

# Quickly Finding the Best Linear Model in High Dimensions

Yahya Sattar

Samet Oymak

**Abstract**—We study the problem of finding the best linear model that can minimize least-squares loss given a dataset. While this problem is trivial in the low-dimensional regime, it becomes more interesting in high-dimensions where the population minimizer is assumed to lie on a manifold such as sparse vectors. We propose projected gradient descent (PGD) algorithm to estimate the population minimizer in the finite sample regime. We establish linear convergence rate and data-dependent estimation error bounds for PGD. Our contributions include: 1) The results are established for heavier tailed sub-exponential distributions besides sub-gaussian. 2) We directly analyze the empirical risk minimization and do not require a realizable model that connects input data and labels. 3) Our PGD algorithm is augmented to learn the bias terms which boosts the performance. The numerical experiments validate our theoretical results.

**Index Terms**—high-dimensional estimation, projected gradient descent, one-bit compressed sensing, gaussian width.

## I. INTRODUCTION

Supervised learning is concerned with finding a relation between the input-output pairs  $(\mathbf{x}_i, y_i)_{i=1}^n \in \mathbb{R}^p \times \mathbb{R}$ . The simplest relations are linear functions where the output  $y_i$  is estimated by a linear function of the input, that is,  $\hat{y}_i = \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle$ . Using quadratic loss, we can find the optimal  $\boldsymbol{\theta}$  with a simple linear regression which minimizes.

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2.$$

If the samples are i.i.d. and input has identity covariance, the population minimizer ( $n \rightarrow \infty$ ) is simply given by

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \mathbb{E}[\mathcal{L}(\boldsymbol{\theta})] = \mathbb{E}[y\mathbf{x}].$$

where  $(\mathbf{x}, y)$  is drawn from same distribution as data. In many applications, we operate in the high-dimensional regime where we have fewer samples than the parameter dimension i.e.  $n \ll p$ . In this case, the problem is ill-posed; however, if  $\boldsymbol{\theta}^*$  lies on a low-dimensional manifold, we can take advantage of this information to solve

the problem. We assume  $\boldsymbol{\theta}^*$  is structured-sparse, for instance, it can be a signal that is sparse in a dictionary or it can be a low-rank matrix. If  $\mathcal{R}$  is a regularization function that promotes this structure, we can solve the regularized empirical risk minimization (ERM)

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_{\ell_2}^2 \quad \text{subject to } \mathcal{R}(\boldsymbol{\theta}) \leq R. \quad (1)$$

where  $\mathbf{y} = [y_1 \dots y_n]^T \in \mathbb{R}^n$  and  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$  are the output labels and data matrix respectively. This problem is well-studied in the statistics and compressed sensing (CS) literature. However, much of the theory literature is concerned with the scenario where the problem is realizable i.e. the outputs are explicitly generated with respect to some ground truth vector  $\mathbf{a}$ . In the simplest scenario, input/output relation can be  $y = \langle \mathbf{x}, \mathbf{a} \rangle + z$  where  $z$  is independent zero-mean noise vector. In this case, one simply has  $\boldsymbol{\theta}^* = \mathbf{a}$ . Such realizability assumption is also common in the single-index models [1], [2]. One contribution of this paper will be analyzing regularized ERM without the realizability assumption.

Bias in the data can negatively affect the estimation quality. Assuming input is zero-mean, instead of solving (1) we can solve a modified problem which accounts for the mean of the output as well. Again, denoting the regularization function by  $\mathcal{R}$ , we will solve the modified problem

$$\hat{\boldsymbol{\theta}}, \hat{\mu} = \arg \min_{\boldsymbol{\theta}, \mu} \mathcal{L}(\boldsymbol{\theta}, \mu) \quad \text{subject to } \mathcal{R}(\boldsymbol{\theta}) \leq R. \quad (2)$$

where the loss is given by  $\mathcal{L}(\boldsymbol{\theta}, \mu) = \frac{1}{2} \|\mathbf{y} - [\mathbf{X} \ \mathbf{1}] \begin{bmatrix} \boldsymbol{\theta} \\ \mu \end{bmatrix}\|_{\ell_2}^2$ . We will show that solving problem (2) is essentially equivalent to solving (1) with debiased output hence it will result in more accurate estimation. The goal of this paper is studying problem (2) under a general algorithmic framework, establishing finite-sample statistical and algorithmic convergence, and addressing practical considerations on the data distribution. In particular, we are interested in how well one can estimate the best linear model (BLM) given by the pair  $(\boldsymbol{\theta}^* = \mathbb{E}[y\mathbf{x}], \mu^* = \mathbb{E}[y])$ .

For estimation, we will utilize the projected gradient descent algorithm given by the iterates

$$\begin{aligned}\boldsymbol{\theta}_{\tau+1} &= \mathcal{P}_{\mathcal{K}}(\boldsymbol{\theta}_{\tau} - \eta \nabla \mathcal{L}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_{\tau}, \mu_{\tau})), \\ \mu_{\tau+1} &= \mu_{\tau} - \eta \nabla \mathcal{L}_{\mu}(\boldsymbol{\theta}_{\tau}, \mu_{\tau}),\end{aligned}\quad (3)$$

where  $\mathcal{P}_{\mathcal{K}}$  projects onto the constraint set  $\mathcal{K} = \{\boldsymbol{\theta} \in \mathbb{R}^p \mid \mathcal{R}(\boldsymbol{\theta}) \leq R\}$  and  $\eta$  is the step size.

#### A. Relation to Prior Work

There is a significant amount of literature on nonlinear (or one-bit) CS [2]–[12]. [4], [13]–[16] study algorithmic and statistical convergence rates for first order methods such as projected/proximal gradient descent. For nonlinear CS, [4], [5], [7], [17] provide statistical analysis of single index estimation with a focus on Gaussian data. Recently, one-bit CS techniques have been extended to sub-gaussian distributions using dithering trick which adds noise before quantization [18]–[21]. Dithering is introduced to guarantee consistent estimation of the ground-truth parameter. The papers [22]–[26] address non-gaussianity by utilizing Stein identity which requires access to the distribution of the input samples. Closer to us [27] studies the constrained empirical risk minimization with linear functions and squared loss with a focus on convex problems. In comparison our analysis applies to a broader class of distributions and focus on first order algorithms. Much of our analysis focuses on addressing subexponential samples, which requires tools from high-dimensional probability [28], [29].

Our results apply to general regularizers and borrow ideas from [4]–[7]. Similar to these, we view the non-linearity between input and output as an additive noise. The convergence analysis of projected gradient descent is a rather well-understood topic and we utilize insights from [13]–[16] for our analysis.

#### B. Contributions

At a high-level our work has three distinguishing features compared to the prior literature.

- Subexponential samples: Most nonlinear CS results apply to Gaussian or subgaussian data when dithering trick is utilized [18]–[21]. We take advantage of the recent techniques for subexponential distributions to provide statistical/computational guarantees for heavier-tailed distributions.

- No realizability assumption: Nonlinear CS literature is typically concerned with a ground-truth vector to be recovered. For instance, one-bit CS aims to learn  $\boldsymbol{\theta}$  from samples of type  $y = \text{sgn}(\boldsymbol{\theta}^T \boldsymbol{x})$ . Unlike these, we do not enforce such relationship to exist between input and output, hence the results apply under much weaker assumptions. Instead of a ground-truth  $\boldsymbol{\theta}$ , we work with

the population BLM  $\boldsymbol{\theta}^*$ . However,  $\boldsymbol{\theta}^*$  can be shown to coincide with ground truth when it exists, if the input distribution is *nice* (e.g. Gaussian) [4]–[7].

- Bias estimation: Our analysis addresses the bias in the output by solving the modified problem (2). We show that (2) can be studied in a similar fashion to (1) by studying the statistical properties of the concatenated data matrix. However, empirically this modification results in a substantial improvement in estimation.

#### C. Paper Organization

We review mathematical background and formulate the problem in Section II. We introduce our main results on statistical and computational convergence guarantees in Section III. Section IV provides numerical experiments to corroborate our theoretical results. Proofs of the main results are provided in Section V and finally the concluding remarks are made in Section VI.

## II. PRELIMINARIES AND PROBLEM FORMULATION

In this section we introduce statistical quantities which are utilized to characterize the benefits of the regularization  $\mathcal{R}$ .

We first set the notation.  $c, c_0, \dots, C$  denote positive absolute constants. For a vector  $\boldsymbol{v}$ , we denote its Euclidean norm by  $\|\boldsymbol{v}\|_{\ell_2}$ . Similarly for a matrix  $\boldsymbol{X}$ , we denote its spectral norm by  $\|\boldsymbol{X}\|$ . Given a set  $S$ , let  $\text{cl}(S)$  and  $\text{clconv}(S)$  be the minimal closed set and minimal closed-convex set containing  $S$  respectively. Let  $\text{rad}(S)$  denote the set radius  $\sup_{\boldsymbol{v} \in S} \|\boldsymbol{v}\|_{\ell_2}$ . For closed sets, let  $\mathcal{P}_S(\cdot)$  be the projection operator defined as  $\mathcal{P}_S(\boldsymbol{a}) = \arg \min_{\boldsymbol{v} \in S} \|\boldsymbol{a} - \boldsymbol{v}\|_{\ell_2}$ .  $\mathcal{N}(\mu, \sigma^2)$  denotes the normal distribution and  $\mathcal{B}^p$  denote the unit ball in  $\mathbb{R}^p$ .  $\mathbf{1}$  is the all ones vector of proper dimension. We will use  $\gtrsim$  and  $\lesssim$  for inequalities that hold up to a constant factor.

Suppose we are given  $n$  i.i.d. samples  $(\boldsymbol{x}_i, y_i)_{i=1}^n \sim (\boldsymbol{x}, y)$ . To keep the exposition clean, we assume that  $\boldsymbol{x}$  is whitened, that is, it has zero-mean and identity covariance. We will aim to find a linear relation between the modified input-output pairs  $([\boldsymbol{x}_i^T \ 1]^T, y_i)_{i=1}^n$ . Let us consider the statistical properties of our modified estimate in the population limit which is given by

$$\begin{aligned}\boldsymbol{\theta}^*, \mu^* &= \arg \min_{\boldsymbol{\theta}, \mu} \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}, \mu)] \\ &= \mathbb{E}[y\boldsymbol{x}], \mathbb{E}[y].\end{aligned}$$

Thus, in the limiting case,  $\mu^*$  captures the mean of the output and  $\boldsymbol{\theta}^*$  is the ideal solution of the problem with debiased output. Our goal is estimating the population minimizer  $\boldsymbol{\theta}^*, \mu^*$ ; which minimizes the expected quadratic loss  $\mathbb{E}[(y - \boldsymbol{\theta}^T \boldsymbol{x} - \mu)^2]$ . As discussed in Section I, assuming  $\boldsymbol{\theta}^*$  is structured sparse, we consider

a non-asymptotic estimation of  $\theta^*, \mu^*$  via problem (2). To proceed with analysis, set

$$\mathcal{K} = \{\theta \in \mathbb{R}^p \mid \mathcal{R}(\theta) \leq R\}, \quad (4)$$

$$\mathcal{K}_{\text{ext}} = \{[\theta^T \ \mu]^T \in \mathbb{R}^{p+1} \mid \mathcal{R}(\theta) \leq R\}. \quad (5)$$

We investigate the PGD algorithm (3) which can be written as

$$\begin{bmatrix} \theta_{\tau+1} \\ \mu_{\tau+1} \end{bmatrix} = \mathcal{P}_{\mathcal{K}_{\text{ext}}} \left( \begin{bmatrix} \theta_{\tau} \\ \mu_{\tau} \end{bmatrix} + \eta [\mathbf{X} \ \mathbf{1}]^T \left( \mathbf{y} - [\mathbf{X} \ \mathbf{1}] \begin{bmatrix} \theta_{\tau} \\ \mu_{\tau} \end{bmatrix} \right) \right), \quad (6)$$

where  $\eta$  is a fixed learning rate and  $[\mathbf{X} \ \mathbf{1}] \in \mathbb{R}^{n \times (p+1)}$  is the modified data matrix constructed as follows

$$[\mathbf{X} \ \mathbf{1}] = \begin{bmatrix} \mathbf{x}_1^T & 1 \\ \vdots & \\ \mathbf{x}_n^T & 1 \end{bmatrix}.$$

Following [4], [30] PGD analysis can be related to the *tangent ball* around the population parameter  $\theta^*$  which is given by

$$\mathcal{C} = \text{cl}(\{\alpha \mathbf{v} \mid \mathbf{v} + \theta^* \in \mathcal{K}, \alpha \geq 0\}) \cap \mathcal{B}^p. \quad (7)$$

Similarly, we define the extended tangent ball as follows

$$\mathcal{C}_{\text{ext}} = \left\{ \begin{bmatrix} \alpha \mathbf{v} \\ \gamma \end{bmatrix} \mid \alpha \geq 0, \mathbf{v} \in \mathcal{C}, \gamma \in \mathbb{R} \right\} \cap \mathcal{B}^{p+1}. \quad (8)$$

The two definitions above are closely related. For any vector  $\mathbf{v} \in \mathcal{C}$ , we have that  $[\sqrt{1-\gamma^2}\mathbf{v}^T \ \gamma]^T \in \mathcal{C}_{\text{ext}}$  for  $0 \leq \gamma \leq 1$ . In the following we will express the convergence rates and residual errors of the PGD algorithm (3) in terms of the statistical properties of the tangent balls.

**Technical approach:** Denoting the parameter estimation error in (6) by  $\mathbf{h}_{\tau} = [\theta_{\tau}^T \ \mu_{\tau}^T]^T - [\theta^{*T} \ \mu^{*T}]^T$  and the effective noise by  $\mathbf{w} = \mathbf{y} - [\mathbf{X} \ \mathbf{1}][\theta^{*T} \ \mu^{*T}]^T$ , the PGD update can be shown to obey [14] (see Eq. (VI.10))

$$\|\mathbf{h}_{\tau+1}\|_{\ell_2} \leq \kappa (\|\mathbf{h}_{\tau}\|_{\ell_2} \rho(\mathcal{C}) + \eta \nu(\mathcal{C})) \quad (9)$$

where  $\kappa$  is a numerical constant which is equal to 1 for convex regularizer  $\mathcal{R}$  and 2 for arbitrary  $\mathcal{R}$  and

$$\rho(\mathcal{C}) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}_{\text{ext}}} |\mathbf{u}^T (\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \mathbf{v}|, \quad (10)$$

$$\nu(\mathcal{C}) = \sup_{\mathbf{v} \in \mathcal{C}_{\text{ext}}} |\mathbf{v}^T [\mathbf{X} \ \mathbf{1}]^T \mathbf{w}|. \quad (11)$$

Here  $\rho$  captures the algorithmic convergence and  $\nu$  captures the statistical accuracy in terms of regularization. To achieve statistical learning bounds, we need to characterize the quantities above in finite sample. Existing literature provides a fairly good understanding of the related terms when  $\mathbf{X}$  has subgaussian rows or  $\mathbf{w}$  is independent of  $\mathbf{X}$ . The technical contributions of this work are i) extending these results to subexponential samples, ii) allowing for nonlinear dependencies between

the noise and data, and iii) addressing the bias term by studying the concatenated matrix  $[\mathbf{X} \ \mathbf{1}]$ . To proceed with statistical analysis, we introduce Gaussian width.

*Definition 2.1 ((Perturbed) Gaussian width [29]):* The Gaussian width of a set  $S \subset \mathcal{B}^p$  is defined as

$$\omega(S) = \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_p)} \left[ \sup_{\mathbf{v} \in S} \mathbf{v}^T \mathbf{g} \right].$$

Let  $C > 0$  be an absolute constant. Given an integer  $n \geq 1$ , the perturbed Gaussian width  $\omega_n(T)$  of  $T \subset \mathcal{B}^d$  is defined as

$$\omega_n(T) = \min_{\substack{\text{clconv}(S) \supseteq T \\ \text{rad}(S) \leq C}} \omega(S) + \frac{\gamma_1(S)}{\sqrt{n}}$$

where  $\gamma_1(S)$  is Talagrand's  $\gamma_1$ -functional (see [28]) with  $\ell_2$ -metric.

Gaussian width helps to quantify the complexity of the regularized problem and determines the sample complexity of the linear inverse problems i.e. high-dimensional problems become manageable in the regime  $n \gtrsim \omega^2(\mathcal{C})$  [30], [31]. Perturbed width is introduced more recently in [29] to address subexponential samples. [29] shows that, for standard regularizers such as  $\ell_0, \ell_1$ , subspace, and rank constraints, we have that

$$\omega^2(\mathcal{C}) \sim \omega_n^2(\mathcal{C}) \quad (12)$$

in the interesting regime  $n \geq \omega^2(\mathcal{C})$ . Hence, perturbed width has the same statistical accuracy of Gaussian width but applies to subexponential samples.

As illustrated in Table I, square of the Gaussian width captures the degrees of freedom for practical regularizers. Table I is obtained by setting  $R = \mathcal{R}(\theta^*)$  in (4). In practice, a good choice for  $R$  can be found by using cross validation. It is also known that the performance of PGD is robust to choice of  $R$  (see Thm 2.6 of [14]).

Constraint	Parameter vector model	$\omega^2(\mathcal{C})$
None	$\theta^* \in \mathbb{R}^p$	$p$
Sparsity $\ \cdot\ _{\ell_0}$	$s$ non-zero entries	$s \log(6p/s)$
$\ell_1$ norm $\ \cdot\ _{\ell_1}$	$s$ non-zero entries	$s \log(6p/s)$
Subspace	$\theta^* \in \mathcal{S}, \dim(\mathcal{S}) = k$	$k$
Matrix rank	$\text{rank}(\text{mat}(\theta^*)) \leq r$	$rp^{1/2}$

TABLE I: List of low-dimensional models and corresponding Gaussian widths (up to a constant factor) for the constraint sets  $\mathcal{K} = \{\theta \mid \mathcal{R}(\theta) \leq \mathcal{R}(\theta^*)\}$ . If constraint is set membership such as subspace,  $\mathcal{R}(\theta) = 0$  inside the set and  $\infty$  outside. Furthermore, we represent the vector  $\theta^* \in \mathbb{R}^p$  in matrix form as  $\text{mat}(\theta^*) \in \mathbb{R}^{\sqrt{p} \times \sqrt{p}}$ .

The next statistical quantity required in our analysis is the Orlicz norm defined as.

*Definition 2.2 (Orlicz norms):* For a scalar random variable Orlicz- $a$  norm is defined as

$$\|X\|_{\psi_a} = \sup_{p \geq 1} p^{-1/a} (\mathbb{E}[|X|^p])^{1/p}$$

Orlicz- $a$  norm of a vector  $\mathbf{x} \in \mathbb{R}^d$  is defined as  $\|\mathbf{x}\|_{\psi_a} = \sup_{\mathbf{v} \in \mathcal{B}^d} \|\mathbf{v}^T \mathbf{x}\|_{\psi_a}$ . Subexponential and subgaussian norms are special cases of Orlicz- $a$  norm given by  $\|\cdot\|_{\psi_1}$  and  $\|\cdot\|_{\psi_2}$  respectively.

Based on perturbed Gaussian width definition, we will show that one can upper bound the critical quantities (10) and (11). In return, this will reveal the statistical and computational performance of the PGD algorithm. This is the topic of the next section which states our main results.

### III. MAIN RESULTS

In this section we estimate the convergence rate and the statistical accuracy of the PGD algorithm as a function of sample size, complexity of the parameter (e.g. sparsity level), and the distribution of the data (whether subgaussian or subexponential). Our main theorem establishes a linear convergence rate of PGD and shows that PGD achieves statistically efficient error rates. We first describe the data model.

*Definition 3.1 (Isotropic vector):*  $\mathbf{x} \in \mathbb{R}^p$  is called an isotropic Orlicz- $a$  vector if it is zero-mean with identity covariance and if its Orlicz- $a$  norm  $\|\mathbf{x}\|_{\psi_a}$  is bounded by an absolute constant.

*Definition 3.2 ( $\sigma$ -noisy datasets):* We assume the samples  $(y_i, \mathbf{x}_i)_{i=1}^n$ . We call a dataset  $\sigma$ -Orlicz- $a$  if the input samples are isotropic Orlicz- $a$  vectors and the residual at the ground truth obeys

$$\|y - \mathbf{x}^T \boldsymbol{\theta}^* - \mu^*\|_{\psi_a} \leq \sigma.$$

We call  $\sigma$ -Orlicz-1 dataset  $\sigma$ -subexponential and  $\sigma$ -Orlicz-2 dataset  $\sigma$ -subgaussian.

Note that residual at the ground truth is the noise in our problem which may be function of the nonlinearity. Our main results capture the PGD performance for different dataset models.

*Theorem 3.3 (Subgaussian):* Suppose  $(\mathbf{x}_i, y_i)_{i=1}^n$  is a  $\sigma$ -subgaussian dataset. Assume  $n \gtrsim (\omega(\mathcal{C}) + t)^2$  and set learning rate  $\eta = 1/n$ . Let  $\mathcal{R}$  be an arbitrary regularizer. Starting from any initial estimate  $[\boldsymbol{\theta}_0^T \mu_0]^T$ , with probability at least  $1 - 6 \exp(-C_0 t^2/2) - 4n^{-100}$ , all PGD iterates (6) obeys

$$\begin{aligned} \left\| \begin{bmatrix} \boldsymbol{\theta}_\tau - \boldsymbol{\theta}^* \\ \mu_\tau - \mu^* \end{bmatrix} \right\|_{\ell_2} &\leq \left( c \frac{\omega(\mathcal{C}) + t}{\sqrt{n}} \right)^\tau \left\| \begin{bmatrix} \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* \\ \mu_0 - \mu^* \end{bmatrix} \right\|_{\ell_2} \\ &\quad + C\sigma \frac{(\omega(\mathcal{C}) + t)\sqrt{\log(n)}}{\sqrt{n}}. \end{aligned}$$

Similarly, for subexponential samples, we have the following theorem which applies to convex regularizers.

*Theorem 3.4 (Subexponential):* Suppose  $(\mathbf{x}_i, y_i)_{i=1}^n$  is a  $\sigma$ -subexponential dataset. Set  $q = (n + p) \log^3(n + p)$ . Set learning rate  $\eta = c_0/q$ , suppose  $\mathcal{R}$  is convex and  $n \gtrsim (\omega_n(\mathcal{C}) + t)^2$ . Starting from initialization  $[\boldsymbol{\theta}_0^T \mu_0]^T$ , with probability at least  $1 - 7 \exp(-C_0 \min\{n, t\sqrt{n}, t^2\}) - 6p(n + p)^{-100}$ , all PGD iterates (6) obey

$$\begin{aligned} \left\| \begin{bmatrix} \boldsymbol{\theta}_\tau - \boldsymbol{\theta}^* \\ \mu_\tau - \mu^* \end{bmatrix} \right\|_{\ell_2} &\leq \left( 1 - \frac{cn}{q} \right)^\tau \left\| \begin{bmatrix} \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* \\ \mu_0 - \mu^* \end{bmatrix} \right\|_{\ell_2} \\ &\quad + C\sigma \frac{(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}}. \end{aligned}$$

Both of these results show that PGD iterates converge to population parameters  $\boldsymbol{\theta}^*$ ,  $\mu^*$  at a linear rate. Subexponential theorem requires a more conservative choice of learning rate. The statistical estimation error grows as  $\omega(\mathcal{C})/\sqrt{n}$  for subgaussian and  $\omega_n(\mathcal{C})/\sqrt{n}$  for subexponential. Since our results apply in the regime  $n \gtrsim \omega^2(\mathcal{C})$ , following (12), statistical errors associated with subgaussian and subexponential are same up to a constant for typical regularizers.

Our main results follow from Theorems 3.5 and 3.6 which are the topics of the following sections.

#### A. Bounding the Error due to Nonlinearity

In this section we provide a bound on the effective noise level  $\nu(\mathcal{C})$ ; which is crucial for assessing statistical accuracy. This term arises from the nonlinearity and noise associated with the relation between input and output. For example, for single-index models, we have  $\mathbb{E}[y | \mathbf{x}] = \phi(\mathbf{x}^T \boldsymbol{\theta}_{\text{GT}})$  for some link function  $\phi$  and ground truth  $\boldsymbol{\theta}_{\text{GT}}$ , and  $\phi$  becomes the source of the nonlinearity. Our approach is similar to [4]–[7], [27] and treats the nonlinearity as a noise. The finite sample noise is captured by the residual vector

$$\mathbf{w} = \mathbf{y} - \mathbf{X} \boldsymbol{\theta}^* - \mathbf{1} \mu^*. \quad (13)$$

Following  $\nu(\mathcal{C})$  term in (11), the contribution of the residual  $\mathbf{w}$  to the estimated parameter is captured by the vector

$$\mathbf{e} = [\mathbf{X} \ \mathbf{1}]^T \mathbf{w} = \sum_{i=1}^n (y_i - \mu^* - \mathbf{x}_i^T \boldsymbol{\theta}^*) \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix}. \quad (14)$$

Our key observation is that the properties of  $\mathbf{e}$  can be characterized under fairly general assumptions compared to the existing literature; which is mostly restricted to zero-mean subgaussian samples.

*Theorem 3.5 (Statistical error):* Suppose  $(\mathbf{x}_i, y_i)_{i=1}^n \sim (\mathbf{x}, y)$  is a  $\sigma$ -subgaussian dataset. Let the tangent balls  $\mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in (7) and (8) respectively.

Assume  $n \gtrsim (\omega(\mathcal{C}) + t)^2$ . Then, with probability at least  $1 - 2 \exp(-t^2/2) - 4n^{-100}$ , we have

$$\frac{\nu(\mathcal{C})}{n} \lesssim \frac{\sigma(\omega(\mathcal{C}) + t) \sqrt{\log(n)}}{\sqrt{n}}.$$

where  $\nu(\mathcal{C})$  is the effective noise given by (11). If  $(\mathbf{x}_i, y_i)_{i=1}^n$  is a  $\sigma$ -subexponential dataset and  $n \gtrsim (\omega_n(\mathcal{C}) + t)^2$ , with probability at least  $1 - 2 \exp(-c \min\{t\sqrt{n}, t^2\}) - 4n^{-100}$ , we have

$$\frac{\nu(\mathcal{C})}{n} \lesssim \frac{\sigma(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}}.$$

This theorem establishes the crucial finite sample upper bounds on  $\nu(\mathcal{C})$  for both subgaussian and subexponential data as a function of Gaussian width of the tangent ball.

### B. Controlling the Convergence Rate of PGD

Next, we study the convergence rate characterized by the  $\rho(\mathcal{C})$  term. The challenges we address are (i) characterizing the *restricted singular values* of the subexponential data matrices and (ii) addressing the concatenated all ones vector.

*Theorem 3.6 (Convergence rate):* Suppose  $(\mathbf{x}_i, y_i)_{i=1}^n$  is a  $\sigma$ -subgaussian dataset and  $[\mathbf{X} \mathbf{1}]$  is the modified-data matrix, where  $\mathbf{1}$  is a vector of all ones. Let  $\mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be the tangent balls as defined in (7) and (8) respectively. Assume  $n \gtrsim (\omega(\mathcal{C}) + t)^2$ . Setting  $\eta = 1/n$ , with probability at least  $1 - 4e^{-ct^2}$  we have

$$\rho(\mathcal{C}) \lesssim \frac{\omega(\mathcal{C}) + t}{\sqrt{n}}.$$

If the dataset is  $\sigma$ -subexponential, then setting  $\eta = c_0/(n+p) \log^3(n+p)$  and assuming  $n \gtrsim (\omega_n(\mathcal{C}) + t)^2$ , with probability  $1 - 5 \exp(-c \min(n, t\sqrt{n}, t^2)) - 3(n+p)^{-100}$ , we have

$$\rho(\mathcal{C}) \leq 1 - C_0 \eta n.$$

Note that, subexponential requires a smaller choice of learning rate which results in slower convergence. Combining our bounds on  $\rho(\mathcal{C})$  and  $\nu(\mathcal{C})$  and utilizing the recursion (9), we can obtain the PGD convergence characteristics and prove the main theorems.

## IV. NUMERICAL EXPERIMENTS

In this section, we discuss experiments that corroborate our theoretical results. We consider a standard single-index model where for some ground truth vector  $\beta$  and link function  $\phi$ , the input/output relation is given by  $y_i = \phi(\beta^T \mathbf{x}_i)$ . We pick  $\beta$  to be a sparse vector with  $s = 20$  nonzeros and  $p = 800$  and set sample size to be  $n = 500$ . Because of sparsity prior, we run PGD as iterative hard thresholding where  $\theta_\tau$  is projected to

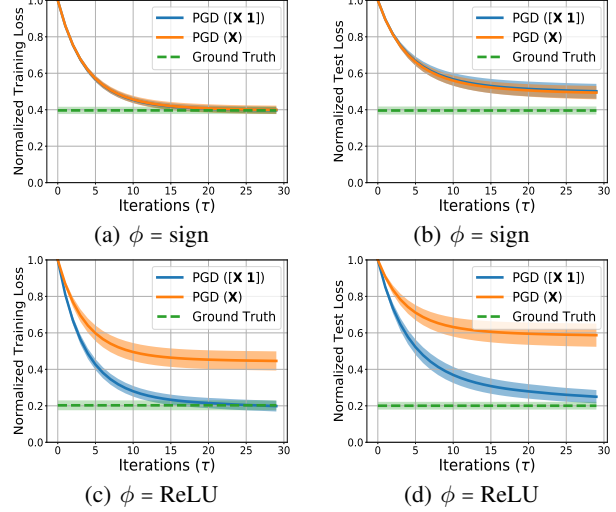


Fig. 1: We run PGD with different activations (ReLU and sign) using  $\mathbf{X}$  and  $[\mathbf{X} \mathbf{1}]$  as data matrices. In both cases, train and test errors decay gracefully to the ground truth baseline (with debiased output). However, PGD using  $[\mathbf{X} \mathbf{1}]$  outperforms PGD using  $\mathbf{X}$  alone for ReLU.

be  $s$ -sparse after every iteration. As link functions, we considered ReLU (i.e.  $\max(x, 0)$ ) and sign functions (maps to  $\pm 1$ ); which are of interest for deep learning and quantization respectively. We generate  $\mathbf{x}_i$ 's with i.i.d. exponentially distributed entries (with parameter  $\lambda = 1$ ) and then remove the mean and normalize the covariance to identity. We pick a learning rate of  $\eta = 1/5n$  in all experiments. The shaded areas in the plots correspond to one standard deviation.

To assess test and training performance of PGD, we use the following three metrics

- the **normalized training error** defined as  $\|\mathbf{y} - \mathbf{X} \theta_\tau - \mu_\tau \mathbf{1}\|_{\ell_2}^2 / \|\mathbf{y}\|_{\ell_2}^2$  and
- the **normalized test error** that is similarly defined but evaluated on a fresh dataset of size  $n$  using the training model  $\theta_\tau$
- **correlation to ground truth** vector  $\beta$  defined as  $\frac{\theta_\tau^T \beta}{\|\theta_\tau\|_{\ell_2} \|\beta\|_{\ell_2}}$ .

We compare two baselines. First one is running PGD with only  $\mathbf{X}$  i.e. without mean estimation. Second one assumes knowledge of ground truth  $\beta$  and fits a model  $\gamma \beta$  by finding  $\gamma$  to minimize the training loss. Numerically, we minimize  $\|\bar{\mathbf{y}} - \gamma \mathbf{X} \beta\|_{\ell_2}^2$  over  $\gamma$  where  $\bar{y}_i = y_i - (1/n) \sum_{i=1}^n y_i$ . This sets  $\gamma = \bar{\mathbf{y}}^T \mathbf{X} \beta / \|\mathbf{X} \beta\|_{\ell_2}^2$ .

Figure 1 plots the loss as a function of the PGD iterations  $\tau$ . Both training and test errors gracefully decays with more iterations for both choices of link functions. The dashed values corresponds to  $\gamma \beta$ 's performance.

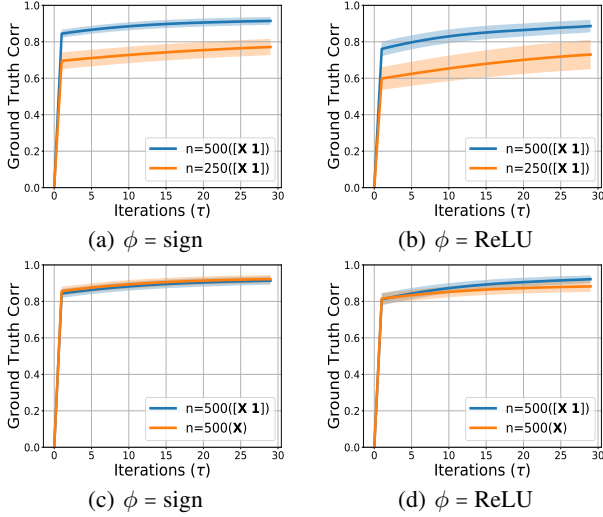


Fig. 2: We estimate correlation of PGD estimate with the ground truth vector  $\beta$ . Correlation increases with more samples. Correlation also improves when we use mean estimation  $[\mathbf{X} \ \mathbf{1}]$  instead of  $\mathbf{X}$ .

While there is a slight mismatch between train/test performances (due to finite samples), high-dimensional estimation via PGD works well and performs on par with ground truth. Observe that for ReLU,  $\mathbb{E}[y]$  is nonzero and estimating mean should be beneficial. Indeed, Figures 1c and 1d demonstrates that  $[\mathbf{X} \ \mathbf{1}]$  substantially outperforms using  $\mathbf{X}$  alone. There is no improvement for sign function since  $\mathbb{E}[y] \approx 0$  (as the sign is symmetric).

In Figure 2 we focus on the parameter estimation question by plotting the correlation between  $\theta_\tau$  and  $\beta$ . Correlation is always between  $-1, 1$  and quantifies how well we can estimate direction of the ground truth vector via PGD. This experiment is conducted with two values of  $n$  namely 250 and 500 while  $p = 800$  in both cases. Observe that, a larger sample size results in more stable estimation (smaller standard deviations) and higher correlation with output. Additionally Figure 2b shows that ReLU problem achieves better correlation once we account for the bias term. Hence, mean estimation is not only beneficial for test performance but also for parameter estimation.

## V. PROOFS OF MAIN THEOREMS

This section proves our main results and outlines the proof of Theorems 3.3, 3.4 3.5 and 3.6. Throughout we use the same notation as described in II.

### A. Proof of Theorem 3.4

We provide our analysis for subexponential samples. The extension to subgaussian samples is accomplished in

an identical fashion. Set the estimation error at iteration  $\tau$  to be  $\mathbf{h}_\tau = \theta_\tau - \theta^*$ . Note that, when  $\rho(\mathcal{C}) < 1$  and  $\mathcal{R}$  is a convex regularizer, then the recursion (9) can be iteratively expanded as

$$\begin{aligned} \|\mathbf{h}_\tau\|_{\ell_2} &\leq \|\mathbf{h}_0\|_{\ell_2} \rho(\mathcal{C})^\tau + \eta \nu(\mathcal{C}) \sum_{k=0}^{\tau-1} \rho(\mathcal{C})^k \\ &\leq \|\mathbf{h}_0\|_{\ell_2} \rho(\mathcal{C})^\tau + \eta \nu(\mathcal{C}) \sum_{k=0}^{\infty} \rho(\mathcal{C})^k \\ &= \|\mathbf{h}_0\|_{\ell_2} \rho(\mathcal{C})^\tau + \frac{\eta \nu(\mathcal{C})}{1 - \rho(\mathcal{C})} \end{aligned} \quad (15)$$

With the advertised probability, subexponential statements of Theorems 3.5 and 3.6 hold. Hence, for some constants, we have that  $\rho(\mathcal{C}) \leq 1 - c_0 \eta n$ ,  $\nu(\mathcal{C}) \leq C \sqrt{n} \sigma(\omega_n(\mathcal{C}) + t) \log(n)$  and  $\eta = c/q$  with  $q = (n + p) \log^3(n + p)$ . Plugging these in (15), we find the following upper bound on the right hand side,

$$\begin{aligned} \|\mathbf{h}_{\tau+1}\|_{\ell_2} &\leq (1 - c_0 \eta n)^\tau \|\mathbf{h}_0\|_{\ell_2} \\ &\quad + \frac{\eta}{c_0 \eta n} C \sqrt{n} \sigma(\omega_n(\mathcal{C}) + t) \log(n) \\ &= \left(1 - \frac{c_0 c n}{q}\right)^\tau \|\mathbf{h}_0\|_{\ell_2} \\ &\quad + \sigma \eta \frac{C (\omega_n(\mathcal{C}) + t) \log(n)}{c_0 \sqrt{n}}, \end{aligned}$$

which is the desired bound. The case of subgaussian samples is again a corollary of Theorems 3.5 and 3.6. This concludes the proof of our main result.

### B. Proof of Theorem 3.6 for subgaussian samples

We start our proof with the following lemma.

*Lemma 5.1:* Let  $(x_i)_{i=1}^n \mathbb{R}^p$  be i.i.d. isotropic subgaussian samples. Let  $\mathbf{X} \in \mathbb{R}^{n \times p}$  be concatenated data and  $[\mathbf{X} \ \mathbf{1}]$  is modified-data matrix, where  $\mathbf{1}$  is a vector of all ones. Let  $\mathcal{T}$  be a closed set with Euclidian radius bounded by a constant and

$$\mathcal{T}_{\text{ext}} = \{\tilde{\mathbf{v}} \mid \tilde{\mathbf{v}} = [\beta \mathbf{v}^T \ \gamma]^T\},$$

where  $\beta \leq C_1$ ,  $\gamma \leq C_2$  for some positive constants  $C_1, C_2$ . Assume  $n \gtrsim (\omega(\mathcal{T}) + t)^2$ . Then, with probability at least  $1 - 2e^{-t^2}$  we have

$$\sup_{\tilde{\mathbf{v}} \in \mathcal{T}_{\text{ext}}} |\tilde{\mathbf{v}}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}}| \lesssim \frac{\omega(\mathcal{T}) + t}{\sqrt{n}}.$$

The proof of Lemma 5.1 is deferred to Section VII-A. Next using the result of Lemma 5.1, we obtain the following lemma which bounds the convergence rate for subgaussian samples.

*Lemma 5.2:* Consider the setup of Lemma 5.1. Furthermore, let the tangent balls  $\mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in (7)

and (8) respectively. With probability at least  $1 - 4e^{-t^2}$ , the following holds

$$\sup_{\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} |\tilde{\mathbf{u}}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}}| \lesssim \frac{\omega(\mathcal{C}) + t}{\sqrt{n}}.$$

The proof of Lemma 5.2 is deferred to Section VII-B. This completes the proof for subgaussian samples.

### C. Proof of Theorem 3.6 for subexponential samples

Let  $(\mathbf{x}_i)_{i=1}^n \sim \mathbf{x}$  be i.i.d. isotropic subexponential vectors and  $\mathbf{X}$  be the associated design matrix as previously. Let  $\mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in 7 and 8 respectively. Assume  $n \gtrsim \omega_n^2(\mathcal{C})$ . Our proof strategy is based on the observation that, we can bound the (restricted) singular values of  $[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]$  with high probability for subexponential data as follows.

1) *Upper bounding the singular values:* In this section we will bound the largest eigenvalue of the matrix  $[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]$  with high probability. Towards this goal, we utilize Matrix Chernoff bound from [33].

*Theorem 5.3 (Matrix Chernoff [33]):* Consider a finite sequence  $\{\mathbf{X}_i\}_{i=1}^n$  of independent, random, Hermitian matrices with common dimension  $d$ . Assume that

$$0 \leq \sigma_{\min}(\mathbf{X}_i) \text{ and } \|\mathbf{X}_i\| \leq L \text{ for } i = 1, \dots, n$$

Define the sum  $\mathbf{M} = \sum_{i=1}^n \mathbf{X}_i$  and  $\zeta_{\max}$  be an upper bound on the spectral norm of the expectation  $\mathbb{E}[\mathbf{M}]$  i.e.  $\zeta_{\max} \geq \|\mathbb{E}[\mathbf{M}]\| = \|\sum_{i=1}^n \mathbb{E}[\mathbf{X}_i]\|$ . We have that

$$\mathbb{P} \{ \|\mathbf{M}\| \geq (1 + \epsilon) \zeta_{\max} \} \leq d \left[ \frac{e^\epsilon}{(1 + \epsilon)^{1 + \epsilon}} \right]^{\frac{\zeta_{\max}}{L}}, \epsilon \geq 0$$

We will use Theorem 5.3 to bound the largest eigenvalue of  $[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]$ . Observe that

$$[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}] = \sum_{i=1}^n \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} [\mathbf{x}_i^T \ 1].$$

Clearly this matrix is positive semidefinite. To bound  $\|[\mathbf{x}_i^T \ 1]^T [\mathbf{x}_i^T \ 1]\|$  we use the following lemma.

*Lemma 5.4 (Spectral norm bound):* Let  $(\mathbf{x}_i)_{i=1}^n$  be i.i.d. isotropic subexponential samples in  $\mathbb{R}^p$ . Then, with probability at least  $1 - 2(n+p)^{-100}$  the spectral norm of all  $\mathbf{x}_i \mathbf{x}_i^T$  matrices can be bounded as

$$\|\mathbf{x}_i \mathbf{x}_i^T\| \leq \|\mathbf{x}_i\|_{\ell_2}^2 \leq cp \log^2(n+p)$$

The proof of lemma 5.4 is deferred to Section VII-C. Lemma 5.4 guarantees that  $\|[\mathbf{x}_i^T \ 1]^T [\mathbf{x}_i^T \ 1]\| \leq \|[\mathbf{x}_i^T \ 1]^T\|_{\ell_2}^2 = \|\mathbf{x}_i\|_{\ell_2}^2 + 1 \leq Cp \log^2(n+p)$ . Hence, we do satisfy the conditions required by Theorem 5.3. Before using Theorem 5.3 we will upper bound the spectral norm of the expectation  $\mathbb{E}[[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]]$  as follows.

*Lemma 5.5 (Spectral norm bound of expectation):* Let  $\mathbf{x}$  be an isotropic subexponential vector,  $\tilde{\mathbf{x}} = [\mathbf{x}^T \ 1]^T$  and let  $B = Cp \log^2(n+p)$  for sufficiently large constant  $C > 0$ . Then we have

$$\mathbb{E} \left[ \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \mid \|\tilde{\mathbf{x}}\|_{\ell_2} \leq \sqrt{B} \right] \leq 2\mathbf{I}_p.$$

The proof of Lemma 5.5 is deferred to Section VII-D. Thus, applying Lemma 5.5 on the set of all  $[\mathbf{x}_i^T \ 1]^T$  satisfying  $\|[\mathbf{x}_i^T \ 1]^T [\mathbf{x}_i^T \ 1]\| \leq Cp \log^2(n+p)$ , we find that with  $1 - 2(n+p)^{-100}$  probability the following holds

$$\begin{aligned} \|\mathbb{E}[[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]]\| &= \|\mathbb{E}[\sum_{i=1}^n \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} [\mathbf{x}_i^T \ 1]]\| \\ &\leq \|\sum_{i=1}^n 2\mathbf{I}_p\| \\ &= 2n. \end{aligned}$$

Hence we can pick  $\zeta_{\max} \geq 2n$  to upper bound the largest eigenvalue of  $\mathbb{E}[[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]]$ . Now, using Theorem 5.3 with  $\zeta_{\max} = C_0 C(n+p) \log^3(n+p)$ ,  $L = Cp \log^2(n+p)$  and  $\epsilon = e - 1$  we get

$$\begin{aligned} \mathbb{P} \{ \|\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]\| \geq eC_0 C(n+p) \log^3(n+p) \} \\ \leq p \left[ \frac{e^{e-1}}{e^e} \right]^{C_0 \frac{n+p}{p} \log(n+p)} \\ = pe^{-C_0 \frac{n+p}{p} \log(n+p)} \leq (n+p)^{-100}. \end{aligned} \quad (16)$$

Union bounding, with probability at least  $1 - 3(n+p)^{-100}$ ,

$$\|[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]\| \lesssim (n+p) \log^3(n+p). \quad (17)$$

2) *Lower bounding the singular values:* In this section we will lower bound the gain of  $[\mathbf{X} \ \mathbf{1}]$  restricted to the tangent ball  $\mathcal{C}_{\text{ext}}$ . We will utilize the notion of restricted singular value (RSV) to proceed.

*Definition 5.6 (Restricted singular value):* Given a matrix  $\mathbf{M}$  and a closed set  $\mathcal{C}$ , the RSV of  $\mathbf{M}$  at  $\mathcal{C}$  is defined as

$$\sigma(\mathbf{M}, \mathcal{C}) = \min_{\mathbf{v} \in \mathcal{C}} \frac{\|\mathbf{M}\mathbf{v}\|_{\ell_2}}{\|\mathbf{v}\|_{\ell_2}}$$

In the following, we will lower bound  $\min_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \|[\mathbf{X} \ \mathbf{1}]\tilde{\mathbf{v}}\|_{\ell_2}$  which is the RSV of  $[\mathbf{X} \ \mathbf{1}]$  at  $\mathcal{C}_{\text{ext}}$ . Recall that any  $\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}$  with unit Euclidian norm obeys  $\tilde{\mathbf{v}} = [\sqrt{1 - \gamma^2} \mathbf{v}^T \ \gamma]^T \in \mathcal{C}_{\text{ext}}$  for  $0 \leq \gamma \leq 1$  and  $\|\mathbf{v}\|_{\ell_2} = 1$ . Consequently

$$\begin{aligned} \|[\mathbf{X} \ \mathbf{1}]\tilde{\mathbf{v}}\|_{\ell_2}^2 &= \|\sqrt{1 - \gamma^2} \mathbf{X}\mathbf{v} + \gamma \mathbf{1}\|_{\ell_2}^2 \\ &= (1 - \gamma^2) \|\mathbf{X}\mathbf{v}\|_{\ell_2}^2 + \gamma^2 \mathbf{1}^T \mathbf{1} + 2\gamma \sqrt{1 - \gamma^2} \mathbf{1}^T \mathbf{X}\mathbf{v} \\ &\geq (1 - \gamma^2) \|\mathbf{X}\mathbf{v}\|_{\ell_2}^2 + \gamma^2 n + 2\gamma \sqrt{1 - \gamma^2} \mathbf{v}^T \sum_{i=1}^n \mathbf{x}_i. \end{aligned}$$

Setting  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$  and taking infimum of both sides, we get

$$\min_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \|[\mathbf{X} \ \mathbf{1}] \tilde{\mathbf{v}}\|_{\ell_2}^2 \geq \quad (18)$$

$$\begin{aligned} &\geq \min_{|\gamma| \leq 1} \left( (1 - \gamma^2) \min_{\mathbf{v} \in \mathcal{C}} \|\mathbf{X} \mathbf{v}\|_{\ell_2}^2 + \gamma^2 n \right) - 2n \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \bar{\mathbf{x}}| \\ &\geq \min_{\mathbf{v} \in \mathcal{C}} \left( \min_{\mathbf{v} \in \mathcal{C}} \|\mathbf{X} \mathbf{v}\|_{\ell_2}^2, n \right) - 2n \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \bar{\mathbf{x}}|. \end{aligned} \quad (19)$$

In essence, (19) bounds RSV of  $[\mathbf{X} \ \mathbf{1}]$  in terms of the RSV of  $\mathbf{X}$  and some simpler terms. The following theorem from [29] (Theorem D.11) gives a lower bound on the RSV of a matrix  $\mathbf{X}$  with i.i.d. subexponential rows.

*Theorem 5.7 (Bounding RSV [29]):* Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be a random matrix with i.i.d. isotropic subexponential rows. Let  $\mathcal{C}$  be a tangent ball as in (7) and suppose the sample size obeys  $n \gtrsim (\omega_n(\mathcal{C}) + t)$ . Then with probability at least  $1 - 3 \exp(-c \min(n, t^2, t\sqrt{n}))$ , we have that

$$\min_{\mathbf{u} \in \mathcal{C}} \|\mathbf{X} \mathbf{u}\|_{\ell_2}^2 \geq c_0 n$$

Next, we shall state a lemma from [29] (Lemma D.7) to upper bound the term involving the sample average  $\bar{\mathbf{x}}$ .

*Lemma 5.8 (Bounding empirical width [29]):* Suppose  $\mathcal{C}$  is a subset of the unit Euclidian ball and  $(\mathbf{x}_i)_{i=1}^n$  are i.i.d. zero-mean vectors with bounded subexponential norm. Define the empirical average vector  $\bar{\mathbf{x}} = n^{-1} \sum_i \mathbf{x}_i$ . We have that

$$\begin{aligned} &\mathbb{P} \left( \sup_{\mathbf{u} \in \mathcal{C}} |\mathbf{u}^T \bar{\mathbf{x}}| \leq C \frac{(\omega_n(\mathcal{C}) + t)}{\sqrt{n}} \right) \\ &\geq 1 - 2 \exp(-c \min(t\sqrt{n}, t^2)) \end{aligned}$$

Combining Theorem 5.7 and Lemma 5.8 into (19) we find that, there exist constants  $c, c_0, C_0 > 0$  such that with probability at least  $1 - 5 \exp(-c \min(n, t\sqrt{n}, t^2))$ , we can lower bound the RSV of  $[\mathbf{X} \ \mathbf{1}]$  following (19)

$$\begin{aligned} \min_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \|[\mathbf{X} \ \mathbf{1}] \tilde{\mathbf{v}}\|_{\ell_2}^2 &\geq c_0 n - C_0 n \frac{\omega_n(\mathcal{C}) + t}{\sqrt{n}} \\ &\geq c_0 n / 2. \end{aligned} \quad (20)$$

where last line follows from the assumption that  $n \gtrsim (\omega_n(\mathcal{C}) + t)^2$ .

3) *Upper bounding the convergence rate:* Union bounding the events (17) and (20), we obtain upper and lower bounds on the singular values of  $[\mathbf{X} \ \mathbf{1}]$  with the desired probability. Hence, we can bound the convergence rate of PGD as follows. Setting  $q = (n + p) \log^3(n + p)$ , we have (17)  $\|[\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]\| \leq Cq$ . Therefore, choosing learning rate  $\eta = 1/Cq$ , the matrix  $\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]$  is positive semidefinite (PSD). Hence,

applying the generalized Cauchy-Schwarz inequality for PSD matrix, we find

$$\begin{aligned} \rho(\mathcal{C}) &= \sup_{\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \tilde{\mathbf{u}}^T (\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}} \\ &\leq \sup_{\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} [(\tilde{\mathbf{u}}^T (\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{u}})^{1/2} \\ &\quad (\tilde{\mathbf{v}}^T (\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}})^{1/2}] \\ &= \sup_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \tilde{\mathbf{v}}^T (\mathbf{I} - \eta [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}} \\ &= 1 - \eta \min_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} \|[\mathbf{X} \ \mathbf{1}] \tilde{\mathbf{v}}\|_{\ell_2}^2 \\ &\leq 1 - c_0 \eta n / 2. \end{aligned}$$

Here the last inequality follows from (20). This completes the proof for subexponential samples.

#### D. Proof of Theorem 3.5 for subgaussian samples

Suppose the dataset  $(\mathbf{x}_i, y_i)_{i=1}^n \sim (\mathbf{x}, y)$  is  $\sigma$ -subgaussian. Let  $\mathbf{X}, [\mathbf{X} \ \mathbf{1}], \mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in Section II, recall  $\mathbf{w}$  from (13) and assume  $n \gtrsim (\omega(\mathcal{C}) + t)^2$ . Representing  $\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}$  as  $\tilde{\mathbf{v}} = [\sqrt{1 - \gamma^2} \mathbf{v}^T \ \gamma]^T$  for  $\mathbf{v} \in \mathcal{C}$  and  $|\gamma| \leq 1$ , we have

$$\begin{aligned} \nu(\mathcal{C}) &= \sup_{\tilde{\mathbf{v}} \in \mathcal{C}_{\text{ext}}} |\tilde{\mathbf{v}}^T [\mathbf{X} \ \mathbf{1}]^T \mathbf{w}| \\ &= \sup_{\mathbf{v} \in \mathcal{C}, |\gamma| \leq 1} |\sqrt{1 - \gamma^2} \mathbf{v}^T \mathbf{X}^T \mathbf{w} + \gamma \mathbf{1}^T \mathbf{w}| \\ &\leq \sup_{\mathbf{v} \in \mathcal{C}, |\gamma| \leq 1} |\sqrt{1 - \gamma^2} \mathbf{v}^T \mathbf{X}^T \mathbf{w}| + \sup_{|\gamma| \leq 1} |\gamma \mathbf{1}^T \mathbf{w}| \\ &\leq \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}| + |\mathbf{1}^T \mathbf{w}| \end{aligned} \quad (21)$$

In the following we will upper bound the terms  $\sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}|$  and  $|\mathbf{1}^T \mathbf{w}|$  separately and will combine them to get an upper bound on the residual error.

1) *Upper bounding the first term in (21):* In order to upper bound the first term in (21), define the clipping function

$$\text{clip}(a, B) = \begin{cases} a & \text{if } |a| \leq B \\ \text{sign}(a)B & \text{else} \end{cases}$$

Following lemma immediately follows from union bounding the large deviations of subgaussian and subexponential variables and shows that  $X = \text{clip}(X, B)$  with high probability.

*Lemma 5.9:* Let  $(w_i)_{i=1}^n$  be i.i.d. subgaussian random variables with  $\|w_i\|_{\psi_2} \leq \sigma$ . There exists a constant  $C > 0$  such that picking  $B = C\sqrt{\log(n)}$ , with probability  $1 - 2n^{-100}$  for all  $i$ , we have

$$w_i = \text{clip}(w_i, \sigma B).$$

If instead  $(w_i)_{i=1}^n$  are i.i.d. subexponential with  $\|w_i\|_{\psi_1} \leq \sigma$ , then picking  $B = C \log(n)$  leads to the same result.



Using Lemma 5.9,  $\|\mathbf{w}\|_\infty \leq \sigma B$  with probability  $1 - 2n^{-100}$ . Conditioned on this event, we have

$$\sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}| = \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n \text{clip}(w_i, \sigma B) \mathbf{x}_i|. \quad (22)$$

Setting  $\mathbf{z}_i = \text{clip}(w_i, \sigma B) \mathbf{x}_i = w_i \mathbf{x}_i$ , (22) can be rewritten as

$$\begin{aligned} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}| &= \frac{1}{n} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n \mathbf{z}_i| \\ &\leq \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n (\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i])| + \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n \mathbb{E}[\mathbf{z}_i]| \\ &\leq \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n (\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i])| + n \|\mathbb{E}[\mathbf{z}_1]\|_{\ell_2}. \end{aligned} \quad (23)$$

Note that  $\mathbf{z}_i = w_i \mathbf{x}_i$  is subgaussian since  $w_i$  is bounded. The subgaussian norm obeys

$$\|\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i]\|_{\psi_2} \lesssim \|\mathbf{z}_i\|_{\psi_2} \lesssim \sigma \sqrt{\log(n)}.$$

Define the average vector  $\bar{\mathbf{z}} = n^{-1/2} \sum_{i=1}^n (\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i])$  which is still subgaussian with same norm (up to a constant). Standard results from functional analysis [28] guarantee

$$\begin{aligned} \frac{1}{n} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T (\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i])| &= n^{-1/2} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \bar{\mathbf{z}}| \\ &\lesssim \frac{\sigma(\omega(\mathcal{C}) + t) \sqrt{\log(n)}}{\sqrt{n}} \end{aligned} \quad (24)$$

with probability at least  $1 - 2e^{-t^2/2}$ . This bounds the first term of (23). Next, we address the expectation term  $\|\mathbb{E}[\mathbf{z}_1]\|_{\ell_2}$  via following lemma.

*Lemma 5.10:* Suppose  $\mathbf{x}$  is an isotropic Orlicz- $\alpha$  vector and  $\|\mathbf{w}\|_{\psi_\alpha} \leq \sigma$ . Let  $B = C \log^{1/\alpha}(n)$  for sufficiently large constant  $C > 0$ . For  $\alpha = 1, 2$ , we have that

$$\|\mathbb{E}[\mathbf{w}\mathbf{x} \mid |w| \leq \sigma B]\|_{\ell_2} \lesssim \sigma p^2 n^{-201}.$$

The proof of Lemma 5.10 is deferred to Section VII-F. Combining (24) and Lemma 5.10 into (23), with probability at least  $1 - 2e^{-t^2/2} - 2n^{-100}$ , we find that,

$$\begin{aligned} \frac{1}{n} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}| &\lesssim \frac{\sigma(\omega(\mathcal{C}) + t) \sqrt{\log(n)}}{\sqrt{n}} \\ &\quad + \sigma p^2 n^{-201} \\ &\lesssim \frac{\sigma(\omega(\mathcal{C}) + t) \sqrt{\log(n)}}{\sqrt{n}} \end{aligned} \quad (25)$$

which is the desired bound for the first term in (21).

2) *Upper bounding the second term in (21):* The vector  $\mathbf{w}$  is zero-mean with  $\|\mathbf{w}\|_{\psi_2} \leq \sigma$ . Hence,  $\|\mathbf{1}^T \mathbf{w}\|_{\psi_2} \leq \sigma \sqrt{n}$  which implies that with probability  $1 - 2n^{-100}$ ,

$$|\mathbf{1}^T \mathbf{w}| \lesssim \sigma \sqrt{n \log n}.$$

Combining the bound above with (25), we get the advertised bound on the residual, namely

$$\frac{1}{n} \nu(\mathcal{C}) \lesssim \frac{\sigma(\omega(\mathcal{C}) + t) \sqrt{\log(n)}}{\sqrt{n}}, \quad (26)$$

with probability at least  $1 - 2 \exp(-t^2/2) - 4n^{-100}$ .

*E. Proof of Theorem 3.5 for subexponential samples*

Suppose the dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  is  $\sigma$ -subexponential. Let  $\mathbf{X}, [\mathbf{X} \ \mathbf{1}], \mathbf{w}, \mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in Section II. Assume  $n \gtrsim (\omega_n(\mathcal{C}) + t)^2$ . Similar to the subgaussian case, we split the residual into two terms via (21) and bound each term separately to get a final bound.

1) *Upper bounding the first term in (21):* Let  $\mathbf{z}_i = w_i \mathbf{x}_i$ . With probability  $1 - 2n^{-100}$ , we have that  $\|\mathbf{w}\|_\infty \lesssim \sigma \log n$ . We continue the analysis conditioned on this event. With bounded  $w_i$ ,  $\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i]$  is subexponential via

$$\|\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i]\|_{\psi_1} \lesssim \|\mathbf{z}_i\|_{\psi_1} \lesssim \sigma \log n \|\mathbf{x}_i\|_{\psi_1} \lesssim \sigma \log n.$$

Lemma 5.8 guarantees that

$$\frac{1}{n} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \sum_{i=1}^n (\mathbf{z}_i - \mathbb{E}[\mathbf{z}_i])| \lesssim \frac{\sigma(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}} \quad (27)$$

with probability at least  $1 - 2 \exp(-\mathcal{O}(\min\{t\sqrt{n}, t^2\}))$ . Using Theorem 5.10, we also upper bound  $\|\mathbb{E}[\mathbf{z}_1]\|_{\ell_2}$  by  $C \sigma p^2 n^{-201}$ . Combining this with (27) and substituting into the (deterministic inequality) (23), with probability at least  $1 - 2 \exp(-\mathcal{O}(\min\{t\sqrt{n}, t^2\})) - 2n^{-100}$  we have,

$$\frac{1}{n} \sup_{\mathbf{v} \in \mathcal{C}} |\mathbf{v}^T \mathbf{X}^T \mathbf{w}| \lesssim \frac{\sigma(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}}. \quad (28)$$

2) *Upper bounding overall error of (21):* Using  $\|\mathbf{w}\|_{\psi_1} \lesssim \sigma$  and applying Lemma 5.8 (over one-dimensional  $\mathbb{R}$ ), we find that  $|\mathbf{1}^T \mathbf{w}| \lesssim \sigma(1+t)\sqrt{n}$  with probability  $1 - 2 \exp(-c \min\{t\sqrt{n}, t^2\})$ . Combining this with (28) and plugging into (21), we get the advertised upper bound

$$\begin{aligned} \frac{1}{n} \nu(\mathcal{C}) &\lesssim \frac{\sigma(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}} + \frac{(1+t)\sigma}{\sqrt{n}} \\ &\lesssim \frac{\sigma(\omega_n(\mathcal{C}) + t) \log(n)}{\sqrt{n}} \end{aligned} \quad (29)$$

which holds with probability at least  $1 - 4 \exp(-c \min\{t\sqrt{n}, t^2\}) - 2n^{-100}$ . This completes the proof for  $\sigma$ -subexponential data.

## VI. CONCLUSION

We studied the problem of finding the best linear model from  $n$  input-output samples under quadratic loss in the high-dimensional regime  $n \ll p$ . For estimation, we utilized the projected gradient descent algorithm and showed its fast convergence as well as statistical accuracy in a data-dependent fashion. Our results are established for subexponential design which is heavier tailed compared to well-studied subgaussian. In both cases, we prove that *nonlinearity of the problem behaves like independent noise* and we establish favorable statistical guarantees as if the problem is linear. We also modified the original regression problem to allow for mean estimation and demonstrated its practical benefit when output labels have nonzero mean via simulations.

It would be desirable to extend our results to general loss function. If a loss function  $\ell$  has the potential to better capture input/output relation, we can solve for

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(y_i, \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle).$$

Specifically this function can still be quadratic but characterized by a nonlinear link function  $\phi$  i.e.  $\ell(y_i, \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) = (y_i - \phi(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle))^2$ . We believe that much of the results presented here extends to strongly-increasing  $\phi$  where the derivative is lower bounded by a constant i.e.  $\phi' \geq \alpha$  for some  $\alpha > 0$ . These functions are shown to behave like linear regression [32]. However, it is not immediately clear if strong statistical and computational guarantees established in this paper (as well as related literature) can be established.

## REFERENCES

- [1] Y. Plan, R. Vershynin, and E. Yudovina, "High-dimensional estimation with geometric constraints," *Information and Inference: A Journal of the IMA*, vol. 6, no. 1, pp. 1–40, 2016.
- [2] P. T. Boufounos and R. G. Baraniuk, "1-bit compressive sensing," in *Information Sciences and Systems, 2008. CISS 2008. 42nd Annual Conference on*. IEEE, 2008, pp. 16–21.
- [3] R. Ganti, N. Rao, R. M. Willett, and R. Nowak, "Learning single index models in high dimensions," *arXiv preprint arXiv:1506.08910*, 2015.
- [4] S. Oymak and M. Soltanolkotabi, "Fast and reliable parameter estimation from nonlinear observations," *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2276–2300, 2017.
- [5] Y. Plan, R. Vershynin, and E. Yudovina, "High-dimensional estimation with geometric constraints," *Information and Inference: A Journal of the IMA*, vol. 6, no. 1, pp. 1–40, 2017.
- [6] Y. Plan and R. Vershynin, "The generalized lasso with non-linear observations," *IEEE Transactions on information theory*, vol. 62, no. 3, pp. 1528–1537, 2016.
- [7] C. Thrampoulidis, E. Abbasi, and B. Hassibi, "Lasso with non-linear measurements is equivalent to one with linear measurements," in *Advances in Neural Information Processing Systems*, 2015, pp. 3420–3428.
- [8] L. Jacques, J. N. Laska, P. T. Boufounos, and R. G. Baraniuk, "Robust 1-bit compressive sensing via binary stable embeddings of sparse vectors," *IEEE Transactions on Information Theory*, vol. 59, no. 4, pp. 2082–2102, 2013.
- [9] R. Vershynin, "Estimation in high dimensions: a geometric perspective," in *Sampling theory, a renaissance*. Springer, 2015, pp. 3–66.
- [10] S. Dirksen, H. C. Jung, and H. Rauhut, "One-bit compressed sensing with partial gaussian circulant matrices," *arXiv preprint arXiv:1710.03287*, 2017.
- [11] S. Dirksen and S. Mendelson, "Robust one-bit compressed sensing with partial circulant matrices," *arXiv preprint arXiv:1812.06719*, 2018.
- [12] Y. Plan and R. Vershynin, "Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach," *Information Theory, IEEE Transactions on*, vol. 59, no. 1, pp. 482–494, 2013.
- [13] A. Agarwal, S. Negahban, and M. J. Wainwright, "Fast global convergence rates of gradient methods for high-dimensional statistical recovery," in *Advances in Neural Information Processing Systems*, 2010, pp. 37–45.
- [14] S. Oymak, B. Recht, and M. Soltanolkotabi, "Sharp time–data tradeoffs for linear inverse problems," *IEEE Transactions on Information Theory*, vol. 64, no. 6, pp. 4129–4158, 2018.
- [15] R. Giryes, Y. C. Eldar, A. M. Bronstein, and G. Sapiro, "Tradeoffs between convergence speed and reconstruction accuracy in inverse problems," *IEEE Transactions on Signal Processing*, vol. 66, no. 7, pp. 1676–1690, 2018.
- [16] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [17] M. Genzel, "High-dimensional estimation of structured signals from non-linear observations with general convex loss functions," *IEEE Transactions on Information Theory*, vol. 63, no. 3, pp. 1601–1619, 2017.
- [18] S. Dirksen and S. Mendelson, "Robust one-bit compressed sensing with non-gaussian measurements," *arXiv preprint arXiv:1805.09409*, 2018.
- [19] C. Thrampoulidis and A. S. Rawat, "The generalized lasso for sub-gaussian measurements with dithered quantization," *arXiv preprint arXiv:1807.06976*, 2018.
- [20] L. Jacques and V. Cambareri, "Time for dithering: fast and quantized random embeddings via the restricted isometry property," *Information and Inference: A Journal of the IMA*, vol. 6, no. 4, pp. 441–476, 2017.
- [21] C. Xu and L. Jacques, "Quantized compressive sensing with rip matrices: The benefit of dithering," *arXiv preprint arXiv:1801.05870*, 2018.
- [22] Z. Yang, Z. Wang, H. Liu, Y. Eldar, and T. Zhang, "Sparse nonlinear regression: Parameter estimation under nonconvexity," in *International Conference on Machine Learning*, 2016, pp. 2472–2481.
- [23] Z. Yang, K. Balasubramanian, and H. Liu, "High-dimensional non-gaussian single index models via thresholded score function estimation," in *International Conference on Machine Learning*, 2017, pp. 3851–3860.
- [24] Z. Yang, K. Balasubramanian, Z. Wang, and H. Liu, "Learning non-gaussian multi-index model via second-order stein's method," *Advances in Neural Information Processing Systems*, 2017.
- [25] Z. Yang, K. Balasubramanian, and H. Liu, "On stein's identity and near-optimal estimation in high-dimensional index models," *arXiv preprint arXiv:1709.08795*, 2017.
- [26] H. L. Yap, M. B. Wakin, and C. J. Rozell, "Stable manifold embeddings with structured random matrices," *IEEE Journal on Selected Topics in Signal Processing*, vol. 7, no. 4, pp. 720–730, 2013.
- [27] M. Genzel and G. Kutyniok, "The mismatch principle: Statistical learning under large model uncertainties," *arXiv preprint arXiv:1808.06329*, 2018.
- [28] M. Talagrand, "Gaussian processes and the generic chaining," in *Upper and Lower Bounds for Stochastic Processes*. Springer, 2014, pp. 13–73.

- [29] S. Oymak, "Learning compact neural networks with regularization," *International Conference on Machine Learning*, 2018.
- [30] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex geometry of linear inverse problems," *Foundations of Computational Mathematics*, vol. 12, no. 6, pp. 805–849, 2012.
- [31] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp, "Living on the edge: Phase transitions in convex programs with random data," *Inform. Inference*, 2014.
- [32] S. Oymak, "Stochastic gradient descent learns state equations with nonlinear activations," *arXiv preprint arXiv:1809.03019*, 2018.
- [33] J. A. Tropp *et al.*, "An introduction to matrix concentration inequalities," *Foundations and Trends® in Machine Learning*, vol. 8, no. 1-2, pp. 1–230, 2015.

## VII. APPENDIX

This section provides the proofs of supporting results.

### A. Proof of Lemma 5.1

We start by expanding the convergence term as follows by substituting  $\tilde{\mathbf{v}} = [\beta \mathbf{v} \ \gamma]^T$

$$\begin{aligned}
& |\tilde{\mathbf{v}}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}}| \\
&= \left| \frac{1}{n} \| [\mathbf{X} \ \mathbf{1}] \tilde{\mathbf{v}} \|_{\ell_2}^2 - \|\tilde{\mathbf{v}}\|_{\ell_2}^2 \right| \\
&= \left| \frac{1}{n} \|\beta \mathbf{X} \mathbf{v} + \gamma \mathbf{1}\|_{\ell_2}^2 - \|[\beta \mathbf{v}^T \ \gamma]^T\|_{\ell_2}^2 \right| \\
&= \left| \frac{1}{n} (\beta^2 \|\mathbf{X} \mathbf{v}\|_{\ell_2}^2 + \gamma^2 \mathbf{1}^T \mathbf{1} + 2\beta\gamma \mathbf{1}^T \mathbf{X} \mathbf{v}) - \beta^2 \|\mathbf{v}\|_{\ell_2}^2 - \gamma^2 \right| \\
&= \left| \frac{1}{n} \beta^2 \|\mathbf{X} \mathbf{v}\|_{\ell_2}^2 - \beta^2 \|\mathbf{v}\|_{\ell_2}^2 + \frac{1}{n} \gamma^2 n - \gamma^2 + 2\beta\gamma \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{v} \right| \\
&\leq \beta^2 \left| \frac{1}{n} \|\mathbf{X} \mathbf{v}\|_{\ell_2}^2 - \|\mathbf{v}\|_{\ell_2}^2 \right| + |2\beta\gamma| \left| \mathbf{v}^T \frac{\sum_{i=1}^n \mathbf{x}_i}{n} \right| \\
&\lesssim |\mathbf{v}^T (\mathbf{I} - \frac{1}{n} \mathbf{X}^T \mathbf{X}) \mathbf{v}| + |\mathbf{v}^T \bar{\mathbf{x}}|, \tag{30}
\end{aligned}$$

where,  $\bar{\mathbf{x}} = n^{-1} \sum_{i=1}^n \mathbf{x}_i$  is the empirical average vector of i.i.d. subgaussian rows  $(\mathbf{x}_i)_{i=1}^n$ . Thus, using (30), we can write

$$\begin{aligned}
& \sup_{\tilde{\mathbf{v}} \in \mathcal{T}_{\text{ext}}} |\tilde{\mathbf{v}}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}}| \\
&\lesssim \sup_{\mathbf{v} \in \mathcal{T}} |\mathbf{v}^T (\mathbf{I} - \frac{1}{n} \mathbf{X}^T \mathbf{X}) \mathbf{v}| + \sup_{\mathbf{v} \in \mathcal{T}} |\mathbf{v}^T \bar{\mathbf{x}}|. \tag{31}
\end{aligned}$$

Given  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is isotropic subgaussian, Lemma 6.14 in [14] guarantees

$$\sup_{\mathbf{v} \in \mathcal{T}} |\mathbf{v}^T (\mathbf{I} - \frac{1}{n} \mathbf{X}^T \mathbf{X}) \mathbf{v}| \lesssim \frac{\omega(\mathcal{T}) + t}{\sqrt{n}} \tag{32}$$

with probability at least  $1 - e^{-t^2}$ . Furthermore, since  $(\mathbf{x}_i)_{i=1}^n$ 's have bounded subgaussian norm,  $\bar{\mathbf{x}}$  is also bounded and standard results from functional analysis guarantee [28]

$$\sup_{\mathbf{v} \in \mathcal{T}} |\mathbf{v}^T \frac{\sum_{i=1}^n \mathbf{x}_i}{n}| = \sup_{\mathbf{v} \in \mathcal{T}} |\mathbf{v}^T \bar{\mathbf{x}}| \lesssim \frac{\omega(\mathcal{T}) + t}{\sqrt{n}} \tag{33}$$

with probability at least  $1 - e^{-t^2}$ . Combining the results (32) and (33) into (31), we find that

$$\sup_{\tilde{\mathbf{v}} \in \mathcal{T}_{\text{ext}}} |\tilde{\mathbf{v}}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \tilde{\mathbf{v}}| \lesssim \frac{\omega(\mathcal{T}) + t}{\sqrt{n}}. \tag{34}$$

holds with probability at least  $1 - 2e^{-ct^2}$ , where  $c > 0$  is a fixed constant.

### B. Proof of Lemma 5.2

Let the tangent cones  $\mathcal{C}$  and  $\mathcal{C}_{\text{ext}}$  be as defined in (7) and (8) respectively. Define the sets

$$\mathcal{T}_- = \mathcal{C}_{\text{ext}} - \mathcal{C}_{\text{ext}} \quad \text{and} \quad \mathcal{T}_+ = \mathcal{C}_{\text{ext}} + \mathcal{C}_{\text{ext}}$$

and note that,

$$\begin{aligned}
\omega(\mathcal{C} - \mathcal{C}) &= \mathbb{E} \left[ \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}} \mathbf{g}^T (\mathbf{u} - \mathbf{v}) \right] \\
&\leq \mathbb{E} \left[ \sup_{\mathbf{u} \in \mathcal{C}} \mathbf{g}^T \mathbf{u} + \sup_{\mathbf{v} \in -\mathcal{C}} \mathbf{g}^T \mathbf{v} \right] = 2\omega(\mathcal{C}).
\end{aligned}$$

Similarly,  $\omega(\mathcal{C} + \mathcal{C}) \leq 2\omega(\mathcal{C})$ . Applying Lemma 5.1 on  $\mathcal{T}_+$  and  $\mathcal{T}_-$ , with advertised probability, we have

$$\sup_{\mathbf{a} \in \mathcal{T}_+ \cup \mathcal{T}_-} |\Lambda(\mathbf{a}, \mathbf{a})| \lesssim \frac{\omega(\mathcal{C}) + t}{\sqrt{n}}.$$

where  $\Lambda(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T (\mathbf{I} - \frac{1}{n} [\mathbf{X} \ \mathbf{1}]^T [\mathbf{X} \ \mathbf{1}]) \mathbf{b}$ . Now, for any  $\mathbf{u}, \mathbf{v} \in \mathcal{C}_{\text{ext}}$ , picking  $\mathbf{u} + \mathbf{v} \in \mathcal{T}_+$ ,  $\mathbf{u} - \mathbf{v} \in \mathcal{T}_-$ , we have

$$|\Lambda(\mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v})|, |\Lambda(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v})| \leq \frac{\omega(\mathcal{C}) + t}{\sqrt{n}}.$$

To proceed, note that

$$\Lambda(\mathbf{u}, \mathbf{v}) = \frac{\Lambda(\mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v}) - \Lambda(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v})}{4}.$$

Hence  $|\Lambda(\mathbf{u}, \mathbf{v})| \lesssim (\omega(\mathcal{C}) + t)/\sqrt{n}$ .

### C. Proof of Lemma 5.4

To start our proof let  $\{\mathbf{x}_i\}_{i=1}^n \sim \mathbf{x}$  be i.i.d. isotropic subgaussian samples in  $\mathbb{R}^p$  and  $\mathbf{X} \in \mathbb{R}^{n \times p}$  be the concatenated design matrix. Let  $x_{ij}$  denotes the  $ij^{\text{th}}$  element of the matrix  $\mathbf{X}$ . Since each  $x_{ij}$  has subexponential norm bounded by a constant, there exists a constant  $C > 0$  such that  $|x_{ij}| \leq C \log(n+p)$  holds with probability at least  $1 - 2(n+p)^{-102}$  using subexponential tail bound. Union bounding over all entries of  $\mathbf{X}$  yields that  $|x_{ij}| \leq C \log(n+p)$  holds for all  $i, j$  with probability at least  $1 - 2(n+p)^{-100}$ . Hence, we can bound each row  $\mathbf{x}_i$  of  $\mathbf{X}$  with probability at least  $1 - 2(n+p)^{-100}$  via

$$\|\mathbf{x}_i\|_{\ell_2} \leq C\sqrt{p} \log(n+p), \tag{35}$$

or equivalently, we have

$$\|\mathbf{x}_i \mathbf{x}_i^T\| \leq \|\mathbf{x}_i\|_{\ell_2}^2 \leq Cp \log^2(n+p)$$

This completes the proof of Lemma 5.4.

### D. Proof of Lemma 5.5

Recall that  $\mathbf{x}_i$ 's are i.i.d. isotropic subexponential. We can estimate the covariance matrix of  $\tilde{\mathbf{x}}$  given  $\|\tilde{\mathbf{x}}\|_{\ell_2} \leq B$  using law of total probability as follows

$$\begin{aligned}
\mathbb{E} \left[ \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \right] &= \mathbb{E} \left[ \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \mid \|\tilde{\mathbf{x}}\|_{\ell_2} \leq B \right] \mathbb{P} \left( \|\tilde{\mathbf{x}}\|_{\ell_2} \leq B \right) \\
&\quad + \mathbb{E} \left[ \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \mid \|\tilde{\mathbf{x}}\|_{\ell_2} > B \right] \mathbb{P} \left( \|\tilde{\mathbf{x}}\|_{\ell_2} > B \right) \tag{36}
\end{aligned}$$

Since a covariance matrix is positive-semidefinite, each term in (36) is individually positive semidefinite. Hence, we will drop

the second term in (36) to get the following lower bound on the covariance matrix

$$\mathbb{E}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T] \geq \mathbb{E}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T \mid \|\tilde{\mathbf{x}}\|_{\ell_2} \leq B] \mathbb{P}(\|\tilde{\mathbf{x}}\|_{\ell_2} \leq B) \quad (37)$$

Using Lemma 5.4, it follows that  $\|\tilde{\mathbf{x}}\|_{\ell_2}^2 = \|\mathbf{x}^T \mathbf{1}\|_{\ell_2}^2 \leq Cp \log^2(n+p) = B^2$  holds with probability at least  $1 - 2(n+p)^{-100}$ . Hence, following (37), we get

$$\begin{aligned} \mathbb{E}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T \mid \|\tilde{\mathbf{x}}\|_{\ell_2} \leq B] &\leq \frac{\mathbb{E}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T]}{\mathbb{P}(\|\tilde{\mathbf{x}}\|_{\ell_2} \leq B)} \\ &\leq \frac{1}{1 - 2(n+p)^{-100}} \mathbf{I}_p \leq 2\mathbf{I}_p. \end{aligned}$$

This completes the proof of Lemma 5.5.

### E. Proof of Lemma 5.9

**Subgaussian case:** Using subgaussian tail, for large enough constant  $C > 0$ , for each  $i$ , we have  $|w_i| \leq C\sigma\sqrt{\log(n)} = \sigma B$  with probability at least  $1 - 2n^{-101}$ . This implies  $\text{clip}(w_i, \sigma B) = w_i$ . Union bounding over all entries of  $\mathbf{w}$ , we find the result which holds with probability at least  $1 - 2n^{-100}$ . **Subexponential case** follows similarly with  $B = C \log(n)$ .

### F. Proof of Lemma 5.10

We prove the result for subexponential samples. Subgaussian case follows similarly. Without losing generality let  $\sigma = 1$  as everything can be scaled accordingly. Defining clip function as previously, set  $\mathbf{z} = \text{clip}(w, B)\mathbf{x}$ . Furthermore, let  $w_{\text{tail}}$  denotes the tail of  $|w|$ , such that,

$$w_{\text{tail}} = \begin{cases} |w| & \text{if } |w| > B \\ 0 & \text{otherwise} \end{cases} \quad (38)$$

$w_{\text{tail}}$  is an upper bound on the error due to clipping i.e.

$$|w - \text{clip}(w, B)| \leq w_{\text{tail}} \quad (39)$$

We proceed by upper bounding  $\|\mathbb{E}[\mathbf{z}]\|_{\ell_2}$  in terms of  $w_{\text{tail}}$ , using subadditive property of  $\ell_2$ -norm and the orthogonality of  $w$  and  $\mathbf{x}$  (i.e.,  $\mathbb{E}[w\mathbf{x}] = 0$ ) as follows

$$\begin{aligned} \|\mathbb{E}[\mathbf{z}]\|_{\ell_2} &= \|\mathbb{E}[\text{clip}(w, B)\mathbf{x}]\|_{\ell_2} \\ &= \|\mathbb{E}[(w - \text{clip}(w, B))\mathbf{x}]\|_{\ell_2} \\ &\leq \mathbb{E}[|w - \text{clip}(w, B)|\|\mathbf{x}\|_{\ell_2}] \\ &\leq \mathbb{E}[w_{\text{tail}} \max(\|\mathbf{x}\|_{\ell_2}, \sqrt{p}B)] \end{aligned} \quad (40)$$

Using subexponentiality, for some constant  $c > 0$ , we have that  $\mathbb{P}(w_{\text{tail}} > \sqrt{ct}) \leq 2e^{-t}$  and  $\mathbb{P}\{\|\mathbf{x}\|_{\ell_2} > \sqrt{cpt}\} \leq 2pe^{-t}$  where latter follows from union bounding entries of  $\mathbf{x}$ . Union bounding, we get the following tail bound on the product,

$$\mathbb{P}\{w_{\text{tail}}\|\mathbf{x}\|_{\ell_2} > c\sqrt{pt}^2\} \leq 4pe^{-t}. \quad (41)$$

For notational convenience, set

$$g = w_{\text{tail}} \max(\|\mathbf{x}\|_{\ell_2}, \sqrt{p}B). \quad (42)$$

and note that  $g$  satisfies the following property due to (38)

$$\begin{cases} \text{either} & g > \sqrt{p}B^2 \\ \text{or} & g = 0 \end{cases}. \quad (43)$$

Furthermore, from (41) we get the following tail distribution

$$Q_g(t) = \mathbb{P}(g > t) \leq 4pe^{-\left[\frac{t}{c\sqrt{p}}\right]^{1/2}}. \quad (44)$$

for  $t \geq \alpha := \sqrt{p}B^2$ . Combining (42), (43) and (44) into (40) and denoting probability density function of  $g$  by  $f_g$ , we get

$$\begin{aligned} \|\mathbb{E}[\mathbf{z}]\|_{\ell_2} &\leq \mathbb{E}[g] = \int_{\alpha}^{\infty} tf_g(t)dt = - \int_{\alpha}^{\infty} tdQ_g(t) \\ &= -tQ_g(t)|_{\alpha}^{\infty} + \int_{\alpha}^{\infty} Q_g(t)dt \\ &= \sqrt{p}B^2 Q_g(\sqrt{p}B^2) + \int_{\alpha}^{\infty} Q_g(t)dt \\ &\stackrel{(a)}{\leq} e^{-B/\sqrt{c}} + 4p \int_{\sqrt{p}B^2}^{\infty} e^{-\left[\frac{t}{c\sqrt{p}}\right]^{1/2}} dt. \end{aligned} \quad (45)$$

where (a) follows from (44). To bound the term on the right hand side, we do a change of variable in (45) by setting  $\tau = [t/(c\sqrt{p})]^{1/2}$  to get,

$$\begin{aligned} 4p \int_{\sqrt{p}B^2}^{\infty} e^{-\left[\frac{t}{c\sqrt{p}}\right]^{1/2}} dt &\leq 8cp^2 \int_{B/\sqrt{c}}^{\infty} \tau e^{-\tau} d\tau \\ &\leq 8cp^2 \left[ -\tau e^{-\tau} \Big|_{B/\sqrt{c}}^{\infty} + \int_{B/\sqrt{c}}^{\infty} e^{-\tau} d\tau \right] \\ &= 8cp^2 \left[ \frac{B}{\sqrt{c}} e^{-B/\sqrt{c}} + e^{-B/\sqrt{c}} \right] \\ &\leq 8cp^2 (B/\sqrt{c} + 1) e^{-B/\sqrt{c}}. \end{aligned} \quad (46)$$

Overall, we found

$$\|\mathbb{E}[\mathbf{z}]\|_{\ell_2} \leq 4p^2 (B^2 + 2c(B/\sqrt{c} + 1)) e^{-B/\sqrt{c}}$$

which is upper bounded by  $\mathcal{O}(n^{-100})$  by picking  $B = C \log n$ . Finally, note that conditioned on  $|w| \leq B$ ,  $\mathbf{z} = w\mathbf{x}$  and

$$\|\mathbb{E}[\mathbf{z}]\|_{\ell_2} \geq \|\mathbb{E}[w\mathbf{x} \mid |w| \leq B]\|_{\ell_2} \mathbb{P}(|w| \leq B).$$

Since  $\mathbb{P}(|w| \leq B) > 1/2$ , this yields  $\|\mathbb{E}[w\mathbf{x} \mid |w| \leq B]\|_{\ell_2} \lesssim p^2 n^{-201}$  which is the advertised result.

Similarly for subgaussian samples, one can show that

$$\|\mathbb{E}[\mathbf{z}]\|_{\ell_2} \lesssim p^2 B^2 e^{-B^2/c}. \quad (47)$$

Picking  $B = C\sqrt{\log n}$ , we conclude with the same result concluding the proof of Lemma 5.10.