic characterization of pseudo-Lipschitz loss between $\hat{\beta}$ and the truth β .

Definition 3.1. Uniformly pseudo-Lipschitz functions [8]: For $k \in \mathbb{N}_{>0}$, a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is pseudo-Lipschitz of order k if there exists a constant L , such that for $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$,

$$\|\phi(\mathbf{a}) - \phi(\mathbf{b})\| \leq L \left(1 + (\|\mathbf{a}\|/\sqrt{d})^{k-1} + (\|\mathbf{b}\|/\sqrt{d})^{k-1} \right) \left(\|\mathbf{a} - \mathbf{b}\|/\sqrt{d} \right). \quad (3.2)$$

A sequence (in p) of pseudo-Lipschitz functions $\{\phi_p\}_{p \in \mathbb{N}_{>0}}$ is uniformly pseudo-Lipschitz of order k if, denoting by L_p the pseudo-Lipschitz constant of ϕ_p , $L_p < \infty$ for each p and $\limsup_{p \rightarrow \infty} L_p < \infty$.

Theorem 3. Under assumptions **(A1)** - **(A5)**, for any uniformly pseudo-Lipschitz sequence of functions $\psi_p : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ and for $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$,

$$\plim_p \psi_p(\hat{\beta}, \beta) = \lim_t \plim_p \mathbb{E}_{\mathbf{Z}}[\psi_p(\text{prox}_{J_{\alpha(p)\tau_t}}(\beta + \tau_t \mathbf{Z}), \beta)],$$

where τ_t is defined in (2.4) and the expectation is taken with respect to \mathbf{Z} .

Theorem 3 tells us that under uniformly pseudo-Lipschitz loss, in the large system limit, distributionally the SLOPE optimizer acts as a ‘denoised’ version of the truth corrupted by additive Gaussian noise where the denoising function is given by the proximal operator, i.e. within uniformly pseudo-Lipschitz loss $\hat{\beta}$ can be replaced with $\text{prox}_{J_{\alpha(p)\tau_t}}(\beta + \tau_t \mathbf{Z})$ for large p, t .

The proof of Theorem 3 can be found in Section 4. We show that Theorem 3 follows from Theorem 2 and recent AMP theory dealing with the state evolution analysis in the case of non-separable denoisers [8], which can be used to demonstrate that the state evolution given in (2.4) characterizes the performance of the SLOPE AMP (1.3) via pseudo-Lipschitz loss functions.

We note that [20, Theorem 1] follows by Theorem 3 and their separability result [20, Proposition 1]. To see this, we use the following lemma that is a simple application of the Law of Large Numbers.

Lemma 3.2. For any function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ that is asymptotically separable, in the sense that there exists some function $\tilde{f} : \mathbb{R} \rightarrow \mathbb{R}$, such that

$$\left| f(\beta) - \frac{1}{p} \sum_{i=1}^n \tilde{f}(\beta_i) \right| \rightarrow 0, \quad \text{as } p \rightarrow \infty,$$

where $\tilde{f}(B)$ is Lebesgue integrable then $\plim_p (f(\beta) - \mathbb{E}_{\mathbf{B}}[\tilde{f}(\mathbf{B})]) = 0$, where $\mathbf{B} \sim \text{i.i.d. } B$.

Now to show the result [20, Theorem 1], consider a special case of Theorem 3 where $\psi_p(\mathbf{x}, \mathbf{y}) = \frac{1}{p} \sum \psi(x_i, y_i)$ for function $\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ that is pseudo-Lipschitz of order $k = 2$. It is easy to show that $\psi_p(\cdot, \cdot)$ is uniformly pseudo-Lipschitz of order $k = 2$. The result of Theorem 3 then says that

$$\plim_p \frac{1}{p} \sum_{i=1}^p \psi(\hat{\beta}_i, \beta_i) = \lim_t \plim_p \frac{1}{p} \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}}[\psi([\text{prox}_{J_{\alpha(p)\tau_t}}(\beta + \tau_t \mathbf{Z})]_i, \beta_i)].$$

Then [20, Theorem 1] follows by [20, Proposition 1], restated below in Lemma 3.3, the Law of Large Numbers, and Theorem 1. Now we restate in Lemma 3.3, the result given in [20, Proposition 1], which says that $\text{prox}_{J_{\alpha\tau_t}}(\cdot)$ becomes asymptotically separable as $p \rightarrow \infty$, for convenience.

Lemma 3.3 (Proposition 1, [20]). *For an input sequence $\{\mathbf{v}(p)\}$, and a sequence of thresholds $\{\boldsymbol{\lambda}(p)\}$, both having empirical distributions that weakly converge to distributions V and Λ , respectively, then there exists a limiting scalar function h (determined by V and Λ) such that as $p \rightarrow \infty$,*

$$\|\text{prox}_{J_{\boldsymbol{\lambda}(p)}}(\mathbf{v}(p)) - h(\mathbf{v}(p); V, \Lambda)\|^2/p \rightarrow 0, \quad (3.3)$$

where h applies $h(\cdot; V, \Lambda)$ coordinate-wise to $\mathbf{v}(p)$ (hence it is separable) and h is Lipschitz(1).

Then [20, Theorem 1] follows from Theorem 3 by using the asymptotic separability of the prox operator. Namely, the result of Lemma 3.3 (using that $\boldsymbol{\alpha}(p)\tau_t$ has an empirical distribution that converges weakly to $A\tau_t$ for A defined by (2.11)), along with Cauchy-Schwarz and the fact that ψ is pseudo-Lipschitz, allow us to apply a dominated convergence argument (see Lemma B.2), from which it follows for some limiting scalar function h^t as specified by Lemma 3.3,

$$\frac{1}{p} \left| \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}}[\psi([\text{prox}_{J_{\boldsymbol{\alpha}(p)\tau_t}}(\boldsymbol{\beta} + \tau_t \mathbf{Z})]_i, \beta_i)] - \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}}[\psi([h^t(\boldsymbol{\beta} + \tau_t \mathbf{Z})]_i, \beta_i)] \right| \rightarrow 0.$$

Then the above allows us to apply Lemma 3.2 and the Law of Large Numbers to show

$$\begin{aligned} \text{plim}_p \frac{1}{p} \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}}[\psi([\text{prox}_{J_{\boldsymbol{\alpha}(p)\tau_t}}(\boldsymbol{\beta} + \tau_t \mathbf{Z})]_i, \beta_i)] &= \lim_p \frac{1}{p} \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}, \mathbf{B}}[\psi(h^t([\mathbf{B} + \tau_t \mathbf{Z}]_i), B_i)] \\ &= \mathbb{E}_{\mathbf{Z}, \mathbf{B}}[\psi(h^t(B + \tau_t Z), B)], \end{aligned}$$

Finally we note that the result of [20, Theorem 1] follows since

$$\lim_t \mathbb{E}_{\mathbf{Z}, \mathbf{B}}[\psi(h^t(B + \tau_t Z), B)] = \mathbb{E}_{\mathbf{Z}, \mathbf{B}}[\psi(h^*(B + \tau_* Z), B)].$$

We highlight that our Theorem 3 allows the consideration of a non-asymptotic case in t . While Theorem 1 motivates an algorithmic way to find a value $\tau_t(p)$ which approximates $\tau_*(p)$ well, Theorem 3 guarantees the accuracy of such approximation for use in practice. One particular use of Theorem 3 is to design the optimal sequence $\boldsymbol{\lambda}$ that achieves the minimum τ_* and equivalently minimum error [20], though a concrete algorithm for doing so is still under investigation.

Finally we show how we use Theorem 3 to study the asymptotic mean-square error between the SLOPE estimator and the truth [12].

Corollary 3.4. *Under assumptions **(A1)** – **(A5)**, $\text{plim}_p \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|^2/p = \delta(\tau_*^2 - \sigma_w^2)$.*

Proof. Applying Theorem 3 to the pseudo-Lipschitz loss function $\psi^1 : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$, defined as $\psi^1(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2/p$, we find $\text{plim}_p \frac{1}{p} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|^2 = \lim_t \text{plim}_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\|\text{prox}_{J_{\boldsymbol{\alpha}\tau_t}}(\boldsymbol{\beta} + \tau_t \mathbf{Z}) - \boldsymbol{\beta}\|^2]$. The desired result follows since $\lim_t \text{plim}_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\|\text{prox}_{J_{\boldsymbol{\alpha}\tau_t}}(\boldsymbol{\beta} + \tau_t \mathbf{Z}) - \boldsymbol{\beta}\|^2] = \delta(\tau_*^2 - \sigma_w^2)$. To see this, note that $\lim_t \delta(\tau_{t+1}^2 - \sigma_w^2) = \delta(\tau_*^2 - \sigma_w^2)$ and

$$\text{plim}_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\|\text{prox}_{J_{\boldsymbol{\alpha}\tau_t}}(\boldsymbol{\beta} + \tau_t \mathbf{Z}) - \boldsymbol{\beta}\|^2] = \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}, \mathbf{B}}[\|\text{prox}_{J_{\boldsymbol{\alpha}\tau_t}}(\mathbf{B} + \tau_t \mathbf{Z}) - \mathbf{B}\|^2] = \delta(\tau_{t+1}^2 - \sigma_w^2),$$

for \mathbf{B} elementwise i.i.d. B independent of $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$. A rigorous argument for the above requires showing that the assumptions of Lemma 3.2 are satisfied and follows similarly to that used to prove property **(P2)** stated in Section 4 and proved in Appendix B. \square

4 Proof for Asymptotic Characterization of the SLOPE Estimate

In this section we prove Theorem 3. To do this, we use a result guaranteeing that the state evolution given in (2.4) characterizes the performance of the SLOPE AMP algorithm (1.3b), given in Lemma 4.1 below. Specifically, Lemma 4.1 relates the state evolution (2.4) to the output of the AMP iteration (1.3b) for pseudo-Lipschitz loss functions. This result follows from [8, Theorem 14], which is a general result relating state evolutions to AMP algorithm with non-separable denoisers. In order to apply [8, Theorem 14], we need to demonstrate that our denoiser, i.e. the proximal operator $\text{prox}_{J_{\alpha\tau_t}}(\cdot)$ defined in (1.4), satisfies two additional properties labeled **(P1)** and **(P2)** below.

Define a sequence of denoisers $\{\eta_p^t\}_{p \in \mathbb{N}_{>0}}$ where $\eta_p^t : \mathbb{R}^p \rightarrow \mathbb{R}^p$ to be those that apply the proximal operator $\text{prox}_{J_{\alpha\tau_t}}(\cdot)$ defined in (1.4), i.e. for a vector $\mathbf{v} \in \mathbb{R}^p$, define

$$\eta_p^t(\mathbf{v}) := \text{prox}_{J_{\alpha\tau_t}}(\mathbf{v}). \quad (4.1)$$

(P1) For each t , denoisers $\eta_p^t(\cdot)$ defined in (4.1) are uniformly Lipschitz (i.e. uniformly pseudo-Lipschitz of order $k = 1$) per Definition 3.1.

(P2) For any s, t with $(\mathbf{Z}, \mathbf{Z}')$ a pair of length- p vectors, where for $i \in \{1, 2, \dots, p\}$, the pair (Z_i, Z'_i) i.i.d. $\sim \mathcal{N}(0, \Sigma)$ with Σ any 2×2 covariance matrix, the following limits exist and are finite.

$$\text{plim}_{p \rightarrow \infty} \frac{1}{p} \|\beta\|, \quad \text{plim}_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\beta^\top \eta_p^t(\beta + \mathbf{Z})], \quad \text{and} \quad \text{plim}_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}[\eta_p^s(\beta + \mathbf{Z}')^\top \eta_p^t(\beta + \mathbf{Z})].$$

We will show that properties **(P1)** and **(P2)** are satisfied for our problem in Appendix B.

Lemma 4.1. [8, Theorem 14] *Under assumptions **(A1)** - **(A4)**, given that **(P1)** and **(P2)** are satisfied, for the AMP algorithm in (1.3b) and for any uniformly pseudo-Lipschitz sequence of functions $\phi_n : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ and $\psi_p : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$, let $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_n)$ and $\mathbf{Z}' \sim \mathcal{N}(0, \mathbb{I}_p)$, then*

$$\begin{aligned} \text{plim}_n \left(\phi_n(\mathbf{z}^t, \mathbf{w}) - \mathbb{E}_{\mathbf{Z}}[\phi_n(\mathbf{w} + \sqrt{\tau_t^2 - \sigma_w^2} \mathbf{Z}, \mathbf{w})] \right) &= 0, \\ \text{plim}_p \left(\psi_p(\beta^t + \mathbf{X}^\top \mathbf{z}^t, \beta) - \mathbb{E}_{\mathbf{Z}'}[\psi_p(\beta + \tau_t \mathbf{Z}', \beta)] \right) &= 0, \end{aligned}$$

where τ_t is defined in (2.4).

We now show that Theorem 3 follows from Lemma 4.1 and Theorem 2.

Proof of Theorem 3. First, for any fixed n and t , the following bound uses that ψ_n is uniformly pseudo-Lipschitz of order k and the Triangle Inequality,

$$\begin{aligned} |\psi_p(\beta^t, \beta) - \psi_p(\hat{\beta}, \beta)| &\leq L \left(1 + \left(\frac{\|(\beta^t, \beta)\|}{\sqrt{2p}} \right)^{k-1} + \left(\frac{\|(\hat{\beta}, \beta)\|}{\sqrt{2p}} \right)^{k-1} \right) \frac{1}{\sqrt{2p}} \|\beta^t - \hat{\beta}\| \\ &\leq L \left(1 + \left(\frac{\|\beta^t\|}{\sqrt{2p}} \right)^{k-1} + \left(\frac{\|\hat{\beta}\|}{\sqrt{2p}} \right)^{k-1} + \left(\frac{\|\beta\|}{\sqrt{2p}} \right)^{k-1} \right) \frac{1}{\sqrt{2p}} \|\beta^t - \hat{\beta}\|. \end{aligned}$$

Now we take limits on either side of the above, first with respect to p and then with respect to t . We note that the term $\frac{1}{\sqrt{n}} \|\beta^t - \hat{\beta}\|$ vanishes by Theorem 2. Then as long as

$$\lim_t \text{plim}_p \left(\|\beta^t\| / \sqrt{p} \right)^{k-1}, \quad \text{plim}_p \left(\|\hat{\beta}\| / \sqrt{p} \right)^{k-1}, \quad \text{and} \quad \text{plim}_p \left(\|\beta\| / \sqrt{p} \right)^{k-1}, \quad (4.2)$$

are all finite, we have $\text{plim}_p \psi_p(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) = \lim_t \text{plim}_p \psi_p(\boldsymbol{\beta}^t, \boldsymbol{\beta})$. But by Theorem 4.1 we also know that

$$\lim_t \text{plim}_p \psi_p(\boldsymbol{\beta}^t, \boldsymbol{\beta}) = \lim_t \text{plim}_p \mathbb{E}[\psi_p(\eta^t(\boldsymbol{\beta} + \tau_t \mathbf{Z}), \boldsymbol{\beta})],$$

giving the desired result.

Finally we convince ourself that the limits in (4.2) are finite. Since k finite, that the third term in (4.2) is finite follows by property **(P2)**. Bounds for the first and second term are demonstrated in Lemma 7.1 found in Appendix 6.

□

5 Proof AMP Finds the SLOPE Solutions

In this section we aim to prove Theorem 2. Define the SLOPE cost function as follows,

$$\mathcal{C}(\mathbf{b}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + J_{\boldsymbol{\lambda}}(\mathbf{b}), \quad (5.1)$$

where $J_{\boldsymbol{\lambda}}(\mathbf{b})$ is the sorted ℓ_1 -norm. The proof of Theorem 2 relies on a technical lemma, Lemma 5.5, stated in Section 5.2 below, that deals carefully with the fact that the SLOPE cost function given in (5.1) is *not* necessarily strongly convex.

In the LASSO case, one works around this challenge by showing that the (LASSO) cost function does have nice properties when considering just the elements of the non-zero support of $\boldsymbol{\beta}^t$ at any (large) iteration t , using that the non-zero support of $\boldsymbol{\beta}$ has size no larger than $n < p$.

In the SLOPE problem, however, it is possible that the support set has size exceeding n , and therefore the LASSO analysis is not immediately applicable. Our proof develops novel techniques that are tailored to the characteristics of the SLOPE solution. Specifically, when considering the SLOPE problem, one can show nice properties (similar to those in the LASSO case) by considering a support-like set, that being the *unique* non-zeros in the estimate $\boldsymbol{\beta}^t$ at any (large) iteration t .

In other words, our strategy is to define an equivalence relation $x \sim y$ when $|x| = |y|$ and partition the entries of the AMP estimate at any iteration t into equivalence classes. This allows us to observe, using (2.10) and the non-negativity of $\boldsymbol{\lambda}$, that the number of equivalence classes is no larger than n . (Recall that $\|\cdot\|_0^*$ counts the unique non-zero magnitudes in a vector.) We see an analogy between SLOPE's equivalence class (or 'maximal atom' as described in Section 5.1) and LASSO's support set. This approach, taken in Lemma 5.5 below, allows us to deal with the fact that we are not guaranteed to have a strongly convex cost. Then Lemma 5.5 is used to prove Theorem 3.

Before we state Lemma 5.5, we include some useful preliminary information on SLOPE that will be needed for the upcoming work. In particular, we introduce in more details the idea of equivalence classes of elements having the same magnitude, a mapping of vector ranking denoted as $\hat{\Pi}$, and a polytope-related mapping whose image is the set of subgradients denoted as \mathcal{P} . These definitions are all given in more detail in Section 5.1.

5.1 Preliminaries on SLOPE

In general, we refer to the function $\mathcal{C}(\cdot)$ stated in (5.1) as the SLOPE cost function and the SLOPE estimator $\hat{\boldsymbol{\beta}}$ is the one that minimizes the SLOPE cost. We note that the SLOPE cost function $\mathcal{C}(\cdot)$ depends on both \mathbf{y} and $\boldsymbol{\lambda}$, so technically a notation like $\mathcal{C}_{(\mathbf{y}, \boldsymbol{\lambda})}(\cdot)$ would be more rigorous, however, we don't think that dropping the explicit dependence on $(\mathbf{y}, \boldsymbol{\lambda})$ will cause any confusion.

For a convex function $f : \mathbb{R}^p \rightarrow \mathbb{R}$, we denote the subgradient of f at a point $\mathbf{x} \in \mathbb{R}^p$ as $\partial f(\mathbf{x})$. We will be interested, particularly, in the subgradient of the SLOPE cost $\partial\mathcal{C}(\mathbf{b})$ which forces us to study the subgradient of the SLOPE norm $\partial J_{\boldsymbol{\lambda}}(\mathbf{b})$. In particular,

Fact 5.1. $\partial\mathcal{C}(\mathbf{b}) = -\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\mathbf{b}) + \partial J_{\boldsymbol{\lambda}}(\mathbf{b})$.

We will now describe explicitly the relevant subgradient, $\partial J_{\boldsymbol{\lambda}} \subset \mathbb{R}^p$. We note that the proximal operator given in (1.4) is linked to the subgradient of the SLOPE norm in the following way.

Fact 5.2. If $\text{prox}_{J_{\boldsymbol{\lambda}}}(\mathbf{v}_1) = \mathbf{v}_2$, then $\mathbf{v}_1 - \mathbf{v}_2 \in \partial J_{\boldsymbol{\lambda}}(\mathbf{v}_2)$.

Define a function $\Pi_{\mathbf{x}} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ to be a mapping (not necessarily unique) that sorts its input by magnitude in descending order according to absolute values of entries in \mathbf{x} . For example, if $\mathbf{x} = (5, 2, -3, -5)$, then there are two possible such mappings $\Pi_{\mathbf{x}}(\mathbf{b}) = (|b_1|, |b_4|, |b_3|, |b_2|)$ or $\Pi_{\mathbf{x}}(\mathbf{b}) = (|b_4|, |b_1|, |b_3|, |b_2|)$. Using this notation, we can rewrite the SLOPE norm as $J_{\boldsymbol{\lambda}}(\mathbf{b}) = \boldsymbol{\lambda} \cdot \Pi_{\mathbf{x}}(\mathbf{b})$. Since such mapping may not be unique, the inverse may not exist and we therefore define a pseudo-inverse mapping, $\hat{\Pi}_{\mathbf{x}}^{-1}$, that is based on the function $\hat{\Pi}_{\mathbf{x}} : \mathbb{R}^p \rightarrow \{\text{maximal atoms}\}$. In words, $\hat{\Pi}_{\mathbf{x}}$ finds the maximal atoms of ranking of the absolute values of \mathbf{x} . Then $\hat{\Pi}_{\mathbf{x}}$ corresponds to the mapping

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ \{1, 2\} & 4 & 3 & \{1, 2\} \end{pmatrix}$$

with $\hat{\Pi}_{\mathbf{x}}(\mathbf{x}) = (\{5, -5\}, \{5, -5\}, -3, 2)$ and $\hat{\Pi}_{\mathbf{x}}^{-1}(\boldsymbol{\lambda}) = (\{\lambda_1, \lambda_2\}, \lambda_4, \lambda_3, \{\lambda_1, \lambda_2\})$. Then it is not hard to see that there exists $\hat{\boldsymbol{\lambda}} \in \hat{\Pi}_{\mathbf{x}}^{-1}(\boldsymbol{\lambda})$ such that $J_{\boldsymbol{\lambda}}(\mathbf{b}) = \boldsymbol{\lambda} \cdot \Pi_{\mathbf{x}}(\mathbf{b}) = \hat{\boldsymbol{\lambda}} \cdot |\mathbf{b}|$. In words, this says there are two equivalent ways to consider the calculation of $J_{\boldsymbol{\lambda}}(\mathbf{b})$ when $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. First $\boldsymbol{\lambda} \cdot \Pi_{\mathbf{x}}(\mathbf{b})$ computes the inner product between $\boldsymbol{\lambda}$ and the *sorted* magnitudes of \mathbf{b} , and in the second case, $\hat{\boldsymbol{\lambda}}^\top |\mathbf{b}|$ computes the inner product between the magnitudes of \mathbf{b} (unsorted), with a rearrangement of the $\boldsymbol{\lambda}$ vector (based on \mathbf{b}) that pairs the values in $\boldsymbol{\lambda}$ with the values of $|\mathbf{b}|$ by magnitude.

Now we define an equivalence relation $x \sim y$ if $|x| = |y|$. Then $\hat{\Pi}_{\mathbf{x}}$ partitions elements in \mathbf{x} into different equivalence classes I . The motivation of using equivalence classes roots from AMP. In calibrating the AMP to the SLOPE problem, we need to calculate ∇prox , which equals the number of non-zero equivalence classes. For example, $\frac{\partial \text{prox}}{\partial \mathbf{v}}|_{\mathbf{v}=(1,0,-1,3)} = (\frac{1}{2}, 0, \frac{1}{2}, 1)$ has a sum of 2.

Now we note that the subgradient of the SLOPE norm can be represented using the idea of the equivalence classes. For a vector $\mathbf{v} \in \mathbb{R}^p$, we use the notation \mathbf{v}_I to be the elements of the vector \mathbf{v} belonging to equivalence class I . Then,

Fact 5.3.

$$\partial J_{\boldsymbol{\lambda}}(\mathbf{s}) = \left\{ \mathbf{v} \in \mathbb{R}^p : \text{for each equivalent class } I, \begin{cases} \text{if } s_I \neq 0 \implies \mathbf{v}_I \in \mathcal{P}([\hat{\Pi}_{\mathbf{s}}^{-1}(\boldsymbol{\lambda})]_I) \text{ sign}(s_I); \\ \text{if } s_I = 0 \implies |\mathbf{v}_I| \in \mathcal{P}_0([\hat{\Pi}_{\mathbf{s}}^{-1}(\boldsymbol{\lambda})]_I) \end{cases} \right\}.$$

In the above, $\mathcal{P}, \mathcal{P}_0$ are polytope-related mappings,

$$\begin{aligned} \mathcal{P}(\mathbf{u}) &:= \{ \mathbf{y} : \mathbf{y} = \mathbf{A}\mathbf{u} \text{ for some doubly stochastic matrix } \mathbf{A} \} \\ \mathcal{P}_0(\mathbf{u}) &:= \{ \mathbf{y} : \mathbf{y} = \mathbf{A}\mathbf{u} \text{ for some doubly sub-stochastic matrix } \mathbf{A} \} \end{aligned}$$

By definition, the doubly stochastic matrix, a.k.a. a Birkhoff polytope, is a square matrix of non-negative real numbers, whose row and column sums equal 1. For example,

$$\mathbf{A} = \begin{pmatrix} 1/3 & 2/3 & 0 \\ 1/6 & 1/3 & 1/2 \\ 1/2 & 0 & 1/2 \end{pmatrix} \tag{5.2}$$

is a doubly stochastic matrix. Similarly, a doubly sub-stochastic matrix is defined as a square matrix of non-negative real numbers, whose row and column sums are at most 1. Note that if all entries of $\boldsymbol{\lambda}$ take the same value, the subgradient in Fact 5.3 gives the usual subgradient of the ℓ_1 norm.

Using the subgradient definition in Fact 5.3, consider $\mathcal{P}((\lambda_1, \lambda_2, \lambda_3))$, relating to a non-zero equivalence class having three entries. Then \mathbf{A} in (5.2) is one possible matrix considered in defining the set $\mathcal{P}((\lambda_1, \lambda_2, \lambda_3))$ and it has the following interpretation. The rows of \mathbf{A} determine how the subgradient \mathbf{v}_I values are calculated by averaging the corresponding threshold values $\boldsymbol{\lambda}$, for example, the first entry of \mathbf{v}_I is a weighted average with 1/3 its weight in λ_1 and 2/3 in λ_2 ; the second entry of \mathbf{v}_I is a weighted average with 1/6 its weight in λ_1 , 1/3 in λ_2 , and 1/2 in λ_3 , etc. You can think of this as determining the threshold each input value s_I receives, as some weighted combination of all the possible threshold values $\boldsymbol{\lambda}$ corresponding to this equivalence class. Similarly, the columns of the doubly-stochastic matrix considered in the mapping \mathcal{P} define how the thresholds $\boldsymbol{\lambda}$ are spread out amongst each element of the subgradient, for example, 1/3 of λ_1 's value goes to the first element of \mathbf{v}_I , 1/6 to the second value, and 1/2 to the third value, etc.

To see why $\partial J_{\boldsymbol{\lambda}}(\mathbf{s})$ takes the form given in Fact 5.3, let's consider again the \mathcal{P} used in the case that $\mathbf{s}_I \neq 0$. Recall the \mathbf{s}_I looks at only the indices of \mathbf{s} appearing in the equivalence class I , so all elements of \mathbf{s}_I have the same absolute value. This means that there are many ways to share the corresponding $\boldsymbol{\lambda}$ threshold values among them. We can think of this as an assignment problem: assign jobs (thresholds $\boldsymbol{\lambda}$) to workers (s_i) where as assignment according to a doubly stochastic matrix is a natural one (all workers take on the same load, and all jobs must be completed). On the other hand, \mathcal{P}_0 does not require that the sharing of the threshold values $\boldsymbol{\lambda}$ amongst the entries of \mathbf{s}_I be strict: row and/or column sums can be smaller than one. This difference is rooted in the subgradient of ℓ_1 norm: i.e. $\partial|x| = \text{sign}(x)$ when $x \neq 0$ and $\partial|x| \in [-1, 1]$ when $x = 0$.

For a rigorous proof of Fact 5.3, we refer the reader to [32, Exercise 8.31], but we give a quick sketch here in the case of $\mathbf{s}_I \neq 0$. The proof uses that $\mathcal{P}(\mathbf{u})$ is a permutohedron, meaning a convex hull with vertices corresponding to permuted entries of \mathbf{u} . Notice that we can rewrite $J_{\boldsymbol{\lambda}}(\mathbf{s})$ as a finite max function $J_{\boldsymbol{\lambda}}(\mathbf{s}) : \max\{\boldsymbol{\lambda}^\top f_1(\mathbf{s}), \dots, \boldsymbol{\lambda}^\top f_m(\mathbf{s})\}$, where $\{f_i(\mathbf{s})\}_{1 \leq i \leq m}$ is the collection of all possible permutations for the entries of $|\mathbf{s}|$. Notice that the permutation that sorts the magnitudes will be chosen by the maximum function. For such a function (see [32, Exercise 8.31]) the subgradient takes the form of a convex hull of the partial derivatives of the maximizing elements:

$$\partial J_{\boldsymbol{\lambda}}(\mathbf{s}) \in \text{conv}\{\nabla_{\mathbf{s}}(\boldsymbol{\lambda}^\top f_i(\mathbf{s})) : i \in A(\mathbf{s})\} \equiv \text{conv}\{f_i^{-1}(\boldsymbol{\lambda}) : i \in A(\mathbf{s})\}, \quad (5.3)$$

where $A(\mathbf{s}) = \{i \in \{1, 2, \dots, m\} : \boldsymbol{\lambda}^\top f_i(\mathbf{s}) = J_{\boldsymbol{\lambda}}(\mathbf{s})\}$ and in our case, the partial derivatives correspond to permutations of the thresholds. Now, without loss of generality, let's consider an input that has only one non-zero equivalence class, i.e. $\mathbf{s} = (s, s, \dots, s) \in \mathbb{R}^d$. Then clearly there are $m = d!$ possible permutations. Therefore,

$$\partial J_{\boldsymbol{\lambda}}(\mathbf{s}) \in \text{conv}\{f_i^{-1}(\boldsymbol{\lambda}) : i \in \{1, 2, \dots, d!\}\} \equiv \text{conv}\{f_i(\boldsymbol{\lambda}) : i \in \{1, 2, \dots, d!\}\}.$$

In other words, the partial derivative lies in the set that is the convex combination of all possible permutations of the threshold $\boldsymbol{\lambda}$. By definition, this is a permutohedron. So, in our case, the subgradient is a convex hull whose vertices are the permuted thresholds, i.e. an image of Birkhoff polytope under the thresholds, which can be characterized by doubly stochastic matrices.

5.2 Main Technical Lemma

Now we state and prove the main technical lemma that will be used to prove Theorem 2. Before we state Lemma 5.5, let us introduce a very important definition:

Definition 5.4. Given a vector $\mathbf{v} \in \mathbb{R}^p$, a set $I \subset \{1, \dots, p\}$ is said to be a maximal atom of indices of \mathbf{v} if $|v_i| = |v_j|$ for all $i, j \in I$ and $|v_i| \neq |v_k|$ for $i \in I$ and all $k \notin I$. With this definition in place, we define the star support of the vector \mathbf{v} as

$$\text{supp}^*(\mathbf{v}) := \{I : I \subset \{1, \dots, p\} \text{ is a maximal atom of indices of } \mathbf{v} \text{ and } \mathbf{v}_I \neq 0\}.$$

For example, if $\mathbf{v} = (1, 1, -1, 0, 2, -1)$, then $\text{supp}^*(\mathbf{v}) = \{\{1, 2, 3, 6\}, \{5\}\}$. Now we state and prove Lemma 5.5.

Lemma 5.5. For constants $c_1, \dots, c_5 > 0$, if the following conditions are satisfied,

- (1) $\frac{1}{\sqrt{p}} \|\boldsymbol{\beta}^t - \hat{\boldsymbol{\beta}}\| \leq c_1$,
- (2) There exists a subgradient $sg(\mathcal{C}, \boldsymbol{\beta}^t) \in \partial \mathcal{C}(\boldsymbol{\beta}^t)$ such that $\frac{1}{\sqrt{p}} \|sg(\mathcal{C}, \boldsymbol{\beta}^t)\| \leq \epsilon$,
- (3) Let $\boldsymbol{\nu}^t := \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t) + sg(\mathcal{C}, \boldsymbol{\beta}^t) \in \partial J_{\boldsymbol{\lambda}}(\boldsymbol{\beta}^t)$ (where $sg(\mathcal{C}, \boldsymbol{\beta}^t)$ is the subgradient from Condition (2)). Denote $s_t(c_2) := \{I \subset [p] : |\boldsymbol{\nu}_I^t| \succeq [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I (1 - c_2)\}$ and $S_t(c_2) := \{i \in I : I \in s(c_2)\}$, where the equivalence classes, I , for both sets are defined via the AMP estimation $\boldsymbol{\beta}^t$, and for a vector $\mathbf{x} \in \mathbb{R}^d$ and a set $\mathbf{A} \subset \mathbb{R}^d$, the notation $\mathbf{x} \succeq \mathbf{A}$ means there exists some $\mathbf{y} \in \mathbf{A}$ such that $\mathbf{x} \geq \mathbf{y}$ elementwise. Then for s' being any set of maximal atoms in $[p]$ with $|s'| \leq c_3 p$ and $S' := \{i \in I : I \in s'\}$, we have $\sigma_{\min}(\mathbf{X}_{S_t(c_2) \cup S'}) \geq c_4$.
- (4) The minimum non-zero and maximum singular value of \mathbf{X} , denoted as $\hat{\sigma}_{\min}^2(\mathbf{X})$ and $\sigma_{\max}^2(\mathbf{X})$, are bounded: i.e. $\hat{\sigma}_{\min}^2(\mathbf{X}) \geq \frac{1}{c_5}$ and $\sigma_{\max}^2(\mathbf{X}) \leq c_5$.
- (5) Define $\mathcal{C}_{\mathbf{x}}(\mathbf{b}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \sum_{i=1}^p \hat{\lambda}_i |b_i|$ for some $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\mathbf{x}}^{-1}(\boldsymbol{\lambda}))$. Then $\mathcal{C}(\boldsymbol{\beta}^t) \geq \mathcal{C}_{\boldsymbol{\beta}^t}(\hat{\boldsymbol{\beta}})$.

then for some function $f(\epsilon) := f(\epsilon, c_1, c_2, c_3, c_4, c_5)$ such that $f(\epsilon) \rightarrow 0$ as $\epsilon \rightarrow 0$,

$$\frac{1}{\sqrt{p}} \|\boldsymbol{\beta}^t - \hat{\boldsymbol{\beta}}\| < f(\epsilon).$$

We wrap up this section by proving Lemma 5.5. Once we have proved Lemma 5.5, we will be able to prove Theorem 2. The major piece of work in proving Theorem 2 is in showing that the five assumptions of Lemma 5.5 are satisfied. Then the result of Theorem 2 is immediate. We show the five assumptions are met in Sections 7.1 - 7.5. Now we prove the Lemma.

Proof of Lemma 5.5. Throughout the proof, we denote ξ_1, ξ_2, \dots as functions of the constants $c_1, \dots, c_5 > 0$ and of ϵ such that $\xi_i(\epsilon) \rightarrow 0$ as $\epsilon \rightarrow 0$ (we omit the dependence of ξ_i on ϵ). We will think of t as a fixed iteration and we denote the residual we are interested in studying as $\mathbf{r} = \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^t$.

The proof strategy is to show that $\frac{1}{p} \|\mathbf{X}\mathbf{r}\|^2 \leq \xi(\epsilon)$ from which a similar result for $\frac{1}{p} \|\mathbf{r}\|^2$ follows when we have control of the singular values of \mathbf{X} as we do with Condition (4). Structurally, the proof is similar to that in the LASSO case (cf. [5, Lemma 3.1]), with the main difference coming through Condition (3), where we need to use star support instead of the support when bounding the minimum singular value of a selection of columns of \mathbf{X} .

For a fixed iteration t , let $S = \{i \in [p] : i \in I \text{ and } I \in \text{supp}^*(\boldsymbol{\beta}^t)\}$, i.e. S is the collection of (unique) indices belonging to the star support of the AMP estimate at iteration t . Then for a vector $\mathbf{v} \in \mathbb{R}^p$ we denote \mathbf{v}_S to mean the vector indexed only over the indices in the set S and we let \bar{S} denote the complement of S . In what follows, we drop the t -dependence on $\boldsymbol{\nu}^t$, writing $\boldsymbol{\nu} = \boldsymbol{\nu}^t$ and for p -length vectors \mathbf{u} and \mathbf{v} , define $\langle \mathbf{u}, \mathbf{v} \rangle := \frac{1}{p} \sum_i u_i v_i$.

First,

$$\begin{aligned} 0 &\stackrel{(a)}{\geq} \frac{1}{p} (\mathcal{C}_{\boldsymbol{\beta}^t}(\hat{\boldsymbol{\beta}}) - \mathcal{C}(\boldsymbol{\beta}^t)) \stackrel{(b)}{=} \frac{1}{2p} (\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 - \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t\|^2) + \langle \hat{\boldsymbol{\lambda}}, |\hat{\boldsymbol{\beta}}| - |\boldsymbol{\beta}^t| \rangle \\ &\stackrel{(c)}{=} \langle \hat{\boldsymbol{\lambda}}_S, |\boldsymbol{\beta}_S^t + \mathbf{r}_S| - |\boldsymbol{\beta}_S^t| \rangle + \langle \hat{\boldsymbol{\lambda}}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle + \frac{1}{2p} (\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t - \mathbf{X}\mathbf{r}\|^2 - \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t\|^2) \\ &\stackrel{(d)}{=} [\langle \hat{\boldsymbol{\lambda}}_S, |\boldsymbol{\beta}_S^t + \mathbf{r}_S| - |\boldsymbol{\beta}_S^t| \rangle - \langle \boldsymbol{\nu}_S, \mathbf{r}_S \rangle] + [\langle \hat{\boldsymbol{\lambda}}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle - \langle \boldsymbol{\nu}_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle] + \langle \boldsymbol{\nu}, \mathbf{r} \rangle - \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t, \mathbf{X}\mathbf{r} \rangle + \frac{\|\mathbf{X}\mathbf{r}\|^2}{2p} \\ &\stackrel{(e)}{=} [\langle \hat{\boldsymbol{\lambda}}_S, |\boldsymbol{\beta}_S^t + \mathbf{r}_S| - |\boldsymbol{\beta}_S^t| \rangle - \langle \boldsymbol{\nu}_S, \mathbf{r}_S \rangle] + [\langle \hat{\boldsymbol{\lambda}}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle - \langle \boldsymbol{\nu}_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle] + \langle \text{sg}(\mathcal{C}, \boldsymbol{\beta}^t), \mathbf{r} \rangle + \frac{\|\mathbf{X}\mathbf{r}\|^2}{2p}. \end{aligned}$$

In the above, step (a) follows immediately from Condition (5) and step (b) holds *for any* $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$ by the definition of $\mathcal{C}_{\boldsymbol{\beta}^t}(\hat{\boldsymbol{\beta}})$, noticing that $J_{\boldsymbol{\lambda}}(\boldsymbol{\beta}^t) = \hat{\boldsymbol{\lambda}}^\top |\boldsymbol{\beta}^t|$ in the SLOPE cost (5.1) since $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$. Below we will select a specific $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$ based on the definition of $\boldsymbol{\nu}$. Step (c) follows by replacing $\hat{\boldsymbol{\beta}}$ with $\boldsymbol{\beta}^t + \mathbf{r}$ and noticing that $\boldsymbol{\beta}_S^t = \mathbf{0}$. Step (d) follows since $\langle \boldsymbol{\nu}, \mathbf{r} \rangle = \langle \boldsymbol{\nu}_S, \mathbf{r}_S \rangle + \langle \boldsymbol{\nu}_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle$ and step (e) from the definition of $\boldsymbol{\nu}$.

Using Conditions (1) and (2), we get by Cauchy-Schwarz

$$[\langle \hat{\boldsymbol{\lambda}}_S, |\boldsymbol{\beta}_S^t + \mathbf{r}_S| - |\boldsymbol{\beta}_S^t| \rangle - \langle \boldsymbol{\nu}_S, \mathbf{r}_S \rangle] + [\langle \hat{\boldsymbol{\lambda}}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle - \langle \boldsymbol{\nu}_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle] + \frac{\|\mathbf{X}\mathbf{r}\|^2}{2p} \leq \frac{\|\text{sg}(\mathcal{C}, \boldsymbol{\beta}^t)\| \|\mathbf{r}\|}{p} \leq c_1 \epsilon \quad (5.4)$$

We now show all three terms on the left side of (5.4) are non-negative. The idea is then: if all three terms are non-negative and their sum tends to 0 as $\epsilon \rightarrow 0$, it must be true that each term tends to 0 too. The third term in (5.4), $\frac{1}{2p} \|\mathbf{X}\mathbf{r}\|^2$, is trivially non-negative, so we focus on the first two.

To show that the other terms are non-negative, we consider choosing a specific vector $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$ such that on the support, $\hat{\boldsymbol{\lambda}}_S = |\boldsymbol{\nu}_S|$, and off the support $\hat{\boldsymbol{\lambda}}_{\bar{S}} \geq |\boldsymbol{\nu}_{\bar{S}}|$, meaning $\hat{\boldsymbol{\lambda}}_I$ is parallel to $|\boldsymbol{\nu}_I|$ for each equivalence class I of $\boldsymbol{\beta}^t$. That such a $\hat{\boldsymbol{\lambda}}$ exists in the set $\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$ follows since $\boldsymbol{\nu}$ is a valid subgradient of $J_{\boldsymbol{\lambda}}(\boldsymbol{\beta}^t)$ (see Fact 5.3).

Using this $\hat{\boldsymbol{\lambda}}$, notice that the sets defined in Condition (3) are equivalent to the following: $s_t(c_2) := \{I \subset [p] : |\boldsymbol{\nu}_I| \geq (1 - c_2)\hat{\boldsymbol{\lambda}}_I\}$ and $S_t(c_2) := \{i : |\nu_i| \geq (1 - c_2)\hat{\lambda}_i\}$, where both use equivalence classes, I , defined for $\boldsymbol{\beta}^t$. To see that this is the case, note that if I is a non-zero equivalence class, by Fact 5.3, since $|\boldsymbol{\nu}_I| \in [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I$, we know that $|\boldsymbol{\nu}_I| \succeq [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I(1 - c_2)$ and similarly, since $\hat{\boldsymbol{\lambda}}_S = |\boldsymbol{\nu}_S|$ we know that $|\boldsymbol{\nu}_I| \geq (1 - c_2)\hat{\boldsymbol{\lambda}}_I$, so I clearly belongs to $s_t(c_2)$ for both definitions. If I is the zero equivalence class, if $|\boldsymbol{\nu}_I| \geq (1 - c_2)\hat{\boldsymbol{\lambda}}_I$ then obviously $|\boldsymbol{\nu}_I| \succeq [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I(1 - c_2)$ since $\hat{\boldsymbol{\lambda}} \in \mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))$. In the other direction, if the non-zero equivalence class I is such that $|\boldsymbol{\nu}_I| \succeq [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I(1 - c_2)$ then there exists a vector $\tilde{\boldsymbol{\nu}}_I \in [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I$ such that $|\boldsymbol{\nu}_I| \geq \tilde{\boldsymbol{\nu}}_I(1 - c_2)$ elementwise. However since $\tilde{\boldsymbol{\nu}}_I \in [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I$, this implies that $|\boldsymbol{\nu}_I| \geq (1 - c_2)\hat{\boldsymbol{\lambda}}_I$ is also true since $\hat{\boldsymbol{\lambda}}_I \in [\mathcal{P}(\hat{\Pi}_{\boldsymbol{\beta}^t}^{-1}(\boldsymbol{\lambda}))]_I$ in the same direction as $|\boldsymbol{\nu}_I|$.

To visualize the choice of $\hat{\boldsymbol{\lambda}}$, we consider an example where $\boldsymbol{\nu}_I = (-1, 2)$ for equivalence class $I = \{1, 2\}$ with $\boldsymbol{\lambda}_I = (4, 1)$ in Figure 3. In the figure, the blue shaded region indicates possible

subgradient values for zero elements and the black line are possible subgradients for zero elements. In this example, the equivalence class is that for zero elements, so we notice that ν_I lies in the blue region. Then λ_I is in the same direction as $|\nu_I|$ but lies on the black line (since $\hat{\lambda} \in \mathcal{P}(\hat{\Pi}_{\beta^t}^{-1}(\lambda))$).

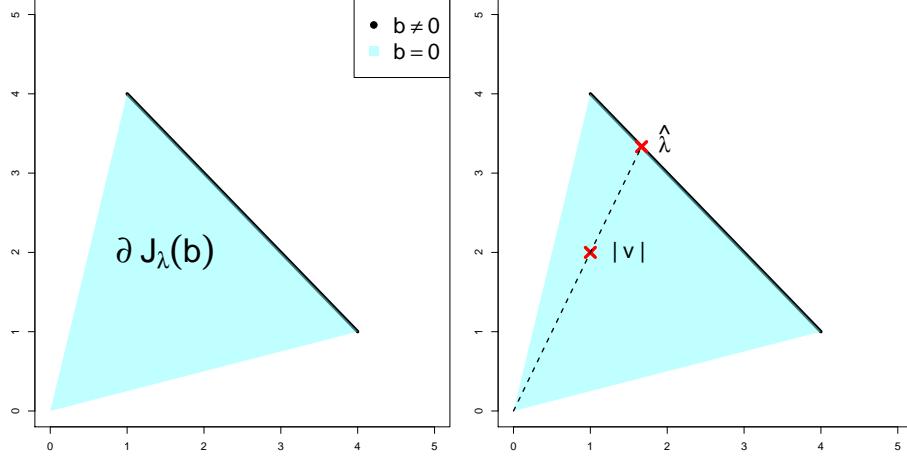


Figure 3: The blue area contained by the black line segment is the set of subgradients; Red crosses are examples of ν_I and $\hat{\lambda}_I$ correspondingly when $b_I = 0$.

Now we would like to show that the first term in (5.4) is non-negative. Specifically, our choice of $\hat{\lambda}$ gives $\nu_i = \text{sign}(\beta_i^t)\hat{\lambda}_i$, for each $i \in S$, and then it suffices, in order to prove the non-negativity of $\langle \hat{\lambda}_S, |\beta_S^t + r_S| - |\beta_S^t| \rangle - \langle \nu_S, r_S \rangle$, to show

$$\begin{aligned} 0 &\leq (|\beta_i^t + r_i| - |\beta_i^t|) - \text{sign}(\beta_i^t)r_i \\ &= (\beta_i^t + r_i)\text{sign}(\beta_i^t + r_i) - \beta_i^t\text{sign}(\beta_i^t) - r_i\text{sign}(\beta_i^t) = (\beta_i^t + r_i)[\text{sign}(\beta_i^t + r_i) - \text{sign}(\beta_i^t)], \end{aligned}$$

which follows since each $(\beta_i^t + r_i)[\text{sign}(\beta_i^t + r_i) - \text{sign}(\beta_i^t)]$ is either equal to 0 (when $\text{sign}(\beta_i^t) = \text{sign}(\beta_i^t + r_i)$) or equal to $2|\beta_i^t + r_i|$ otherwise.

Finally, the second term in (5.4) is also non-negative. It suffices to show for each $i \in \bar{S}$, we have $0 \leq \hat{\lambda}_i|r_i| - \nu_i r_i$, or equivalently $0 \leq \hat{\lambda}_i - \nu_i \text{sign}(r_i) = \hat{\lambda}_i(1 - \text{sign}(\beta_i^t)\text{sign}(r_i))$ which is clearly true. Since all three terms in (5.4) are non-negative and their sum tends to 0 as $\epsilon \rightarrow 0$, it must be true that each term tends to 0,

$$\langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle \leq \xi_1(\epsilon), \quad (5.5)$$

$$\|\mathbf{X}\mathbf{r}\|^2 \leq p\xi_1(\epsilon). \quad (5.6)$$

We now make use of these inequalities to construct the bound for $\frac{1}{p}\|\mathbf{r}\|^2$.

Decompose \mathbf{r} as $\mathbf{r} = \mathbf{r}^\perp + \mathbf{r}^\parallel$, with $\mathbf{r}^\parallel \in \ker(\mathbf{X})$ and $\mathbf{r}^\perp \in \ker^\perp(\mathbf{X})$ so that $\mathbf{X}\mathbf{r} = \mathbf{X}\mathbf{r}^\perp$. We will now use (5.5) and (5.6) to obtain bounds for $\|\mathbf{r}^\perp\|^2$ and $\|\mathbf{r}^\parallel\|^2$. First notice that by (5.6) and Condition (4) we have $\frac{1}{c_5}\|\mathbf{r}^\perp\|^2 \leq \hat{\sigma}_{\min}^2(\mathbf{X})\|\mathbf{r}^\perp\|^2 \leq \|\mathbf{X}\mathbf{r}^\perp\|^2 = \|\mathbf{X}\mathbf{r}\|^2 \leq p\xi_1(\epsilon)$.

In the case $\ker(\mathbf{X}) = \{0\}$, the proof is concluded. Otherwise, we prove a similar bound for $\|\mathbf{r}^\parallel\|^2$. To bound $\|\mathbf{r}^\parallel\|^2$, we use the fact that this can be done if there exists sets $Q \in [p]$ and $\bar{Q} \in [p]/Q$ such that we can bound $\|\mathbf{r}_{\bar{Q}}^\parallel\|^2$ and show a high probability lower bound for $\sigma_{\min}^2(\mathbf{X}_Q)$.

In (5.5), decompose $\mathbf{r}_{\bar{S}} = \mathbf{r}_{\bar{S}}^\perp + \mathbf{r}_{\bar{S}}^{\parallel}$ and observe that by Cauchy Schwarz inequality and the bound just obtained,

$$\langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\parallel}| \rangle \leq \frac{1}{p} \|\hat{\lambda}_{\bar{S}}\| \|\mathbf{r}_{\bar{S}}^{\parallel}\| \leq \frac{1}{p} \|\hat{\lambda}\| \|\mathbf{r}^{\perp}\| \leq \frac{1}{\sqrt{p}} \|\hat{\lambda}\| \sqrt{c_5 \xi_1(\epsilon)}. \quad (5.7)$$

Then we use the fact that

$$\begin{aligned} \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\parallel}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}}^{\parallel} \rangle &= \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}} - \mathbf{r}_{\bar{S}}^{\perp}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}} - \mathbf{r}_{\bar{S}}^{\perp} \rangle \leq \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle + \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\perp}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle + \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}}^{\perp} \rangle, \\ &= \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}| \rangle + \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\perp}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle + \langle \hat{\lambda}_{\bar{S}} \text{ sign}(\beta_{\bar{S}}^t), \mathbf{r}_{\bar{S}}^{\perp} \rangle \leq \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\perp}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}} \rangle + 2\langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\perp}| \rangle, \end{aligned}$$

to get from (5.5) and (5.7) that

$$\langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\parallel}| \rangle - \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}}^{\parallel} \rangle \leq \xi_2(\epsilon). \quad (5.8)$$

Next we would like to show

$$\langle \hat{\lambda}_{\bar{S}(c_2)}, |\mathbf{r}_{\bar{S}(c_2)}^{\parallel}| \rangle - \langle \nu_{\bar{S}(c_2)}, \mathbf{r}_{\bar{S}(c_2)}^{\parallel} \rangle (1 - c_2)^{-1} \geq 0. \quad (5.9)$$

Note that it suffices again to prove this elementwise for each $i \in \bar{S}(c_2)$. Specifically, note that $(1 - c_2)^{-1} |\nu_i| < \hat{\lambda}_i$ for each $i \in \bar{S}(c_2)$ by the set's definition and therefore $\hat{\lambda}_i |r_i^{\parallel}| - \nu_i r_i^{\parallel} (1 - c_2)^{-1} \geq |\nu_i| |r_i^{\parallel}| (1 - c_2)^{-1} - \nu_i r_i^{\parallel} (1 - c_2)^{-1} \geq 0$. Therefore,

$$\begin{aligned} \langle \hat{\lambda}_{\bar{S}(c_2)}, |\mathbf{r}_{\bar{S}(c_2)}^{\parallel}| \rangle &\stackrel{(a)}{\leq} \frac{1}{c_2} \langle \lambda_{\bar{S}(c_2)}, |\mathbf{r}_{\bar{S}(c_2)}^{\parallel}| \rangle - \frac{1}{c_2} \langle \nu_{\bar{S}(c_2)}, \mathbf{r}_{\bar{S}(c_2)}^{\parallel} \rangle = \frac{1}{c_2} \langle \hat{\lambda}_{\bar{S}(c_2)} - \nu_{\bar{S}(c_2)} \text{ sign}(\mathbf{r}_{\bar{S}(c_2)}^{\parallel}), |\mathbf{r}_{\bar{S}(c_2)}^{\parallel}| \rangle \\ &\stackrel{(b)}{\leq} \frac{1}{c_2} \langle \hat{\lambda}_{\bar{S}} - \nu_{\bar{S}} \text{ sign}(\mathbf{r}_{\bar{S}}^{\parallel}), |\mathbf{r}_{\bar{S}}^{\parallel}| \rangle = \frac{1}{c_2} \langle \hat{\lambda}_{\bar{S}}, |\mathbf{r}_{\bar{S}}^{\parallel}| \rangle - \frac{1}{c_2} \langle \nu_{\bar{S}}, \mathbf{r}_{\bar{S}}^{\parallel} \rangle \stackrel{(c)}{\leq} c_2^{-1} \xi_2(\epsilon). \end{aligned} \quad (5.10)$$

In particular, step (a) follows by (5.9), step (b) since $S \subseteq S_t(c_2)$ implies $\bar{S}_t(c_2) \subseteq \bar{S}$ along with the fact that $\hat{\lambda}_{\bar{S}} - \nu_{\bar{S}} \text{ sign}(\mathbf{r}_{\bar{S}}^{\parallel}) \geq 0$ elementwise (for each $i \in \bar{S}$, we have $\hat{\lambda}_i - \nu_i \text{ sign}(r_i^{\parallel}) > 0$ by $\hat{\lambda}_i \geq |\nu_i|$). Finally step (c) holds by (5.8). We now use the bound in (5.10) to bound components of \mathbf{r}^{\parallel} .

In order to bound $\|\mathbf{r}^{\parallel}\|^2$, we would like to exploit a relationship between the ℓ_1 and ℓ_2 norms. To do this, we consider an ordering of the elements of the vector \mathbf{r}^{\parallel} by magnitude. Recall that $\bar{S}_t(c_2) \subseteq \bar{S}$ and we first assume $|\bar{S}_t(c_2)| \geq pc_3/2$. Now we partition $\bar{S}_t(c_2) = \cup_{\ell=1}^K S_\ell$, where $(pc_3/2) \leq |S_\ell| \leq pc_3$, and such that for each $i \in S_\ell$ and $j \in S_{\ell+1}$, it follows that $|r_i^{\parallel}| \geq |r_j^{\parallel}|$. Finally, define $\bar{S}_+ := \cup_{\ell=2}^K S_\ell \subseteq \bar{S}_t(c_2)$, i.e. the set union of all the partitions except the first one corresponding to the indices containing the largest elements in \mathbf{r}^{\parallel} . Now we note for any $i \in S_\ell$, we have $|r_i^{\parallel}| \leq \|\mathbf{r}_{S_{\ell-1}}^{\parallel}\| / |S_{\ell-1}|$, that is, in terms of absolute value, for any i in group ℓ , it should be smaller than the average of all the elements in the previous group $\ell - 1$.

Then,

$$\begin{aligned} \|\mathbf{r}_{\bar{S}_+}^{\parallel}\|^2 &\stackrel{(a)}{=} \sum_{\ell=2}^K \|\mathbf{r}_{S_\ell}^{\parallel}\|^2 \stackrel{(b)}{\leq} \sum_{\ell=2}^K |S_\ell| \frac{\|\mathbf{r}_{S_{\ell-1}}^{\parallel}\|_1^2}{|S_{\ell-1}|^2} \stackrel{(c)}{\leq} \frac{4}{pc_3} \sum_{\ell=2}^K \|\mathbf{r}_{S_{\ell-1}}^{\parallel}\|_1^2 \leq \frac{4}{pc_3} \left[\sum_{\ell=2}^K \|\mathbf{r}_{S_{\ell-1}}^{\parallel}\|_1 \right]^2 \\ &\stackrel{(d)}{\leq} \frac{4}{pc_3} \|\mathbf{r}_{\bar{S}(c_2)}^{\parallel}\|_1^2 \stackrel{(e)}{\leq} \frac{4\xi_2(\epsilon)^2 p}{c_2^2 c_3 (\min \hat{\lambda}_{\bar{S}(c_2)})^2} =: p\xi_3(\epsilon). \end{aligned} \quad (5.11)$$

In the above, step (a) follows from the definition of \bar{S}_+ , step (b) from the fact that for $i \in S_\ell$, we have $|r_i^{\parallel}| \leq \|r_{S_{\ell-1}}^{\parallel}\| / |S_{\ell-1}|$, step (c) since $(pc_3/2) \leq |S_\ell| \leq pc_3$, and step (d) since $\sum_{\ell=2}^K S_\ell \subset \sum_{\ell=1}^K S_\ell = \bar{S}_t(c_2)$. Finally step (e) follows using that $\frac{1}{p} \min\{\hat{\lambda}_{\bar{S}(c_2)}\} \|r_{\bar{S}(c_2)}^{\parallel}\|_1 \leq \langle \hat{\lambda}_{\bar{S}(c_2)}, |r_{\bar{S}(c_2)}^{\parallel}\rangle$.

Now, recalling $S_+ = S_t(c_2) \cup S_1$ and $|S_1| \leq pc_3$, by Condition (3), $\sigma_{\min}(\mathbf{X}_{S_+}) \geq c_4$ and therefore,

$$c_4^2 \|r_{S_+}^{\parallel}\|^2 \leq \sigma_{\min}^2(\mathbf{X}_{S_+}) \|r_{S_+}^{\parallel}\|^2 \leq \|\mathbf{X}_{S_+} r_{S_+}^{\parallel}\|^2 \stackrel{(a)}{=} \|\mathbf{X}_{\bar{S}_+} r_{\bar{S}_+}^{\parallel}\|^2 \stackrel{(b)}{\leq} 2c_5 \|r_{\bar{S}_+}^{\parallel}\|^2. \quad (5.12)$$

In the above, in step (a) we use that $\mathbf{0} = \mathbf{X} r^{\parallel} = \mathbf{X}_{S_+} r_{S_+}^{\parallel} + \mathbf{X}_{\bar{S}_+} r_{\bar{S}_+}^{\parallel}$. In step (b) we use Condition (4) and the fact that $\|\mathbf{X}_{\bar{S}_+} r_{\bar{S}_+}^{\parallel}\|^2 \leq \sigma_{\max}^2(\mathbf{X}) \|r_{\bar{S}_+}^{\parallel}\|^2$. Therefore, to conclude the proof, it is sufficient to prove a bound for $\|r_{S_+}^{\parallel}\|^2$.

Decomposing $\|r^{\parallel}\|^2 = \|r_{S_+}^{\parallel}\|^2 + \|r_{\bar{S}_+}^{\parallel}\|^2$, we find from (5.11) and (5.11) the desired bound:

$$\|r^{\parallel}\|^2 \leq \|r_{S_+}^{\parallel}\|^2 + \|r_{\bar{S}_+}^{\parallel}\|^2 \leq \left(\frac{2c_5}{c_4^2} + 1 \right) \|r_{\bar{S}_+}^{\parallel}\|^2 \leq \left(\frac{2c_5}{c_4^2} + 1 \right) p \xi_3(\epsilon).$$

This finishes the proof when $|\bar{S}_t(c_2)| \geq pc_3/2$. When $|\bar{S}_t(c_2)| < pc_3/2$, we can take $\bar{S}_+ = \emptyset$ and $S_+ = [p]$. Hence, the result holds as a special case of the above inequality. \square

6 Expansion of the AMP State Evolution Ideas

In this section, we develop ideas and notation specifically for the SLOPE AMP algorithm given in (1.3). Most are adapted from the work in [8] that studies general non-separable AMP algorithms. These results relate to the performance analysis of the AMP algorithm and will be useful in proving Lemma 5.5. Throughout this section, we use the $\{\eta_p^t\}_{p \in \mathbb{N}_{>0}}$ notation introduced in Section 4 and defined in (4.1). Namely, we consider a sequence of denoisers $\eta_p^t : \mathbb{R}^p \rightarrow \mathbb{R}^p$ to be those that apply the proximal operator $\text{prox}_{J_{\alpha\tau_t}}(\cdot)$ defined in (1.4), i.e. $\eta_p^t(\mathbf{v}) := \text{prox}_{J_{\alpha\tau_t}}(\mathbf{v})$ for a vector $\mathbf{v} \in \mathbb{R}^p$.

Given $\mathbf{w} \in \mathbb{R}^n$ and $\beta \in \mathbb{R}^p$, define sequences of column vectors $\mathbf{h}^{t+1} \in \mathbb{R}^p$ and $\mathbf{m}^t \in \mathbb{R}^n$ for $t \geq 0$. At each iteration t , the sequence \mathbf{h}^{t+1} measures the difference between the truth β and the pseudo-data $\mathbf{X}^\top \mathbf{z}^t + \beta^t$, that is the input to the denoiser, and the sequence \mathbf{m}^t measures the difference between the noise \mathbf{w} and the AMP residual \mathbf{z}^t . Namely, define $\mathbf{m}^t, \mathbf{h}^{t+1}$: for $t \geq 0$,

$$\mathbf{h}^{t+1} = \beta - (\mathbf{X}^\top \mathbf{z}^t + \beta^t) \quad \text{and} \quad \mathbf{m}^t = \mathbf{w} - \mathbf{z}^t. \quad (6.1)$$

We next introduce a generalization to the state evolution given in (2.4), that will be useful in studying the limiting properties of functions of the AMP estimates β^s and β^t at different iterations s and t . To do this, we will recursively define covariances $\{\Sigma_{s,t}\}_{s,t \geq 0}$: for \mathbf{B} elementwise i.i.d. $\sim B$, set $\Sigma_{0,0} = \sigma_w^2 + \frac{1}{\delta} \mathbb{E}[B^2]$ and

$$\Sigma_{0,t+1} = \sigma_w^2 + \lim_p \frac{1}{\delta p} \mathbb{E}\{-\mathbf{B}^\top [\eta_p^t(\mathbf{B} + \tau_t \mathbf{Z}_t) - \mathbf{B}]\}, \quad (6.2)$$

for $\mathbf{Z}_t \sim \mathcal{N}(0, \mathbb{I})$ independent of \mathbf{B} . Then for each $t \geq 0$, given $(\Sigma_{s,r})_{0 \leq s,r \leq t}$, define

$$\Sigma_{s+1,t+1} = \sigma_w^2 + \lim_p \frac{1}{\delta p} \mathbb{E}\left\{ [\eta_p^s(\mathbf{B} + \tau_s \mathbf{Z}_s) - \mathbf{B}]^\top [\eta_p^t(\mathbf{B} + \tau_t \mathbf{Z}_t) - \mathbf{B}] \right\}, \quad (6.3)$$

where \mathbf{Z}_s and \mathbf{Z}_r are length- p jointly Gaussian vectors, independent of $\mathbf{B} \sim B$ i.i.d. elementwise, with $\mathbb{E}[\mathbf{Z}_s] = \mathbb{E}[\mathbf{Z}_r] = \mathbf{0}$, $\mathbb{E}\{([\mathbf{Z}_s]_i)^2\} = \mathbb{E}\{([\mathbf{Z}_r]_i)^2\} = 1$ for any element $i \in [p]$, and $\mathbb{E}\{[\mathbf{Z}_s]_i[\mathbf{Z}_r]_j\} = \frac{\Sigma_{s,r}}{\tau_r \tau_s} \mathbb{I}\{i = j\}$. Note that $\Sigma_{t,t} = \tau_t^2$ defined in (2.4).

Using the above covariances, we have the following result that characterizes the asymptotic empirical distributions of the difference vectors defined in (F.1) and generalizes Lemma (4.1). This result follows by [8, Theorem 1].

Lemma 6.1. [8, Theorem 1] *Assuming that $\Sigma_{0,0}, \dots, \Sigma_{t+1,t+1} > \sigma_w^2$, then for any deterministic sequence $\phi_p : (\mathbb{R}^p \times \mathbb{R}^n)^t \times \mathbb{R}^p \rightarrow \mathbb{R}$ of uniformly pseudo-Lipschitz functions of order k ,*

$$\text{plim}_p \left(\phi_p(\boldsymbol{\beta}, \mathbf{m}^0, \mathbf{h}^1, \dots, \mathbf{m}^t, \mathbf{h}^{t+1}) - \mathbb{E}[\phi_p(\boldsymbol{\beta}, \sqrt{\tau_0^2 - \sigma_w^2} \mathbf{Z}'_0, \tau_0 \mathbf{Z}_0, \dots, \sqrt{\tau_t^2 - \sigma_w^2} \mathbf{Z}'_t, \tau_t \mathbf{Z}_t)] \right) = 0,$$

for $(\mathbf{Z}_0, \mathbf{Z}_1, \dots, \mathbf{Z}_t)$ defined in (6.3) independent of $(\mathbf{Z}'_0, \mathbf{Z}'_1, \dots, \mathbf{Z}'_t)$ and the expectation is taken with respect to the collection $(\mathbf{Z}_0, \mathbf{Z}'_0, \mathbf{Z}_1, \mathbf{Z}'_1, \dots, \mathbf{Z}'_t, \mathbf{Z}_t)$. We note that \mathbf{Z}'_s and \mathbf{Z}'_r are length- n jointly Gaussian vectors, with $\mathbb{E}[\mathbf{Z}'_s] = \mathbb{E}[\mathbf{Z}'_r] = \mathbf{0}$, $\mathbb{E}\{([\mathbf{Z}'_s]_i)^2\} = \mathbb{E}\{([\mathbf{Z}'_r]_i)^2\} = 1$ for any element $i \in [n]$, and $\mathbb{E}\{[\mathbf{Z}'_s]_i[\mathbf{Z}'_r]_j\} = (\Sigma_{s,r} - \sigma_w^2)((\tau_r^2 - \sigma_w^2)(\tau_s^2 - \sigma_w^2))^{-1/2} \mathbb{I}\{i = j\}$.

We use Lemma 6.1 to explicitly state asymptotic characterizations of AMP quantities that will be useful in our analysis.

Lemma 6.2. *Under the condition of Theorem 3, for \mathbf{z}^t and $\boldsymbol{\beta}^{t+1}$ defined in (1.3) and the generalized state evolution sequence defined in (6.3),*

$$\text{plim}_n \left(\frac{1}{n} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|^2 - (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2) \right) = 0, \quad (6.4)$$

$$\text{plim}_p \left(\frac{1}{\delta p} \|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t\|^2 - (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2) \right) = 0. \quad (6.5)$$

Proof. The major tools in proving (6.4)-(6.5) are first recognizing that we can write the differences $\mathbf{z}^t - \mathbf{z}^{t-1}$ and $\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t$ as a function of the values $(\boldsymbol{\beta}, \mathbf{m}^0, \mathbf{h}^1, \dots, \mathbf{m}^t, \mathbf{h}^{t+1})$ defined in (F.1) and finally making an appeal to the Law of Large Numbers. We prove (6.5) and (6.4) follows similarly.

By (1.3a), $\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t = \eta_p^t(\boldsymbol{\beta}^t + \mathbf{X}^\top \mathbf{z}^t) - \eta_p^{t-1}(\boldsymbol{\beta}^{t-1} + \mathbf{X}^\top \mathbf{z}^{t-1}) = \eta_p^t(\boldsymbol{\beta} - \mathbf{h}^{t+1}) - \eta_p^{t-1}(\boldsymbol{\beta} - \mathbf{h}^t)$. Therefore, we will appeal to Lemma 6.1 for the uniformly pseudo-Lipschitz function

$$\phi_p(\boldsymbol{\beta}, \mathbf{m}^0, \mathbf{h}^1, \dots, \mathbf{m}^t, \mathbf{h}^{t+1}) = \frac{1}{\delta p} \|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t\|^2 = \frac{1}{\delta p} \|\eta_p^t(\boldsymbol{\beta} - \mathbf{h}^{t+1}) - \eta_p^{t-1}(\boldsymbol{\beta} - \mathbf{h}^t)\|^2.$$

We note that it easy to show that the above function is uniformly pseudo-Lipschitz, though we don't do this here. Then by Lemma 6.1,

$$\text{plim}_p \left(\frac{1}{\delta p} \|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t\|^2 - \frac{1}{\delta p} \mathbb{E} \|\eta_p^t(\boldsymbol{\beta} - \tau_t \mathbf{Z}_t) - \eta_p^{t-1}(\boldsymbol{\beta} - \tau_{t-1} \mathbf{Z}_{t-1})\|^2 \right) = 0. \quad (6.6)$$

Now to prove result (6.4), we note that by Lemma 3.2,

$$\text{plim}_{\delta p} \frac{1}{p} \mathbb{E} \|\eta_p^t(\boldsymbol{\beta} - \tau_t \mathbf{Z}_t) - \eta_p^{t-1}(\boldsymbol{\beta} - \tau_{t-1} \mathbf{Z}_{t-1})\|^2 = \lim_p \frac{1}{\delta p} \mathbb{E} \|\eta_p^t(\mathbf{B} - \tau_t \mathbf{Z}_t) - \eta_p^{t-1}(\mathbf{B} - \tau_{t-1} \mathbf{Z}_{t-1})\|^2,$$

where $\mathbf{B} \sim B$ i.i.d. elementwise independent of \mathbf{Z}_t and \mathbf{Z}_{t-1} . The argument for showing that the assumptions of Lemma 3.2 are met follows like that used in Appendix B in the proof of Proposition (P2) introduced in Section 4. Then, $\lim_p \frac{1}{\delta p} \mathbb{E} \|\eta_p^t(\mathbf{B} - \tau_t \mathbf{Z}_t) - \eta_p^{t-1}(\mathbf{B} - \tau_{t-1} \mathbf{Z}_{t-1})\|^2 = \Sigma_{t,t} - 2\Sigma_{t,t-1} + \Sigma_{t-1,t-1}$.

□

We finally state a lemma that characterizes the asymptotic value of the normalized ℓ_2 norm of the residuals in AMP algorithm (1.3b) following from Lemma 4.1.

Lemma 6.3. *For \mathbf{z}^t defined in (1.3b) and τ_t^2 given in (2.4),*

$$\text{plim}_n (\|\mathbf{z}^t\|^2/n - \tau_t^2) = 0. \quad (6.7)$$

Proof. This follows from Lemma 4.1, using the uniformly pseudo-Lipschitz (of order 2) sequence of functions $\phi_n(\mathbf{a}, \mathbf{b}) = \frac{1}{n}\|\mathbf{a}\|^2$ to get, $\text{plim}_n \|\mathbf{z}^t\|^2/n = \text{plim}_n \mathbb{E}_{\mathbf{Z}}[\|\mathbf{w} + \sqrt{\tau_t^2 - \sigma_w^2} \mathbf{Z}\|^2]/n$ for $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I})$. Then the final result follows by noticing that $\mathbb{E}_{\mathbf{Z}}\|\mathbf{w} + \sqrt{\tau_t^2 - \sigma_w^2} \mathbf{Z}\|^2 = \|\mathbf{w}\|^2 + (\tau_t^2 - \sigma_w^2)\mathbb{E}_{\mathbf{Z}}\|\mathbf{Z}\|^2 = \|\mathbf{w}\|^2 + n(\tau_t^2 - \sigma_w^2)$, and therefore, using that $\text{plim}_n \|\mathbf{w}\|^2/n = \sigma_w^2$ by the Law of Large Numbers,

$$\text{plim}_n \frac{1}{n} \mathbb{E}_{\mathbf{Z}}\|\mathbf{w} + \sqrt{\tau_t^2 - \sigma_w^2} \mathbf{Z}\|^2 = (\tau_t^2 - \sigma_w^2) + \text{plim}_n \frac{1}{n} \|\mathbf{w}\|^2 = \tau_t^2.$$

□

7 Verification of Main Technical Lemma Conditions

We now verify that the Lemma 5.5 conditions 1-5 are met for the SLOPE cost function and the associated AMP algorithm. We note that conditions 1, 4, and 5 are straightforward, so their proof is presented first. On the other hand, condition 2 and condition 3 are quite technical. Their proofs are given in Section 7.4 and Section 7.5 below.

7.1 Condition (4)

This follows by standard limit theorems about the singular values of Wishart matrices (see Appendix G, Theorem H.2).

7.2 Condition (5)

Recall, $\mathcal{C}_{\mathbf{x}}(\mathbf{b}) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \sum_{i=1}^p \hat{\lambda}_i |b_i|$ for some $\hat{\lambda} \in \mathcal{P}(\hat{\Pi}_{\mathbf{x}}^{-1}(\boldsymbol{\lambda}))$, and by definition, $\mathcal{C}_{\mathbf{x}}(\mathbf{x}) = \mathcal{C}(\mathbf{x})$ for all \mathbf{x} . Since $\hat{\beta}$ is the minimizer of $\mathcal{C}(\cdot)$ we have $\mathcal{C}(\beta^t) \geq \mathcal{C}(\hat{\beta})$ and by the rearrangement inequality, $\mathcal{C}_{\hat{\beta}}(\hat{\beta}) \geq \mathcal{C}_{\beta^t}(\hat{\beta})$. Therefore, $\mathcal{C}(\beta^t) \geq \mathcal{C}(\hat{\beta}) = \mathcal{C}_{\hat{\beta}}(\hat{\beta}) \geq \mathcal{C}_{\beta^t}(\hat{\beta})$.

7.3 Condition (1)

Condition (1) follows, for large enough p , from Lemma 7.1, stated below, which proves the asymptotic boundedness of the norms of the AMP estimates β^t and the SLOPE estimate $\hat{\beta}$.

Lemma 7.1. *For any parameter vector $\boldsymbol{\lambda} \in \mathbb{R}^p$ defining a SLOPE cost as in (1.2), let $\boldsymbol{\alpha} = \boldsymbol{\alpha}(\boldsymbol{\lambda})$, then for $t \geq 0$,*

$$\text{plim}_p \frac{1}{p} \|\beta^t\|^2 = \text{plim}_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\|\eta_p^t(\beta + \tau_t \mathbf{Z})\|^2] \leq 2\sigma_{\beta}^2 + 2\tau_t^2, \quad (7.1)$$

for $\eta_p^t(\cdot)$ defined in (4.1) with $\sigma_{\beta}^2 := \mathbb{E}[B^2] < \infty$ and $\sigma_{\beta}^2 + \tau_*^2 < \infty$ and

$$\text{plim}_p \frac{1}{p} \|\hat{\beta}\|^2 \leq C, \quad (7.2)$$

where $C := C(\delta, \sigma_{\beta}^2, \sigma_w^2, B_{max}, B_{min}, \lambda_{min})$ is a positive constant depending on $\delta, \sigma_{\beta}^2, \sigma_w^2$, along with the singular values of \mathbf{X} through $B_{max} \geq \lim_p \sigma_{max}^2(\mathbf{X})$, and $B_{min} \leq \lim_p \hat{\sigma}_{min}^2(\mathbf{X})$, and a lower bound on the parameter values $\lambda_{min} := \lim_p \min(\boldsymbol{\lambda})$.

Proof. The proof is included in Appendix D. \square

7.4 Condition (2)

Condition (2) follows from Lemma 7.2 stated below, for ϵ arbitrarily small when t is large enough.

Lemma 7.2. *Under the conditions of Theorem 3, for every iteration t , there exists a subgradient $sg(C, \boldsymbol{\beta}^t)$ of C defined in (5.1) at point $\boldsymbol{\beta}^t$ such that almost surely,*

$$\lim_t \plim_p \frac{1}{p} \|sg(C, \boldsymbol{\beta}^t)\|^2 = 0.$$

The proof is an adaption of [5, Lemma 3.3], though, the subgradient for the SLOPE cost function (studied extensively in Section 5.1) is quite different than that of the LASSO cost and our analysis requires handling this carefully. Before we prove Lemma 7.2, we state and prove a result which tells us that the asymptotic difference between the AMP output at any two iterations t and $t - 1$ goes to zero in ℓ_2 norm as the algorithm runs. This result is crucial to the proof of Lemma 7.2.

Lemma 7.3. *Under the condition of Theorem 3, the estimates $\{\boldsymbol{\beta}^t\}_{t \geq 0}$ and residuals $\{\mathbf{z}^t\}_{t \geq 0}$ of AMP almost surely satisfy*

$$\lim_t \plim_p \frac{1}{\delta p} \|\boldsymbol{\beta}^t - \boldsymbol{\beta}^{t-1}\|^2 = 0, \quad \text{and} \quad \lim_t \plim_p \frac{1}{n} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|^2 = 0$$

Proof of Lemma 7.3. This result uses Lemma 6.2, which characterizes the large system limit of $\frac{1}{n} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|^2$ and $\frac{1}{\delta p} \|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t\|^2$ as both being equal to $\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2$ where $\Sigma_{t,t-1}$ is the generalized state evolution sequence defined in (6.3). Then Lemma E.1 (which is stated and proved in Appendix E) shows that $\lim_t (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2) = 0$. \square

Proof of Lemma 7.2. For any vector $\boldsymbol{\nu}^t \in \partial J_{\boldsymbol{\lambda}}(\boldsymbol{\beta}^t)$, note that $\boldsymbol{\nu}^t - \mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t)$ is a valid subgradient belonging to the set $\partial \mathcal{C}(\boldsymbol{\beta}^t)$ as defined in Fact 5.1. Moreover, by AMP (1.3b), $\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t = \mathbf{z}^t - \omega^t \mathbf{z}^{t-1}$ with $\omega^t := \frac{1}{\delta p} [\nabla \eta^{t-1}(\boldsymbol{\beta}^{t-1} + \mathbf{X}^\top \mathbf{z}^{t-1})]$. Therefore we can write,

$$\begin{aligned} \boldsymbol{\nu}^t - \mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t) &= \boldsymbol{\nu}^t - \mathbf{X}^\top(\mathbf{z}^t - \omega^t \mathbf{z}^{t-1}) = \boldsymbol{\nu}^t - \mathbf{X}^\top(\mathbf{z}^t - \mathbf{z}^{t-1}) - (1 - \omega^t) \mathbf{X}^\top \mathbf{z}^{t-1} \\ &= (\boldsymbol{\nu}^t - \mu_t \mathbf{X}^\top \mathbf{z}^{t-1}) - \mathbf{X}^\top(\mathbf{z}^t - \mathbf{z}^{t-1}) + (\mu_t - (1 - \omega^t)) \mathbf{X}^\top \mathbf{z}^{t-1}, \end{aligned} \tag{7.3}$$

where we define $\mu_t := \langle \boldsymbol{\lambda}, \boldsymbol{\theta}_{t-1} \rangle / \|\boldsymbol{\theta}_{t-1}\|^2$ as the ratio of $\boldsymbol{\lambda}$ to $\boldsymbol{\theta}_{t-1}$ so that $\boldsymbol{\lambda} = \mu_t \boldsymbol{\theta}_{t-1}$ (here $\boldsymbol{\theta}_{t-1} := \boldsymbol{\alpha} \tau_{t-1}$ and recall that $\boldsymbol{\alpha}$ is calibrated to be parallel to $\boldsymbol{\lambda}$). It follows that $\partial J_{\boldsymbol{\lambda}}(\mathbf{x}) = \mu_t \partial J_{\boldsymbol{\theta}_{t-1}}(\mathbf{x})$.

Now, by the definition of the proximal operator used in (1.3a) and by Fact 5.2, we have that $(\mathbf{X}^\top \mathbf{z}^{t-1} + \boldsymbol{\beta}^{t-1}) - \boldsymbol{\beta}^t \in \partial J_{\boldsymbol{\theta}^{t-1}}(\boldsymbol{\beta}^t)$. Hence we choose $\boldsymbol{\nu}^t$ to be the specific subgradient defined by

$$\boldsymbol{\nu}^t = \mu_t (\mathbf{X}^\top \mathbf{z}^{t-1} + \boldsymbol{\beta}^{t-1} - \boldsymbol{\beta}^t) \in \partial J_{\boldsymbol{\lambda}}(\boldsymbol{\beta}^t), \tag{7.4}$$

which leads to $\boldsymbol{\nu}^t - \mu_t \mathbf{X}^\top \mathbf{z}^{t-1} = \mu_t (\boldsymbol{\beta}^{t-1} - \boldsymbol{\beta}^t)$. Plugging into (7.3),

$$\boldsymbol{\nu}^t - \mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t) = \mu_t (\boldsymbol{\beta}^{t-1} - \boldsymbol{\beta}^t) - \mathbf{X}^\top(\mathbf{z}^t - \mathbf{z}^{t-1}) + (\mu_t - (1 - \omega^t)) \mathbf{X}^\top \mathbf{z}^{t-1}. \tag{7.5}$$

Then taking the norm, dividing by \sqrt{p} , and using the triangular inequality, we have

$$\frac{1}{\sqrt{p}} \|\boldsymbol{\nu}^t - \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t)\| \leq \frac{\mu_t}{\sqrt{p}} \|\boldsymbol{\beta}^{t-1} - \boldsymbol{\beta}^t\| + \frac{1}{\sqrt{p}} \|\mathbf{X}^\top (\mathbf{z}^t - \mathbf{z}^{t-1})\| + \frac{(\mu_t - (1 - \omega^t))}{\sqrt{p}} \|\mathbf{X}^\top \mathbf{z}^{t-1}\|.$$

Using Lemma 6.2, that $\sigma_{\max}(\mathbf{X})$ is almost surely bounded as $p \rightarrow \infty$ (cf. Theorem 2), and that $\lim_t \lim_p \mu_t = 1 - \lim_p \frac{1}{\delta p} \mathbb{E} \| \text{prox}_{J_{\mathbf{A}(p)\tau_*}}(\mathbf{B} + \tau_* \mathbf{Z}) \|_0^*$ as in (2.11) is finite, the first two terms on the right side of the above $\rightarrow 0$. Finally, for the third term, Lemma 6.3 gives $\lim_t \text{plim}_p \|z^t\|/\sqrt{p} = \tau_*$, and together with the calibration formula (2.11), that $\sigma_{\max}(\mathbf{X})$ is almost surely bounded as $p \rightarrow \infty$, and the definition of ω in the proof of Lemma 2.2, we find $\lim_t \lim_p (\mu_t - (1 - \omega^t)) = 0$, and thus the third term $\rightarrow 0$. As $\boldsymbol{\nu}^t - \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^t) \in \partial \mathcal{C}(\boldsymbol{\beta}^t)$, the proof is complete. \square

7.5 Condition (3)

We take $\boldsymbol{\nu}^t$ to be the subgradient defined in (7.4) and since t is fixed, we drop the superscript t writing $\boldsymbol{\nu} := \boldsymbol{\nu}^t$. Recall the sets $s_t(c_2)$ and $S_t(c_2)$ defined in Condition (3). Then for s' being any set of maximal atoms in $[p]$ with $|s'| \leq c_3 p$ and $S' := \{i \in I : I \in s'\}$, we would like to show $\sigma_{\min}(\mathbf{X}_{S_t(c_2) \cup S'}) \geq c_4$. This holds by Proposition 7.4, stated below, whose proof is the main challenge. We state the proposition and then we identify two auxiliary lemmas, Lemma 7.5 and 7.6, that will be used to ultimately prove Proposition 7.4.

Proposition 7.4. *There exist constants $c_2 \in (0, 1)$, $c_3, c_4 > 0$ and $t_{\min} < \infty$ such that, for any $t \geq t_{\min}$, and set S_t defined in Condition (3)*

$$\min_{s'} \{\sigma_{\min}(\mathbf{X}_{S_t(c_2) \cup S'}) : S' \subseteq [p], |s'| \leq c_3 p, S' = \{i \in I : I \in s'\}\} \geq c_4$$

eventually almost surely as $p \rightarrow \infty$.

The proof of Proposition 7.4 will use two auxiliary lemmas, Lemma 7.5 and 7.6, stated below.

Lemma 7.5. *Let the set s_t be measurable on the σ -algebra \mathfrak{S}_t generated by $\{\mathbf{z}^0, \dots, \mathbf{z}^{t-1}\}$ and $\{\boldsymbol{\beta}^0 + \mathbf{X}^* \mathbf{z}^0, \dots, \boldsymbol{\beta}^{t-1} + \mathbf{X}^* \mathbf{z}^{t-1}\}$ and assume $|s_t| \leq p(\delta - c)$ for some $c > 0$. Define $S_t \subseteq [p]$ as $\{i \in I \text{ for some } I \in s_t\}$. Then there exists $a_1 = a_1(c) > 0$ (independent of t) and $a_2 = a_2(c, t) > 0$ (depending on t and c) such that*

$$\min_{s'} \{\sigma_{\min}(\mathbf{X}_{S_t \cup S'}) : S' \subseteq [p], |s'| \leq a_1 p, S' = \{i \in I : I \in s'\}\} \geq a_2,$$

eventually almost surely as $p \rightarrow \infty$.

Proof. The proof of Lemma 7.5 is given in Appendix F. The key difference in SLOPE case (Lemma 7.5) and LASSO case (cf. [5, Lemma 3.4]) is the concept of equivalence classes of indices. On a high level, the set s describes some structure in the support space S and such structure restricts the dimension of some linear spaces in the proof of Lemma 7.5. \square

Lemma 7.6. *[5, Lemma 3.5] Fix $\gamma \in (0, 1)$ and let the sequence $\{S_t(\gamma)\}_{t \geq 0}$ be defined as before. For any $\xi > 0$ there exists $t_* = t_*(\xi, \gamma) < \infty$ such that, for all $t_2 \geq t_1 \geq t_*$ fixed, we have*

$$\frac{1}{p} |S_{t_2}(\gamma) \setminus S_{t_1}(\gamma)| < \xi, \tag{7.6}$$

eventually almost surely as $p \rightarrow \infty$.

Proof. For LASSO, this result was given in [5, Lemma 3.5], and for SLOPE, the proof stays largely the same so we don't repeat it here. The major difference is that where the work in [5] can appeal to AMP analysis in [4], for SLOPE, we appeal to similar results given in [8] (e.g. Lemma 6.1). \square

Proof of Proposition 7.4. The subgradient in Condition (2) is given by $sg(\mathcal{C}, \beta^t) := \nu^t - \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta^t)$ where $\nu^t \in \partial J_\lambda(\beta^t)$ is the subgradient defined in the Condition (2) proof at Eq. (7.4). Recall, $S_t(c_2) = \{i \in I : |\nu_I^t| \geq \mathcal{P}([\hat{\Pi}_{\beta^t}^{-1}(\lambda)]_I)(1 - c_2)\}$. We include a simple visualization for the set $S_t(c_2)$ in Figure 4. We have plotted the subgradient $\nu_I^t = (-1, 2)$ for (zero) equivalence class $I = \{1, 2\}$ when $\lambda = (4, 1)$ and $\beta^t = (0, 0)$. Then indices of $|\nu_I^t|$, namely $(1, 2)$ are in $S_t(c_2)$ unless $c_2 < 0.4$.

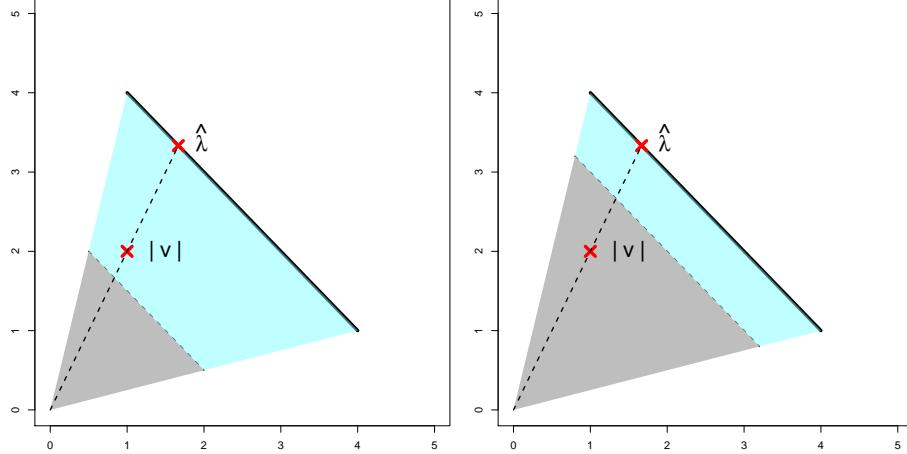


Figure 4: Left: $c_2 = 0.5$; Right: $c_2 = 0.2$; Blue area is $\{\nu \in \partial J_\lambda(0, 0) : |\nu| \geq (1 - c_2)\mathcal{P}(\lambda_1, \lambda_2)\}$ and grey area is complement of blue area in $\partial J_\lambda(0, 0)$.

We know from the proof of Lemma 7.2 Eq. (7.4) that $\nu^t = \mu_t(\mathbf{X}^\top \mathbf{z}^{t-1} + \beta^{t-1} - \beta^t) \in \mu_t J_{\theta^t}(\beta^t)$ where $\mu_t := \langle \lambda, \theta_{t-1} \rangle / \|\theta_{t-1}\|^2$ and $\lambda = \mu_t \theta^{t-1}$. Therefore, summing over all equivalence classes I ,

$$\begin{aligned} |s_t(c_2)| &= \sum_I \mathbb{I}\{|\nu_I^t| \geq \mathcal{P}([\hat{\Pi}_{\beta^t}^{-1}(\lambda)]_I)(1 - c_2)\} \\ &= \sum_I \mathbb{I}\{|\beta^t - [\mathbf{X}^\top \mathbf{z}^{t-1}] - \beta^{t-1}|_I \geq \mathcal{P}([\hat{\Pi}_{\beta^t}^{-1}(\theta^{t-1})]_I)(1 - c_2)\}. \end{aligned} \quad (7.7)$$

As detailed in the proof of Lemma 5.5, for non-zero equivalence classes, let $\hat{\lambda}_I = |\nu_I|$, and for the zero equivalence class, let $\hat{\lambda}_I \geq |\nu_I|$, meaning $\hat{\lambda}_I$ is parallel to $|\nu_I|$ for each equivalence class I of β^t . That such a $\hat{\lambda}$ exists in the set $\mathcal{P}(\hat{\Pi}_{\beta^t}^{-1}(\lambda))$ follows since ν is a valid subgradient of $J_\lambda(\beta^t)$ (see Fact 5.3). We can then simplify the set definitions of $s_t(c_2)$ and $S_t(c_2)$ to be $s_t(c_2) := \{I \subset [p] : |\nu_I| \geq (1 - c_2)\hat{\lambda}_I\}$ and $S_t(c_2) := \{i : |\nu_i| \geq (1 - c_2)\hat{\lambda}_i\}$, where both use equivalence classes, I , defined for β^t . Then since $\lambda = \mu_t \theta^{t-1}$, we also let $\hat{\theta}^{t-1}$ be defined such that $\hat{\lambda} = \mu_t \hat{\theta}^{t-1}$.

Therefore, by (7.7), $|s_t(c_2)| = \sum_I \mathbb{I}\{|\beta^t - [\mathbf{X}^\top \mathbf{z}^{t-1}] - \beta^{t-1}|_I \geq \hat{\theta}_I^{t-1}(1 - c_2)\}$. In the notation of (F.1), $\beta^t - [\mathbf{X}^\top \mathbf{z}^{t-1}] - \beta^{t-1} = \mathbf{h}^t + \eta^{t-1}(\beta - \mathbf{h}^t) - \beta$ and $\beta^t = \eta^{t-1}(\beta - \mathbf{h}^t)$ and therefore by (7.7),

$$|s_t(c_2)| = \sum_I \mathbb{I}\{|\mathbf{h}^t + \eta^{t-1}(\beta - \mathbf{h}^t) - \beta|_I \geq \hat{\theta}_I^{t-1}(1 - c_2)\}.$$

Now, we note that Lemma 6.1 implies weak convergence of the empirical distribution of \mathbf{h}^t to $\tau_{t-1}\mathbf{Z}_{t-1}$ for \mathbf{Z}_{t-1} a vector of i.i.d. standard Gaussian and τ_{t-1} given by the state evolution (2.4). Therefore a careful argument using continuous approximations to indicators gives,

$$\begin{aligned} & \text{plim}_p \frac{1}{p} \sum_I \mathbb{I}\left\{|\mathbf{h}^t + \eta^{t-1}(\boldsymbol{\beta} - \mathbf{h}^t) - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right\} \\ &= \lim_p \mathbb{E}_{\mathbf{Z}_{t-1}} \left\{ \frac{1}{p} \sum_I \mathbb{I}\left\{|\tau_{t-1}\mathbf{Z}_{t-1} + \eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}) - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right\} \right\}, \end{aligned} \quad (7.8)$$

where in the right side of the above, the equivalence classes I are taken with respect to $\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})$ and $\hat{\boldsymbol{\theta}}_I^{t-1}$ as equal to or larger than $|\tau_{t-1}\mathbf{Z}_{t-1} + \eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}) - \boldsymbol{\beta}|_I$ depending on whether I is the zero equivalence class or not. We justify the substitution of $\tau_{t-1}\mathbf{Z}_{t-1}$ for \mathbf{h}^t by approximating the sum of indicators with a function that counts the number of elements in $\eta^{t-1}(\boldsymbol{\beta} - \mathbf{h}^t)$ that are strictly greater than its neighbour. Then this function converges to a continuous and bounded function, the function that measures the proportion of η^{t-1} that is non-flat, to which we apply the Portmanteau Theorem (cf. [20], Lemma 1(b) in [4] and Lemma F.3(b) in [5]).

Now, using (7.8), we can simplify:

$$\text{plim}_p \frac{1}{p} |s_t(c_2)| = \lim_p \frac{1}{p} \sum_I \mathbb{P}_{\mathbf{Z}_{t-1}} \left(|\tau_{t-1}\mathbf{Z}_{t-1} - \eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}) - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2) \right), \quad (7.9)$$

and we study the probability on the right side of the above, for a fixed equivalence class I , writing $\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})$ to be η^{t-1} , dropping the input.

$$\begin{aligned} & \mathbb{P}\left(|\tau_{t-1}\mathbf{Z}_{t-1} + \eta^{t-1} - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right) \\ &= \mathbb{P}\left(|\tau_{t-1}\mathbf{Z}_{t-1} + \eta^{t-1} - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2), \eta_I^{t-1} = \mathbf{0}\right) \\ &\quad + \mathbb{P}\left(|\tau_{t-1}\mathbf{Z}_{t-1} + \eta^{t-1} - \boldsymbol{\beta}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2), \eta_I^{t-1} \neq \mathbf{0}\right) \\ &\stackrel{(a)}{=} \mathbb{P}\left(\hat{\boldsymbol{\theta}}_I^{t-1} \geq |\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right) + \mathbb{P}\left(\hat{\boldsymbol{\theta}}_I^{t-1} \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right) \mathbb{P}(\eta_I^{t-1} \neq \mathbf{0}). \\ &= \mathbb{P}\left(\hat{\boldsymbol{\theta}}_I^{t-1} \geq |\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}|_I \geq \hat{\boldsymbol{\theta}}_I^{t-1}(1 - c_2)\right) + \mathbb{P}(\eta_I^{t-1} \neq \mathbf{0}). \end{aligned} \quad (7.10)$$

In the above, step (a) follows when $\eta_I^{t-1} = [\text{prox}_{J_{\boldsymbol{\theta}^{t-1}}}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})]_I = \mathbf{0}$, since we must have $|\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}|_I \leq \hat{\boldsymbol{\theta}}_I^{t-1}$, and when $\eta_I^{t-1} \neq \mathbf{0}$, by Fact 5.2 and Fact 5.3, we know that $|\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1}) - (\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})|_I \in \mathcal{P}([\hat{\Pi}_{\eta^{t-1}}^{-1}(\boldsymbol{\theta}^{t-1})]_I)$.

It obvious that one can make the first probability arbitrarily small by bringing c_2 to 0. To see this, say $1 \in I$ and notice that $\mathcal{P}([\hat{\Pi}_{\eta^{t-1}}^{-1}(\boldsymbol{\theta}^{t-1})]_I)$ always has Lebesgue measure 0 because it is a subset of the hyperplane $\{\mathbf{x} \in \mathbb{R}^p : \sum_{j \in I} x_j = \sum_{j \in I} \theta_j^{t-1}\}$.

On the other hand, notice that

$$\sum_I \mathbb{P}([\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})]_I \neq \mathbf{0}) = \sum_I \mathbb{E}\{\mathbb{I}([\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})]_I \neq \mathbf{0})\} = \mathbb{E}_{\mathbf{Z}_{t-1}} \|\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1}\mathbf{Z}_{t-1})\|_0^*,$$

and that η^{t-1} is asymptotically separable by Lemma 3.3. Define $h^{t-1}(x) = h(x; B + \tau_{t-1}Z, \Theta^{t-1})$ with Θ^{t-1} being the distribution to which the empirical distribution of $\boldsymbol{\theta}^{t-1}$ converges, and also define

$$\mathbf{W}_{t-1} := \left\{ x \mid h^{t-1}(x) \neq 0 \text{ and } m\{z \mid |h^{t-1}(z)| = |h^{t-1}(x)|\} = 0 \right\}$$

similarly to (2.12), where m is the Lebesgue measure. Then,

$$\begin{aligned} \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}_{t-1}} \|\eta^{t-1}(\boldsymbol{\beta} - \tau_{t-1} \mathbf{Z}_{t-1})\|_0^* &= \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}_{t-1}} \|h^{t-1}(\boldsymbol{\beta} - \tau_{t-1} \mathbf{Z}_{t-1})\|_0^* \\ &= \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}_{t-1}} \sum_{i=1}^p \mathbb{I}\{(\beta_i - \tau_{t-1} Z_{t-1,i}) \in \mathbf{W}_{t-1}\} = \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}_{t-1}, \mathbf{B}} \|\eta^{t-1}(\mathbf{B} - \tau_{t-1} \mathbf{Z}_{t-1})\|_0^*, \end{aligned}$$

where the last equality holds by Lemma 3.2.

Then (2.10) gives this term is smaller than δ for large t . Hence, by (7.9) and (7.10),

$$\begin{aligned} \text{plim}_p \frac{1}{p} |s_t(c_2)| \\ = \lim_p \frac{1}{p} \sum_I \mathbb{P}(\hat{\theta}_I^{t-1} \geq |\boldsymbol{\beta} - \tau_{t-1} \mathbf{Z}_{t-1}|_I \geq \hat{\theta}_I^{t-1}(1 - c_2)) + \lim_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}_{t-1}, \mathbf{B}} \|\eta^{t-1}(\mathbf{B} - \tau_{t-1} \mathbf{Z}_{t-1})\|_0^*, \end{aligned}$$

Therefore, for some $c > 0$, choose $c_2 \in (0, 1)$ such that the first term on the right side of the above is arbitrarily small along with $t_{\min,1}(c)$ such that the second term is arbitrarily close to δ , meaning

$$\lim_p \mathbb{P}\left(\frac{1}{p} |s_t(c_2)| < \delta - c\right) = 1,$$

for all fixed t larger than some $t_{\min,1}(c)$.

For any $t \geq t_{\min,1}(c)$ we can apply Lemma 7.5 for some $a_1(c)$, $a_2(c, t)$. Note this doesn't immediately give the result we use since the lower bound, a_2 , depends on t . To get around this we additionally appeal to Lemma 7.6 that tells us after some time t_* , the supports of the AMP estimates don't change appreciably. Now we fix $c > 0$ and consequently $a_1 = a_1(c)$ is fixed. Define $t_{\min} = \max(t_{\min,1}, t_*(a_1/2, c_2))$ with $t_*(\cdot)$ defined as in Lemma 7.6 and let $a_2 = a_2(c, t_{\min})$. Then, by Lemma 7.5 and the fact that $a_2(c, t)$ is non-increasing in t ,

$$\min \{\sigma_{\min}(\mathbf{X}_{S_{t_{\min}}(c_2) \cup S'}) : S' \subseteq [p], |s'| \leq a_1 p\} \geq a_2.$$

In addition, by Lemma 7.6, $|S_t(c_2) \setminus S_{t_{\min}}(c_2)| \leq p a_1/2$. Both events hold eventually almost surely as $p \rightarrow \infty$. The proof completes with $c_3 = a_1(c)/2$ and $c_4 = a_2(c, t_{\min})$, fixed with respect to t . \square

8 Discussion and Future Work

This work develops and analyzes the dynamics of an approximate message passing (AMP) algorithm with the purpose of solving the SLOPE convex optimization procedure for high-dimensional linear regression. By employing recent theoretical analysis of AMP when the non-linearities used in the algorithm are non-separable [8], as is the case for the SLOPE problem, we provide rigorous proof that the proposed AMP algorithm finds the SLOPE solution asymptotically. Moreover empirical evidence suggests that the AMP estimate is already very close to the SLOPE solution even in few iterations. By leveraging our analysis showing AMP provably solves SLOPE, we provide an exact asymptotic characterization of the ℓ_2 risk of the SLOPE estimator from the underlying truth and insight into other statistical properties of the SLOPE estimator. Though this asymptotic analysis of the SLOPE solution has been demonstrated in other recent work [20] using a different proof strategy, we believe that our AMP-based approach offers a more concrete and algorithmic understanding of the finite-sample behavior of the SLOPE estimator.

A limitation of this approach is that the theory assumes an i.i.d. Gaussian measurement matrix, and moreover, the AMP algorithm can become unstable when the measurement matrix is far from i.i.d., creating the need for heuristic techniques to provide convergence in applications where the measurement matrix is generated by nature (i.e., a real-world experiment or observational study). Additionally, the asymptotical regime studied here, $n/p \rightarrow \delta \in (0, \infty)$, requires that the number of columns of the measurement matrix p grow at the same rate as the number of rows n . It is of practical interest to extend the results to high-dimensional settings where p grows faster than n .

References

- [1] Z. Bai and Y. Yin. Limit of the smallest eigenvalue of a large dimensional sample covariance matrix. In *Advances In Statistics*, pages 108–127. World Scientific, 2008.
- [2] R. F. Barber and E. J. Candès. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [3] M. Bayati, M. A. Erdogdu, and A. Montanari. Estimating lasso risk and noise level. In *Advances in Neural Information Processing Systems*, pages 944–952, 2013.
- [4] M. Bayati and A. Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Trans. on Inf. Theory*, 57(2):764–785, 2011.
- [5] M. Bayati and A. Montanari. The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2011.
- [6] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [7] P. C. Bellec, G. Lecué, and A. B. Tsybakov. SLOPE meets lasso: improved oracle bounds and optimality. *The Annals of Statistics*, 46(6B):3603–3642, 2018.
- [8] R. Berthier, A. Montanari, and P.-M. Nguyen. State evolution for approximate message passing with non-separable functions. *arXiv preprint arXiv:1708.03950*, 2017.
- [9] M. Bogdan, E. Van Den Berg, C. Sabatti, W. Su, and E. J. Candès. SLOPE—adaptive variable selection via convex optimization. *The Annals of Applied Statistics*, 9(3):1103, 2015.
- [10] H. D. Bondell and B. J. Reich. Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with oscar. *Biometrics*, 64(1):115–123, 2008.
- [11] D. Brzyski, A. Gossman, W. Su, and M. Bogdan. Group SLOPE—adaptive selection of groups of predictors. *Journal of the American Statistical Association*, pages 1–15, 2018.
- [12] M. Celentano and A. Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *arXiv preprint arXiv:1903.10603*, 2019.
- [13] A. Chambolle, R. A. De Vore, N.-Y. Lee, and B. J. Lucier. Nonlinear wavelet image processing: variational problems, compression, and noise removal through wavelet shrinkage. *IEEE Transactions on Image Processing*, 7(3):319–335, 1998.
- [14] I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 57(11):1413–1457, 2004.
- [15] D. Donoho and A. Montanari. High dimensional robust m-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3-4):935–969, 2016.

- [16] D. L. Donoho, A. Maleki, and A. Montanari. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- [17] D. L. Donoho, A. Maleki, and A. Montanari. The noise-sensitivity phase transition in compressed sensing. *IEEE Transactions on Information Theory*, 57(10):6920–6941, 2011.
- [18] J. L. Doob. *Stochastic processes*, volume 101. New York Wiley, 1953.
- [19] M. Figueiredo and R. Nowak. Ordered weighted l₁ regularized regression with strongly correlated covariates: Theoretical aspects. In *Artificial Intelligence and Statistics*, pages 930–938, 2016.
- [20] H. Hu and Y. M. Lu. Asymptotics and optimal designs of SLOPE for sparse linear regression. *arXiv preprint arXiv:1903.11582*, 2019.
- [21] A. Javanmard and A. Montanari. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- [22] B. S. Kashin. Diameters of some finite-dimensional sets and classes of smooth functions. *Izvestiya Rossiiskoi Akademii Nauk. Seriya Matematicheskaya*, 41(2):334–351, 1977.
- [23] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová. Probabilistic reconstruction in compressed sensing: algorithms, phase diagrams, and threshold achieving matrices. *J. Stat. Mech. Theory Exp.*, (8), 2012.
- [24] M. Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2001.
- [25] A. E. Litvak, A. Pajor, M. Rudelson, and N. Tomczak-Jaegermann. Smallest singular value of random matrices and geometry of random polytopes. *Advances in Mathematics*, 195(2):491–523, 2005.
- [26] F. J. MacWilliams and N. J. A. Sloane. *The theory of error-correcting codes*, volume 16. Elsevier, 1977.
- [27] A. Montanari. Graphical models concepts in compressed sensing. In Y. C. Eldar and G. Kutyniok, editors, *Compressed Sensing*, pages 394–438. Cambridge University Press, 2012.
- [28] A. Mousavi, A. Maleki, R. G. Baraniuk, et al. Consistent parameter estimation for lasso and approximate message passing. *The Annals of Statistics*, 46(1):119–148, 2018.
- [29] N. Parikh, S. Boyd, et al. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, 2014.
- [30] L. D. Pitt. Positively correlated normal variables are associated. *The Annals of Probability*, pages 496–499, 1982.
- [31] S. Rangan. Generalized approximate message passing for estimation with random linear mixing. In *Proc. IEEE Int. Symp. Inf. Theory*, pages 2168–2172, 2011.
- [32] R. T. Rockafellar and R. J.-B. Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- [33] H. L. Royden. *Real analysis*. Krishna Prakashan Media, 1968.
- [34] W. Rudin et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- [35] C. Rush and R. Venkataramanan. Finite sample analysis of approximate message passing algorithms. *IEEE Trans. on Inf. Theory*, 64(11):7264–7286, 2018.
- [36] W. Su, M. Bogdan, and E. Candès. False discoveries occur early on the lasso path. *The Annals of Statistics*, 45(5):2133–2150, 2017.
- [37] W. Su and E. Candès. SLOPE is adaptive to unknown sparsity and asymptotically minimax. *The Annals of Statistics*, 44(3):1038–1068, 2016.

- [38] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 58(1):267–288, 1996.
- [39] X. Zeng and M. A. Figueiredo. Decreasing weighted sorted ℓ_1 regularization. *IEEE Signal Processing Letters*, 21(10):1240–1244, 2014.

A State Evolution Analysis

We first prove Theorem 1 and then provide a proof of Proposition 2.6.

A.1 Proving Theorem 1

Proof of Theorem 1. To begin with, we prove that $F(\tau^2, \alpha\tau)$ defined in (2.8) is concave with respect to τ^2 . The proof follows along the same lines as the proof of [5, Proposition 1.3], however, whereas the proof of [5, Proposition 1.3] proceeds by explicitly expressing the first derivative of the corresponding function F , and then differentiating on the explicit form to get the second derivative, in SLOPE case, because of the averaging that occurs within the proximal operation, it is extremely difficult to similarly derive an explicit form. To work around this, we keep all differentiation implicit. First,

$$\begin{aligned} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) &= \frac{\partial}{\partial \tau^2} [\sigma_w^2 + \frac{1}{\delta p} \mathbb{E} \|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2] \stackrel{(a)}{=} \frac{1}{\delta} \mathbb{E} \left\{ \frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2 \right\} \\ &= \frac{2}{\delta p} \sum_{i=1}^p \mathbb{E} \{ ([\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i - B_i) \frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i \}. \end{aligned} \quad (\text{A.1})$$

We note that the interchange between the derivative (a limit) and the expectation in step (a) of the above holds due to a dominated convergence argument that relies on the following lemma. First we introduce a bit of notation that will be used throughout the proof. Define an equivalence classes I_i for each index $i = \{1, 2, \dots, p\}$, defined as

$$I_i := \{j : |[\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_j| = |[\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i|\}.$$

For any $j \in I_i$, with the above definition, $I_j = I_i$. In general, we use I , without any specific index, to represent an entire equivalence class and let \mathcal{I} indicate the collection of unique equivalence classes.

Lemma A.1.

$$\left| \frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z}) - \mathbf{B}\|^2 \right| \leq \frac{1}{p} \sum_{I \in \mathcal{I}} \frac{1}{|I|} \left(\sum_{i \in I} |\text{sign}(B_i + \tau Z_i) Z_i - \alpha_i| \right)^2. \quad (\text{A.2})$$

Lemma A.1 will be proved below, after we solve $\frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i$.

Now we describe how the bound in Lemma A.1 can be used to produce the dominated convergence result needed in step (a) of (A.1). First note,

$$\begin{aligned} \frac{1}{p} \mathbb{E} \left\{ \sum_{I \in \mathcal{I}} \frac{1}{|I|} \left(\sum_{i \in I} |\text{sign}(B_i + \tau Z_i) Z_i - \alpha_i| \right)^2 \right\} &\leq \frac{1}{p} \mathbb{E} \left\{ \sum_{I \in \mathcal{I}} \sum_{i \in I} \left(|\text{sign}(B_i + \tau Z_i) Z_i - \alpha_i| \right)^2 \right\} \\ &\leq \frac{2}{p} \mathbb{E} \left\{ \sum_{I \in \mathcal{I}} \sum_{i \in I} (Z_i^2 + \alpha_i^2) \right\} = \frac{2}{p} \mathbb{E} \left\{ \sum_{i \in [p]} (Z_i^2 + \alpha_i^2) \right\} = 2 + 2\|\boldsymbol{\alpha}\|^2/p < \infty \end{aligned}$$

The first and second inequalities follow from $(\sum_{i=1}^n x_i)^2 \leq n \sum_i x_i^2$. The last inequality comes from entries of $\boldsymbol{\alpha}$ being finite and then $\|\boldsymbol{\alpha}\|^2/p \leq \max_i \alpha_i^2 < \infty$. Therefore we can invoke the dominated convergence theorem that allows the exchange of the derivative and expectation in step (a) of (A.1).

Now we want to further simplify (A.1). For each $1 \leq i \leq p$, we would like to study $\frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i$. We first note that the mapping $\tau^2 \mapsto [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i$ can be considered

as $f(g(\tau^2))$, where $g : \mathbb{R} \rightarrow \mathbb{R}^{2p}$ is defined as $y \mapsto g(y) := (\mathbf{B} + \mathbf{Z}\sqrt{y}, \boldsymbol{\alpha}\sqrt{y})$ and $f : \mathbb{R}^{2p} \rightarrow \mathbb{R}$ is defined as $(\mathbf{a}, \mathbf{b}) \mapsto f(\mathbf{a}, \mathbf{b}) := [\text{prox}_{J_b}(\mathbf{a})]_i$. Hence,

$$\frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i = \mathbf{J}_{f \circ g}(\tau^2) \stackrel{(a)}{=} \mathbf{J}_f(g(\tau^2))\mathbf{J}_g(\tau^2) = \left[\nabla_{\mathbf{a}} f(g(\tau^2)), \nabla_{\mathbf{b}} f(g(\tau^2)) \right] \left[\frac{\mathbf{Z}}{2\tau}, \frac{\boldsymbol{\alpha}}{2\tau} \right]^\top, \quad (\text{A.3})$$

where $\mathbf{J}_h \in \mathbb{R}^{m \times n}$ is the Jacobian matrix of a function $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and step (a) follows by the chain rule. We denote the proximal operator using a function $\eta : \mathbb{R}^{2p} \rightarrow \mathbb{R}^p$ as $\eta(\mathbf{a}, \mathbf{b}) := \text{prox}_{J_b}(\mathbf{a})$ and consider the partial derivatives of η with respect to its first and second arguments. Denote

$$\partial_1 \eta(\mathbf{a}, \mathbf{b}) := \text{diag} \left[\frac{\partial}{\partial a_1}, \frac{\partial}{\partial a_2}, \dots, \frac{\partial}{\partial a_p} \right] \eta(\mathbf{a}, \mathbf{b}), \quad \text{and} \quad \partial_2 \eta(\mathbf{a}, \mathbf{b}) := \text{diag} \left[\frac{\partial}{\partial b_1}, \frac{\partial}{\partial b_2}, \dots, \frac{\partial}{\partial b_p} \right] \eta(\mathbf{a}, \mathbf{b}). \quad (\text{A.4})$$

Recall that the derivatives computed in $\partial_1 \eta(\mathbf{a}, \mathbf{b})$ are defined in (2.2), and by anti-symmetry between two arguments, $\frac{d}{db_j}[\eta(\mathbf{a}, \mathbf{b})]_i = -\text{sign}([\eta(\mathbf{a}, \mathbf{b})]_j) \frac{d}{da_j}[\eta(\mathbf{a}, \mathbf{b})]_i$. Then using the result of (2.2):

$$\frac{\partial [\text{prox}_{J_{\lambda}}(\mathbf{v})]_i}{\partial v_j} = \frac{\partial [\eta(\mathbf{v}, \boldsymbol{\lambda})]_i}{\partial v_j} = \frac{\mathbb{I}\{|\eta(\mathbf{v}, \boldsymbol{\lambda})|_i = |\eta(\mathbf{v}, \boldsymbol{\lambda})|_j|\} \text{sign}([\eta(\mathbf{v}, \boldsymbol{\lambda})]_i[\eta(\mathbf{v}, \boldsymbol{\lambda})]_j)}{\#\{1 \leq k \leq p : |\eta(\mathbf{v}, \boldsymbol{\lambda})|_k = |\eta(\mathbf{v}, \boldsymbol{\lambda})|_i\}}$$

we have

$$\frac{d}{da_j} f(\mathbf{a}, \mathbf{b}) = \frac{d}{da_j} [\eta(\mathbf{a}, \mathbf{b})]_i = \mathbb{I}\{|\eta(\mathbf{a}, \mathbf{b})|_i = |\eta(\mathbf{a}, \mathbf{b})|_j|\} \text{sign}([\eta(\mathbf{a}, \mathbf{b})]_i[\eta(\mathbf{a}, \mathbf{b})]_j) [\partial_1 \eta(\mathbf{a}, \mathbf{b})]_i, \quad (\text{A.5})$$

and similarly,

$$\frac{d}{db_j} f(\mathbf{a}, \mathbf{b}) = \frac{d}{db_j} [\eta(\mathbf{a}, \mathbf{b})]_i = -\mathbb{I}\{|\eta(\mathbf{a}, \mathbf{b})|_i = |\eta(\mathbf{a}, \mathbf{b})|_j|\} \text{sign}([\eta(\mathbf{a}, \mathbf{b})]_i) [\partial_1 \eta(\mathbf{a}, \mathbf{b})]_i.$$

Now plugging the above into (A.3), we have

$$\begin{aligned} & \frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})]_i \\ &= \frac{1}{2\tau} [\partial_1 \eta(\mathbf{B} + \tau\mathbf{Z}, \alpha\tau)]_i \text{sign}([\eta(\mathbf{B} + \tau\mathbf{Z}, \alpha\tau)]_i) \sum_{j \in I_i} (\text{sign}([\eta(\mathbf{B} + \tau\mathbf{Z}, \alpha\tau)]_j) Z_j - \alpha_j) \end{aligned} \quad (\text{A.6})$$

In what follows, we drop the explicit statement of the $\eta(\cdot, \cdot)$ input to save space, writing η_i to mean $[\eta(\mathbf{B} + \tau\mathbf{Z}, \alpha\tau)]_i$ or $[\partial_1 \eta]_i$ to mean $[\partial_1 \eta(\mathbf{B} + \tau\mathbf{Z}, \alpha\tau)]_i$ for example. Using (A.6) in (A.1),

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \tau^2}(\tau^2, \alpha\tau) &= \frac{1}{\delta p \tau} \sum_{i=1}^p \sum_{j \in I_i} \mathbb{E} \left\{ (\eta_i - B_i) [\partial_1 \eta]_i \text{sign}(\eta_i) (\text{sign}(\eta_j) Z_j - \alpha_j) \right\} \\ &= \frac{1}{\delta p} \sum_{i=1}^p \sum_{j \in I_i} \mathbb{E} \left\{ ([\partial_1 \eta]_i)^2 + (\eta_i - B_i) [\partial_1^2 \eta]_i \right\} - \frac{1}{\delta p \tau} \sum_{i=1}^p \sum_{j \in I_i} \mathbb{E} \left\{ (\eta_i - B_i) [\partial_1 \eta]_i \text{sign}(\eta_i) \alpha_j \right\}. \end{aligned} \quad (\text{A.7})$$

where the second equality follows by Stein's lemma for a fixed i and $j \in I_i$, namely, for standard Gaussian Z we have $\mathbb{E}\{f(Z)Z\} = \mathbb{E}\{f'(Z)\}$ and therefore,

$$\begin{aligned} \frac{1}{\tau} \mathbb{E} \left\{ [\partial_1 \eta]_i \text{sign}(\eta_i) (\eta_i - B_i) \text{sign}(\eta_j) Z_j \right\} &= \mathbb{E} \left\{ \text{sign}(\eta_i) \text{sign}(\eta_j) \left[(\eta_i - B_i) \frac{d}{da_j} [\partial_1 \eta]_i + [\partial_1 \eta]_i \frac{d}{da_j} [\eta]_i \right] \right\} \\ &= \mathbb{E} \left\{ (\eta_i - B_i) [\partial_1^2 \eta]_i + ([\partial_1 \eta]_i)^2 \right\}. \end{aligned}$$

where the last step uses the definition of $\frac{d}{da_j}[\eta(\mathbf{a}, \mathbf{b})]_i$ given in (A.5) and the fact that $\frac{d}{da_j}[\partial_1 \eta(\mathbf{a}, \mathbf{b})]_i = \text{sign}(\eta_i) \text{sign}(\eta_j) [\partial_1^2 \eta(\mathbf{a}, \mathbf{b})]_{ij}$.

Therefore, simplifying (A.7), we have shown

$$(\delta p\tau) \times \frac{\partial \mathcal{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau) = \sum_{i=1}^p \mathbb{E} \left\{ \tau |I_i| \left([\partial_1 \eta]_i^2 + (\eta_i - B_i) [\partial_1^2 \eta]_i \right) - [\partial_1 \eta]_i \text{sign}(\eta_i) (\eta_i - B_i) \sum_{j \in I_i} \alpha_j \right\}. \quad (\text{A.8})$$

We now have the tools to prove Lemma A.1.

Proof of Lemma A.1. First,

$$\frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z}) - \mathbf{B}\|^2 = \frac{2}{p} \sum_{i=1}^p \left([\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z})]_i - B_i \right) \frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z})]_i.$$

As in the work above, we denote the proximal operator using a function $\eta : \mathbb{R}^{2p} \rightarrow \mathbb{R}^p$ as $\eta(\mathbf{a}, \mathbf{b}) := \text{prox}_{J_b}(\mathbf{a})$. Now from (A.6), denoting $I_i := \{j : |[\eta(\mathbf{a}, \mathbf{b})]_j| = |[\eta(\mathbf{a}, \mathbf{b})]_i|\}$, again dropping the explicit statement of the $\eta(\cdot, \cdot)$ input to save space,

$$\frac{\partial}{\partial \tau^2} [\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z})]_i = \frac{1}{2\tau} [\partial_1 \eta]_i \text{sign}(\eta_i) \sum_{j \in I_i} (\text{sign}(\eta_j) Z_j - \alpha_j).$$

Therefore,

$$\left| \frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z}) - \mathbf{B}\|^2 \right| = \frac{1}{\tau p} \left| \sum_{i=1}^p (\eta_i - B_i) [\partial_1 \eta]_i \text{sign}(\eta_i) \sum_{j \in I_i} (\text{sign}(\eta_j) Z_j - \alpha_j) \right|.$$

Since the averaging operation reduces the dot product (meaning informally that for a vector $\mathbf{v} \in \mathbb{R}^p$, $(\text{mean}(\mathbf{v}), \dots, \text{mean}(\mathbf{v})) \cdot \mathbf{v} \leq \|\mathbf{v}\|^2$), we have for any $i \in \{1, 2, \dots, p\}$ that $[\eta(\mathbf{B} + \tau \mathbf{Z}, \boldsymbol{\alpha}\tau)]_i - B_i$ can be replaced with $B_i + \tau Z_i - \text{sign}(\eta_i) \alpha_i \tau - B_i$. Using this in the above,

$$\begin{aligned} \left| \frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z}) - \mathbf{B}\|^2 \right| &\leq \frac{1}{p} \left| \sum_{i=1}^p \sum_{j \in I_i} (Z_i - \text{sign}(\eta_i) \alpha_i) [\partial_1 \eta]_i \text{sign}(\eta_i) (\text{sign}(\eta_j) Z_j - \alpha_j) \right| \\ &= \frac{1}{p} \left| \sum_{i=1}^p \sum_{j \in I_i} (\text{sign}(\eta_i) Z_i - \alpha_i) (\text{sign}(\eta_j) Z_j - \alpha_j) [\partial_1 \eta]_i \right|. \end{aligned} \quad (\text{A.9})$$

Next, using that $0 \leq |[\partial_1 \eta]_i| \leq 1/|I_i|$,

$$\left| \sum_{i=1}^p \sum_{j \in I_i} (\text{sign}(\eta_i) Z_i - \alpha_i) (\text{sign}(\eta_j) Z_j - \alpha_j) [\partial_1 \eta]_i \right| \leq \sum_{i=1}^p \frac{1}{|I_i|} \sum_{j \in I_i} \left| (\text{sign}(\eta_i) Z_i - \alpha_i) (\text{sign}(\eta_j) Z_j - \alpha_j) \right|.$$

Finally we make the following observation. Any equivalence class I_i is a collection of indices $j \in \{1, 2, \dots, p\}$ such that $|[\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z})]_j| = |[\text{prox}_{J_{\boldsymbol{\alpha}\tau}}(\mathbf{B} + \tau \mathbf{Z})]_i|$, so for any $j \in I_i$, it follows $I_j = I_i$. Recall, \mathbb{I} indicates the collection of unique equivalence classes, and we have

$$\sum_{i=1}^p \frac{1}{|I_i|} \sum_{j \in I_i} \left| (\text{sign}(\eta_i) Z_i - \alpha_i) (\text{sign}(\eta_j) Z_j - \alpha_j) \right| = \sum_{I \in \mathbb{I}} \frac{1}{|I|} \sum_{i,j \in I} \left| (\text{sign}(\eta_i) Z_i - \alpha_i) (\text{sign}(\eta_j) Z_j - \alpha_j) \right|.$$

Now plugging back into (A.9),

$$\begin{aligned} \left| \frac{\partial}{\partial \tau^2} \frac{1}{p} \|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau \mathbf{Z}) - \mathbf{B}\|^2 \right| &\leq \frac{1}{p} \sum_{I \in \mathcal{I}} \frac{1}{|I|} \sum_{i,j \in I} |(\text{sign}(\eta_i) Z_i - \alpha_i)(\text{sign}(\eta_j) Z_j - \alpha_j)| \\ &= \frac{1}{p} \sum_{I \in \mathcal{I}} \frac{1}{|I|} \left(\sum_{j \in I} |\text{sign}(\eta_j) Z_j - \alpha_j| \right)^2. \end{aligned}$$

□

Now considering (A.8), for simplicity in our future calculations, we suppress $|I_i|$ to 1 without loss of generality. To see this, recall that $I_i := \{j : |[\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_j| = |[\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_i|\}$ and note that when $|[\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_j|$ equals $|[\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_i|$, the terms will remain equal after small changes in τ . Therefore $|I_i|$ is treated as a constant in the derivative and since all operations below preserves linearity, it can safely be assumed to be equal to 1. Note that similarly, $\sum_{j \in I_i} \alpha_j$, will pass through future calculations as a constant. Therefore (A.8) becomes

$$(\delta p\tau) \times \frac{\partial \mathsf{F}}{\partial \tau^2}(\tau^2, \alpha\tau) = \sum_{i=1}^p \left[\mathbb{E}\{\tau([\partial_1 \eta]_i)^2 + \tau(\eta_i - B_i)[\partial_1^2 \eta]_i - \alpha_i \text{sign}(\eta_i)(\eta_i - B_i)[\partial_1 \eta]_i\} \right]. \quad (\text{A.10})$$

In what follows we will need to take care with the points (\mathbf{x}, \mathbf{y}) such that $[\partial_1^2 \eta(\mathbf{x}, \mathbf{y})]_i$ is not equal to 0. We refer to such points as ‘kink’ points, since these are points where the partial derivative jumps (and the second partial gradient acts like Dirac delta function $\delta(x)$), or in other words the points where the two (sorted, averaged) arguments in η are equal to each other. Informally, define a ‘kink’ point as an index where the sorted vector \mathbf{x} matches the corresponding threshold \mathbf{y} exactly. In LASSO, for example, the correspond to the ‘kinks’ of the soft-thresholding function. We have

$$[\partial_1^2 \eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_i = \delta(B_i + \tau Z_i - \alpha_i \tau) - \delta(B_i + \tau Z_i + \alpha_i \tau) \quad (\text{A.11})$$

and

$$\begin{aligned} &\mathbb{E}_{\mathbf{Z}, \mathbf{B}} \left\{ ([\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_i - B_i)[\partial_1^2 \eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)]_i \right\} \\ &= -\mathbb{E}_{\mathbf{B} \mid \mathbf{Z}} \left\{ B_i [\delta(B_i + \tau Z_i - \alpha_i \tau) - \delta(B_i + \tau Z_i + \alpha_i \tau)] \right\} \\ &= -\frac{1}{\tau} \mathbb{E}_{\mathbf{B}} \left\{ B_i [\phi(\alpha_i - \frac{1}{\tau} B_i) - \phi(-\alpha_i - \frac{1}{\tau} B_i)] \right\}. \end{aligned} \quad (\text{A.12})$$

Therefore, denoting \odot as elementwise multiplication of vectors, by (A.10) and (A.12),

$$\begin{aligned} &(\delta p\tau) \times \frac{\partial \mathsf{F}}{\partial \tau^2}(\tau^2, \alpha\tau) \\ &= \tau \mathbb{E} \|\partial_1 \eta\|^2 - \mathbb{E}_{\mathbf{B}} \left\{ \mathbf{B}^\top [\phi(\alpha - \frac{1}{\tau} \mathbf{B}) - \phi(-\alpha - \frac{1}{\tau} \mathbf{B})] \right\} - \mathbb{E} \{ [\alpha \odot \text{sign}(\eta) \odot (\eta - \mathbf{B})]^\top \partial_1 \eta \}. \end{aligned} \quad (\text{A.13})$$

Now we have shown the first derivative, so we consider the second derivative to prove concavity.

Notice, however, that in order to prove concavity of $\mathsf{F}(\tau^2, \alpha\tau)$ it suffices to show $\frac{\partial}{\partial \tau} [\frac{\partial \mathsf{F}}{\partial \tau^2}(\tau^2, \alpha\tau)] \leq 0$ because $\frac{\partial}{\partial \tau^2} (\frac{\partial \mathsf{F}}{\partial \tau^2}) = \frac{\partial \tau}{\partial \tau^2} [\frac{\partial}{\partial \tau} (\frac{\partial \mathsf{F}}{\partial \tau^2})] = \frac{1}{2\tau} [\frac{\partial}{\partial \tau} (\frac{\partial \mathsf{F}}{\partial \tau^2})]$.

We now show $\frac{\partial}{\partial \tau} [\frac{\partial \mathbf{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau)] \leq 0$. First,

$$(\delta p) \times \frac{\partial}{\partial \tau} \left[\frac{\partial \mathbf{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau) \right] = \frac{\partial}{\partial \tau} \mathbb{E} \|\partial_1 \eta\|^2 - \frac{\partial}{\partial \tau} \frac{1}{\tau} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B}) - \phi(-\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B})] \} \\ - \frac{\partial}{\partial \tau} \frac{1}{\tau} \mathbb{E} \{ [\boldsymbol{\alpha} \odot \text{sign}(\eta) \odot (\eta - \mathbf{B})]^\top \partial_1 \eta \}. \quad (\text{A.14})$$

To show that (A.14) is ≤ 0 , we find simplified representations of the three terms on the right side. This requires the same techniques as were used to find the first derivative above and so aren't given in full detail.

The first term on the right side of (A.14) can be simplified to the following:

$$\frac{\partial}{\partial \tau} \mathbb{E} \|\partial_1 \eta\|^2 = -\frac{1}{\tau^2} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B}) - \phi(\boldsymbol{\alpha} + \frac{1}{\tau} \mathbf{B})] \}. \quad (\text{A.15})$$

Doing so requires smart uses of the chain rule, a dominated convergence argument, the partials in (A.6), and special care for the 'kink' points as discussed above. Similarly, using (A.12), one can easily show for the third term on the right side of (A.14),

$$\frac{\partial}{\partial \tau} \frac{1}{\tau} \mathbb{E} \{ [\boldsymbol{\alpha} \odot \text{sign}(\eta) \odot (\eta - \mathbf{B})]^\top \partial_1 \eta \} \geq \frac{1}{\tau^3} \mathbb{E}_{\mathbf{B}} \{ [\boldsymbol{\alpha} \odot \mathbf{B}^2]^\top [\phi(\boldsymbol{\alpha} + \frac{1}{\tau} \mathbf{B}) + \phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B})] \}. \quad (\text{A.16})$$

Finally, using $\phi'(u) = -u\phi(u)$ and a dominated convergence argument, the second term on the right side of (A.14) equals

$$-\frac{\partial}{\partial \tau} \frac{1}{\tau} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B}) - \phi(-\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B})] \} \\ = \frac{1}{\tau^2} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B}) - \phi(-\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B})] \} \\ - \frac{1}{\tau^3} \mathbb{E}_{\mathbf{B}} \{ (\mathbf{B}^2)^\top [(\frac{1}{\tau} \mathbf{B} - \boldsymbol{\alpha}) \odot \phi(\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B}) - (\boldsymbol{\alpha} + \frac{1}{\tau} \mathbf{B}) \odot \phi(-\boldsymbol{\alpha} - \frac{1}{\tau} \mathbf{B})] \}. \quad (\text{A.17})$$

Now we plug (A.15), (A.16), and (A.17) back into (A.14) to show that $\frac{\partial}{\partial \tau} [\frac{\partial \mathbf{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau)] \leq 0$.

$$(\delta p) \times \frac{\partial}{\partial \tau} \left[\frac{\partial \mathbf{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau) \right] \\ \leq -\frac{1}{\tau^2} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \mathbf{B}/\tau) - \phi(\boldsymbol{\alpha} + \mathbf{B}/\tau)] \} + \frac{1}{\tau^2} \mathbb{E}_{\mathbf{B}} \{ \mathbf{B}^\top [\phi(\boldsymbol{\alpha} - \mathbf{B}/\tau) - \phi(-\boldsymbol{\alpha} - \mathbf{B}/\tau)] \} \\ - \frac{1}{\tau^3} \mathbb{E}_{\mathbf{B}} \{ (\mathbf{B}^2)^\top [(\mathbf{B}/\tau - \boldsymbol{\alpha}) \odot \phi(\boldsymbol{\alpha} - \mathbf{B}/\tau) - (\boldsymbol{\alpha} + \mathbf{B}/\tau) \odot \phi(-\boldsymbol{\alpha} - \mathbf{B}/\tau)] \} \\ - \frac{1}{\tau^3} \mathbb{E}_{\mathbf{B}} \{ [\boldsymbol{\alpha} \odot \mathbf{B}^2]^\top [\phi(\boldsymbol{\alpha} + \mathbf{B}/\tau) + \phi(\boldsymbol{\alpha} - \mathbf{B}/\tau)] \} \\ = -\frac{1}{\tau^4} \mathbb{E}_{\mathbf{B}} \{ [\mathbf{B}^3]^\top [\phi(\boldsymbol{\alpha} - \mathbf{B}/\tau) - \phi(\boldsymbol{\alpha} + \mathbf{B}/\tau)] \}. \quad (\text{A.18})$$

We justify non-positivity of (A.18) by showing that the elementwise term inside the expectation is less than or equal to 0. First assume $B_i \geq 0$, then $\alpha_i - B_i/\tau \leq \alpha_i + B_i/\tau$ and $\phi(\alpha_i - B_i/\tau) \geq \phi(\alpha_i + B_i/\tau)$. The other case $B_i \leq 0$ follows similarly.

Now (A.18), implies $\frac{\partial}{\partial \tau} [\frac{\partial \mathbf{F}}{\partial \tau^2}(\tau^2, \boldsymbol{\alpha}\tau)] \leq 0$ and therefore, we have shown that $\mathbf{F}(\tau^2, \boldsymbol{\alpha}\tau)$ defined in (2.8), is concave with respect to τ^2 .

Next we show that $\tau^2 \mapsto F(\tau^2, \alpha\tau)$ is strictly increasing. To do so, it is sufficient to show that $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ is positive as $\tau \rightarrow \infty$ because the concavity implies that $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ is non-increasing. Define $f(\alpha) := \delta \lim_{\tau \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$. First recall that $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ is given in (A.8). In particular,

$$\delta \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) = \frac{1}{p} \sum_{i=1}^p \mathbb{E} \left\{ |I_i| \left([\partial_1 \eta]_i^2 + (\eta_i - B_i)[\partial_1^2 \eta]_i \right) - \frac{1}{\tau} [\partial_1 \eta]_i \operatorname{sign}(\eta_i)(\eta_i - B_i) \sum_{j \in I_i} \alpha_j \right\}, \quad (\text{A.19})$$

Then taking $\tau \rightarrow \infty$ in the above, it is easy to see that $f(\alpha)$ is equivalent to setting $\mathbf{B} = \mathbf{0}$ in $\eta(\mathbf{B} + \tau \mathbf{Z}, \alpha\tau)$ and using that $\eta(\tau \mathbf{Z}, \alpha\tau) = \tau \eta(\mathbf{Z}, \alpha)$ (implying that $\partial_1 \eta(\tau \mathbf{Z}, \alpha\tau) = \partial_1 \eta(\mathbf{Z}, \alpha)$). We note that using a simplification of $[\partial_1^2 \eta]_i$ as in (A.11)-(A.12), means that this term will go to zero as $\tau \rightarrow \infty$. Therefore, using $\operatorname{sign}(\eta(\mathbf{Z}, \alpha)) \odot \eta(\mathbf{Z}, \alpha) = |\eta(\mathbf{Z}, \alpha)|$,

$$f(\alpha) = \frac{1}{p} \sum_{i=1}^p \mathbb{E} \left\{ [D(\eta(\mathbf{Z}, \alpha))]_i ([\partial_1 \eta(\mathbf{Z}, \alpha)]_i)^2 - [\partial_1 \eta(\mathbf{Z}, \alpha)]_i |\eta(\mathbf{Z}, \alpha)|_i \sum_{j: |\eta(\mathbf{Z}, \alpha)|_j = |\eta(\mathbf{Z}, \alpha)|_i} \alpha_j \right\}.$$

In the above we have used the following definition: for a vector $\mathbf{v} \in \mathbb{R}^p$, define \mathbf{D} elementwise as $[\mathbf{D}(\mathbf{v})]_i := \#\{j : |v_j| = |v_i|\} = |I_i|$ if $v_i \neq 0$ and ∞ otherwise. Using that $\partial_1 \eta(\mathbf{Z}, \alpha) = \frac{1}{D(\eta(\mathbf{Z}, \alpha))}$,

$$f(\alpha) = \frac{1}{p} \sum_{i=1}^p \mathbb{E} \left\{ \left(1 - |\eta(\mathbf{Z}, \alpha)|_i \sum_{j: |\eta(\mathbf{Z}, \alpha)|_j = |\eta(\mathbf{Z}, \alpha)|_i} \alpha_j \right) \frac{1}{[D(\eta(\mathbf{Z}, \alpha))]_i} \right\} \quad (\text{A.20})$$

This simplification can be efficiently computed because only $|\eta(\mathbf{Z}, \alpha)|$ and α need to be memorized.

Now considering (A.20), let $\alpha \rightarrow \infty$ and note that since $|\mathbf{Z}| < \alpha$ almost surely as $\alpha \rightarrow \infty$, it follows that $\eta(\mathbf{Z}, \alpha) = \partial_1 \eta(\mathbf{Z}, \alpha) = \mathbf{0}$. Therefore $\lim_{\alpha \rightarrow \infty} f(\alpha) = 0$. By a very similar argument to the proof of concavity, it is easy to see $f'(\alpha) < 0$, and together these facts imply $f(\alpha) > 0$ for all α . The monotonicity of F is now obvious: since F is concave (implying $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ is non-increasing) and strictly increasing for τ^2 large enough, it is increasing everywhere. Moreover, the monotonicity of F implies the monotonicity of the sequence $\{\tau_t^2(p)\}_{t \geq 0}$.

Finally we show that there exists a unique τ_* such that $F(\tau_*^2, \alpha\tau_*) = \tau_*^2$, from which it follows that the monotone sequence $\{\tau_t^2(p)\}_{t \geq 0}$ converges to $\tau_*^2(p)$ as $t \rightarrow \infty$. First, by (A.20), we know $f(\mathbf{0}) = \mathbb{E} \|\partial_1 \eta(\tau \mathbf{Z}, \mathbf{0})\|^2 / p = \mathbb{E} \|\mathbf{1}\|^2 / p = 1$. This, along with the fact that $f'(\alpha) < 0$, tells us that $0 < f(\alpha) < 1$ for all α . Recall the definition of the set \mathbf{A}_{\min} , namely $\mathbf{A}_{\min} := \{\alpha : f(\alpha) = \delta\}$. We know that this set is non-empty since the LASSO case shows $\alpha = (\alpha_{\min}, \dots, \alpha_{\min})$ belongs to \mathbf{A}_{\min} where α_{\min} is the unique non-negative solution of $(1 + \alpha^2)\Phi(-\alpha) - \alpha\phi(\alpha) = \delta/2$. We write $\alpha \succeq \mathbf{A}_{\min}$ to mean α is larger than at least one element in \mathbf{A}_{\min} , where we consider one vector \mathbf{v} to be larger than another vector \mathbf{u} if $v_i \geq u_i$ for all i and $v_j > u_j$ for some j .

To complete the proof, we show that $F(\tau^2, \alpha\tau) > \tau^2$ for small enough τ^2 and $F(\tau^2, \alpha\tau) < \tau^2$ for large enough τ^2 . Therefore, there is at least one τ_* such that $F(\tau_*^2, \alpha\tau_*) = \tau_*^2$ since F is continuous in τ . It follows from the concavity of F that the solution is unique and the sequence of iterates $\tau_t^2(p)$ converge to $\tau_*^2(p)$. We first show that $F(\tau^2, \alpha\tau) > \tau^2$ for small enough τ^2 . Consider the function $G(\tau^2) := F(\tau^2, \alpha\tau) - \tau^2$. Recalling the definition of $F(\tau^2, \alpha\tau)$ in (2.8), namely, $F(\tau^2, \alpha\tau) = \sigma_w^2 + \mathbb{E} \|\operatorname{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau \mathbf{Z}) - \mathbf{B}\|^2 / (\delta p)$, clearly $F(0, \mathbf{0}) = \sigma_w^2 \geq 0$ and therefore $G(0) = \sigma_w^2 \geq 0$ (with equality only if $\sigma_w^2 = 0$). Now we show that $F(\tau^2, \alpha\tau) < \tau^2$ for large enough τ^2 . Since $f(\alpha)$ is decreasing in α , for $\alpha \succeq \mathbf{A}_{\min}$, it must be that $f(\alpha) < \delta$. Moreover, $\lim_{\tau \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) = \frac{1}{\delta} f(\alpha) \leq 1$ for $\alpha \succeq \mathbf{A}_{\min}$. Therefore, $\lim_{\tau \rightarrow \infty} \frac{\partial G}{\partial \tau^2}(\tau^2) \leq 0$ meaning G is

eventually decreasing (as τ^2 grows) for any $\alpha \succeq A_{\min}$. Also, $G(\tau^2)$ is concave and therefore for τ^2 large enough we will have $G(\tau^2) < 0$, in which case $F(\tau^2, \alpha\tau) < \tau^2$.

Finally, $|\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)|$ evaluated at $\tau^2 = \tau_*^2$ is upper bounded by 1 when $\alpha \succeq A_{\min}$, as the concavity of F implies that $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ is strictly decreasing in τ^2 along with $\lim_{\tau \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) = \frac{1}{\delta} f(\alpha) \leq 1$ when $\alpha \succeq A_{\min}$. If this were not the case then there would be multiple fixed points. \square

A.2 Proving Proposition 2.6

Proof of Proposition 2.6. This proof is a generalized result of [5, Proposition 1.4] (originally proved in [17]) and [5, Corollary 1.7]. Here we fixed p and denote $\tau(p)$ as τ .

Recall in the proof of Theorem 1 we have shown the following facts: **(A)** $0 < \lim_{\tau^2 \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) < 1$; **(B)** $\tau^2 \mapsto F(\tau^2, \alpha\tau)$ is concave; **(C)** $\tau^2 \mapsto F(\tau^2, \alpha\tau)$ is strictly increasing; and **(D)** $\frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ evaluated at $\tau = \tau_*$, which we denote $\frac{\partial F}{\partial \tau^2}(\tau_*, \alpha\tau_*)$ is such that $0 < \frac{\partial F}{\partial \tau^2}(\tau_*, \alpha\tau_*) < 1$.

First we claim $\alpha \mapsto \tau_*^2(\alpha)$ is continuously differentiable on \mathbb{R}_+^p . This follows from the implicit function theorem on function $G(\alpha, \tau^2) := \tau^2 - F(\tau^2, \alpha\tau)$ and from Fact **(D)**: G is continuously differentiable and $0 < \frac{\partial G}{\partial \tau^2} < 1$. Hence τ^2 can be written as $\tau^2(\alpha)$ which is continuously differentiable. Defining $g(\alpha, \tau^2) := \alpha\tau[1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{\alpha\tau}}(\mathbf{B} + \tau\mathbf{Z})\|_0^*]$, notice that $\lambda(\alpha) = g(\alpha, \tau_*^2(\alpha))$. Clearly g is continuously differentiable in α and so is $\alpha \mapsto \lambda(\alpha)$.

In the next step, we consider $\alpha \succeq A_{\min}(\delta)$ such that $\alpha \rightarrow a_{\min}$ for some $a_{\min} \in A_{\min}(\delta)$ (denote as $\alpha \downarrow A_{\min}(\delta)$). We claim $\tau_*^2(\alpha) \rightarrow +\infty$ as $\alpha \downarrow A_{\min}(\delta)$. Recall, $f(\alpha) := \delta \lim_{\tau \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ (cf. Theorem 1). Then by concavity of $F(\tau^2, \alpha\tau)$ in τ ,

$$\tau_*^2 = F(\tau_*, \alpha\tau_*) \geq F(0, \mathbf{0}) + \tau_*^2 \lim_{\tau^2 \rightarrow \infty} \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau) = F(0, \mathbf{0}) + \frac{1}{\delta} \tau_*^2 f(\alpha) \Rightarrow \tau_*^2 \geq \frac{F(0, \mathbf{0})}{1 - f(\alpha)/\delta}$$

Recall $F(0, \mathbf{0}) = \sigma_w^2$ and $f(a_{\min}) = \delta$ for any $a_{\min} \in A_{\min}(\delta)$. Hence $\tau_*^2(\alpha) \rightarrow +\infty$ as $\alpha \downarrow A_{\min}(\delta)$.

Define $\ell(\alpha) := 1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{\alpha\tau_*}}(\mathbf{B} + \tau_*\mathbf{Z})\|_0^*$. Then when $\tau_*^2(\alpha) \rightarrow +\infty$ as $\alpha \downarrow A_{\min}(\delta)$,

$$\ell_* := \lim_{\alpha \rightarrow a_{\min}} \ell(\alpha) = \lim_{\alpha \rightarrow a_{\min}} \left(1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{\alpha\tau_*}}(\tau_*\mathbf{Z})\|_0^*\right) = 1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{a_{\min}}}(\mathbf{Z})\|_0^*.$$

We claim that $\ell_* < 0$. Using the definition of the vector \mathbf{D} and the set $A_{\min}(\delta)$ in (2.7),

$$\begin{aligned} \ell_* &= 1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{a_{\min}}}(\mathbf{Z})\|_0^* = 1 - \frac{1}{\delta}\mathbb{E}\left\langle \frac{1}{\mathbf{D}(\text{prox}_{J_{a_{\min}}}(\mathbf{Z}))} \right\rangle \\ &< 1 - \frac{1}{\delta p} \sum_i \mathbb{E} \left\{ \frac{1}{[\mathbf{D}(\text{prox}_{J_{a_{\min}}}(\mathbf{Z}))]_i} \left(1 - \sum_{j \in I_i} [\mathbf{a}_{\min}]_j \cdot |[\text{prox}_{J_{a_{\min}}}(\mathbf{Z})]_i|\right) \right\} = 0, \end{aligned}$$

where (writing $\boldsymbol{\eta}$ to mean $\text{prox}_{J_{a_{\min}}}(\mathbf{Z})$ and \mathbf{D} to mean $\mathbf{D}(\boldsymbol{\eta})$) the inequality in the above uses the fact that

$$\frac{1}{\mathbf{D}_i} - \frac{1}{\mathbf{D}_i} \left(1 - \sum_{j \in I_i} [\mathbf{a}_{\min}]_j |\boldsymbol{\eta}_i|\right) = \frac{1}{\mathbf{D}_i} \sum_{j \in I_i} [\mathbf{a}_{\min}]_j |\boldsymbol{\eta}_i| \geq 0.$$

Notice in the above, the equality only holds when $\boldsymbol{\eta}_i = \mathbf{0}$ but $\boldsymbol{\eta} \neq \mathbf{0}$ almost surely. Therefore, using that $\lambda(\alpha) = g(\alpha, \tau_*^2(\alpha)) = \alpha\tau_*(\alpha)[1 - \frac{1}{n}\mathbb{E}\|\text{prox}_{J_{\alpha\tau_*}(\alpha)}(\mathbf{B} + \tau_*(\alpha)\mathbf{Z})\|_0^*]$,

$$\lim_{\alpha \downarrow A_{\min}(\delta)} \lambda(\alpha) = \ell_* \cdot \lim_{\alpha \downarrow A_{\min}(\delta)} \alpha\tau_*(\alpha) = -\infty. \tag{A.21}$$

Finally we consider the case $\alpha \rightarrow \infty$ and observe $\tau_*^2(\alpha) \rightarrow \sigma_w^2 + \mathbb{E}\{B^2\}/\delta$. To see this, notice that $F(\tau^2, \alpha\tau) \rightarrow \sigma_w^2 + \mathbb{E}\{B^2\}/\delta$ as $\alpha \rightarrow \infty$ since $\tau_*^2(\alpha) = F(\tau_*^2(\alpha), \alpha\tau_*(\alpha))$ is bounded above. Moreover, since $\tau_*(\alpha)$ is bounded, $\alpha\tau_*(\alpha)$ is unbounded as $\alpha \rightarrow \infty$ and we have $\lim_{\alpha \rightarrow \infty} \ell(\alpha) = 1$ whence

$$\lim_{\alpha \rightarrow \infty} \lambda(\alpha) = 1 \cdot \lim_{\alpha \rightarrow \infty} \alpha\tau_*(\alpha) = \infty. \quad (\text{A.22})$$

We pause here to summarize that $\alpha \mapsto \lambda(\alpha)$ is continuously differentiable on the domain $\{\alpha : \alpha \geq A_{\min}(\delta)\}$ with $\lambda(A_{\min}(\delta)) = -\infty$ and $\lim_{\alpha \rightarrow \infty} \lambda(\alpha) = +\infty$.

Now to prove the inverse mapping $\lambda \mapsto \alpha(\lambda)$ is continuous and non-decreasing when $p \rightarrow \infty$, we claim that the invertibility of $\alpha \mapsto \lambda(\alpha)$ is sufficient. Precisely, (1) invertibility implies strict monotonicity; (2) monotonicity plus (A.21) and (A.22) implies both $\alpha \mapsto \lambda(\alpha)$ and $\lambda \mapsto \alpha(\lambda)$ are increasing; and (3) continuity of $\alpha \mapsto \lambda(\alpha)$ implies continuity of $\lambda \mapsto \alpha(\lambda)$.

Now we prove the invertibility by contradiction. Assume that there are two distinct such values α_1, α_2 satisfying $\tilde{\lambda} = \lambda(\alpha_1) = \lambda(\alpha_2)$. Apply Theorem 3 to both $\alpha(\tilde{\lambda}) = \alpha_1, \alpha_2$ with $\psi(\mathbf{x}, \mathbf{y}) = \langle (\mathbf{x} - \mathbf{y})^2 \rangle$. Then, together with Corollary 3.4,

$$\operatorname{plim}_{p \rightarrow \infty} \|\hat{\beta} - \beta\|^2/p = \operatorname{plim}_{p \rightarrow \infty} \mathbb{E}\langle \|\operatorname{prox}_{J_{\alpha\tau_*}}(\beta + \tau_* \mathbf{Z}; \alpha\tau_*) - \beta\|_2^2 \rangle = \delta(\tau_*^2 - \sigma_w^2).$$

Since $\operatorname{plim}_{p \rightarrow \infty} \|\hat{\beta} - \beta\|^2/p$ is independent of α , the right side gives $\tau_*(\alpha_1) = \tau_*(\alpha_2)$. Next apply Theorem 3 with $\psi(\mathbf{x}, \mathbf{y}) = \langle |\mathbf{x}| \rangle$, giving $\operatorname{plim}_{p \rightarrow \infty} \|\hat{\beta}\|_1/p = \operatorname{plim}_{p \rightarrow \infty} \mathbb{E}\langle \|\operatorname{prox}_{J_{\alpha\tau_*}}(\beta + \tau_* \mathbf{Z}; \alpha\tau_*)\|_1 \rangle$. Obviously, for τ_* and p fixed, $\theta \mapsto \mathbb{E}\langle \|\operatorname{prox}_{J_{\alpha\tau_*}}(\beta + \tau_* \mathbf{Z}; \theta)\|_1 \rangle$ is strictly decreasing in θ . Therefore $\alpha_1\tau_*(\alpha_1) = \alpha_2\tau_*(\alpha_2)$ implying $\alpha_1 = \alpha_2$, since $\tau_*(\alpha_1) = \tau_*(\alpha_2)$, which is a contradiction. \square

B Verifying Properties (P1) and (P2)

In this appendix we demonstrate that the properties (P1) and (P2) given in Section 4 and relating to the denoiser $\eta_p^t(\cdot)$ defined in (4.1) are true.

Verifying Properties (P1) and (P2). Property (P1) follows since $\eta_p^t(\cdot) = \operatorname{prox}_{J_{\alpha\tau_t}}(\cdot)$, as it is easy to show that proximal operators are Lipschitz continuous with Lipschitz constant one. Namely

$$\|\eta_p^t(\mathbf{v}_1) - \eta_p^t(\mathbf{v}_2)\| = \|\operatorname{prox}_{J_{\alpha\tau_t}}(\mathbf{v}_1) - \operatorname{prox}_{J_{\alpha\tau_t}}(\mathbf{v}_2)\| \leq \|\mathbf{v}_1 - \mathbf{v}_2\|.$$

Next we show that property (P2) is true. We restate property (P2) for convenience: for any s, t with $(\mathbf{Z}, \mathbf{Z}')$ a pair of length- p vectors such that (Z_i, Z'_i) are i.i.d. $\sim \mathcal{N}(0, \Sigma)$ for $i \in [p]$ where Σ is any 2×2 covariance matrix, the following limits exist and are finite.

$$\operatorname{plim}_{p \rightarrow \infty} \frac{1}{p} \|\beta\|, \quad \operatorname{plim}_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}_{\mathbf{Z}}[\beta^\top \eta_p^t(\beta + \mathbf{Z})], \quad \operatorname{plim}_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}[\eta_p^s(\beta + \mathbf{Z}')^\top \eta_p^t(\beta + \mathbf{Z})]. \quad (\text{B.1})$$

We first note that the first limit in (B.1) exists by Assumption (A2) and the strong law of large numbers. We focus on the other two limits. These results follow by [20, Proposition 1] given in Lemma 3.3 and the following lemma, which is a classic result in probability theory.

Lemma B.1 (Doob's L^1 maximal inequality, [18] Chapter VII, Theorem 3.4). *Let X_1, X_2, \dots, X_p be a sequence of nonnegative i.i.d. random variables such that $\mathbb{E}[X_1 \max\{0, \log(X_1)\}] < \infty$. Then,*

$$\mathbb{E}\left[\sup_{p \geq 1} \left\{\frac{1}{p}(X_1 + X_2 + \dots + X_p)\right\}\right] \leq \frac{e}{e-1}(1 + \mathbb{E}[X_1 \max\{0, \log(X_1)\}]).$$

Proof. Let $M_p = \frac{1}{p}(X_1 + X_2 + \dots + X_p)$. Then the sequence $\{M_p\}$ is a submartingale and hence by Doob's maximal inequality,

$$\mathbb{E}\left[\sup_{p' \geq p \geq 1} M_p\right] \leq \frac{e}{e-1}(1 + \mathbb{E}[M_{p'} \max\{0, \log(M_{p'})\}]).$$

Note the mapping $x \mapsto x \max\{0, \log x\}$ is convex and hence $\mathbb{E}[M_{p'} \max\{0, \log(M_{p'})\}] \leq \mathbb{E}[X_1 \max\{0, \log(X_1)\}]$. The result follows by Fatou's lemma and by noting that $\sup_{p' \geq p \geq 1} M_p \uparrow \sup_{p \geq 1} M_p$ as $p' \rightarrow \infty$. \square

Before we prove that the second and third limits in (B.1) exist and are finite, we state one more result that will be helpful in the proof. This result uses Lemma B.1 along with a Dominated Convergence argument to study expectations taken with respect to $(\mathbf{Z}, \mathbf{Z}')$ like those in (B.1).

Lemma B.2. *Consider a function $\psi_p : \mathbb{R}^p \times \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ such that for iterations $s, t \geq 0$,*

$$\frac{1}{p} \left| \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}')) - \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right| \rightarrow 0, \quad \text{as } p \rightarrow \infty, \quad (\text{B.2})$$

where h^s, h^t are the unspecified functions of Lemma 3.3, and $(\mathbf{Z}, \mathbf{Z}')$ are independent Gaussian vectors having zero-mean and independent entries with finite variance. Assume, for some constant $L > 0$ not depending on p ,

$$\frac{1}{p} \left| \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}')) - \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right| \leq L \left(1 + \frac{\|\boldsymbol{\beta}\|^2}{p} + \frac{\|\mathbf{Z}\|^2}{p} + \frac{\|\mathbf{Z}'\|^2}{p} \right). \quad (\text{B.3})$$

Then, as $p \rightarrow \infty$,

$$\frac{1}{p} \left| \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})) \right\} - \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right\} \right| \rightarrow 0. \quad (\text{B.4})$$

Proof. We begin by showing that $\mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \sup_{p \geq 1} \frac{1}{p} |\psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}'))| \right\} < \infty$. Using (B.3), it is clear that this expectation is finite almost surely if

$$\mathbb{E}\left[\sup_{p \geq 1} \left\{ \frac{1}{p} \|\mathbf{Z}(p)\|^2 \right\}\right] < \infty, \quad \mathbb{E}\left[\sup_{p \geq 1} \left\{ \frac{1}{p} \|\mathbf{Z}'(p)\|^2 \right\}\right] < \infty, \quad \text{and} \quad \mathbb{E}\left[\sup_{p \geq 1} \left\{ \frac{1}{p} \|\boldsymbol{\beta}(p)\|^2 \right\}\right] < \infty,$$

where we have made the dependence of the vectors on the dimension p explicit. But Lemma B.1 immediately implies the above since $\mathbb{E}[B^2 \max\{0, \log B\}] < \infty$ by assumption **(A2)**.

Now by dominated convergence we have,

$$\begin{aligned} & \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \text{plim}_p \frac{1}{p} \left| \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}')) - \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right| \right\} \\ &= \text{plim}_p \frac{1}{p} \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left| \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}')) - \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right| \\ &\geq \text{plim}_p \frac{1}{p} \left| \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \psi_p(\boldsymbol{\beta}, \eta_p^s(\boldsymbol{\beta} + \mathbf{Z}), \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})) \right\} - \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \left\{ \psi_p(\boldsymbol{\beta}, h^s(\boldsymbol{\beta} + \mathbf{Z}), h^t(\boldsymbol{\beta} + \mathbf{Z}')) \right\} \right|. \end{aligned}$$

Then the above implies the desired result (B.4) from assumption (B.2). \square

First consider the second limit in (B.1). By Cauchy-Schwarz, (3.3) of Lemma 3.3 implies that $|\boldsymbol{\beta}^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - \boldsymbol{\beta}^\top h^t(\boldsymbol{\beta} + \mathbf{Z})|/p \rightarrow 0$, as $p \rightarrow \infty$. This follows because

$$|\boldsymbol{\beta}^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - \boldsymbol{\beta}^\top h^t(\boldsymbol{\beta} + \mathbf{Z})|/p \leq \|\boldsymbol{\beta}\| \|\eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^t(\boldsymbol{\beta} + \mathbf{Z})\|/p.$$

Then the right side of the above $\rightarrow 0$ with growing p because $\|\boldsymbol{\beta}\|/\sqrt{p}$ limits to a constant as justified above (this is the limit in (B.1)), and the other term $\rightarrow 0$ by (3.3) of Lemma 3.3. This means that assumption (B.2) of Lemma B.2 is satisfied. Assumption (B.3) of Lemma B.2 is also satisfied since both η_p^t and h^t are Lipschitz(1), by Cauchy-Schwarz inequality. Therefore Lemma B.2 implies $|\mathbb{E}_{\mathbf{Z}}\{\boldsymbol{\beta}^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})\} - \mathbb{E}_{\mathbf{Z}}\{\boldsymbol{\beta}^\top h^t(\boldsymbol{\beta} + \mathbf{Z})\}|/p \rightarrow 0$, as $p \rightarrow \infty$. Therefore,

$$\text{plim}_{p \rightarrow \infty} \mathbb{E}_{\mathbf{Z}}[\boldsymbol{\beta}^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})]/p = \text{plim}_{p \rightarrow \infty} \sum_{i=1}^p \beta_{0,i} \mathbb{E}_{\mathbf{Z}}\{h^t(\beta_{0,i} + Z_i)\}/p = \mathbb{E}[Bh^t(B + Z)],$$

where B, Z are univariate. By the Cauchy-Schwarz inequality, $\mathbb{E}[Bh^t(B + Z)] < \infty$ if $\mathbb{E}[B^2] < \infty$ and $\mathbb{E}[h^t(B + Z)^2] < \infty$. Since $\mathbb{E}[B^2] = \sigma_{\boldsymbol{\beta}}^2 < \infty$ is given by our assumption, it suffices to show $\mathbb{E}[h^t(B + Z)^2] < \infty$. But this follows from the fact that $h^t(\cdot)$ is Lipschitz(1) and therefore $\mathbb{E}[h^t(B + Z)^2] < \mathbb{E}[(B + Z)^2] \leq \mathbb{E}[B^2] + \mathbb{E}[Z^2] = \sigma_{\boldsymbol{\beta}}^2 + \Sigma_{11} < \infty$.

Finally consider the third limit in (B.1). Similarly to the work in studying the second limit in (B.1), we will appeal to Lemma B.2. First we will show that

$$|\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}')^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^s(\boldsymbol{\beta} + \mathbf{Z}')^\top h^t(\boldsymbol{\beta} + \mathbf{Z})|/p \rightarrow 0, \quad \text{as } p \rightarrow \infty, \quad (\text{B.5})$$

meaning that assumption (B.2) of Lemma B.2 is satisfied. Then, again, assumption (B.3) of Lemma B.2 is satisfied since both $\eta_p^t(\cdot)$ and $h^t(\cdot)$ are Lipschitz(1), using Cauchy-Schwarz.

Now we want to prove (B.5). By repeated applications of Cauchy-Schwarz it is not hard to show,

$$\begin{aligned} & \text{plim}_p |\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}')^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^s(\boldsymbol{\beta} + \mathbf{Z}')^\top h^t(\boldsymbol{\beta} + \mathbf{Z})|/p \\ & \leq \text{plim}_p \|h^s(\boldsymbol{\beta} + \mathbf{Z}')\| \|\eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^t(\boldsymbol{\beta} + \mathbf{Z})\|/p + \text{plim}_p \|h^t(\boldsymbol{\beta} + \mathbf{Z})\| \|\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}') - h^s(\boldsymbol{\beta} + \mathbf{Z}')\|/p \\ & \quad + \text{plim}_p \|\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}') - h^s(\boldsymbol{\beta} + \mathbf{Z}')\| \|\eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^t(\boldsymbol{\beta} + \mathbf{Z})\|/p. \end{aligned}$$

Now, (B.5) follows since the right side of the above goes to 0 as p grows. This follows since, by (3.3) of Lemma 3.3, as $p \rightarrow \infty$,

$$\|\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}') - h^s(\boldsymbol{\beta} + \mathbf{Z}')\|/\sqrt{p} \rightarrow 0 \quad \text{and} \quad \|\eta_p^t(\boldsymbol{\beta} + \mathbf{Z}) - h^t(\boldsymbol{\beta} + \mathbf{Z})\|/\sqrt{p} \rightarrow 0.$$

Moreover, since $h^s(\cdot)$ and $h^t(\cdot)$ are separable, by the Law of Large Numbers,

$$\begin{aligned} & \text{plim}_p \|h^s(\boldsymbol{\beta} + \mathbf{Z}')\|^2/p = \text{plim}_p \sum_{i=1}^p [h^s(\beta_i + Z'_i)]^2/p = \mathbb{E}[(h^s(B + Z'))^2] < \infty, \\ & \text{plim}_p \|h^t(\boldsymbol{\beta} + \mathbf{Z})\|^2/p = \text{plim}_p \sum_{i=1}^p [h^t(\beta_i + Z_i)]^2/p = \mathbb{E}[(h^t(B + Z))^2] < \infty, \end{aligned}$$

where the inequalities follow since $\mathbb{E}[(h^s(B + Z'))^2] \leq \mathbb{E}[(B + Z')^2] \leq \sigma_{\boldsymbol{\beta}}^2 + \Sigma_{22} < \infty$ and $\mathbb{E}[(h^t(B + Z))^2] \leq \mathbb{E}[(B + Z)^2] \leq \sigma_{\boldsymbol{\beta}}^2 + \Sigma_{11} < \infty$. This proves (B.5) and therefore we can apply Lemma B.2.

Then Lemma B.2 implies,

$$|\mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}\{\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}')^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})\} - \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}\{h^s(\boldsymbol{\beta} + \mathbf{Z}')^\top h^t(\boldsymbol{\beta} + \mathbf{Z})\}|/p \rightarrow 0, \quad \text{as } p \rightarrow \infty.$$

But now, using the above, we find that

$$\begin{aligned} \operatorname{plim}_{p \rightarrow \infty} \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}\{\eta_p^s(\boldsymbol{\beta} + \mathbf{Z}')^\top \eta_p^t(\boldsymbol{\beta} + \mathbf{Z})\}/p &= \operatorname{plim}_{p \rightarrow \infty} \sum_{i=1}^p \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'}\{h^s(\beta_i + Z'_i)h^t(\beta_i + Z_i)\}/p \\ &= \mathbb{E}[h^s(B + Z')h^t(B + Z)], \end{aligned}$$

where B, Z' , and Z are univariate and $\mathbb{E}[h^s(B + Z')h^t(B + Z)] < \infty$ by Cauchy-Schwarz and the fact that $h^s(\cdot)$ and $h^t(\cdot)$ are Lipschitz(1). Namely, this gives the bound

$$\begin{aligned} (\mathbb{E}[h^s(B + Z')h^t(B + Z)])^2 &\leq \mathbb{E}[(h^s(B + Z'))^2]\mathbb{E}[(h^t(B + Z))^2] \leq \mathbb{E}[(B + Z')^2]\mathbb{E}[(B + Z)^2] \\ &= (\mathbb{E}[B^2] + \mathbb{E}[Z'^2])(\mathbb{E}[B^2] + \mathbb{E}[Z^2]) = (\sigma_{\boldsymbol{\beta}}^2 + \Sigma_{22})(\sigma_{\boldsymbol{\beta}}^2 + \Sigma_{11}) < \infty. \end{aligned}$$

We have now shown that property **(P2)** is true. \square

C Proof of Fact 2.7

Proof. The fact follows from the asymptotic separability of the proximal operator [20, Proposition 1] (restated in Lemma 3.3) and the dominated convergence theorem [33] allowing for interchange of limit and expectation. We sketch the proof of the existence of the limit in (2.4) (and the result for the limit in (2.11) follows similarly). By Lemma 3.3, the weak convergence of $\boldsymbol{\alpha}(p)$ to A , and the Weak Law of Large Numbers, one can argue that

$$\lim_p \|\operatorname{prox}_{J_{\boldsymbol{\alpha}(p)\tau_*}}(\mathbf{B} + \tau_* \mathbf{Z}) - \mathbf{B}\|^2 / (\delta p) = \mathbb{E}\{(h(B + \tau_* Z) - B)^2\}/\delta, \quad (\text{C.1})$$

where $h(\cdot) := h(\cdot; B + \tau_* Z, A\tau_*)$ is the unspecified, separable function of Lemma 3.3. This is consistent with [Lemma 29, [20]]. The limit in (2.4) exists if $\mathbb{E}\{(h(B + \tau_* Z) - B)^2\}/\delta < \infty$ and

$$\begin{aligned} \mathbb{E}\{(h(B + \tau_* Z) - B)^2\} &\leq 2\mathbb{E}\{h(B + \tau_* Z)^2 + B^2\} \leq 2\mathbb{E}\{(B + \tau_* Z)^2 + B^2\} \\ &\leq 2\mathbb{E}\{2B^2 + 2\tau_*^2 Z^2 + B^2\} = 6\mathbb{E}\{B^2\} + 4\tau_*^2 < \infty. \end{aligned}$$

Here the first and third inequalities follow from $(x - y)^2 \leq 2(x^2 + y^2)$ and the second inequality follows from h being Lipschitz(1): $|h(x)| = |h(x) - h(0)| \leq |x - 0| = |x|$. \square

D Proof of Lemma 7.1

Proof. First, the proof of (7.1) follows from Theorem 4.1. To see this, note that by (1.3a), we have $\boldsymbol{\beta}^{t+1} = \operatorname{prox}_{J_{\theta_t}}(\mathbf{X}^\top \mathbf{z}^t + \boldsymbol{\beta}^t) = \eta_p^t(\mathbf{X}^\top \mathbf{z}^t + \boldsymbol{\beta}^t)$, and therefore we apply Theorem 4.1 with uniformly pseudo-Lipschitz function $\psi_p(\boldsymbol{\beta}^t + \mathbf{X}^\top \mathbf{z}^t, \boldsymbol{\beta}) = \|\eta_p^t(\boldsymbol{\beta}^t + \mathbf{X}^\top \mathbf{z}^t)\|^2/p$ to get

$$\operatorname{plim}_p \|\boldsymbol{\beta}^t\|^2/p \stackrel{p}{=} \operatorname{plim}_p \mathbb{E}_{\mathbf{Z}}[\|\eta_p^t(\boldsymbol{\beta} + \tau_t \mathbf{Z})\|^2]/p, \quad (\text{D.1})$$

for $\mathbf{Z} \sim \mathcal{N}(0, \mathbb{I}_p)$. By the Lipschitz property of η_p^t (Assumption **(A4)**), we have $\mathbb{E}_{\mathbf{Z}}[\|\eta_p^t(\boldsymbol{\beta} + \tau_t \mathbf{Z})\|^2] \leq \mathbb{E}_{\mathbf{Z}}[\|\boldsymbol{\beta} + \tau_t \mathbf{Z}\|^2] \leq 2\|\boldsymbol{\beta}\|^2 + 2p\tau_t^2$. Plugging into (D.1), we find $\text{plim}_p \|\boldsymbol{\beta}^t\|^2/p \stackrel{p}{=} 2 \text{plim}_p \|\boldsymbol{\beta}\|^2/p + 2\tau_t^2 = 2\sigma_{\boldsymbol{\beta}}^2 + 2\tau_t^2$, where the final inequality follows by Assumption **(A2)**.

Now consider the $\hat{\boldsymbol{\beta}}$ result in (7.2). First, note that by definition $\mathcal{C}(\hat{\boldsymbol{\beta}}) \leq \mathcal{C}(\mathbf{0})$ where the cost function $\mathcal{C}(\cdot)$ is defined in (1.2). Using that

$$\mathcal{C}(\mathbf{0}) = \frac{1}{2}\|\mathbf{y}\|^2 = \frac{1}{2}\|\mathbf{X}\boldsymbol{\beta} + \mathbf{w}\|^2 \leq \|\mathbf{X}\boldsymbol{\beta}\|^2 + \|\mathbf{w}\|^2 \leq \sigma_{\max}^2(\mathbf{X})\|\boldsymbol{\beta}\|^2 + \|\mathbf{w}\|^2, \quad (\text{D.2})$$

where $\sigma_{\max}(\mathbf{X})$ is the maximum singular value of \mathbf{X} . We note that this value, $\sigma_{\max}(\mathbf{X})$, is bounded almost surely as $p \rightarrow \infty$ using standard estimates on the singular values of random matrices since \mathbf{X} has i.i.d. Gaussian entries by Assumption **(A1)** (see, for example, [8, Lemma F.2]). Therefore,

$$\text{plim}_p \mathcal{C}(\hat{\boldsymbol{\beta}})/p \leq \text{plim}_p \sigma_{\max}^2(\mathbf{X})\|\boldsymbol{\beta}\|^2/p + \text{plim}_p \|\mathbf{w}\|^2/p \leq \mathsf{B}_{\max}\sigma_{\boldsymbol{\beta}}^2 + \sigma_w^2, \quad (\text{D.3})$$

where we've defined B_{\max} to be a bound on the limit of the maximum singular value, i.e. $\lim_p \sigma_{\max}^2(\mathbf{X}) \leq \mathsf{B}_{\max}$, and the final inequality holds by Assumptions **(A2)** and **(A3)**.

Now we will relate $\frac{1}{p}\|\hat{\boldsymbol{\beta}}\|^2$ to $\frac{1}{p}\mathcal{C}(\hat{\boldsymbol{\beta}})$ and other terms lower-bounded by a constant with high probability. We write $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^\perp + \hat{\boldsymbol{\beta}}^\parallel$ where $\hat{\boldsymbol{\beta}}^\perp \in \ker(\mathbf{X})^\perp$ and $\hat{\boldsymbol{\beta}}^\parallel \in \ker(\mathbf{X})$. Since $\hat{\boldsymbol{\beta}}^\parallel \in \ker(\mathbf{X})$ and $\ker(\mathbf{X})$ is a random subspace of size $p - n = p(1 - \delta)$, by Kashin Theorem (Theorem H.1.), we have that for some constant $\nu_1 = \nu_1(\delta)$, with high probability

$$\|\hat{\boldsymbol{\beta}}^\parallel\|_2^2 \leq \nu_1 \|\hat{\boldsymbol{\beta}}^\parallel\|_1^2/p. \quad (\text{D.4})$$

Then we have the following bound

$$\|\hat{\boldsymbol{\beta}}\|^2 = \|\hat{\boldsymbol{\beta}}^\parallel\|^2 + \|\hat{\boldsymbol{\beta}}^\perp\|^2 \stackrel{(a)}{\leq} \nu_1 \|\hat{\boldsymbol{\beta}}^\parallel\|_1^2/p + \|\hat{\boldsymbol{\beta}}^\perp\|^2 \stackrel{(b)}{\leq} 2\nu_1 \|\hat{\boldsymbol{\beta}}\|_1^2/p + (2\nu_1 + 1)\|\hat{\boldsymbol{\beta}}^\perp\|^2, \quad (\text{D.5})$$

where step (a) holds by (D.4) and step (a) by the Triangle Inequality and Cauchy-Schwarz as follows

$$\|\hat{\boldsymbol{\beta}}^\parallel\|_1^2 = \|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^\perp\|_1^2 \leq (\|\hat{\boldsymbol{\beta}}\|_1 + \|\hat{\boldsymbol{\beta}}^\perp\|_1)^2 \leq 2\|\hat{\boldsymbol{\beta}}\|_1^2 + 2\|\hat{\boldsymbol{\beta}}^\perp\|_1^2 \leq 2\|\hat{\boldsymbol{\beta}}\|_1^2 + 2p\|\hat{\boldsymbol{\beta}}^\perp\|^2.$$

Now we bound the second term on the right side of (D.5). Define $\hat{\sigma}_{\min}(\mathbf{X})$ as the minimum non-zero singular value of \mathbf{X} . By standard results in linear algebra, $\hat{\sigma}_{\min}^2(\mathbf{X})\|\hat{\boldsymbol{\beta}}^\perp\|^2 \leq \|\mathbf{X}\hat{\boldsymbol{\beta}}^\perp\|^2$. Therefore,

$$\hat{\sigma}_{\min}^2(\mathbf{X})\|\hat{\boldsymbol{\beta}}^\perp\|^2 \leq \|\mathbf{X}\hat{\boldsymbol{\beta}}^\perp\|^2 \leq \|\mathbf{X}\hat{\boldsymbol{\beta}}^\perp - \mathbf{y} + \mathbf{y}\|^2 \leq 2\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^\perp\|^2 + 2\|\mathbf{y}\|^2 \leq 2\mathcal{C}(\hat{\boldsymbol{\beta}}) + 2\mathcal{C}(\mathbf{0}) \leq 2\mathcal{C}(\hat{\boldsymbol{\beta}}).$$

Therefore, using (D.2) and (D.3), we have

$$\text{plim}_p \frac{1}{p}\|\hat{\boldsymbol{\beta}}^\perp\|^2 \leq \text{plim}_p \frac{\frac{2}{p}\mathcal{C}(\mathbf{0})}{\hat{\sigma}_{\min}^2(\mathbf{X})} \leq \frac{2(\mathsf{B}_{\max}\sigma_{\boldsymbol{\beta}}^2 + \sigma_w^2)}{\mathsf{B}_{\min}}. \quad (\text{D.6})$$

where we've defined B_{\min} to be a bound on the limit of the minimum non-zero singular value, i.e. $\lim_p \hat{\sigma}_{\min}^2(\mathbf{X}) \geq \mathsf{B}_{\min}$.

Now we bound the first term on the right side of (D.5). Recall the definition of the sort-ed ℓ_1 norm, i.e. $J_{\boldsymbol{\lambda}}(\mathbf{b}) = \sum \lambda_i |\mathbf{b}|_{(i)}$, then using $\lambda_{\min} := \lim_p \min(\boldsymbol{\lambda})$ to lower bound the threshold values,

$$\lambda_{\min} \|\hat{\boldsymbol{\beta}}\|_1 = \sum \lambda_{\min} |\hat{\boldsymbol{\beta}}_i| = \sum \lambda_{\min} |\hat{\boldsymbol{\beta}}|_{(i)} \leq \sum \lambda_i |\hat{\boldsymbol{\beta}}|_{(i)} = J_{\boldsymbol{\lambda}}(\hat{\boldsymbol{\beta}}) \leq \mathcal{C}(\hat{\boldsymbol{\beta}}) \leq \mathcal{C}(\mathbf{0}).$$

Then, using (D.2) and (D.3), we see

$$\plim_p \frac{1}{p} \|\hat{\beta}\|_1 \leq \plim_p \frac{1}{\lambda_{\min}} \left(\frac{1}{p} \mathcal{C}(\mathbf{0}) \right) \leq \frac{1}{\lambda_{\min}} (\mathsf{B}_{\max} \sigma_{\beta}^2 + \sigma_w^2). \quad (\text{D.7})$$

By (D.7), along with the upper bound in (D.5), we have

$$\plim_p \frac{\|\hat{\beta}\|^2}{p} \leq 2\nu_1 \plim_p \frac{\|\hat{\beta}\|_1^2}{p^2} + (2\nu_1 + 1) \plim_p \frac{\|\hat{\beta}^\perp\|^2}{p} \leq \left[\frac{2\nu_1 (\mathsf{B}_{\max} \sigma_{\beta}^2 + \sigma_w^2)}{\lambda_{\min}} \right]^2 + \frac{2(2\nu_1 + 1)(\mathsf{B}_{\max} \sigma_{\beta}^2 + \sigma_w^2)}{\mathsf{B}_{\min}}.$$

□

E Proof of Lemma 7.3

The proof of Lemma 7.3 relies on the following result, Lemma E.1, about the exponential rate of the convergence of the state evolution sequence defined in (6.3). We state and prove Lemma E.1, and Lemma 7.3 is proved afterward.

Lemma E.1. *Assume $\alpha > \mathbf{A}_{\min}(\delta)$ and let $\{\Sigma_{s,t}\}_{s,t \geq 0}$ be defined by the recursion (6.3) with initial condition (6.2). Then there exists constants $B_1, r_1 > 0$ such that for all $t \geq 0$, letting $\tau_* := \lim_t \tau_t$,*

$$|\Sigma_{t,t} - \tau_*^2| \leq B_1 e^{-r_1 t}, \quad \text{and} \quad |\Sigma_{t,t+1} - \tau_*^2| \leq B_1 e^{-r_1 t}.$$

Proof. Throughout the proof, we use the $\{\eta_p^t\}_{p \in \mathbb{N}_{>0}}$ notation introduced in Section 4 and defined in (4.1) with a slight modification to explicitly state the thresholds. Namely, we consider a sequence of denoisers $\eta_p : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^p$ to be those that apply the proximal operator $\text{prox}_{J_{\alpha\tau_t}}(\cdot)$ defined in (1.4), i.e. $\eta_p(\mathbf{v}; \alpha\tau_t) := \text{prox}_{J_{\alpha\tau_t}}(\mathbf{v})$ for a vector $\mathbf{v} \in \mathbb{R}^p$.

Then, per the definition in (6.3), we have

$$\Sigma_{s+1,t+1} = \sigma_w^2 + \lim_p \mathbb{E}\{\eta_p(\mathbf{B} + \tau_s \mathbf{Z}_s; \alpha\tau_s) - \mathbf{B}\}^\top [\eta_p(\mathbf{B} + \tau_t \mathbf{Z}_t; \alpha\tau_t) - \mathbf{B}]/(\delta p),$$

where $\mathbf{B} \sim \mathcal{B}$ i.i.d. elementwise, independent of length- p jointly Gaussian vectors \mathbf{Z}_s and \mathbf{Z}_r having $\mathbb{E}[\mathbf{Z}_s] = \mathbb{E}[\mathbf{Z}_r] = \mathbf{0}$, with covariance $\mathbb{E}\{(\mathbf{Z}_s)_i^2\} = \mathbb{E}\{(\mathbf{Z}_r)_i^2\} = 1$ for any element $i \in [p]$, and $\mathbb{E}\{(\mathbf{Z}_s)_i (\mathbf{Z}_r)_j\} = \frac{\Sigma_{s,r}}{\tau_r \tau_s} \mathbb{I}\{i = j\}$. Recall, $\Sigma_{t,t} = \tau_t^2$ defined in (2.4) and by Theorem 1 we know that $\{\Sigma_{t,t}\}_{t \geq 0}$ is monotone and converges to τ_*^2 as $t \rightarrow \infty$. To prove exponential convergence of $\{\Sigma_{t-1,t}\}_{t \geq 0}$ as claimed in the lemma statement, we construct a discrete dynamical system below.

For $t \geq 1$, define the vector $\mathbf{y}_t = (y_{t,1}, y_{t,2}, y_{t,3}) \in \mathbb{R}^3$ as

$$y_{t,1} \equiv \Sigma_{t-1,t-1} = \tau_{t-1}^2, \quad y_{t,2} \equiv \Sigma_{t,t} = \tau_t^2, \quad y_{t,3} \equiv \Sigma_{t-1,t-1} - 2\Sigma_{t,t-1} + \Sigma_{t,t}. \quad (\text{E.1})$$

A careful argument shows that the vector $\mathbf{y}_t = (y_{t,1}, y_{t,2}, y_{t,3})$ belongs to \mathbb{R}_+^3 . Essentially this requires showing that a matrix $R_T :=$ as in [5, Lemma 5.8] is strictly positive definite. Using the definition of the Σ recursion in (6.3), it is immediate to see that this sequence is updated according to the mapping $\mathbf{y}_{t+1} = G(\mathbf{y}_t)$ where

$$G_1(\mathbf{y}_t) \equiv y_{t,2}, \quad (\text{E.2})$$

$$G_2(\mathbf{y}_t) \equiv \sigma_w^2 + \lim_p \mathbb{E}\{\|\eta_p(\mathbf{B} + \sqrt{y_{t,2}} \mathbf{Z}_t; \alpha\sqrt{y_{t,2}}) - \mathbf{B}\|^2\}/(\delta p), \quad (\text{E.3})$$

$$G_3(\mathbf{y}_t) \equiv \lim_p \mathbb{E}\{\|\eta_p(\mathbf{B} + \sqrt{y_{t,2}} \mathbf{Z}_t; \alpha\sqrt{y_{t,2}}) - \eta_p(\mathbf{B} + \sqrt{y_{t,1}} \mathbf{Z}_{t-1}; \alpha\sqrt{y_{t,1}})\|^2\}/(\delta p), \quad (\text{E.4})$$

where $(\mathbf{Z}_t, \mathbf{Z}_{t-1})$ are length- p jointly Gaussian vectors, independent of $\mathbf{B} \sim B$ i.i.d. elementwise, having $\mathbb{E}[\mathbf{Z}_t] = \mathbb{E}[\mathbf{Z}_{t-1}] = \mathbf{0}$ and with covariance $\mathbb{E}\{([\mathbf{Z}_t]_i)^2\} = \mathbb{E}\{([\mathbf{Z}_{t-1}]_i)^2\} = 1$ for any element $i \in [p]$, and $\mathbb{E}\{[\mathbf{Z}_t]_i [\mathbf{Z}_{t-1}]_j\} = \frac{\Sigma_{t,t-1}}{\tau_t \tau_{t-1}} \mathbb{I}\{i = j\}$. Notice that $\mathbb{E}\{\|\sqrt{y_{t,2}} \mathbf{Z}_t - \sqrt{y_{t,1}} \mathbf{Z}_{t-1}\|^2\} = y_{t,3}$, where we emphasize that $G_3(\mathbf{y}_t)$ depends on $y_{t,3}$ through the covariance of \mathbf{Z}_t and \mathbf{Z}_{t-1} . Moreover, if $\sigma_w^2 > 0$, then $y_{t,1}$ and $y_{t,2}$ are both strictly positive and by the map defined above it is easy to see that $y_{t,3}$ for all $t \geq 0$. This mapping is defined for $y_{t,3} \leq 2(y_{t,1} + y_{t,2})$.

In the following, we will show by induction on t , for $t \geq 1$, that the stronger inequality $y_{t,3} < (y_{t,1} + y_{t,2})$ holds. The initial condition implied by Eq. (6.2) is

$$\begin{aligned} y_{1,1} &= \sigma_w^2 + \mathbb{E}[B^2]/\delta, & y_{1,2} &= \sigma_w^2 + \lim_p \mathbb{E}\{\|\eta_p(\mathbf{B} + \tau_0 \mathbf{Z}_0; \boldsymbol{\alpha} \tau_0) - \mathbf{B}\|^2\}/(\delta p), \\ y_{1,3} &= \lim_p \mathbb{E}\{\|\eta_p(\mathbf{B} + \tau_0 \mathbf{Z}_0; \boldsymbol{\alpha} \tau_0)\|^2\}/(\delta p), \end{aligned}$$

It follows that

$$\begin{aligned} y_{1,1} + y_{1,2} - y_{1,3} &= 2\sigma_w^2 + 2 \lim_p \mathbb{E}\{\mathbf{B}^\top (\mathbf{B} - \eta_p(\mathbf{B} + \tau_0 \mathbf{Z}_0; \boldsymbol{\alpha} \tau_0))\}/(\delta p) \\ &= 2\sigma_w^2 + 2 \lim_p \mathbb{E}_{\mathbf{B}}\{\mathbf{B}^\top (\mathbf{B} - \mathbb{E}_{\mathbf{Z}_0}\{\eta_p(\mathbf{B} + \tau_0 \mathbf{Z}_0; \boldsymbol{\alpha} \tau_0)\})\}/(\delta p). \end{aligned}$$

Using the above, it is easy to show $y_{1,3} < y_{1,1} + y_{1,2}$. This follows since $\mathbb{E}_{\mathbf{B}}\{\mathbf{B}^\top (\mathbf{B} - \mathbb{E}_{\mathbf{Z}_0}\{\eta_p^0(\mathbf{B} + \tau_0 \mathbf{Z}_0)\})\}$ is asymptotically separable using Lemma 3.3 and because the function $x \mapsto x - \mathbb{E}_Z h^0(x + \tau_0 Z)$ is monotone increasing. It follows that $\lim_p \mathbb{E}_{\mathbf{B}}\{\mathbf{B}^\top (\mathbf{B} - \mathbb{E}_{\mathbf{Z}_0}\{\eta_p^0(\mathbf{B} + \tau_0 \mathbf{Z}_0)\})\}/(\delta p) > 0$.

Suppose that $y_{t,3} < y_{t,1} + y_{t,2}$, we want to show $y_{t+1,3} < y_{t+1,1} + y_{t+1,2}$. By the induction hypothesis, $\mathbb{E}\{[\mathbf{Z}_t]_i [\mathbf{Z}_{t-1}]_i\} = \frac{y_{t,1} + y_{t,2} - y_{t,3}}{2\sqrt{y_{t,1}y_{t,2}}} > 0$, so elementwise \mathbf{Z}_t and \mathbf{Z}_{t-1} are positively correlated.

$$\begin{aligned} y_{t+1,1} + y_{t+1,2} - y_{t+1,3} \\ = 2\sigma_w^2 + \lim_p 2\mathbb{E}\{[\eta_p(\mathbf{B} + \sqrt{y_{t,2}} \mathbf{Z}_t; \boldsymbol{\alpha} \sqrt{y_{t,2}}) - \mathbf{B}]^\top [\eta_p(\mathbf{B} + \sqrt{y_{t,1}} \mathbf{Z}_{t-1}; \boldsymbol{\alpha} \sqrt{y_{t,1}}) - \mathbf{B}]\}/(\delta p). \end{aligned} \quad (\text{E.5})$$

Notice that $x \mapsto \eta(b+c \cdot x; \theta) - b$ is monotone for any constants b and $c > 0$ and consider the following result: for g , a monotone function, and X_1 and X_2 , two positively correlated standard Gaussians, $\mathbb{E}[g(X_1)g(X_2)] \geq 0$. This is a special case of a theorem in [30], which shows $\mathbb{E}[g(X_1)g(X_2)] \geq \mathbb{E}[g(X_1)]\mathbb{E}[g(X_2)] = (\mathbb{E}[g(X_1)])^2 > 0$. Then since \mathbf{Z}_t and \mathbf{Z}_{t-1} are positively correlated, $\mathbb{E}\{[\eta_p(\mathbf{B} + \sqrt{y_{t,2}} \mathbf{Z}_t; \boldsymbol{\alpha} \sqrt{y_{t,2}}) - \mathbf{B}]^\top [\eta_p(\mathbf{B} + \sqrt{y_{t,1}} \mathbf{Z}_{t-1}; \boldsymbol{\alpha} \sqrt{y_{t,1}}) - \mathbf{B}]\} \geq 0$, which yields $y_{t+1,3} < (y_{t+1,1} + y_{t+1,2})$.

We can hereafter therefore assume $y_{t,3} < y_{t,1} + y_{t,2}$ for all t .

We will consider the above iteration for arbitrary initialization y_0 (satisfying $y_{0,3} < y_{0,1} + y_{0,2}$) and will show the following three facts:

Fact (i). $y_{t,1}, y_{t,2} \rightarrow \tau_*^2$ as $t \rightarrow \infty$. Further the convergence is monotone.

Fact (ii). If $y_{0,1} = y_{0,2} = \tau_*^2$ and $y_{0,3} \leq 2\tau_*^2$, then $y_{t,1} = y_{t,2} = \tau_*^2$ for all t and $y_{t,3} \rightarrow 0$.

Fact (iii). The Jacobian $J = J_G(y_*)$ of G at $y_* = (\tau_*^2, \tau_*^2, 0)$ has spectral radius $\sigma(J) < 1$.

By simple compactness arguments, Facts (i) and (ii) imply $y_t \rightarrow y_*$ as $t \rightarrow \infty$. (Notice that $y_{t,3}$ remains bounded since $y_{t,3} \leq (y_{t,1} + y_{t,2})$ and by the convergence of $y_{t,1}, y_{t,2}$.) Fact (iii) implies that convergence is exponentially fast.

Proof of Fact (i). Notice that $y_{t,2}$ evolves independently by $y_{t+1,2} = G_2(y_t) = F(y_{2,t}, \boldsymbol{\alpha}\sqrt{y_{2,t}})$, with $F(\cdot, \cdot)$ the state evolution mapping introduced in (2.8). It follows from Proposition 1.3 that $y_{t,2} \rightarrow \tau_*^2$ monotonically for any initial condition. Since $y_{t+1,1} = y_{t,2}$, the same happens for $y_{t,1}$.

Proof of Fact (ii). Consider the function

$$G_*(x) = G_3(\tau_*^2, \tau_*^2, x) = \lim_p \mathbb{E}\{\|\eta_p(\mathbf{B} + \tau_* \mathbf{Z}_t; \boldsymbol{\alpha}\tau_*) - \eta_p(\mathbf{B} + \tau_* \mathbf{Z}_{t-1}; \boldsymbol{\alpha}\tau_*)\|^2\}/(\delta p),$$

where

$$\mathbb{E}\{\mathbf{Z}_t[i] \mathbf{Z}_{t-1}[i]\} = \frac{y_{t,1} + y_{t,2} - y_{t,3}}{2\sqrt{y_{t,1}y_{t,2}}} = \frac{2\tau_*^2 - x}{2\tau_*^2}$$

is no longer time-dependent. This function is defined for $x \in [0, 2\tau_*^2]$. Further G_* can be represented as follows in terms of the independent random vectors $\mathbf{Z}, \mathbf{W} \sim N(0, \mathbb{I})$:

$$G_*(x) = \lim_p \frac{1}{\delta p} \mathbb{E}\{\|\eta_p(\mathbf{B} + \mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} + \mathbf{W}(\frac{1}{2}\sqrt{x}); \boldsymbol{\alpha}\tau_*) - \eta_p(\mathbf{B} + \mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} - \mathbf{W}(\frac{1}{2}\sqrt{x}); \boldsymbol{\alpha}\tau_*)\|^2\},$$

where

$$(\tau_* \mathbf{Z}_{t-1}, \tau_* \mathbf{Z}_t) \stackrel{d}{=} \left(\mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} - \mathbf{W}(\frac{1}{2}\sqrt{x}), \mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} + \mathbf{W}(\frac{1}{2}\sqrt{x}) \right).$$

Obviously $G_*(0) = 0$. A simple Taylor expansion about the first argument around \mathbf{B} yields (recall higher derivatives of η are 0 almost everywhere)

$$\begin{aligned} G_*(x) &= \lim_p \mathbb{E}\left\{\|\eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*) + \left(\mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} + \mathbf{W}(\frac{1}{2}\sqrt{x})\right) \odot \partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*) \right. \\ &\quad \left. - \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*) - \left(\mathbf{Z}\sqrt{\tau_*^2 - \frac{1}{4}x} - \mathbf{W}(\frac{1}{2}\sqrt{x})\right) \odot \partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\| \|^2\right\}/(\delta p) \\ &= \lim_p x \mathbb{E}\{\|\mathbf{W} \odot \partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\}/(\delta p) = \lim_p x \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\}/(\delta p). \end{aligned}$$

Using the above, we study $G'_*(x)$. First, we can exchange the limit and differentiation because $f_p(x) := x \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\}/(\delta p)$ converges uniformly to $f(x) := \lim_p x \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\}/(\delta p)$. To see this, notice f_p, f are linear in x and defined on $[0, 2\tau_*^2]$. Hence for every $\epsilon > 0$, there exists p_0 such that

$$\begin{aligned} |f_{p_0}(x) - f(x)| &= x \left| \frac{1}{\delta p_0} \mathbb{E}\{\|\partial_1 \eta_{p_0}(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\} - \lim_p \frac{1}{\delta p} \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\} \right| \\ &\leq 2\tau_*^2 \left| \frac{1}{\delta p_0} \mathbb{E}\{\|\partial_1 \eta_{p_0}(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\} - \lim_p \frac{1}{\delta p} \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\} \right| < \epsilon. \end{aligned}$$

By uniform convergence we have,

$$G'_*(x) = \lim_p \frac{1}{\delta p} \mathbb{E}\{\|\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)\|^2\} = G'_*(0) \leq \lim_p \frac{1}{\delta p} \sum_{i=1}^p \mathbb{E}\{[\partial_1 \eta_p(\mathbf{B}; \boldsymbol{\alpha}\tau_*)]_i\}.$$

Hence $G'_*(0) < 1$, using (2.10) since $\boldsymbol{\lambda} > \mathbf{0}$. Then $y_{t,3} = [G'_*(0)]^t y_{0,3} \rightarrow 0$ as $t \rightarrow \infty$ as claimed.

Proof of Fact (iii). By the definition of G , the Jacobian is given by

$$J_G(y_*) = \begin{pmatrix} 0 & 1 & 0 \\ 0 & F'(\tau_*^2) & 0 \\ a & G'_*(0) & b \end{pmatrix}$$

denoting $F'(\tau_*^2) \equiv \frac{\partial F}{\partial \tau^2}(\tau^2, \alpha\tau)$ evaluated at $\tau^2 = \tau_*^2$ with a and b constants whose values are not important to the proof. Computing the eigenvalues of the Jacobian, we get $\sigma(J) = \max \{ F'(\tau_*^2), G'_*(0) \}$. Since $G'_*(0) < 1$ proved above and $F(\tau_*^2) < 1$ by Theorem 1, the claim follows. \square

Proof of Lemma 7.3. We show that Lemma 7.3 follows by Lemmas E.1 and 6.2. By Lemma 6.2,

$$\begin{aligned} \text{plim}_n (\|z^t - z^{t-1}\|^2/n - (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2)) &= 0, \\ \text{plim}_p (\|\beta^{t+1} - \beta^t\|^2/(\delta p) - (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2)) &= 0, \end{aligned}$$

and so it is sufficient to show that $\lim_t (\tau_t^2 - 2\Sigma_{t,t-1} + \tau_{t-1}^2) = 0$. Note that this follows from Lemma E.1 since $\tau_t^2 = \Sigma_{t,t}$ and $\tau_{t-1}^2 = \Sigma_{t-1,t-1}$ both converge to τ_*^2 as does $\Sigma_{t,t-1}$. \square

F Technical Details for the Condition (3) Proof

We first introduce some notation and ideas that will be used throughout the proof. The proof is similar to [5, Section 5.3], with the key difference being the concept of equivalence classes as described in Section 5.1.

We now introduce a more general recursion than the AMP algorithm in (1.3a)-(1.3b). Given $w \in \mathbb{R}^n$ and $\beta \in \mathbb{R}^p$, define the column vectors $h^{t+1}, q^{t+1} \in \mathbb{R}^p$ and $b^t, m^t \in \mathbb{R}^n$, recursively, for $t \geq 0$ as follows, starting with initial condition $\beta^0 = 0$ and $z^0 = y$.

$$h^{t+1} = \beta - (X^\top z^t + \beta^t), \quad q^t = \beta^t - \beta, \quad b^t = w - z^t, \quad m^t = -z^t. \quad (\text{F.1})$$

Note that these definitions of h^t and m^t match those used in Section 6.

Denoting $[u|v]$ to mean the matrix of concatenating vectors u, v horizontally, we define

$$\begin{aligned} \underbrace{[h^1 + q^0 | \cdots | h^t + q^{t-1}]}_{A_t} &= X^\top \underbrace{[m^0 | \cdots | m^{t-1}]}_{M_t}, \\ \underbrace{[b^0 | b^1 + \kappa_1 m^0 | \cdots | b^{t-1} + \kappa_{t-1} m^{t-2}]}_{Y_t} &= X \underbrace{[q^0 | \cdots | q^{t-1}]}_{Q_t}, \end{aligned} \quad (\text{F.2})$$

where the scalars κ_t are defined as $\kappa_t := -[\nabla \eta^{t-1}(\beta - h^{t-1})]/n$.

Define the σ -algebra generated by $b^0, \dots, b^{t-1}, m^0, \dots, m^{t-1}, h^1, \dots, h^t, q^0, \dots, q^t$ as \mathfrak{S}_t . Then [4, 8], says that the conditional distribution of the random matrix X given \mathfrak{S}_t is

$$X|_{\mathfrak{S}_t} \stackrel{d}{=} E_t + P_{M_t}^\perp \tilde{X} P_{Q_t}^\perp, \quad (\text{F.3})$$

where $\tilde{X} \stackrel{d}{=} X$ is independent of the conditioning sigma-algebra \mathfrak{S}_t and $E_t = \mathbb{E}(X|\mathfrak{S}_t)$ is given by:

$$E_t := Y_t(Q_t^\top Q_t)^{-1} Q_t^\top + M_t(M_t^\top M_t)^{-1} A_t^\top + M_t(M_t^\top M_t)^{-1} M_t^\top Y_t(Q_t^\top Q_t)^{-1} Q_t^\top.$$

In (F.3), we use the notation $P_{M_t}^\perp = \mathbb{I} - P_{M_t}$ and $P_{Q_t}^\perp = \mathbb{I} - P_{Q_t}$ where P_{Q_t} and P_{M_t} are orthogonal projectors onto column spaces of Q_t, M_t respectively. From now on, since t is fixed, we will drop the subscript t when it is clear. A proof of (F.3) can be found in [4, Lemma 11]. We note that there are no differences in this conditional distribution in the nonseparable case, since the analysis (in both cases) is just that of an i.i.d. Gaussian matrix conditional on linear constraints.

Given the above notations, we claim that Lemma 7.5 is implied by the following statement.

Lemma F.1. Let s be a set of maximal atoms in $[p]$ such that $|s| \leq p(\delta - \gamma)$, for some $\gamma > 0$. Then there exists $\alpha_1 = \alpha_1(\gamma) > 0$ (independent of t) and $\alpha_2 = \alpha_2(\gamma, t) > 0$ (depending on t and γ) with

$$\mathbb{P} \left\{ \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s} \|\mathbf{E}\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \leq \alpha_2 \middle| \mathfrak{S}_t \right\} \leq e^{-p\alpha_1},$$

eventually almost surely as $p \rightarrow \infty$, with $\mathbf{E}\mathbf{v} = \mathbf{Y}(\mathbf{Q}^*\mathbf{Q})^{-1}\mathbf{Q}^*\mathbf{P}_Q\mathbf{v} + \mathbf{M}(\mathbf{M}^*\mathbf{M})^{-1}\mathbf{X}^*\mathbf{P}_Q^\perp\mathbf{v}$.

We prove such implication in the next section now.

Proof of Lemma 7.5. The proof is adapted from [5, Section 5.3.1]. First note that by Borel-Cantelli, it is sufficient to show that, for s measurable on \mathfrak{S}_t and $|s| \leq p(\delta - c)$ there exist $a_1 = a_1(c) > 0$ and $a_2 = a_2(c, t) > 0$, such that

$$\mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s \cup s'} \|\mathbf{X}\mathbf{v}\| < a_2 \right\} \leq 1/p^2,$$

for all p large enough, using $\sigma_{\min}(\mathbf{X}_{S_t \cup S'}) = \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s \cup s'} \|\mathbf{X}\mathbf{v}\|$. To shorten notation, the set $\{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s \cup s'\}$ is denoted $\mathbf{v}(s')$. Now, conditioning on \mathfrak{S}_t , by a union bound,

$$\begin{aligned} \mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\} &\leq \sum_{|s'| \leq a_1 p} \mathbb{P} \left\{ \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\} \\ &\leq \left[\sum_{k=1}^{a_1 p} \binom{p}{k} \right] \max_{|s'| \leq p a_1} \mathbb{P} \left\{ \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\} \leq e^{ph(a_1)} \max_{|s'| \leq a_1 p} \mathbb{P} \left\{ \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\}, \end{aligned} \tag{F.4}$$

where $h(a) = -a \log a - (1-a) \log(1-a)$ is the binary entropy function (cf. [26, Chapter 10, Corollary 9]). Therefore, using iterated expectation and (F.4),

$$\begin{aligned} \mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \right\} &= \mathbb{E} \left\{ \mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\} \right\} \\ &\leq e^{ph(a_1)} \mathbb{E} \left\{ \max_{|s'| \leq a_1 p} \mathbb{P} \left\{ \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \middle| \mathfrak{S}_t \right\} \right\}, \end{aligned}$$

Now, we fix $a_1 < c/2$ in such a way that $h(a_1) \leq \frac{1}{2}\alpha_1(\frac{c}{2})$ and let $a_2 = \frac{1}{2}\alpha_2(\frac{c}{2}, t)$ where α_1 and α_2 are defined by Lemma F.1. Then,

$$\begin{aligned} \mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \right\} &\leq e^{\frac{1}{2}p\alpha_1(\frac{c}{2})} \mathbb{E} \left\{ \max_{|s'| \leq a_1 p} \mathbb{P} \left\{ \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s \cup s'} \|\mathbf{X}\mathbf{v}\| < \frac{1}{2}\alpha_2(\frac{c}{2}, t) \middle| \mathfrak{S}_t \right\} \right\} \\ &\leq e^{\frac{1}{2}p\alpha_1(\frac{c}{2})} \mathbb{E} \left\{ \max_{|s''| \leq p(\delta - \frac{c}{2})} \mathbb{P} \left\{ \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s''} \|\mathbf{X}\mathbf{v}\| < \frac{1}{2}\alpha_2(\frac{c}{2}, t) \middle| \mathfrak{S}_t \right\} \right\}. \end{aligned}$$

Finally, using (cf. [5, Lemma 5.1]),

$$\mathbf{X}\mathbf{v}|_{\mathfrak{S}} \stackrel{d}{=} \mathbf{Y}(\mathbf{Q}^*\mathbf{Q})^{-1}\mathbf{Q}^*\mathbf{P}_Q\mathbf{v} + \mathbf{M}(\mathbf{M}^*\mathbf{M})^{-1}\mathbf{X}^*\mathbf{P}_Q^\perp\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp\mathbf{v}. \tag{F.5}$$

to estimate $\mathbf{X}\mathbf{v}$ and applying Lemma F.1, we get, for all p large enough,

$$\mathbb{P} \left\{ \min_{|s'| \leq a_1 p} \min_{\mathbf{v}(s')} \|\mathbf{X}\mathbf{v}\| < a_2 \right\} \leq e^{\frac{1}{2}p\alpha_1} \mathbb{E} \left\{ \max_{|s''| \leq p(\delta - \frac{c}{2})} e^{-p\alpha_1} \right\} \leq 1/p^2.$$

□

Now we prove Lemma F.1, using a proof that is similar to that of [5, Section 5.3.2]. We first state some lemmas that will be used in the proof, but we will not migrate the full proofs from [5] for the sake of brevity. Instead, we describe the key points of proofs with an emphasis on the technical differences for the SLOPE problem and provide pointers to the original proofs.

The concept of maximal atoms are reflected in these lemmas via the sets s and correspondingly \mathbf{P}_s , where \mathbf{P}_s is the $p \times p$ projector matrix onto the subspace of vectors whose supp^* equals s . In the LASSO case where $\text{supp}^* \equiv \text{supp}$ and $s \equiv S$, the projector is orthogonal, but in general, we must define $\mathbf{P}_s[\cdot, j] = \frac{1}{|I|} \sum_{i \in I} \mathbf{e}_i$ for $j \in I$ where $\mathbf{P}_s[\cdot, j]$ is the j^{th} column of \mathbf{P}_s for $1 \leq j \leq p$ and \mathbf{e}_i is the i^{th} vector of the standard basis. For example, when $p = 4$ and $s = \{\{1\}, \{2, 4\}\}$,

$$\mathbf{P}_s = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 1/2 \end{pmatrix}.$$

Such a projector is not necessarily orthogonal and its rank is described via $|s|$ (the number of equivalence classes), not via $|S|$ (the number of non-zero elements) as for the LASSO. We may view this projector as an orthogonal projector onto the subspace of maximal atoms: for a maximal atom $I \in s$, the projector maps elements whose indices belong to I onto their average value.

We begin with the auxiliary lemmas.

Lemma F.2. *[Adapted from [5, Lemma 5.4]] Let s be a set of maximal atoms in $[p]$ such that $|s| \leq p(\delta - \gamma)$, for some $\gamma > 0$. Recall that $\mathbf{E}\mathbf{v} = \mathbf{Y}(\mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{Q}^\top \mathbf{P}_Q \mathbf{v} + \mathbf{M}(\mathbf{M}^\top \mathbf{M})^{-1} \mathbf{A}^\top \mathbf{P}_Q^\perp \mathbf{v}$ and consider the event*

$$\varepsilon_1 :=$$

$$\left\{ \|\mathbf{E}\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 \geq \frac{\gamma}{4\delta} \|\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 + \frac{\gamma}{4\delta} \|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 \quad \forall \mathbf{v} \text{ s.t. } \|\mathbf{v}\| = 1 \text{ and } \text{supp}^*(\mathbf{v}) \subseteq s \right\}.$$

Then there exists $a = a(\gamma) > 0$ such that $\mathbb{P}\{\varepsilon_1 | \mathfrak{S}_t\} \geq 1 - e^{-pa}$.

Sketch proof. Define an event $\tilde{\varepsilon}_1$ as follows:

$$\tilde{\varepsilon}_1 = \left\{ |(\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v})^\top (\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v})| \leq \left(1 - \frac{\gamma}{2\delta}\right)^{1/2} \|\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \right\}, \quad (\text{F.6})$$

where the event $\tilde{\varepsilon}_1$ is meant to hold for all \mathbf{v} such that $\|\mathbf{v}\| = 1$ and $\text{supp}^*(\mathbf{v}) \subseteq s$. We claim that $\mathbb{P}\{\tilde{\varepsilon}_1 | \mathfrak{S}_t\} \geq 1 - e^{-pa}$. To prove the claim, we use that for any \mathbf{v} , the unit vector $\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v} / \|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|$ belongs to the random linear space $\text{im}(\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{P}_s)$ with dimension at most $p(\delta - \gamma)$. Also, $\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}$ belongs to space spanned by the column space of the matrices \mathbf{M} and of \mathbf{B} where $\mathbf{B}_t = [\mathbf{b}^0 | \dots | \mathbf{b}^{t-1}]$ defined in (F.1) and (F.2), having dimension at most $2t$. Applying Proposition G.1 using $m = n, m\lambda = p(\delta - \gamma)$, $d = 2t$ and $\varepsilon = (1 - \frac{\gamma}{2\delta})^{1/2} (1 - \frac{\gamma}{\delta})^{1/2}$ gives that the event

$$\left(\frac{\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}}{\|\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|} \right)^\top \frac{\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}}{\|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|} \leq \sqrt{\lambda} + \varepsilon = \left(1 - \frac{\gamma}{2\delta}\right)^{1/2},$$

holds with the desired probability, proving the claim. Conditional on event (F.6), one can show

$$\|\mathbf{E}\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 \geq \left(1 - \left(1 - \frac{\gamma}{2\delta}\right)^{1/2}\right) \left\{ \|\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 + \|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\|^2 \right\}.$$

Finally observe that $1 - (1 - \frac{\gamma}{2\delta})^{1/2} \geq \frac{\gamma}{4\delta}$ and therefore since event $\tilde{\varepsilon}_1$ occurring implies ε_1 occurs, giving the desired probability of ε_1 as well. \square

Next we estimate the term $\|\tilde{\mathbf{X}}\mathbf{P}_Q^\perp \mathbf{v}\|^2$ in the above lower bound.

Lemma F.3. [Adapted from [5, Lemma 5.5]] Let s be a set of maximal atoms in $[p]$ such that $|s| \leq p(\delta - \gamma)$, for some $\gamma > 0$. Then there exists constant $c_1 = c_1(\gamma)$, $c_2 = c_2(\gamma)$ such that the event

$$\varepsilon_2 := \left\{ \|\tilde{\mathbf{X}}\mathbf{P}_Q^\perp \mathbf{v}\| \geq c_1(\gamma) \|\mathbf{P}_Q^\perp \mathbf{v}\| \quad \forall \mathbf{v} \text{ such that } \text{supp}^*(\mathbf{v}) \subseteq s \right\}$$

holds with probability $\mathbb{P}\{\varepsilon_2 | \mathfrak{S}_t\} \geq 1 - e^{-pc_2}$.

Sketch proof. Let V be the linear space $V = \text{im}(\mathbf{P}_Q^\perp \mathbf{P}_s)$ having dimension at most $p(\delta - \gamma)$. For all \mathbf{v} with $\text{supp}^*(\mathbf{v}) \subseteq s$,

$$\|\tilde{\mathbf{X}}\mathbf{P}_Q^\perp \mathbf{v}\| \geq \sigma_{\min}(\tilde{\mathbf{X}}|_V) \|\mathbf{P}_Q^\perp \mathbf{v}\|, \quad (\text{F.7})$$

where $\tilde{\mathbf{X}}|_V$ refers to the restriction of $\tilde{\mathbf{X}}$ to V . Then $\sigma_{\min}(\tilde{\mathbf{X}}|_V)$ is distributed as the minimum singular value of a Gaussian matrix of dimensions $p\delta \times \dim(V)$, which is almost surely bounded away from 0 as $p \rightarrow \infty$ (see Theorem G. 2). Large deviation estimates [25] imply that the probability that σ_{\min} is smaller than a constant $c_1(\gamma)$ is exponentially small. \square

In the next step we estimate the norm $\mathbf{E}\mathbf{v}$ by quoting the following result.

Lemma F.4. [5, Lemma 5.6] There exists a constant $c = c(t) > 0$ such that, defining the event,

$$\mathcal{E}_3 := \left\{ \|\mathbf{E}\mathbf{P}_Q \mathbf{v}\| \geq c(t) \|\mathbf{P}_Q \mathbf{v}\|, \|\mathbf{E}\mathbf{P}_Q^\perp \mathbf{v}\| \leq c(t)^{-1} \|\mathbf{P}_Q^\perp \mathbf{v}\|, \text{ for all } \mathbf{v} \in \mathbb{R}^p \right\}, \quad (\text{F.8})$$

we have that \mathcal{E}_3 holds eventually almost surely as $p \rightarrow \infty$.

Finally, we can now prove Lemma F.1 with the ingredients given in Lemmas F.2-F.4. We restate the proof from [5, Lemma 5.3] with minor changes.

Proof of Lemma F.1. We start with Lemma F.4 by which we assume that event \mathcal{E}_3 holds for some function $c = c(t)$ (without loss of generality $c < 1/2$). For $\alpha_2(t) > 0$ small enough, let \mathcal{E} be the event

$$\mathcal{E} := \left\{ \min_{\|\mathbf{v}\|=1, \text{supp}^*(\mathbf{v}) \subseteq s} \|\mathbf{E}\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \leq \alpha_2(t) \right\}. \quad (\text{F.9})$$

First assume $\|\mathbf{P}_Q^\perp \mathbf{v}\| \leq c^2/10$, from which it follows,

$$\begin{aligned} \|\mathbf{E}\mathbf{v} - \mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| &\geq \|\mathbf{E}\mathbf{P}_Q \mathbf{v}\| - \|\mathbf{E}\mathbf{P}_Q^\perp \mathbf{v}\| - \|\mathbf{P}_M \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \\ &\geq c \|\mathbf{P}_Q \mathbf{v}\| - (c^{-1} + \|\tilde{\mathbf{X}}\|_2) \|\mathbf{P}_Q^\perp \mathbf{v}\| \geq \frac{c}{2} - \frac{c}{10} - \|\tilde{\mathbf{X}}\|_2 \frac{c^2}{10} = \frac{2c}{5} - \|\tilde{\mathbf{X}}\|_2 \frac{c^2}{10}, \end{aligned}$$

where the last inequality uses $\|\mathbf{P}_Q \mathbf{v}\| = \sqrt{1 - \|\mathbf{P}_Q^\perp \mathbf{v}\|^2} \geq 1/2$ under the assumption $\|\mathbf{P}_Q^\perp \mathbf{v}\| \leq c^2/10$. Therefore, using Lemma F.2, we get

$$\mathbb{P}\{\mathcal{E} | \mathfrak{S}_t\} \leq \mathbb{P}\left\{ \frac{2c}{5} - \|\tilde{\mathbf{X}}\|_2 \frac{c^2}{10} \leq \left(\frac{4\delta}{\gamma}\right)^{1/2} \alpha_2(t) \middle| \mathfrak{S}_t \right\} + e^{-pa},$$

and the thesis follows from large deviation bounds on the norm $\|\tilde{\mathbf{X}}\|_2$ (see [24]) by first taking c small enough, and then choosing $\alpha_2(t) < \frac{c}{5} \sqrt{\frac{\gamma}{4\delta}}$.

Next assume $\|\mathbf{P}_Q^\perp \mathbf{v}\| \geq c^2/10$. By Lemma F.2 and F.3, we can assume events \mathcal{E}_1 and \mathcal{E}_2 hold. Therefore $\|\mathbf{E}\mathbf{v} + \mathbf{P}_M^\perp \tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \geq (\frac{\gamma}{4\delta})^{1/2} \|\tilde{\mathbf{X}} \mathbf{P}_Q^\perp \mathbf{v}\| \geq (\frac{\gamma}{4\delta})^{1/2} c_1(\gamma) \|\mathbf{P}_Q^\perp \mathbf{v}\|$, proving our thesis. \square

G Some Useful Auxiliary Material

We collect some auxiliary results that are necessary in our proof. Most of these are results that were initially stated in [5] that we repeat here for the reader.

The following proposition is used in the proof of Lemma F.2. The proof is identical to that of [5, Proposition E.1] and it follows from a standard concentration of measure argument in [24]. For this reason, we don't repeat it here.

Proposition G.1. *Let $V \subseteq \mathbb{R}^m$ a uniformly random linear space of dimension d . For $\lambda \in (0, 1)$, let \mathbf{P}_λ denote the projector onto the first $m\lambda$ maximal atoms in $[m]$: assume that $s = \{I_1, \dots, I_d\}$, is the set of maximal atoms, then the j^{th} column, $\mathbf{P}_\lambda[:, j] = \frac{1}{|I_r|} \sum_{i \in I_r} \mathbf{e}_i$ if $j \in I_r$ for some $r \leq m\lambda$; otherwise $\mathbf{P}_\lambda[:, j] = \mathbf{0}$. Define $Z(\lambda) := \sup\{\|\mathbf{P}_\lambda \mathbf{v}\| : \mathbf{v} \in V, \|\mathbf{v}\| = 1\}$. Then, for any $\varepsilon > 0$ there exists $c(\varepsilon) > 0$ such that, for all m large enough (and d fixed) $\mathbb{P}\{|Z(\kappa) - \sqrt{\lambda}| \geq \varepsilon\} \leq e^{-mc(\varepsilon)}$.*

We next state a result due to Kashin [22] relating to the equivalence of ℓ^2 and ℓ^1 norms on random vector spaces (cf. also [5, Theorem F.1]).

Theorem G.1. [22] *For any positive number v there exist a universal constant c_v such that for any $n \geq 1$, with probability at least $1 - 2^{-n}$, for a uniformly random subspace $V_{n,v}$ of dimension $\lfloor n(1-v) \rfloor$, for all $x \in V_{n,v}$, we have $c_v \|x\|_2 \leq \|x\|_1 / \sqrt{n}$.*

Finally we state a general result about the limit behavior of extreme singular values of random matrices, as proved in [1] (cf. also [5, Theorem F.2]).

Theorem G.2. [1] *Let $\mathbf{A} \in \mathbb{R}^{n \times p}$ have i.i.d. entries with $\mathbb{E}\{A_{ij}\} = 0$, $\mathbb{E}\{A_{ij}^2\} = 1/n$, and $n/p = \delta$. Let $\sigma_{\max}(\mathbf{A})$ be its largest singular value, and $\hat{\sigma}_{\min}(\mathbf{A})$ be its smallest non-zero singular value. Then,*

$$\lim_{p \rightarrow \infty} \sigma_{\max}(\mathbf{A}) \stackrel{a.s.}{=} 1/\sqrt{\delta} + 1, \quad \text{and} \quad \lim_{p \rightarrow \infty} \hat{\sigma}_{\min}(\mathbf{A}) \stackrel{a.s.}{=} 1/\sqrt{\delta} - 1.$$