

Moving Beyond Sub-Gaussianity in High Dimensional Statistics: Applications in Covariance Estimation and Linear Regression*

Arun Kumar Kuchibhotla^{1†} Abhishek Chakraborty^{1‡}

May 11, 2022

Abstract

Concentration inequalities form an essential toolkit in the study of high dimensional statistical methods. Most of the relevant statistics literature in this regard is, however, based on the assumptions of sub-Gaussian or sub-exponential random variables/vectors. In this paper, we first bring together, through a unified exposition, various probabilistic inequalities for sums of independent random variables under much more general exponential type (namely sub-Weibull) tail assumptions. These results extract a part sub-Gaussian tail behavior of the sum in finite samples, matching the asymptotics governed by the central limit theorem, and are compactly represented in terms of a new Orlicz quasi-norm – the Generalized Bernstein-Orlicz norm – that typifies such kind of tail behaviors.

We illustrate the usefulness of these inequalities through the analysis of four fundamental problems in high dimensional statistics. In the first two problems, we study the rate of convergence of the sample covariance matrix in terms of the maximum elementwise norm and the maximum k -sub-matrix operator norm which are key quantities of interest in bootstrap procedures and high dimensional structured covariance matrix estimation, as well as in high dimensional and post-selection inference. The third example concerns the restricted eigenvalue condition, required in high dimensional linear regression, which we verify for all sub-Weibull random vectors through a unified analysis, and also prove a more general result related to restricted strong convexity in the process. In the final example, we consider the Lasso estimator for linear regression and establish its rate of convergence to be generally $\sqrt{k \log p/n}$, for k -sparse signals, under much weaker than usual tail assumptions (on the errors as well as the covariates), while also allowing for misspecified models and both fixed and random design. To our knowledge, these are the first such results for Lasso obtained in this generality. The common feature in all our results over all the examples is that the convergence rates under most exponential tails match the usual (optimal) ones obtained under sub-Gaussian assumptions. Finally, we also establish some complementary results on analogous tail bounds for the suprema of empirical processes indexed by sub-Weibull variables. All our results are finite sample.

Keywords: Concentration Inequalities, Orlicz Norms, Sub-Weibull Random Variables, Structured Covariance Matrix Estimation, Restricted Eigenvalue Condition, High Dimensional Linear Regression and Lasso, Empirical Processes.

*To appear in *Information and Inference: A Journal of the IMA*.

[†]Department of Statistics and Data Science, Carnegie Mellon University. (Email: arunku@cmu.edu)

[‡]Department of Statistics, Texas A&M University. (Email: abhishek@stat.tamu.edu; Corresponding author)

¹The authors contributed equally to this work.

1 Introduction and Motivation

In the current era of big data, with an abundance of information often available for a large number of variables, there has been a burst of statistical methods dealing with high dimensional data. In particular, estimation and inference methods are being developed for settings with a huge number of variables often larger than the number of observations available. In these settings, classical statistical methods such as the least squares or the maximum likelihood principle usually do not lead to meaningful estimators, and regularization methods have been widely used as an alternative; see, e.g., [Wainwright \(2019\)](#) for an overview. These methods typically penalize the original loss function, e.g. squared error loss or the negative log-likelihood function, with a penalty on the parameter vector that reduces the “effective” number of parameters being estimated. The theoretical analyses of most of these methods, despite all their diversities, generally obey a *common unifying theme* wherein a key quantity to control is the maximum of a (high dimensional) vector of averages of mean zero random variables. Since the dimension is potentially larger than the sample size, it is important to analyze the behavior of the maximum in a non-asymptotic way. Concentration inequalities and probabilistic tail bounds form a major part of the toolkit required for such analyses.

Some of the most commonly used probabilistic tail bounds are of the exponential type, including, in particular, Hoeffding’s and Bernstein’s inequalities; see Section 3.1 of [Giné and Nickl \(2016\)](#) for a review. In the classical versions of these inequalities, the random variables are assumed to be bounded, but this assumption can be relaxed to sub-Gaussian and sub-exponential random variables, respectively; see Sections 2.6 and 2.8 of [Vershynin \(2018\)](#), and also [Wainwright \(2019\)](#). A random variable is called sub-Gaussian if its survival function is bounded by that of a Gaussian distribution. A sub-exponential random variable is defined similarly (see Section 2). Note that in both these cases, the moment generating function (MGF) exists in a neighborhood around zero. Most of the high dimensional statistics literature is based on the assumption of sub-Gaussian or sub-exponential random variables/vectors. But in many applications, these assumptions may not be appropriate. For instance, consider the following two simple examples that exemplify the main issues.

- Suppose $(X_1, Y_1), \dots, (X_n, Y_n)$ are independent and identically distributed (i.i.d.) observations of a random vector $(X, Y) \in \mathbb{R}^2$ and let $\hat{\beta} = \sum X_i Y_i / \sum X_i^2$ denote the linear regression slope estimator for regressing Y on X . Under a possibly misspecified linear model, the estimation of the asymptotic variance of $\hat{\beta}$ involves $\sum X_i^2 (Y_i - X_i \beta)^2$, where β is the limit of $\hat{\beta}$; see [Buja et al. \(2019\)](#) for details. It is clear that if the initial random variables X and Y are sub-exponential, then the random variables $X_i^2 (Y_i - X_i \beta)^2$ do not have a finite MGF. The same holds even when the ingredient random variables X and Y are further assumed to be sub-Gaussian.
- Let Y be a response variable, and X_1, X_2 be two covariates, all having a finite MGF in a neighborhood of zero. In many applications, it is important to consider regression models with interaction effects among the covariates, and more generally, second (or higher) order effects such as $X_1^2, X_1 X_2$ etc. The presence of such second order effects clearly implies that the summands involved in the analyses of these linear regression estimators may not necessarily have a finite MGF anymore.

These examples are not high dimensional in nature, but are mainly presented here as some basic examples where the core problem becomes apparent. The requirement of controlling averages defined by higher order or product-type terms, as in the second example, also arises inevitably in

the case of high dimensional regression and covariance estimation; see also the recent work of [Yu et al. \(2019\)](#) on problems of a similar flavor. The first example, apart from its relevance in inference for linear regression estimators, also appears in the problem of testing for the existence of active predictors in linear regression. This problem can be reduced to a simultaneous significance testing problem based on all the marginal regressions, as shown in [McKeague and Qian \(2015\)](#). For this type of marginal testing problems, uniform consistency of the estimators of the variance of all the marginal regression coefficient estimators is required, thus creating the need for a non-asymptotic analysis.

1.1 Our Contributions

Tail bounds for sums of independent random variables play an important role in probability and statistics. In statistics, especially in the high dimensional statistics literature, most of the applications are studied only under strong (or light) exponential tail assumptions such as sub-Gaussian or sub-exponential. Although tail bounds do exist for sums of independent random variables with “heavy” exponential tails (scattered mostly in the probability literature), the impact of moving from sub-Gaussian/sub-exponential (i.e. light-tailed) variables to those with heavy exponential tails on the rates of convergence and the dependence on the dimension does not seem to be well-studied in the statistics literature. These heavy exponential tailed random variables are what we call *sub-Weibull* variables (see Definition 2.2).

The first goal of our article is to provide a clear (and unified) exposition of concentration inequalities related to sub-Weibull random variables, which constitutes the *first part* of our paper. It provides in one place a user-friendly off-the-shelf toolset that can be readily used in the analysis of a variety of modern statistical problems (and yet under much weaker tail conditions than those typically assumed).

In the *second part*, we provide applications of these concentration inequalities for four such fundamental problems in high dimensional statistics. A detailed account of our contributions along both these lines is provided next.

Exposition of Tail Bounds. The outline of our probability exposition is as follows. We first propose a new Orlicz quasi-norm called the *Generalized Bernstein-Orlicz (GBO) norm* that allows for a compact representation of the results regarding sub-Weibull random variables. [van de Geer and Lederer \(2013\)](#) introduced its predecessor, the Bernstein-Orlicz norm, that provides a formal understanding of the nature of the tail bound given by Bernstein’s Inequality (see Section 2 for details). The recent paper [Wellner \(2017\)](#) extends the results of [van de Geer and Lederer \(2013\)](#) to capture the tail behavior given by Bennett’s inequality. Although it was not stressed in [van de Geer and Lederer \(2013\)](#), one of the main features of Bernstein’s inequality is that even for sub-exponentials, it provides a part sub-Gaussian tail behavior for the sum. This, in turn, plays a key role in proving the rate of convergence of a maximum of several such sums to be the same as that in the case of sub-Gaussian variables. The GBO norm is constructed with the aim of capturing a similar tail behavior for the general case of sub-Weibulls. The results on unbounded empirical processes from [Adamczak \(2008\)](#), along with a maximal inequality (Theorem 5.2) of [Chernozhukov et al. \(2014\)](#) and the results of [Latała \(1997\)](#), will be exploited to provide a sequence of ready-to-use results (Theorems 3.1–3.4) about sub-Weibull random variables. This is essentially the probability contribution of the current article. The results of [Adamczak \(2008\)](#) are derived based on Chapter 6 of [Ledoux and Talagrand \(1991\)](#), and those of [Chernozhukov et al. \(2014\)](#) are based on the maximal

inequality of [van der Vaart and Wellner \(2011\)](#). All of our results are derived under the assumption of independence only and allow for *non-identically distributed* ingredient variables. The extensions for the supremum of empirical processes with sub-Weibull envelope functions are further discussed in [Appendix B](#).

Lastly, we would also like to point out that we mainly focus on exponential-type tails in this paper, since in all our high dimensional applications, a logarithmic dependence on the dimension is desired (our proof techniques, however, also apply equally to polynomial-type tails). The initial version of the current paper was available in ArXiv since 2018 ([Kuchibhotla and Chakraborty, 2018](#)). Recently, [Bakhshizadeh et al. \(2020\)](#) further explored some refinements of our tail bound results there. But a unified exposition of concentration results, coupled with a thorough demonstration of their usefulness in various important statistical applications as discussed below, is still lacking in the literature to the best of our knowledge.

High Dimensional Statistical Applications. Following the exposition of concentration inequalities, we apply these probabilistic tools to four fundamental problems in high dimensional statistics. In all these examples, we establish precise tail bounds and rates of convergence, under the assumption of sub-Weibull random variables/vectors only. The results, apart from being seamlessly unified and general in terms of the underlying tail assumptions, also exhibit several interesting features and provide some key insights into the behavior of these problems. In particular, *a common outcome of all our analyses is that the rates of convergence generally match those obtained under the sub-Gaussian assumption.*

Furthermore, most results in high dimensional statistics that involve *random vectors* are only derived under tail assumptions on the *joint* distribution of the random vector (for example, a random vector X is sub-Gaussian if $\theta^\top X$ is uniformly sub-Gaussian over all θ of unit Euclidean norm). Although commonly adopted in the literature, such a condition imposes (often implicitly) certain strong restrictions on the joint distribution, as discussed at the beginning of [Section 4](#). Throughout this paper, we make a formal *distinction* between such a ‘*joint*’ assumption on the tail behavior of a random vector versus a much weaker ‘*marginal*’ assumption on the tail behaviors of its coordinates only; see [Definitions 2.4](#) and [2.5](#). All of our applications are also studied under such an assumption only on the marginal distributions, and often with nearly (if not exactly) similar results and convergence rates.

The description and the main implications of our results for each of the four high dimensional statistical applications we consider in this paper are enlisted below. (In all examples, p denotes the ambient dimension of the random vectors and n denotes the sample size.)

1. *Covariance Estimation (Maximum Elementwise Norm).* A central part of high dimensional inference hinges on an application of the central limit theorem through a bootstrap procedure. The consistency of the bootstrap in this case requires consistent estimation of the covariance matrix in terms of the maximum elementwise norm. This norm also appears in the coupling inequality for maxima of sums of random vectors; see [Theorem 4.1](#) of [Chernozhukov et al. \(2014\)](#). In [Section 4.1](#), we prove a finite sample tail bound (via [Theorems 4.1](#) and [4.2](#)) for the error of the sample covariance matrix in terms of this norm under the assumption of (marginally) sub-Weibull (α) ingredient random vectors. The rate of convergence is shown to be $\sqrt{\log p/n}$ if $\log p = o(n^{\alpha/(4-\alpha)})$; see [Remark 4.1](#). This rate of convergence can be easily shown to be optimal in case the random vectors are standard multivariate Gaussian. Furthermore, the tail bounds presented in this section also play a central role in sparse covariance matrix estimation, as

shown in [Bickel and Levina \(2008\)](#) and [Cai and Liu \(2011\)](#). Both these papers deal with jointly sub-Gaussian random vectors, while the second paper additionally deals with fixed polynomial moments. Using our results in [Section 4.1](#), the problem of sparse covariance matrix estimation can be analyzed under weaker assumptions with logarithmic dependence on the dimension. Finally, the results in this section also establish the consistency of bootstrap procedures when applied to (high dimensional) marginally sub-Weibull random vectors.

2. *Covariance Estimation (Maximum k -Sub-Matrix Operator Norm)*. Covariance matrices play an important role in statistical analyses through principal component analysis, factor analysis and so on. Clearly, for most of these methods, consistency of the covariance matrix estimator in terms of the operator norm is important. In high dimensions, however, the sample covariance matrix is known to be not consistent in the operator norm. Under such settings, in practice, one often selects a (random) subset of variables and focuses on the spectral properties of the corresponding covariance (sub)-matrix. In [Section 4.2](#), we study the consistency of the sample covariance matrix of (marginal or joint) sub-Weibull (α) ingredient random vectors, in terms of the maximum sub-matrix operator norm with sub-matrix size $k \leq p$. We show through [Theorem 4.3](#) that the rate of convergence is $\sqrt{k \log(ep/k)/n}$ for most values of $\alpha > 0$. This rate was previously obtained for the joint sub-Gaussian case by [Loh and Wainwright \(2012\)](#); see [Lemma 15](#) therein. This norm was possibly first studied by [Rudelson and Vershynin \(2008\)](#) for bounded random variables. The convergence rate of this norm plays a key role in studying post-Lasso least squares linear regression estimators and in structured covariance matrix estimation. The post-Lasso linear regression estimator was studied in [Belloni and Chernozhukov \(2013\)](#), and more generally, in [Kuchibhotla et al. \(2018\)](#) for post-selection inference. Lastly, for adaptive estimation of so-called bandable covariance matrices, a thresholding mechanism was introduced by [Cai and Yuan \(2012\)](#), where a result about maximum sub-matrix operator norm is also required. [Cai and Yuan \(2012\)](#) deal with Gaussian random vectors, and using our results this method can be thus extended to sub-Weibull random vectors.
3. *Restricted Eigenvalues*. [Bickel et al. \(2009\)](#) introduced the restricted eigenvalue (RE) condition to analyze the Lasso and the Dantzig selector. The RE condition concerns the minimum eigenvalue of the sample covariance matrix when the directions are restricted to lie in a specific cone (see [Section 4.3](#) for a precise definition), and its verification forms a key step in high dimensional linear regression. A well known result in this regard is that of [Rudelson and Zhou \(2013\)](#) who verified the RE condition for the covariance matrices of jointly sub-Gaussian random vectors. Some extensions under weaker tail assumptions (e.g. sub-exponentials) have also been considered by [Lecué and Mendelson \(2017\)](#), among others; see [Section 4.3](#) for further details. Based on our results in [Section 4.2](#), we prove in [Section 4.3](#) that covariance matrices of both jointly and marginally sub-Weibull random vectors satisfy the RE condition with probability tending to one. In fact, we prove a more general result (in [Theorem 4.4](#)) related to *restricted strong convexity* from which the RE condition's verification follows as a consequence. To our knowledge, such unified results regarding the RE condition are not so easily accessible in the core statistics literature.
4. *Linear Regression via Lasso*. One of the most popular and possibly the first high dimensional linear regression technique is the Lasso introduced by [Tibshirani \(1996\)](#). The general results of [Negahban et al. \(2012\)](#) provide an easy recipe for studying the rate of convergence of the

Lasso estimator. Based on this general recipe and equipped with the verification of the RE condition, we prove in Section 4.4 (via Theorems 4.5 and 4.6) the rate of convergence of the Lasso estimator to be $\sqrt{k \log p/n}$ (the near minimax optimal rate) under sub-Weibull covariates and sub-Weibull/polynomial-tailed errors when the “true” regression parameter is assumed to be k -sparse. We also *allow* for both fixed and random designs, as well as for misspecified models. Apart from admitting several other extensions (see Remark 4.16), our results *only* assume a marginal sub-Weibull property of the covariates, thus making them stronger than most existing results for Lasso which usually provide the rates under jointly sub-Gaussian/sub-exponential covariate vectors. To our knowledge, these are the first such results for the Lasso obtained in this generality.

1.2 Organization

The rest of this paper is organized as follows. In Section 2, we define the class of sub-Weibull random variables and introduce the Generalized Bernstein-Orlicz norm. A detailed discussion of several useful and basic properties of the GBO norm is deferred to Appendix A. Section 3 provides several ready-to-use bounds for sums of independent mean zero sub-Weibull random variables. Using the results of Section 3, the fundamental statistical applications discussed above are studied in Section 4 (via Sections 4.1-4.4 dedicated respectively to these four problems). We conclude with a summary and directions for future research in Section 5.

In the [Supplementary Material](#) (Appendices B–F), we include additional results and technical materials that could not be accommodated in the main article. In Appendix B, we provide some supplementary results on tail bounds for suprema of empirical processes with sub-Weibull envelopes, and maximal inequalities based on uniform and bracketing entropy. Proofs of all the results in Section 2 (along with those in Appendix A) and Section 3 are presented in Appendices C and D, respectively. The results of Section 4, as well as Appendix B, are proved in Appendices E and F, respectively.

2 The Generalized Bernstein-Orlicz (GBO) Norm

We first recall the general definition of an Orlicz norm for random variables. For a historical account of Orlicz norms, and sub-Gaussian, sub-exponential (and sub-Weibull) variables, we refer to Section 1 of [Wellner \(2017\)](#) and the references therein.

Definition 2.1 (Orlicz Norms). *Let $g : [0, \infty) \rightarrow [0, \infty)$ be a non-decreasing function with $g(0) = 0$. The “ g -Orlicz norm” of a real-valued random variable X is given by*

$$\|X\|_g := \inf\{\eta > 0 : \mathbb{E}[g(|X|/\eta)] \leq 1\}. \quad (2.1)$$

The function $\|\cdot\|_g$ on the space of real-valued random variables is not a norm unless g is additionally a convex function. We define the g -Orlicz norm here under the only assumption of monotonicity of g , since in the following, convexity is not satisfied and is also not required. It readily follows from (2.1) that

$$\mathbb{P}(|X| \geq \eta g^{-1}(t)) \leq \frac{1}{t} \quad \text{for all } t \geq 0. \quad (2.2)$$

Two very important special cases of g are given by $\psi_2(x) := \exp(x^2) - 1$ and $\psi_1(x) := \exp(x) - 1$, which correspond to *sub-Gaussian* and *sub-exponential* random variables, respectively. As a generalization, we now define *sub-Weibull* random variables as follows.

Definition 2.2 (Sub-Weibull Variables). *A random variable X is said to be sub-Weibull of order $\alpha > 0$, denoted as sub-Weibull (α) , if*

$$\|X\|_{\psi_\alpha} < \infty, \quad \text{where } \psi_\alpha(x) := \exp(x^\alpha) - 1 \quad \text{for } x \geq 0.$$

Based on this definition, it follows that if X is sub-Weibull (α) , then

$$\mathbb{P}(|X| \geq t) \leq 2 \exp\left(-\frac{t^\alpha}{\|X\|_{\psi_\alpha}^\alpha}\right), \quad \text{for all } t \geq 0. \quad (2.3)$$

The right hand side here resembles the survival function of a Weibull random variable of order $\alpha > 0$, and hence the name sub-Weibull random variable. It is also clear from inequality (2.3) that the smaller the α is, the more heavy-tailed the random variable is.

A simple calculation implies that a converse of the tail bound result in (2.3) also holds. It can further be shown that X is sub-Weibull of order α , if and only if, its moments satisfy

$$\sup_{r \geq 1} r^{-1/\alpha} \|X\|_r < \infty,$$

where $\|X\|_r := (\mathbb{E}[|X|^r])^{1/r}$; see Propositions 2.5.2 and 2.7.1 of [Vershynin \(2018\)](#) for similar results. Clearly, sub-exponential and sub-Gaussian random variables are sub-Weibull of orders 1 and 2 respectively, while bounded variables are sub-Weibulls of order ∞ . Also, X is sub-exponential if and only if $|X|^{1/\alpha}$ is sub-Weibull of order α ; this follows readily from Definition 2.2.

Next, to define the Generalized Bernstein-Orlicz norm, we first recall the classical Bernstein inequality for sub-exponential random variables. Suppose X_1, \dots, X_n are independent mean zero sub-exponential random variables, then

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) \leq 2 \times \begin{cases} \exp(-t^2/(4\sigma_n^2)), & \text{if } t < \sigma_n^2/C_n, \\ \exp(-t/(4C_n)), & \text{otherwise,} \end{cases} \quad (2.4)$$

where $\sigma_n^2 := 2 \sum_{i=1}^n \|X_i\|_{\psi_1}^2$ and $C_n := \max\{\|X_i\|_{\psi_1} : 1 \leq i \leq n\}$; see Proposition 3.1.8 of [Giné and Nickl \(2016\)](#). Clearly, the tail of the sum behaves like a Gaussian for smaller values of t and behaves like an exponential for larger t .

An equivalent way of writing inequality (2.4) that leads to the Bernstein-Orlicz norm is

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i\right| \geq \eta_1 \sqrt{\sigma_n^2 \log(1+t)} + \eta_2 C_n \log(1+t)\right) \leq \frac{1}{t},$$

for some constants $\eta_1, \eta_2 > 0$. Comparing this inequality with (2.2), one can define an Orlicz norm through a function $g_\eta(\cdot)$ whose inverse is given by:

$$g_\eta^{-1}(t) := \sqrt{\log(1+t)} + \eta \log(1+t),$$

parametrized by $\eta > 0$. The corresponding Orlicz norm $\|\cdot\|_{g_\eta}$ is exactly the Bernstein-Orlicz norm introduced by [van de Geer and Lederer \(2013\)](#). The Generalized Bernstein-Orlicz (GBO) norm is now defined analogously as follows.

Definition 2.3 (Generalized Bernstein-Orlicz Norm). Fix $\alpha > 0$ and $L \geq 0$. Define the function $\Psi_{\alpha,L}(\cdot)$ based on the inverse function

$$\Psi_{\alpha,L}^{-1}(t) := \sqrt{\log(1+t)} + L(\log(1+t))^{1/\alpha} \quad \text{for all } t \geq 0. \quad (2.5)$$

The Generalized Bernstein-Orlicz (GBO) norm of a random variable X is then given by $\|X\|_{\Psi_{\alpha,L}}$ as in Definition 2.1.

Remark 2.1 It is easy to verify from (2.5) that $\Psi_{\alpha,L}(\cdot)$ is monotone and $\Psi_{\alpha,L}(0) = 0$ and so, Definition 2.1 is applicable. The function $\Psi_{\alpha,L}(\cdot)$ does not have a closed form expression in general, and is not convex for $\alpha < 1$. But $\|\cdot\|_{\Psi_{\alpha,L}}$ is a quasi-norm; see Proposition A.5 in Appendix A. The properties proved for the Bernstein-Orlicz norm in van de Geer and Lederer (2013) also hold for the GBO norm $\|\cdot\|_{\Psi_{\alpha,L}}$ even though the function $\Psi_{\alpha,L}(\cdot)$ is not convex for $\alpha < 1$. Several basic properties of the GBO norm, along with equivalent tail and moment bound properties and some maximal inequalities, are presented in Appendix A. \diamond

The ready-to-use concentration inequality results in Section 3 are presented in terms of the $\|\cdot\|_{\Psi_{\alpha,L}}$ norm and for this reason, we briefly mention here the precise nature of the tail behavior captured by the GBO norm. If $\|X\|_{\Psi_{\alpha,L}} < \infty$, then

$$\mathbb{P}\left(|X| \geq \|X\|_{\Psi_{\alpha,L}} \left\{ \sqrt{t} + Lt^{1/\alpha} \right\}\right) \leq 2 \exp(-t) \quad \text{for all } t \geq 0.$$

So, for t small enough, the survival function of X behaves like a Gaussian, and for t larger, the survival function behaves like a Weibull of order α . Hence, the results from Section 3 will imply that the tail of a sum of independent sub-Weibull random variables behaves like a combination of a Gaussian tail and a Weibull tail.

2.1 Sub-Weibull Random Vectors

For our applications, we consider the following two definitions of sub-Weibull random vectors. For any vector $x \in \mathbb{R}^q$, let $x(j)$ represent the j -th coordinate of x for all $1 \leq j \leq q$, and let $\|x\|_r := \left(\sum_{j=1}^q |x(j)|^r\right)^{1/r}$ denote the vector L_r -norm of x for any $r \geq 1$. (For $r = 2$, we sometimes also refer to the vector L_2 -norm simply as the Euclidean norm.)

Definition 2.4 (Joint Sub-Weibull Vectors). A random vector $X \in \mathbb{R}^q$ is said to be jointly sub-Weibull of order $\alpha > 0$ if for every $\theta \in \mathbb{R}^q$ of unit Euclidean norm, $X^\top \theta$ is sub-Weibull of order α , and the joint sub-Weibull (α) norm of X , $\|X\|_{J,\psi_\alpha}$ (where the subscript “J” stands for “joint”), is given by

$$\|X\|_{J,\psi_\alpha} := \sup_{\theta \in \mathbb{R}^q, \|\theta\|_2=1} \left\| X^\top \theta \right\|_{\psi_\alpha}.$$

This is one of the most commonly adopted type of tail assumptions on random vectors (especially with $\alpha = 2$); see Section 3.4 of Vershynin (2018). As with random variables, the cases $\alpha = 1, 2$ correspond to sub-exponential and sub-Gaussian random vectors, respectively.

Definition 2.5 (Marginal Sub-Weibull Vectors). A random vector $X \in \mathbb{R}^q$ is said to be marginally sub-Weibull of order $\alpha > 0$ if for every $1 \leq j \leq q$, $X(j)$ is sub-Weibull of order α , and the marginal

sub-Weibull (α) norm of X , $\|X\|_{M,\psi_\alpha}$ (where the subscript “M” stands for “marginal”), is given by

$$\|X\|_{M,\psi_\alpha} := \sup_{1 \leq j \leq q} \|X(j)\|_{\psi_\alpha}.$$

Clearly, $\|X\|_{M,\psi_\alpha} \leq \|X\|_{J,\psi_\alpha}$ for any random vector X , and hence, a marginal sub-Weibull property is (much) *weaker* than a joint sub-Weibull property. A detailed comparison of the marginal and joint sub-Weibull properties is deferred to the beginning of Section 4.

3 Norms of Sums of Independent Random Variables

The following sequence of results show the use of the $\Psi_{\alpha,L}$ -norm in representing the part sub-Gaussian tail behavior in finite samples for sums of independent random variables when the ingredient random variables are sub-Weibull (α). All results in this section are stated for independent random variables that are possibly non-identically distributed. Extensions to the case of dependent random variables also exist in the literature; see [Merlevède et al. \(2011\)](#) and Appendix B of [Kuchibhotla et al. \(2018\)](#). The proofs of all the results in this section are given in Appendix D.

The following result can be derived from Theorem 2 of [Latała \(1997\)](#). (Note that the constants here are explicit, but they are not optimized and could possibly be improved.)

Theorem 3.1. *If X_1, \dots, X_n are independent mean zero random variables with $\|X_i\|_{\psi_\alpha} < \infty$ for all $1 \leq i \leq n$ and some $\alpha > 0$, then for any vector $a = (a_1, \dots, a_n) \in \mathbb{R}^n$, the following bounds hold true:*

$$\left\| \sum_{i=1}^n a_i X_i \right\|_{\Psi_{\alpha, L_n(\alpha)}} \leq 2eC(\alpha) \|b\|_2,$$

and

$$\mathbb{P} \left(\left| \sum_{i=1}^n a_i X_i \right| \geq 2eC(\alpha) \|b\|_2 \sqrt{t} + 2eL_n^*(\alpha) t^{1/\alpha} \|b\|_{\beta(\alpha)} \right) \leq 2e^{-t} \quad \text{for all } t \geq 0, \quad (3.1)$$

where $b = (a_1 \|X_1\|_{\psi_\alpha}, \dots, a_n \|X_n\|_{\psi_\alpha}) \in \mathbb{R}^n$,

$$C(\alpha) := \max\{\sqrt{2}, 2^{1/\alpha}\} \times \begin{cases} \sqrt{8}e^3(2\pi)^{1/4}e^{1/24}(e^{2/e}/\alpha)^{1/\alpha}, & \text{if } \alpha < 1, \\ 4e + 2(\log 2)^{1/\alpha}, & \text{if } \alpha \geq 1, \end{cases}$$

and for $\beta(\alpha) = \infty$ when $\alpha \leq 1$ and $\beta(\alpha) = \alpha/(\alpha - 1)$ when $\alpha > 1$,

$$L_n(\alpha) := \frac{4^{1/\alpha}}{\sqrt{2} \|b\|_2} \times \begin{cases} \|b\|_{\beta(\alpha)}, & \text{if } \alpha < 1, \\ 4e \|b\|_{\beta(\alpha)} / C(\alpha), & \text{if } \alpha \geq 1, \end{cases}$$

and for (3.1), the quantity $L_n^*(\alpha) = L_n(\alpha)C(\alpha)\|b\|_2/\|b\|_{\beta(\alpha)}$.

Remark 3.1 (Sharpness of Theorem 3.1). Theorem 3.1 provides a useful generalization of Theorem 2.8.1 of [Vershynin \(2018\)](#) for $\alpha \neq 1$. The transition in our result at $\alpha = 1$ is due to the fact that Weibull random variables are log-convex for $\alpha \leq 1$ and log-concave for $\alpha \geq 1$. It is worth noting that the conclusion of Theorem 3.1 cannot be improved in terms of dependence on $a = (a_1, \dots, a_n)$ and are optimal in the sense that there exists distributions for X_i satisfying

$\|X_i\|_{\psi_\alpha} \leq 1$ for which there is a lower bound matching the upper bound; see Theorem 2 and Examples 3.2 and 3.3 of [Latała \(1997\)](#). In particular, Examples 3.2 and 3.3 of [Latała \(1997\)](#) show that the moment bounds implied by Theorem 3.1 (via Proposition A.3) have matching lower bounds when the random variables X_1, \dots, X_n satisfy $\mathbb{P}(|X_i| \geq t) = \exp(-t^\alpha)$ for all $t \geq 0$. It should also be noted that these optimality results were also derived earlier by [Gluskin and Kwapien \(1995\)](#) and [Hitczenko et al. \(1997\)](#). In particular, for $\alpha \geq 1$, the corollary on page 307 of [Gluskin and Kwapien \(1995\)](#) shows that the probability tail bound implied by Theorem 3.1 (via Proposition A.3) is optimal in that there is a lower bound on the tail probability that only differs from the upper bound by a universal constant. We are not aware of a similar result for $\alpha < 1$. It is worth stressing here that the lower bounds mentioned in this remark should be understood in a minimax sense: there exists a distribution setting for independent random variables X_1, \dots, X_n for which the bound implied by Theorem 3.1 is sharp (i.e., Theorem 3.1 cannot be improved without further assumptions). \diamond

Tail Bounds Scaling with Variance. The bound provided by Theorem 3.1 is solely in terms of $\|X_i\|_{\psi_\alpha}$. It is clear, however, from the classical central limit theorem (CLT) that asymptotically the distribution of the sum (properly scaled) is determined by the variance of the sum. Although it is impossible to prove an exponential tail bound solely in terms of the variance, we expect at least the Gaussian part of the tail to depend on the variance only. This is the content of the next three results - Theorems 3.2–3.3 (on norm bounds) and Theorem 3.4 (on tail bounds). The proofs are based on the techniques of [Adamczak \(2008\)](#).

Theorem 3.2 (Bounds Scaling with Variance – the Case $\alpha \leq 1$). *If X_1, \dots, X_n are independent mean zero random variables with $\|X_i\|_{\psi_\alpha} < \infty$ for all $1 \leq i \leq n$ and some $0 < \alpha \leq 1$, then*

$$\left\| \sum_{i=1}^n X_i \right\|_{\Psi_{\alpha, L_n(\alpha)}} \leq 2e\sqrt{6} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{1/2},$$

with

$$L_n(\alpha) = \frac{4^{1/\alpha} K_\alpha C_\alpha}{2\sqrt{6}} (\log(n+1))^{1/\alpha} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{-1/2} \max_{1 \leq i \leq n} \|X_i\|_{\psi_\alpha},$$

for some constants $C_\alpha, K_\alpha > 0$ depending only on α .

The following result is the analogue of Theorem 3.2 for the case $\alpha \geq 1$.

Theorem 3.3 (Bounds Scaling with Variance – the Case $\alpha \geq 1$). *If X_1, \dots, X_n are independent mean zero random variables with $\|X_i\|_{\psi_\alpha} < \infty$ for all $1 \leq i \leq n$ and some $\alpha \geq 1$, then*

$$\left\| \sum_{i=1}^n X_i \right\|_{\Psi_{1, L_n(\alpha)}} \leq 2e\sqrt{6} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{1/2},$$

with

$$L_n(\alpha) := \frac{4^{1/\alpha} C_\alpha}{2\sqrt{6}} (\log(n+1))^{1/\alpha} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{-1/2} \max_{1 \leq i \leq n} \|X_i\|_{\psi_\alpha},$$

for some constant $C_\alpha > 0$ depending only on α .

Optimality of Theorems 3.2 and 3.3. Theorem 3.3 proves a bound on the $\Psi_{1,L_n(\alpha)}$ -norm irrespective of how light-tailed the initial random variables are (or in other words, how large $\alpha > 1$ is). Observe that this result reduces to the usual Bernstein’s inequality for bounded random variables by taking $\alpha = \infty$. Bennet’s inequality, which is a slight improvement of Bernstein’s inequality (Wellner, 2017), is known to be optimal for bounded random variables, as shown in Major (2005, Example 2.4). In light of this, it seems not possible to prove Theorem 3.3 for a $\Psi_{\alpha,L}$ -norm with $\alpha > 1$ as long as the bound is needed in terms of the variance. Note further that even though the result uses the $\Psi_{1,L}$ -norm, the parameter L behaves as $(\log n)^{1/\alpha}/\sqrt{n}$ with the exponent of $\log n$ being $1/\alpha$ instead of 1. So, this result cannot be obtained by simply applying Theorem 3.2 with $\alpha = 1$.

Section 2.2 of Adamczak (2008) provides a counterexample proving that it is not possible to replace the factor $(\log(n+1))^{1/\alpha}$ by anything of smaller order with only the hypothesis of $\|X_i\|_{\psi_\alpha} < \infty$ if the norm bound is desired to be in terms of the variance itself. Formally, if we assume a bound of the form

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i\right| \geq C \sqrt{t \sum_{i=1}^n \mathbb{E}[X_i^2] + Ct^{1/\alpha^*}(\log n)^u}\right) \leq 3e^{-t} \quad \text{for all } t \geq 0,$$

holds true with some $u \geq 0$ for all independent mean zero random variables X_1, \dots, X_n satisfying $\|X_i\|_{\psi_\alpha} \leq 1$, then $u \geq 1/\alpha$. This, again, should be understood in the minimax sense: for the result to hold for all distributions of X_1, \dots, X_n , then u must be at least $1/\alpha$. This follows from Section 2.2 of Adamczak (2008) by considering (as $r \rightarrow \infty$) i.i.d. random variables $X_1 = \varepsilon_1 Y_1, \dots, X_n = \varepsilon_n Y_n$ with $\mathbb{P}(Y_i = r^{1/\alpha}) = e^{-r} = 1 - \mathbb{P}(Y_i = 0)$ and $\varepsilon_1, \dots, \varepsilon_n$ are independent Rademacher random variables; we refer the reader to Adamczak (2008) for more details. Furthermore, Theorems 3.2 and 3.3 can be considered optimal in light of the large deviation results from Bakhshizadeh et al. (2020, Section III).

The main advantage of Theorems 3.2 and 3.3 over Theorem 3.1 is the appearance of the variance in the bound, as opposed to the $\|\cdot\|_{\psi_\alpha}$ norm, at the cost of the log factor in $L_n(\alpha)$ (which also explains the gain in the logarithmic factor mentioned after Theorem 8 of van de Geer and Lederer (2013)). This distinction can impact the convergence rate if $\mathbb{E}(X_i^2)$ is of much smaller order than $\|X_i\|_{\psi_\alpha}^2$; see Remark 3.3 for an example involving kernel smoothing estimators where this is indeed the case. Once again, we stress here that Theorem 3.1 can lead to a better tail bound if one does not care about dependence of the tail probability of the sum on the variance or if $\sqrt{\text{var}(X_i)}$ and $\|X_i\|_{\psi_\alpha}$ are of the same order.

Most of our examples in Section 4 involve the maximum of many averages. For this reason, we present a generally useful tail bound result for such maximums explicitly as a theorem below, although it is in fact a simple corollary of Theorems 3.2 and 3.3 (depending on whether $\alpha \leq 1$ or $\alpha \geq 1$). For a vector $v \in \mathbb{R}^q$, let $\|v\|_\infty$ denote $\max_{1 \leq j \leq q} |v(j)|$.

Theorem 3.4 (Tail Bounds for Maximums using Theorems 3.2–3.3). *Suppose X_1, \dots, X_n are independent mean zero random vectors in \mathbb{R}^q , for any $q \geq 1$, such that for some $\alpha > 0$ and $K_{n,q} > 0$,*

$$\max_{1 \leq i \leq n} \max_{1 \leq j \leq q} \|X_i(j)\|_{\psi_\alpha} \leq K_{n,q}, \quad \text{and define } \Gamma_{n,q} := \max_{1 \leq j \leq q} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^2(j)].$$

Then for any $t \geq 0$, with probability at least $1 - 3e^{-t}$,

$$\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_{\infty} \leq 7 \sqrt{\frac{\Gamma_{n,q}(t + \log q)}{n}} + \frac{C_{\alpha} K_{n,q} (\log(2n))^{1/\alpha} (t + \log q)^{1/\alpha^*}}{n},$$

where $\alpha^* := \min\{\alpha, 1\}$ and $C_{\alpha} > 0$ is some constant depending only on α .

Remark 3.2 (Comparison with Existing Maximal Inequalities). One of the most important conclusions of Theorem 3.4 is a bound on the expectation of the maximum, which are usually referred to as *maximal inequalities*. In particular, Theorem 3.4 yields

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_{\infty} \right] \leq C_1 \sqrt{\frac{\Gamma_{n,q} \log(eq)}{n}} + C_2(\alpha, K_{n,q}) \frac{(\log(2n))^{1/\alpha} (\log(eq))^{1/\alpha^*}}{n}, \quad (3.2)$$

for a universal constant $C_1 \geq 0$ and a constant $C_2(\alpha, K_{n,q}) \geq 0$ is a constant depending only on α and $K_{n,q}$. (Here $e \approx 2.71$ represents the natural logarithm constant.) This bound compares favorably with the existing maximal inequalities applicable for this case. For instance, Lemma E.1 of Chernozhukov et al. (2017) yields the bound

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_{\infty} \right] \leq C'_1 \sqrt{\frac{\Gamma_{n,q} \log(eq)}{n}} + C'_2(\alpha, K_{n,q}) \frac{(\log(eqn))^{1/\alpha} \log(eq)}{n}, \quad (3.3)$$

where the constants C'_1 and $C'_2(\alpha, K_{n,q})$ are similar to $C_1, C_2(\alpha, K_{n,q})$ in (3.2); see also Lemmas 3.4.2 and 3.4.3 of van der Vaart and Wellner (1996) for similar maximal inequalities. (3.3) is the best possible inequality obtained from Lemma E.1 of Chernozhukov et al. (2017) since the quantity $\sqrt{\mathbb{E}[M^2]}$ (in the referenced paper) can be bounded only by $(\log(eqn))^{1/\alpha}$ under the assumption of Theorem 3.4. Comparing (3.2) and (3.3), we note that the requirements for the average to converge to zero in the respective displays are given by:

$$\max\{\log(eq), (\log(2n))^{1/\alpha} (\log(eq))^{1/\alpha^*}\} = o(n) \quad \text{for (3.2),}$$

and

$$\max\{\log(eq), (\log(eqn))^{1/\alpha} (\log(eq))\} = o(n) \quad \text{for (3.3).}$$

The former is *strictly better* than the latter, especially if $\log(eq) = O(n^{\gamma})$ for some γ . These two conditions match only when the random vectors are uniformly bounded vis-a-vis $\alpha = \infty$. Lemma E.1 of Chernozhukov et al. (2017) is improved by Proposition B.1 in Kuchibhotla and Patra (2019) which, in fact, is built on an earlier version of the current paper. \diamond

Remark 3.3 (Tail Bounds for Linear Kernel Averages: An Illustration of Theorem 3.4). An important illustration of some of the main features of our results is in the derivation of (pointwise) deviation bounds for linear kernel average estimators (LKAEs) involving sub-Weibull variables. Such estimators are encountered in kernel smoothing based methods for non-parametric regression and density estimation.

Let $\{(Y_i, X_i) : i = 1, \dots, n\}$ denote n i.i.d. realizations of a random vector (Y, X) having finite second moments, where $Y \in \mathbb{R}$ and $X \in \mathbb{R}^p$. Assume for simplicity that X has a Lebesgue density

$f(\cdot)$. Let $m(x) := \mathbb{E}(Y|X = x)$ and $\psi(x) := m(x)f(x)$. Let $K(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$ denote any kernel function (e.g. the Gaussian kernel on \mathbb{R}^p). Consider the following LKAE of $\psi(x)$, given by

$$\widehat{\psi}(x) := \frac{1}{nh^p} \sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h}\right), \quad \text{where } h \equiv h_n > 0 \text{ is the bandwidth.}$$

Suppose $\|Y\|_{\psi_\alpha} \leq C_Y$ for some $\alpha, C_Y > 0$ and $g(x) := \mathbb{E}(Y^2|X = x)f(x)$ is bounded, i.e. $0 \leq g(x) \leq M_Y$ for all x , for some constant $M_Y \geq 0$. Assume further that $K(\cdot)$ is bounded and square integrable, i.e. for some constants $C_K, R_K \geq 0$, $|K(x)| \leq C_K$ for all x and $\int_{\mathbb{R}^p} K^2(x)dx \leq R_K$. Then, for any fixed $x \in \mathbb{R}^p$ and any $t \geq 0$, we have with probability at least $1 - 3e^{-t}$,

$$\left| \widehat{\psi}(x) - \mathbb{E}\{\widehat{\psi}(x)\} \right| \leq \frac{7\Gamma_{Y,K}}{\sqrt{nh^p}} \sqrt{t} + \frac{C_\alpha \Upsilon_{Y,K} (\log(2n))^{1/\alpha}}{nh^p} t^{1/\alpha^*}, \quad (3.4)$$

where $\Gamma_{Y,K} := (M_Y R_K)^{\frac{1}{2}}$, $\Upsilon_{Y,K} := C_Y C_K$, $\alpha^* := \min\{\alpha, 1\}$ and $C_\alpha > 0$ is some constant depending only on α . (3.4) provides a ready-to-use deviation bound for sub-Weibull LKAEs with a convergence rate of $(nh^p)^{-1/2}$ for any $\alpha > 0$, assuming $nh^p \rightarrow \infty$ as $n \rightarrow \infty$. Note that to extract this (sharp) rate, it is *necessary* to exploit that $h^{-p} Y K\{(X - x)/h\}$ has a variance of much *smaller* order than its squared $\|\cdot\|_{\psi_\alpha}$ norm. The proof of (3.4) is given in Appendix D. Under standard smoothness conditions and a q -th order kernel $K(\cdot)$, for some $q \geq 2$, it can be shown that $|\mathbb{E}\{\widehat{\psi}(x)\} - \psi(x)| \leq O(h^q)$ uniformly in x (see, for instance, Hansen (2008) and references therein) and hence, a tail bound for $|\widehat{\psi}(x) - \psi(x)|$ can also be obtained. The result provided here is mostly for illustration purposes and can possibly be extended in several directions; see Section 5 for further discussion. \diamond

Orlicz Norms of Products of Random Variables. In all our results, the random variables are only required to be sub-Weibull of some order $\alpha > 0$. In many applications, one may need to deal with products of two or more such sub-Weibull variables. The following result (proved as Proposition D.2 in Appendix D) provides a Hölder type inequality establishing a bound on the $\|\cdot\|_{\psi_\alpha}$ norm of such product variables. The two examples mentioned in the introduction can also be easily dealt with using this result. If W_i , $1 \leq i \leq k$, are (possibly dependent) random variables satisfying $\|W_i\|_{\psi_{\alpha_i}} < \infty$ for some $\alpha_i > 0$, then

$$\left\| \prod_{i=1}^k W_i \right\|_{\psi_\beta} \leq \prod_{i=1}^k \|W_i\|_{\psi_{\alpha_i}} \quad \text{where} \quad \frac{1}{\beta} := \sum_{i=1}^k \frac{1}{\alpha_i}. \quad (3.5)$$

See also Lemma 2.7.7 of Vershynin (2018) for a similar result.

Tail Bounds for Powers of Sub-Gaussians. As a simple application of the above discussion on products, coupled with our general results in this section, one can obtain tail bounds for powers of sub-Gaussians which are often useful in practice. For example, consider $X_i = \varepsilon_i |G_i|^r$, $1 \leq i \leq n$ with $r \geq 0$, Rademacher ε_i and sub-Gaussian G_i . Then, using (3.5), X_i 's satisfy $\|X_i\|_{\psi_{2/r}} \leq \mathfrak{C} < \infty$ for some constant \mathfrak{C} whenever $\|G_i\|_{\psi_2} \leq \mathfrak{C}$. For such random variables, one can apply Theorem 3.1 or Theorems 3.2–3.3 to obtain a tail bound. Note that $\|X_i\|_{\psi_{2/r}} = \|G_i\|_{\psi_2} \leq \mathfrak{C}$ and $\mathbb{E}[X_i^2] = \mathbb{E}[G_i^{2r}] \leq \mathfrak{C}r^r$. This implies that the standard deviation and the $\psi_{2/r}$ -norm of the random variables

are of the same order if r is treated as a constant and the $2r$ -th moment of G_i is of the same order as $\|G_i\|_{\psi_2}^{2r}$. Then, Theorem 3.1 implies

$$\mathbb{P}\left(\left|\sum_{i=1}^n \varepsilon_i |G_i|^r\right| \geq \mathfrak{C}_r (nt)^{1/2} + \mathfrak{C}_r t^{r/2} n^{(1-r/2)_+}\right) \leq 2e^{-t} \quad \text{for all } t \geq 0. \quad (3.6)$$

Here, \mathfrak{C}_r is a constant depending only on r and $(u)_+ = \max\{u, 0\}$. In this case, Theorems 3.2–3.3 may yield a sub-optimal result because it does not account for the fact that the standard deviation and the $\psi_{2/r}$ -norm are of the same order. If, however, the $2r$ -th moments of G_i 's are not of the same order as $\|G_i\|_{\psi_2}^{2r}$, then Theorem 3.2 or 3.3 (as the case may be) yields a *better* tail bound. Finally, note that we consider the symmetrized form $\varepsilon_i |G_i|^r$ involving the Rademacher ε_i 's here to ensure the random variables are all mean zero. A similar bound as (3.6) continues to hold if $\varepsilon_i |G_i|^r$ is replaced by $|G_i|^r - \mathbb{E}(|G_i|^r)$. Furthermore, the form $|G_i|^r$ with absolute value is considered to ensure it is well defined for any $r \geq 0$. A similar bound as (3.6) continues to hold for $G_i^r - \mathbb{E}(G_i^r)$ whenever r is any positive integer.

4 Applications in High Dimensional Statistics

Outline. In this section, we study in detail the four fundamental statistical applications mentioned in the introduction, through Sections 4.1–4.4. Below we first provide a *high-level organization* – in terms of the problems considered in each sub-section, and pointers to the corresponding main results and key discussions. A more detailed outline for each is provided within the respective sub-sections themselves.

1. Section 4.1 – *Covariance matrix estimation in maximum elementwise norm.* (Main results: Theorems 4.1 and Theorem 4.2 (in Section 4.1.1); key discussions: Remarks 4.1 and 4.4.)
2. Section 4.2 – *Covariance matrix estimation in maximum k -sub-matrix operator norm, and the restricted isometry property (RIP).* (Main result: Theorem 4.3; key discussions: Remarks 4.5–4.9.)
3. Section 4.3 – *The restricted eigenvalue (RE) and restricted strong convexity (RSC) conditions.* (Main result: Theorem 4.4; key discussions: Remarks 4.12–4.13, as well as the results and associated discussions in Section 4.3.1 on verification of the RE condition for general sub-Weibulls.)
4. Section 4.4 – *High dimensional linear regression via Lasso.* (Main results: Theorems 4.5 and 4.6; key discussions: Remarks 4.14 and 4.15, as well as the general oracle inequality in Remark 4.16.)

A Discussion on Sub-Weibull Random Vectors: Joint vs. Marginal. Before proceeding to these applications, we provide a brief discussion that suggests that for random vectors the joint sub-Weibull property (Definition 2.4), although commonly adopted in the literature (especially for the sub-Gaussian case; e.g., see Vershynin (2018)), is a much more restrictive assumption than the marginal one (Definition 2.5). A careful examination of the joint sub-Weibull property implies an “almost independence” restriction on the coordinates for a dimension-free bound on the joint sub-Weibull norm.

As a simple (albeit a bit extreme) example, consider the random vector $X \in \mathbb{R}^q$ where all the coordinates are exactly the same $X(1) = \dots = X(q)$. In this case, it is clear that

$$\|X\|_{J,\psi_\alpha} = \sup_{\theta \in \mathbb{R}^q, \|\theta\|_2=1} \|\theta\|_1 \|X(1)\|_{\psi_\alpha} = \sqrt{q} \|X(1)\|_{\psi_\alpha}. \quad (4.1)$$

Although this is a pathological example, it shows that if the coordinates of X are highly dependent, then the random vector *cannot* have a “small” joint sub-Weibull norm; see Section 3.4 of [Vershynin \(2018\)](#) for a similar discussion. For all the high dimensional applications we consider, the (polynomial) dependence on the dimension in (4.1) can render the rates useless. Note that even though a Gaussian vector $X \in \mathbb{R}^q$ is jointly sub-Gaussian, $\|X\|_{J,\psi_2}$ will depend on the maximum eigenvalue of $\Sigma := \text{Cov}(X)$, which may not be dimension-free if X has correlated components (e.g., if Σ is an equicorrelation matrix).

The “almost independence” restriction implied by the joint sub-Weibull property may not necessarily be satisfied in practice and it is also hard to find results for high dimensional statistical methods in the literature under *marginal* sub-Gaussian/sub-exponential tails. So, we consider both the marginal and the joint sub-Weibull assumptions in deriving the tail bounds as well as the rates of convergence in all the following statistical applications.

4.1 Covariance Matrix Estimation: Maximum Elementwise Norm

Outline. In this section, we consider concentration properties of covariance matrices for sub-Weibulls under the maximum elementwise norm, which plays a crucial role in various high dimensional inference problems as well as in bootstrap. Our main result here is Theorem 4.1 (along with Theorem 4.2 in Section 4.1.1 that further allows for data dependent centering). It proves a finite sample tail bound under the assumption of only *marginally* sub-Weibull (α) ingredient random vectors. Remark 4.1 provides useful discussions on its implications and shows, in particular, the rate of convergence to be $\sqrt{\log p/n}$ if $\log p = o(n^{\alpha/(4-\alpha)})$. This rate can be easily shown to be optimal in case the random vectors are standard multivariate Gaussian. Finally, we discuss applications of these results in sparse covariance matrix estimation (Remark 4.3) and in establishing consistency of bootstrap (Remark 4.4) for (high dimensional) marginally sub-Weibull random vectors. Below we introduce the problem setup, followed by our results.

Suppose X_1, \dots, X_n are independent random vectors in \mathbb{R}^p . Define the (gram) matrices

$$\hat{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n X_i X_i^\top \quad \text{and} \quad \Sigma_n := \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X_i X_i^\top]. \quad (4.2)$$

Note that $\hat{\Sigma}_n$ is unbiased for Σ_n . Assuming that X_i 's have mean 0, Σ_n is also the covariance matrix of the \sqrt{n} -scaled sample mean, $\sqrt{n}\bar{X}_n$, and $\hat{\Sigma}_n$ is a natural estimator of Σ_n .

Define the *elementwise maximum norm* of $\hat{\Sigma}_n - \Sigma_n$ as

$$\Delta_n := \|\hat{\Sigma}_n - \Sigma_n\|_\infty = \max_{1 \leq j \leq k \leq p} \left| \frac{1}{n} \sum_{i=1}^n \{X_i(j)X_i(k) - \mathbb{E} [X_i(j)X_i(k)]\} \right|.$$

As shown in Remark 4.1 of [Chernozhukov et al. \(2017\)](#), it is necessary to control Δ_n , the elementwise maximum norm between the empirical and population covariance matrices, to establish consistency of the multiplier bootstrap.

Theorem 4.1 below (proved in Appendix E.1), the main result of this section, controls Δ_n under only a *marginal* sub-Weibull (α) assumption. Only the case $\alpha \leq 2$ is considered here (the case $\alpha > 2$ can be derived similarly from Theorems 3.3 and 3.4). Recall Definition 2.5.

Theorem 4.1. *Let X_1, \dots, X_n be independent marginally sub-Weibull random vectors in \mathbb{R}^p satisfying*

$$\max_{1 \leq i \leq n} \|X_i\|_{M, \psi_\alpha} \leq K_{n,p} < \infty \quad \text{for some } 0 < \alpha \leq 2. \quad (4.3)$$

Fix $n, p \geq 1$. Then for any $t \geq 0$, with probability at least $1 - 3e^{-t}$,

$$\Delta_n \leq 7A_{n,p} \sqrt{\frac{t + 2 \log p}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + 2 \log p)^{2/\alpha}}{n},$$

where $C_\alpha > 0$ is a constant depending only on α , and $A_{n,p}^2$ is given by

$$A_{n,p}^2 := \max_{1 \leq j \leq k \leq p} \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i(j)X_i(k)).$$

Remark 4.1 (Rate of Convergence). Firstly, we reiterate that Theorem 4.1 *only* requires a marginal sub-Weibull assumption on the X_i 's, as in (4.3). Next, it is clear from Theorem 4.1 that the rate of convergence of Δ_n is given by

$$\Delta_n = O_p \left(\max \left\{ A_{n,p} \sqrt{\frac{\log p}{n}}, K_{n,p}^2 \frac{(\log n)^{2/\alpha} (\log p)^{2/\alpha}}{n} \right\} \right).$$

Thus if $(\log p)^{2/\alpha-1/2} = o(\sqrt{n}(\log n)^{-2/\alpha})$, then $\Delta_n = O_p \left(A_{n,p} \sqrt{\log p/n} \right)$. It is easy to verify under assumption (4.3) that $A_{n,p} \leq C_\alpha K_{n,p}^2$; see Proposition 2.5.2 of Vershynin (2018) for a proof. Note that if $\alpha = 2$, i.e. X_i 's are marginally sub-Gaussian, then the (known) rate of convergence is $\sqrt{\log p/n}$. Thus, the key implication of the above calculations is that the *rate of convergence can match that of the sub-Gaussian case* for a wide range of $\alpha > 0$. This is the main importance of the tail bounds stated in Section 3 and the *same phenomenon is observed in all subsequent applications in Sections 4.2–4.4 too*. Also, it is clear that the same result continues to hold under a (stronger) joint sub-Weibull assumption. \diamond

Remark 4.2 (Application to Coupling Inequality). The quantity Δ_n also appears in a coupling inequality for the maximum of a sum of random vectors. The coupling inequality refers to bounding

$$\left| \max_{1 \leq j \leq p} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i(j) - \max_{1 \leq j \leq p} \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i(j) \right|,$$

where $X_i \in \mathbb{R}^p$ are mean zero and $Z_i \sim N_p(0, \mathbb{E}[X_i X_i^\top])$ constructed on the same probability space as X_i 's. For this quantity to converge in probability to zero, Theorem 4.1 of Chernozhukov et al. (2014) requires Δ_n to converge to zero, among other terms. \diamond

4.1.1 Gram Matrix to Covariance Matrix (Accounting for Centering)

The quantity Δ_n only measures the difference between the sample and the population *gram matrices* that involve the *uncentered* X_i 's, and this is important in applications involving linear regression since only the gram matrix directly appears there and not the covariance matrix. In some applications, however, it is of interest to deal with the *covariance matrices*

$$\begin{aligned}\hat{\Sigma}_n^* &:= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n) (X_i - \bar{X}_n)^\top, \quad \text{and} \\ \Sigma_n^* &:= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[(X_i - \bar{\mu}_n) (X_i - \bar{\mu}_n)^\top \right],\end{aligned}\tag{4.4}$$

where $\bar{X}_n := \sum_{i=1}^n X_i/n$ and $\bar{\mu}_n := \mathbb{E}[\bar{X}_n] = \sum_{i=1}^n \mathbb{E}[X_i]/n$. Note, however, that Σ_n^* is *not* the variance of \bar{X}_n unless $\mu_i = \bar{\mu}_n$ for all i . Define the maximum elementwise norm error between the sample and population covariance matrices $\hat{\Sigma}_n^*$ and Σ_n^* , respectively, as

$$\Delta_n^* := \|\hat{\Sigma}_n^* - \Sigma_n^*\|_\infty.$$

Theorems 4.1 and 3.4 together imply the following result (proved in Appendix E.1) for Δ_n^* .

Theorem 4.2. *Under the setting of Theorem 4.1, for any $t \geq 0$, with probability at least $1 - 6e^{-t}$,*

$$\Delta_n^* \leq 7A_{n,p}^* \sqrt{\frac{t + 2 \log p}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + 2 \log p)^{2/\alpha}}{n},$$

where

$$A_{n,p}^* := \max_{1 \leq j \leq k \leq p} \left(\frac{1}{n} \sum_{i=1}^n \text{Var} [(X_i(j) - \bar{\mu}_n(j))(X_i(k) - \bar{\mu}_n(k))] \right)^{1/2}.$$

In comparison to Theorem 4.1 which applied to gram matrices, the only change with covariance matrices is the replacement of $A_{n,p}$ therein with $A_{n,p}^*$ as above.

Remark 4.3 (Applications in Sparse Covariance Matrix Estimation). The basic technique of sparse covariance matrix estimation is thresholding. For simplicity, consider the case of identically distributed random vectors. Recall the definition of the usual covariance matrix $\hat{\Sigma}_n^*$ from (4.4) and define for $\lambda > 0$, the matrix $\check{\Sigma}_{n,\lambda}$ given by

$$\check{\Sigma}_{n,\lambda}(j, k) := \begin{cases} \hat{\Sigma}_n^*(j, k), & \text{if } |\hat{\Sigma}_n^*(j, k)| \geq \lambda, \\ 0, & \text{otherwise,} \end{cases}$$

for $1 \leq j \leq k \leq p$. This estimator essentially sets to zero those elements of $\hat{\Sigma}_n^*$ that are ‘‘small’’. This is referred to sometimes as universal hard thresholding since λ does not depend on (j, k) . The parameter λ is called the thresholding parameter. It is easy to verify that

$$\mathbb{P} \left(\Sigma_n^*(j, k) = 0 \text{ and } \check{\Sigma}_{n,\lambda}(j, k) \neq 0 \text{ for some } j, k \right) \leq \mathbb{P}(\Delta_n^* > \lambda).$$

So, the right cut-off λ for consistent support recovery would be of the same order as the rate of convergence of Δ_n^* which is $\sqrt{\log p/n}$, as shown in Theorem 4.2 (under additional conditions, as

in Remark 4.1). So, for a wide range of α , the cut-off used for Gaussians works for marginally sub-Weibull random vectors too. For a more careful study of the properties of $\check{\Sigma}_{n,\lambda}$ in terms of the operator norm and extensions to weakly sparse matrices, see [Bickel and Levina \(2008\)](#), [Cai and Liu \(2011\)](#) and [Fan et al. \(2016\)](#). As can be seen from the analysis there, a result similar to Theorem 4.2 plays a key role. It should be noted here that most of the literature about covariance matrix estimation is based on a joint sub-Gaussian assumption on the ingredient random vectors. Our setting above is clearly more general. \diamond

Remark 4.4 (Bootstrap Consistency). From Remark 4.1 and Theorem 4.2 of [Chernozhukov et al. \(2017\)](#), it follows that the consistency of either the multiplier bootstrap or Efron’s empirical bootstrap for high dimensional averages requires the convergence of Δ_n^* to zero. In fact, the multiplier bootstrap error is bounded by a multiple of $(\Delta_n^*)^{1/3}$. Hence, our results in this section prove the bootstrap consistency under weaker tail assumptions. \diamond

4.2 Covariance Matrix Estimation: Maximum k-Sub-Matrix Operator Norm

This section focuses on estimation of covariance matrices of sub-Weibull random vectors under the so-called sub-matrix operator norm. In the previous sub-section, a bound on the elementwise maximum norm for such covariance matrices was provided. It is clear that the maximum norm only deals with the elements of the matrix. In many applications and practical data exploration, it is of much more importance to study functionals of the covariance matrix such as the eigenvalues and eigenvectors. A key ingredient in studying these functionals is consistency of the covariance matrix in the operator norm.

As expected, if the dimension of the random vectors X_i is larger than the sample size n , then the covariance matrix is *not* consistent in the operator norm. Also, in high-dimensions it is a common practice to select a subset of “significant” group of coordinates of X_i ’s and explore the properties of that subset. Motivated by this discussion, we study the *maximum k-sparse sub-matrix operator norm* of the gram matrix, for any $1 \leq k \leq p$. This norm is also of importance in high dimensional linear regression due to its connections to the *restricted isometry property* (RIP) ([Candes and Tao, 2007](#)) and the restricted eigenvalue (RE) condition ([Bickel et al., 2009](#)). Define, for $k \leq p$,

$$\text{RIP}_n(k) := \sup_{\substack{\theta \in \mathbb{R}^p, \\ \|\theta\|_0 \leq k, \|\theta\|_2 \leq 1}} |\theta^\top (\hat{\Sigma}_n - \Sigma_n) \theta|, \quad (4.5)$$

with $\hat{\Sigma}_n$ and Σ_n as defined in (4.2). Here, $\|\theta\|_0$ denotes the number of non-zero entries (i.e. the *sparsity*) of θ . Note further that $\text{RIP}_n(k)$ is actually a norm for $k \geq 2$.

The quantity $\text{RIP}_n(k)$ also plays an important role in post-Lasso linear regression asymptotics (see condition RSE(m) in [Belloni and Chernozhukov \(2013\)](#)) and more generally, in post-selection inference (see [Kuchibhotla et al. \(2018\)](#) for details). This norm was possibly first studied (with Σ_n being the identity matrix) in [Rudelson and Vershynin \(2008\)](#) under the assumption of marginally bounded random vectors or equivalently, assumption (4.3) with $\alpha = \infty$. Also see Appendix C of [Belloni and Chernozhukov \(2013\)](#) for similar results.

An Easier but Sub-Optimal Bound for $\text{RIP}_n(k)$. Our main results on tail bounds for $\text{RIP}_n(k)$ are presented in Theorem 4.3. However, using the results of Section 4.1, *an easier but generally sub-optimal bound* on $\text{RIP}_n(k)$ may also be obtained which we present below for the sake of com-

pleteness. Note that

$$\text{RIP}_n(k) \leq \left(\sup_{\|\theta\|_0 \leq k, \|\theta\|_2 \leq 1} \|\theta\|_1^2 \right) \|\hat{\Sigma}_n - \Sigma_n\|_\infty \leq k \|\hat{\Sigma}_n - \Sigma_n\|_\infty.$$

This is a deterministic inequality and using the bounds on Δ_n derived previously in Section 4.1, it is easy to derive bounds for $\text{RIP}_n(k)$. For simplicity, we only present here an expectation bound instead of general tail bounds (or moment bounds) for $\text{RIP}_n(k)$. Under the hypothesis of Theorem 4.1 in Section 4.1, we have

$$\mathbb{E}[\text{RIP}_n(k)] \leq C_\alpha \left(A_{n,p} k \sqrt{\frac{\log p}{n}} + K_{n,p}^2 \frac{k(\log p \log(2n))^{2/\alpha}}{n} \right), \quad (4.6)$$

for some constant $C_\alpha > 0$ depending only on α . This bound provides the rate of $k\sqrt{\log p/n}$ only for $\text{RIP}_n(k)$ using the arguments of Remark 4.1. Note that this is derived *only* under a marginal ψ_α -bound, and the factor k here is, in a sense, optimal under the marginal ψ_α -bound hypothesis as can be seen from the pathological example discussed before Section 4.1. (For this example, the factor $\sqrt{\log p}$ disappears from the rate.) A bound alternative to (4.6) can be derived under the hypothesis of a joint ψ_α assumption on X_i . Under this joint hypothesis, the dominating term becomes $\sqrt{k \log p/n}$ which is the more familiar rate.

Main Result (Outline). We next derive a bound on $\text{RIP}_n(k)$ in a unified way, using a different approach, that always presents the dominating term of the (optimal) order $\sqrt{k \log p/n}$ (upto a distributional constant factor) under either of these assumptions (i.e., marginal or joint sub-Weibull). This is presented in Theorem 4.3 below (proved in Appendix E.2), the main result of this section. Once again, we only present the result for $0 < \alpha \leq 2$ and a similar result for $\alpha > 2$ can be derived using Theorem 3.3. The result is presented in two parts: (a) *marginal* case and (b) *joint* case. The implications, including the behavior of the bound and its rate of convergence, as well as the sample complexity requirements under either cases, are discussed in detail in Remarks 4.5–4.8, followed by a thorough comparison with the existing literature on RIP in Remark 4.9. Overall, to our knowledge, Theorem 4.3(a) is the first result on $\text{RIP}_n(k)$ for the *marginal* case, while for the *joint* case, Theorem (4.3)(b) matches existing (and optimal) results for the special case of sub-Gaussians (i.e., $\alpha = 2$), and also extends these to general sub-Weibulls.

Theorem 4.3 (Unified Bounds for RIP). *Let X_1, \dots, X_n be independent random vectors in \mathbb{R}^p . Define*

$$\Theta_k := \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq k, \|\theta\|_2 \leq 1\} \quad \text{and} \quad \Upsilon_{n,k} := \sup_{\theta \in \Theta_k} \frac{1}{n} \sum_{i=1}^n \text{Var} \left[\left(X_i^\top \theta \right)^2 \right].$$

Fix $0 < \alpha \leq 2$. Then, for every $1 \leq k \leq p$, the following bounds hold true for $\text{RIP}_n(k)$ as in (4.5):

- (a) (*Marginal Sub-Weibull Case*). *If $\|X_i\|_{M, \psi_\alpha} \leq K_{n,p}$ for all $1 \leq i \leq n$, then for any $t > 0$, with probability at least $1 - 3e^{-t}$,*

$$\begin{aligned} \text{RIP}_n(k) \leq & 14 \sqrt{\frac{\Upsilon_{n,k}(t + k \log(36p/k))}{n}} \\ & + \frac{C_\alpha K_{n,p}^2 k (\log(2n))^{2/\alpha} (t + k \log(36p/k))^{2/\alpha}}{n}. \end{aligned} \quad (4.7)$$

(b) (*Joint Sub-Weibull Case*). If $\|X_i\|_{J,\psi_\alpha} \leq K_{n,p}$ for all $1 \leq i \leq n$, then for any $t > 0$, with probability at least $1 - 3e^{-t}$,

$$\begin{aligned} \text{RIP}_n(k) \leq & 14 \sqrt{\frac{\Upsilon_{n,k}(t + k \log(36p/k))}{n}} \\ & + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + k \log(36p/k))^{2/\alpha}}{n}. \end{aligned} \quad (4.8)$$

Here, in both cases, $C_\alpha > 0$ represents a constant depending only on α .

Comparing between the two bounds from parts (a) and (b) above, the only difference is an extra factor of k in the second term (for part (a) – which uses the weaker assumption of marginal sub-Weibull) which is usually of lower order than the first term.

Remark 4.5 (Rate of Convergence). The bounds (4.7) and (4.8) – obtained for the marginal and joint sub-Weibull cases, respectively – both provide the same rate of $(\Upsilon_{n,k} k \log p/n)^{1/2}$ for a wide range of k following the arguments of Remark 4.1, and this is actually what is expected from the central limit theorem as well. More details on the sample complexity requirements for both results are provided in Remarks 4.7 and 4.8, for the joint and marginal cases, respectively. It is also worth mentioning that to the best of our knowledge, Theorem 4.3(a) is the *first* such result in the literature for the *marginal* case (that too for general sub-Weibulls), while for the *joint* case, Theorem (4.3)(b) *matches* existing results for the special case of sub-Gaussians (i.e., $\alpha = 2$), while also *extending* them to general sub-Weibulls. Further discussions on comparison with the existing literature are provided in Remark 4.9. \diamond

Remark 4.6 (Growth of $\Upsilon_{n,k}$). The leading term in the bounds of Theorem 4.3 depends on $\Upsilon_{n,k}$ which relates to the fourth moment of linear combinations. Such quantities have also appeared in several other problems, including likelihood methods with diverging number of parameters (Portnoy, 1988), sub-Gaussian estimation of means (Joly et al., 2017), tail bounds for lower eigenvalues of covariance matrices (Oliveira, 2013) and verification of so-called small-ball conditions (Lecué and Mendelson, 2017). In some of these works, the fourth moment of linear combinations is assumed to be bounded by the square of the second moment. Such an assumption coupled with a bounded operator norm of Σ_n implies that $\Upsilon_{n,k}$ is of constant order. In general, $\Upsilon_{n,k}$ can grow with k and it is not clear the rate at which it can grow for arbitrary distributions. However, under a joint sub-Weibull assumption as in part (b) of Theorem 4.3, it is at most a constant multiple of $K_{n,p}^4$. \diamond

Remark 4.7 (Conditions for Theorem 4.3(b) – the *Joint Sub-Weibull Case*). For the joint sub-Weibull case, the bound (4.8) for $\text{RIP}_n(k)$ converges to zero whenever

$$n \gg \max\{k \log(ep/k), k^{2/\alpha} (\log(2n))^{2/\alpha} (\log(ep/k))^{2/\alpha}\}. \quad (4.9)$$

In particular, for the special case of joint sub-Gaussian (i.e., $\alpha = 2$), the sample complexity requirement (4.9) simplifies to: $n \gg k \log(ep/k) \log(n)$. We clarify that the appearance of the $\log(n)$ factor here is due to our usage of Theorem 3.3 (and ultimately Theorem 3.4) and *can* be avoided if instead one directly uses Theorem 3.1 – this essentially relates to our earlier discussion on the optimality of Theorems 3.2–3.3 (see Section 3). It is worth noting that the sample complexity $n \gg k \log(ep/k) \log(n)$ *matches* (possibly upto a $\log(n)$ factor) the scaling requirements of most results known in the literature, including those of Candes and Tao (2005, 2007); Baraniuk et al. (2008) and Loh and Wainwright (2012, Appendix G.1), among several others. In fact, most settings

considered in the existing literature are included as special cases under our *joint* sub-Weibull setting for the choice $\alpha = 2$; see Remark 4.9 for further details. \diamond

Remark 4.8 (Conditions for Theorem 4.3(a) – the *Marginal* Sub-Weibull Case). Firstly, before we discuss the bound (4.7), we note that for the initial bound of $\text{RIP}_n(k)$ provided in (4.6), although the rate obtained there is generally sub-optimal, convergence to zero of the bound therein requires

$$n \gg \max\{k^2 \log p, k(\log p)^{2/\alpha}\}, \quad (4.10)$$

whenever $X_i \in \mathbb{R}^p$ are marginally sub-Weibull, i.e., satisfy $\|X_i\|_{M,\alpha} < \infty$. However, for Theorem 4.3(a), which also requires only a marginal sub-Weibull property and provides a bound with a much sharper (and optimal) rate, convergence to zero of $\text{RIP}_n(k)$ requires

$$n \gg \max\{k \log(ep/k), k^{1+2/\alpha}(\log(2n))^{2/\alpha}(\log(ep/k))^{2/\alpha}\}. \quad (4.11)$$

For α considerably smaller than 1, the requirement (4.11) for Theorem 4.3(a) thus appears more stringent in terms of k , compared to (4.10). This deficiency can be explained by the fact that the proof of Theorem 4.3 uses the bound $\|\max_{\theta \in \Theta_k} \theta^\top X_i\|_{\psi_\alpha} \leq (k \log(ep/k))^{1/\alpha}$ (for applying Theorem 3.4), but using $\max_{\theta \in \Theta_k} \theta^\top X_i \leq \sqrt{k} \|X_i\|_\infty$, we can get a sharper bound: $\|\max_{\theta \in \Theta_k} \theta^\top X_i\|_{\psi_\alpha} \leq K_{n,p} \sqrt{k} (\log(ep))^{1/\alpha}$. Formally, Lemma E.1 of Chernozhukov et al. (2017) implies that

$$\begin{aligned} \mathbb{E} \left[\sup_{\theta \in \Theta_k} \left| \frac{1}{n} \sum_{i=1}^n \{(\theta^\top X_i)^2 - \mathbb{E}[(\theta^\top X_i)^2]\} \right| \right] &\lesssim \sqrt{\frac{k \log(ep/k)}{n}} \sup_{\theta \in \Theta_k} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\theta^\top X_i)^4] \right)^{1/2} \\ &\quad + \frac{k \log(ep/k)}{n} \left(\mathbb{E} \left[\max_{1 \leq i \leq n} \sup_{\theta \in \Theta_k} |\theta^\top X_i|^4 \right] \right)^{1/2} \\ &\lesssim \sqrt{\frac{k \log(ep/k)}{n}} \sup_{\theta \in \Theta_k} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\theta^\top X_i)^4] \right)^{1/2} \\ &\quad + \frac{k^2 (\log(epn))^{1+2/\alpha}}{n}. \end{aligned} \quad (4.12)$$

The second inequality above follows from $\max_{\theta \in \Theta_k} \theta^\top X_i \leq \sqrt{k} \|X_i\|_\infty$ and the marginal sub-Weibull assumption. The right hand side of (4.12) converges to zero whenever

$$\max\{k \log(ep/k), k^2 (\log(epn))^{1+2/\alpha}\} = o(n). \quad (4.13)$$

(4.13) therefore improves our original requirement (4.11) for Theorem 4.3(a), and matches (4.10) upto a log factor. However, we do not follow this approach further, mostly because we want to present the statistical applications as direct corollaries of the results in Section 3. In any case, (4.10) at least shows $n \gg k^2 (\log p)^{2/\alpha}$ suffices for $\text{RIP}_n(k)$ to go to zero under only a marginal sub-Weibull assumption. \diamond

Remark 4.9 (Comparison of Theorem 4.3 with Existing Literature). Theorem 4.3 succinctly provides a unified set of results on $\text{RIP}_n(k)$ as in (4.5) under very general conditions – both in terms of the tail behavior (i.e., sub-Weibull) as well as its nature (marginal vs. joint). To the best of our knowledge, the results in Theorem 4.3(a) for the *marginal* case are the first such results in the literature obtained under a (much) weaker assumption than most in the existing literature

on verifying $\text{RIP}_n(k)$, and should therefore be of substantial interest in the future. Secondly, most of the literature on $\text{RIP}_n(k)$ has focused on specific cases of our *joint* sub-Weibull setting in Theorem 4.3(b), and our convergence rates as well as sample complexity requirements, as discussed in Remarks 4.5 and 4.7, match these results – these include the well known works of Candes and Tao (2005, 2007); Baraniuk et al. (2008); Loh and Wainwright (2012), among many others. Rudelson and Zhou (2013) provides a comprehensive review of the existing literature on verification of $\text{RIP}_n(k)$ as we consider; see in particular their discussion in Section I (pg. 3434) and the references cited therein. To our knowledge, apart from the obvious flexibility (and novelty) of allowing for the marginal case, our results also enjoy the benefits of extension to the general sub-weibull case (i.e., for a general α) *even* in the joint case, where most of the existing literature can be summarized as special cases in some form of our joint sub-Weibull setting with the choice $\alpha = 2$.

It is worth noting that there has certainly been some work on the joint sub-Weibull setting (for a general α) as well, but for a *different* RIP problem. This includes, in particular, the works of Adamczak et al. (2011) and Guédon et al. (2014, 2015). However, there are some important differences in their setting versus ours. Their definition of RIP, translated in our notation, is given by

$$\text{RIP}_n^*(k) := \sup_{\substack{\alpha \in \mathbb{R}^n, \\ \|\alpha\|_0 \leq k, \|\alpha\|_2 \leq 1}} \left| \frac{1}{p} \sum_{j=1}^p \left| \sum_{i=1}^n \alpha(i) X_i(j) \right|^2 - 1 \right|. \quad (4.14)$$

Here, as in Theorem 4.3, $X_i \in \mathbb{R}^p$ are independent random vectors. The main difference between our $\text{RIP}_n(k)$ in (4.5) and $\text{RIP}_n^*(k)$ above is that for the former, the sparse linear combination involved in (4.5) is being taken over coordinates of X_i 's, while for $\text{RIP}_n^*(k)$, the linear combination involved in (4.14) is over X_i 's themselves. Thus, the maximizing space in (4.14) is *not* our Θ_k ; in fact, it is not even a subspace of \mathbb{R}^p (rather it is in \mathbb{R}^n). Under the definition (4.14) of $\text{RIP}_n^*(k)$, Adamczak et al. (2011) proves that $n \gg k \log^{2/\alpha}(p/k)$ suffices for controlling $\text{RIP}_n^*(k)$ when X_i 's are jointly sub-Weibull(α) with $\alpha \in [1, 2]$ ($\|X_i\|_{J,\alpha} < \infty$) and Guédon et al. (2015) extends this result to the case $\alpha \in (0, 1]$. Although this condition is better in terms of dependence on k compared to our requirement (4.9) under the joint sub-Weibull case when $\alpha < 2$, it does *not* apply for our definition (4.5) of RIP. Further, we would also like to point out that these works require the joint sub-Weibull assumption, while we allow for the marginal case as well.

We remark here that the goal of Adamczak et al. (2011); Guédon et al. (2015) is to derive RIP constants for the reconstruction of sparse signals, and their definition of RIP as in (4.14) works for this purpose. However, they focus on the RIP of a very different type of matrices, ones with column spaces in \mathbb{R}^n , not \mathbb{R}^p , which makes their setting fundamentally *different* from ours, and their results not directly comparable to ours either. Our main motivation for studying $\text{RIP}_n(k)$, as in (4.5), stems from considering the approximation error between the sample and population covariance matrices, which is very much needed in post-selection inference applications (Kuchibhotla et al., 2018) as well as in high dimensional linear regression (Candes and Tao, 2007; Negahban et al., 2012; Wainwright, 2019). \diamond

Remark 4.10 (Gram Matrix to Covariance Matrix). Using the results in Section 4.1.1, the results of this section can be easily modified to bound $\text{RIP}_n(k)$ when the gram matrices $\hat{\Sigma}_n$ and Σ_n are replaced by the covariance matrices $\hat{\Sigma}_n^*$ and Σ_n^* respectively; see Remark 4.5 of Kuchibhotla et al. (2018) for more details. Similar comments also apply for the results in the next section on the restricted eigenvalue condition and will not be repeated. \diamond

Remark 4.11 (Applications in Adaptive Covariance Matrix Estimation). Concentration in-

equalities for $\text{RIP}_n(k)$ are also needed in adaptive estimation of a *bandable* covariance matrix. A matrix $\Sigma_n \in \mathbb{R}^{p \times p}$ is said to be *k-bandable*, for some $k \geq 1$, if

$$\Sigma_n(i, j) = 0 \quad \text{for all } |i - j| \geq k, \quad \text{for some } k \geq 1.$$

An adaptive estimator was proposed in [Cai and Yuan \(2012\)](#) based on the idea of block thresholding. Similar to the thresholding used in sparse covariance matrix estimation (see [Remark 4.3](#)), block thresholding sets to zero a sub-matrix if its operator norm is smaller than a threshold. The actual procedure is more complicated than this and is described in [Section 2.2 of Cai and Yuan \(2012\)](#). Theoretical study of such a block thresholding procedure requires a result similar to [Theorem 4.3](#); see [Theorem 3.3 in Section 3.2 of Cai and Yuan \(2012\)](#) for more details. The main difference in comparison with our result is that we do not require sub-Gaussian tails whereas the proof of [Theorem 3.3](#) there relies heavily on the normality of the random vectors; see also [Cai et al. \(2016\)](#) for a survey about high dimensional structured covariance matrix estimation. Using our results from this section, the performance of the adaptive estimator can be studied under much weaker assumptions of marginal sub-Weibull tail behaviors. \diamond

4.3 Restricted Eigenvalue (RE) Condition

One of the most well known estimators for high dimensional linear regression is the Lasso ([Tibshirani, 1996](#)). A crucial assumption in the proof of the oracle inequalities for Lasso is the *restricted eigenvalue* (RE) condition introduced by [Bickel et al. \(2009\)](#) for the matrix $\hat{\Sigma}_n$ as defined in [\(4.2\)](#); see [Section 4.4](#) for further details on its application to the theoretical analysis of Lasso. This section focuses on the RE condition and its bounds for sub-Weibulls. For any $1 \leq k \leq p$, the *RE(k) condition* on $\hat{\Sigma}_n$ is given by:

$$\inf_{\substack{S \subseteq \{1, \dots, p\}, \\ |S| \leq k}} \inf_{\theta \in \mathcal{C}(S; \delta)} \frac{\theta^\top \hat{\Sigma}_n \theta}{\theta^\top \theta} \geq \gamma_n > 0, \quad (4.15)$$

for some constant γ_n , where for any subset $S \subseteq \{1, 2, \dots, p\}$ and any $\delta \geq 1$,

$$\mathcal{C}(S; \delta) := \{\theta \in \mathbb{R}^p : \|\theta(S^c)\|_1 \leq \delta \|\theta(S)\|_1\}, \quad \text{where}$$

$\|v\|_1$ denotes the L_1 norm of any vector $v \in \mathbb{R}^p$, and $\theta(S)$ represents the sub-vector of θ with indices in S ; see [Equation \(11.10\) of Hastie et al. \(2015\)](#). Note, however, that for the specific application of RE conditions in the analysis of Lasso type estimators, the first infimum in [\(4.15\)](#) over all S with $|S| \leq k$ is *not* needed. Instead it *only* needs to be verified for S being the true support of the regression parameter β_0 (as in [Section 4.4](#)) with $\|\beta_0\|_0 \leq k$. [Rudelson and Zhou \(2013\)](#) verified [assumption \(4.15\)](#) for covariance matrices of sub-Gaussian random vectors, extending the work of [Raskutti et al. \(2010\)](#) for Gaussians. It is worth mentioning that the assumption of [Rudelson and Zhou \(2013\)](#) is that of jointly sub-Gaussian random vectors. Some extensions under weaker tail behavior, including sub-exponentials have also been considered in [Adamczak et al. \(2011\)](#) and [Lecué and Mendelson \(2017\)](#), for instance, although the latter's result applies more generally (see [Remark 4.13](#) for more discussion).

A general result proving this assumption based on a bound on the maximum elementwise norm is given in [Lemma 10.1 of van de Geer and Bühlmann \(2009\)](#). This result, coupled with our bounds on Δ_n in [Section 4.1](#), implies that if the random vectors X_i are (marginally) sub-Weibull as in

(4.3), then $\hat{\Sigma}_n$ satisfies the RE(k) condition (4.15) with probability converging to 1 as long as Σ_n satisfies its own corresponding RE(k) condition and the following holds:

$$kA_{n,p}\sqrt{\frac{\log p}{n}} + K_{n,p}^2 \frac{k(\log n)^{2/\alpha}(\log p)^{2/\alpha}}{n} = o(1). \quad (4.16)$$

This result does not allow for the optimal largest size for k , as noted in Raskutti et al. (2010, Section 3.2) as well, but it *does* relax the sub-Gaussianity assumption largely. Further, it is possible to get better rates using the bounds on $\text{RIP}_n(k)$ from Section 4.2, as shown below.

Main Result (Outline). In the following, we prove that gram matrices obtained from marginal/joint sub-Weibull random vectors satisfy the RE condition with high probability. The main result is Theorem 4.4 below which is proved (in Appendix E.3) using Theorem 4.3, and Lemma 12 of Loh and Wainwright (2012). Theorem 4.4 actually proves a *stronger* result – a sufficient condition regarding *restricted strong convexity* (RSC), a notion introduced by Negahban et al. (2012). As shown later in Section 4.3.1, the RE condition’s verification follows directly from this result. For simplicity, we again only consider the case $0 < \alpha \leq 2$ (the case $\alpha > 2$ is similar). Similar to Theorem 4.3, the result has two cases: (a) *marginal* and (b) *joint*. For each case, the corresponding bounds (4.17) and (4.18) in Theorem 4.4 establish the RSC property under appropriate conditions. Furthermore, Section 4.3.1 verifies the RE condition under *both* marginal and joint sub-Weibull assumptions, and also provides the *optimized* sample complexities required in each case; see in particular (4.24) for the joint case and (4.26)–(4.27) for the marginal case. In particular, these match the existing (optimal) scaling known for the joint sub-Gaussian case. More details on the implications of the results, as well as comparisons with the existing literature on the RE condition are discussed in Remarks 4.12–4.13. To our knowledge, a unified set of results obtained in this generality for the RE condition is not available (at least not easily) within the core statistics literature.

Theorem 4.4 (RSC: Unified Bounds for Sub-Weibulls). *Under the setting of Theorem 4.3 and recalling $\Upsilon_{n,s}$ as defined therein, the following high probability statements hold true: for every $1 \leq s \leq p$,*

(a) *(Marginal Sub-Weibull Case). If $\|X_i\|_{M,\psi_\alpha} \leq K_{n,p}$ for all $1 \leq i \leq n$, then setting*

$$\Xi_{n,s}^{(M)} := 14\sqrt{2}\sqrt{\frac{\Upsilon_{n,s}s \log(36np/s)}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (s \log(36np/s))^{2/\alpha}}{n},$$

we have with probability at least $1 - 3s(np)^{-1}$, simultaneously for all $\theta \in \mathbb{R}^p$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_{n,s}^{(M)} \right) \|\theta\|_2^2 - \frac{54\Xi_{n,s}^{(M)}}{s} \|\theta\|_1^2. \quad (4.17)$$

(b) *(Joint Sub-Weibull Case). If $\|X_i\|_{J,\psi_\alpha} \leq K_{n,p}$ for all $1 \leq i \leq n$, then setting*

$$\Xi_{n,s}^{(J)} := 14\sqrt{2}\sqrt{\frac{\Upsilon_{n,s}s \log(36np/s)}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (s \log(36np/s))^{2/\alpha}}{n},$$

we have with probability at least $1 - 3s(np)^{-1}$, simultaneously for all $\theta \in \mathbb{R}^p$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_{n,s}^{(J)} \right) \|\theta\|_2^2 - \frac{54\Xi_{n,s}^{(J)}}{s} \|\theta\|_1^2. \quad (4.18)$$

Here, in both cases, $C_\alpha > 0$ represents a constant depending only on α (but possibly different in the two cases), and $\lambda_{\min}(\Sigma_n)$ denotes the minimum eigenvalue of Σ_n .

Note: Bounds of the type (4.17)–(4.18) were discussed in Negahban et al. (2012) as sufficient conditions for verifying their general RSC condition (Definition 2); see Eqns. (20) and (31) therein. With a slight abuse of terminology, we ignore the distinction between their original RSC condition and these sufficient conditions, and call the latter by the same name here.

Remark 4.12 (Implications of Theorem 4.4) The “parameter” s in Theorem 4.4 is *not* directly related to the sparsity k of the regression parameter β_0 (as in Section 4.4). It is a free parameter that can be *chosen* (or optimized) suitably over $1 \leq s \leq p$. E.g., if we take $s = 1$, then $\Xi_{n,s}^{(M)} = \Xi_{n,s}^{(J)}$ and both these quantities converge to 0 if $n \gg (\log(np))^{2/\alpha}(\log(2n))^{2/\alpha}$. This implies that with probability at least $1 - 3/(np)$, simultaneously for all $\theta \in \mathbb{R}^p$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_{n,1}^{(J)} \right) \|\theta\|_2^2 - 54\Xi_{n,1}^{(J)} \|\theta\|_1^2. \quad (4.19)$$

Note that $\Xi_{n,1}^{(J)} = C_1 \sqrt{\log(np)/n} + C_2 (\log(n) \log(np))^{2/\alpha}/n$ (treating $K_{n,p}$ and $\Upsilon_{n,s}$ as constants). The inequality (4.19) can be compared to Proposition 8 in the recent work of Wong et al. (2020) where a similar result is derived under a joint sub-Weibull assumption only, but allowing for dependence through β -mixing of the observations. In comparison, our result’s sample complexity is similar to theirs, while having a better (faster) coefficient for $\|\theta\|_1^2$. \diamond

4.3.1 Verification of the RE(k) Condition

As mentioned earlier, Theorem 4.4 proves a stronger sufficient condition regarding RSC (as we call it; see the note below Theorem 4.4). We now show that this indeed implies the RE(k) condition (4.15), where k is set to denote the true sparsity of the regression parameter β_0 (as in Section 4.4). In our application for Lasso, we only need the RE(k) condition (4.15) with $\delta = 3$. Hence, for simplicity, we only prove (4.15) with $\delta = 3$. To this end, first note that for any $S \subseteq \{1, 2, \dots, p\}$ with $|S| \leq k$, and for any $\theta \in \mathcal{C}(S; 3)$, we have

$$\|\theta(S^c)\|_1 \leq 3 \|\theta(S)\|_1 \leq 3\sqrt{k} \|\theta(S)\|_2 \quad \Rightarrow \quad \|\theta\|_1 \leq 4\sqrt{k} \|\theta\|_2.$$

Now, let Ξ_s be either $\Xi_{n,s}^{(M)}$ or $\Xi_{n,s}^{(J)}$, as in Theorem 4.4, for *any* $1 \leq s \leq p$. The inequality above and Theorem 4.4 then together imply that for any given k and for all $1 \leq s \leq p$, with probability at least $1 - 3s/(np)$, simultaneously for all S with $|S| \leq k$ and for all $\theta \in \mathcal{C}(S; 3)$,

$$\begin{aligned} \theta^\top \hat{\Sigma}_n \theta &\geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_s \right) \|\theta\|_2^2 - \frac{54\Xi_s}{s} \|\theta\|_1^2 \\ &\geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_s - \frac{864k\Xi_s}{s} \right) \|\theta\|_2^2. \end{aligned} \quad (4.20)$$

The inequality (4.20) holds for *every* $1 \leq s \leq p$. If we choose $s = k$, then $\lambda_{\min}(\Sigma_n) \geq 1782\Xi_k$ is needed to conclude the RE(k) condition (4.15) with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$; see footnote 4 of Negahban et al. (2012) for a related calculation. Under a joint sub-Weibull assumption, this corresponds to requiring $n \gg (k \log(np/k) \log n)^{2/\alpha}$ (treating $\Upsilon_{n,k}$ and $K_{n,p}$ as constants). For $\alpha = 2$, this reduces to the familiar requirement of $n \gg k \log(np/k)$ upto a $\log n$ factor.

Similarly, if one chooses $s = 1$ in (4.20), then $\lambda_{\min}(\Sigma_n) \geq 27(32k + 1)\Xi_1 \gtrsim k\Xi_1$ is needed to conclude the RE(k) condition (4.15) with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$. Under either a marginal or a joint sub-Weibull assumption, this corresponds to a sample complexity similar to (4.16).

In general, one can *optimize* the right hand side of (4.20) over $1 \leq s \leq p$ in order to derive a better sample complexity. Note that $\Xi_s + k\Xi_s/s$ is non-random, and so, if $s^{(o)}$ minimizes $\Xi_s + k\Xi_s/s$, then (4.20) implies that with probability at least $1 - 3s^{(o)}/(np) \geq 1 - 3/n$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 864 \min_{1 \leq s \leq p} \left\{ \Xi_s + \frac{k\Xi_s}{s} \right\} \right) \|\theta\|_2^2 \quad \text{for all } \theta \in \bigcup_{|S| \leq k} \mathcal{C}(S; 3). \quad (4.21)$$

It is, however, hard to minimize $\Xi_s + k\Xi_s/s$ exactly because it involves four different powers of s , and hence, we find a simpler upper bound on the minimum, separately, under the joint and marginal sub-Weibull assumptions, considering different sub-cases for k in each case. (In the following calculations, we treat the quantities $\Upsilon_{n,s}$, for all $1 \leq s \leq p$, and $K_{n,p}$ as constants, and also ignore any multiplicative constants for rate optimization purposes.)

Optimized Sample Complexity in the Joint Sub-Weibull Case. In light of (4.21) with $\Xi_s = \Xi_{n,s}^{(J)}$, we consider minimizing $\Xi_{n,s}^{(J)} + k\Xi_{n,s}^{(J)}/s$, over $1 \leq s \leq p$. To this end, define

$$s_J^* := \frac{n^{\alpha/(4-\alpha)}}{(\log n)^{4/(4-\alpha)} \log(np)}. \quad (4.22)$$

This is the value¹ of s obtained by minimizing $k\Xi_{n,s}^{(J)}/s$ which consists of two terms that behave antagonistically with s (i.e. one increases while the other decreases). To find the best sample complexity for the RE(k) condition to hold, we now consider three cases:

$$\text{Case (i): } \alpha = 2, \text{ or Case (ii): } \alpha < 2, k \leq s_J^*, \text{ or Case (iii): } \alpha < 2, k > s_J^*. \quad (4.23)$$

For Cases (i) and (ii), we take $s = k$ in (4.20), with $\Xi_s \equiv \Xi_{n,s}^{(J)}$, to obtain: with probability at least $1 - 3k/(np)$, simultaneously for all $\theta \in \mathcal{C}(S; 3)$ and all S with $|S| \leq k$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - (27 + 864)\Xi_{n,k}^{(J)} \right) \|\theta\|_2^2.$$

Now, under Case (i) in (4.23), i.e. if $\alpha = 2$, we have:

$$\Xi_{n,k}^{(J)} \lesssim \sqrt{\frac{k \log(np)}{n}} + \frac{(\log n)(k \log(np))}{n} = o(1) \text{ whenever } n \gg k \log(np) \log n.$$

Under Case (ii) in (4.23), i.e. if $\alpha < 2$ and $k \leq s_J^*$, and with $\Xi_{n,s}^{(J)}$ monotone in s , we have:

$$\Xi_{n,k}^{(J)} \leq \Xi_{n,s_J^*}^{(J)} \lesssim \left(\frac{n^{\alpha/(4-\alpha)}}{n(\log n)^{4/(4-\alpha)}} \right)^{1/2} + \frac{1}{n} \left(\frac{n^{\alpha/(4-\alpha)}}{(\log n)^{4/(4-\alpha)}} \right)^{2/\alpha} = o(1) \text{ whenever } n \gg 1.$$

¹Formally, s_J^* should be defined as the smallest integer that is larger than (or equal to) the right hand side in (4.22) above. But this minor adjustment is irrelevant for the purpose of rate or sample complexity calculations, and therefore, we disregard this technicality in our calculations here.

Finally, under Case (iii) in (4.23), i.e. if $\alpha < 2$ and $k > s_J^*$, appealing to (4.21), bounding $\min_{1 \leq s \leq p} \{\Xi_{n,s}^{(J)} + k\Xi_{n,s}^{(J)}/s\}$ suffices. We do so using the following inequalities:

$$\begin{aligned} \min_{1 \leq s \leq p} \left\{ \Xi_{n,s}^{(J)} + \frac{k}{s} \Xi_{n,s}^{(J)} \right\} &\leq \Xi_{n,s_J^*}^{(J)} + \frac{k}{s_J^*} \Xi_{n,s_J^*}^{(J)} \leq \frac{2k}{s_J^*} \Xi_{n,s_J^*}^{(J)} \\ &\lesssim \frac{k}{\sqrt{s_J^*}} \sqrt{\frac{\log(np)}{n}} + \frac{k(s_J^*)^{(2-\alpha)/\alpha}}{n} (\log n)^{2/\alpha} (\log(np))^{2/\alpha} \\ &= \frac{2k \log(np) (\log n)^{2/(4-\alpha)}}{n^{2/(4-\alpha)}} = o(1) \text{ if } n \gg (k \log(np))^{2-\alpha/2} \log n. \end{aligned}$$

Summarizing, we conclude that under a joint sub-Weibull assumption, the $\text{RE}(k)$ condition (4.15) holds (with probability converging to 1 and with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$) over all the cases in (4.23) with the following corresponding *sample complexity requirements*:

$$(i): n \gg k \log(np) \log n, \quad (ii): n \gg 1, \quad (iii): n \gg (k \log(np))^{2-\alpha/2} \log n. \quad (4.24)$$

Combining all cases in (4.24), we only require $n \gg (k \log(np))^{2-\alpha/2} \log n$. \square

Optimized Sample Complexity in the Marginal Sub-Weibull Case. Again, in light of (4.21) with $\Xi_s = \Xi_{n,s}^{(M)}$, bounding $\min_{1 \leq s \leq p} \{\Xi_{n,s}^{(M)} + k\Xi_{n,s}^{(M)}/s\}$ is sufficient. Define

$$s_M^* := \frac{n^{\alpha/(4+\alpha)}}{(\log(np))^{(4-\alpha)/(4+\alpha)} (\log n)^{4/(4+\alpha)}}.$$

This is the value² of s obtained by minimizing $k\Xi_{n,s}^{(M)}/s$ which consists of two terms that behave antagonistically with s . We now consider two different cases for k :

$$\text{Case (i): } k \leq s_M^*, \quad \text{or} \quad \text{Case (ii): } k > s_M^*. \quad (4.25)$$

Under Case (i) in (4.25), i.e. when $k \leq s_M^*$, noting that $\Xi_{n,s}^{(M)}$ is monotone in s , we can use

$$\begin{aligned} \min_{1 \leq s \leq p} \left\{ \Xi_{n,s}^{(M)} + \frac{k}{s} \Xi_{n,s}^{(M)} \right\} &\leq \Xi_{n,k}^{(M)} + \frac{k}{k} \Xi_{n,k}^{(M)} = 2\Xi_{n,k}^{(M)} \leq 2\Xi_{n,s_M^*}^{(M)} \\ &\lesssim \left(\frac{n^{\alpha/(4+\alpha)} (\log(np))^{2\alpha/(4+\alpha)}}{n (\log n)^{4/(4+\alpha)}} \right)^{1/2} + \frac{1}{n} \left(\frac{n^{\alpha/(4+\alpha)} (\log(np))^{2\alpha/(4+\alpha)}}{(\log n)^{4/(4+\alpha)}} \right)^{2/\alpha} \\ &= o(1) \text{ if } n \log n \gg (\log(np))^{\alpha/2} \text{ and } n (\log n)^{8/(2\alpha+\alpha^2)} \gg (\log(np))^{4/(2+\alpha)}. \end{aligned} \quad (4.26)$$

²See Footnote 1 earlier regarding s_J^* . The same comments apply here for s_M^* as well and are not repeated.

Similarly, under Case (ii) in (4.25), i.e. if $k > s_M^*$, we can use

$$\begin{aligned}
\min_{1 \leq s \leq p} \left\{ \Xi_{n,s}^{(M)} + \frac{k}{s} \Xi_{n,s}^{(M)} \right\} &\leq \Xi_{n,s_M^*}^{(M)} + \frac{k}{s_M^*} \Xi_{n,s_M^*}^{(M)} \leq \frac{2k}{s_M^*} \Xi_{n,s_M^*}^{(M)} \\
&\lesssim \frac{k}{\sqrt{s_M^*}} \sqrt{\frac{\log(np)}{n}} + \frac{k(s_M^*)^{2/\alpha}}{n} (\log n)^{2/\alpha} (\log(np))^{2/\alpha} \\
&= 2k \frac{\sqrt{\log(np)}}{\sqrt{n}} \frac{(\log(np))^{(4-\alpha)/(8+2\alpha)} (\log n)^{4/(8+2\alpha)}}{n^{\alpha/(8+2\alpha)}} = 2k \frac{(\log(np))^{4/(4+\alpha)} (\log n)^{2/(4+\alpha)}}{n^{(2+\alpha)/(4+\alpha)}} \\
&= o(1) \text{ if } n \gg k^{(4+\alpha)/(2+\alpha)} (\log(np))^{4/(2+\alpha)} (\log n)^{2/(2+\alpha)}. \tag{4.27}
\end{aligned}$$

Thus, (4.26) and (4.27) provide the required sample complexities under Cases (i) and (ii) in (4.25). Combining both the cases, the requirement in (4.27) suffices for the RE(k) condition to hold (with probability converging to 1) under a marginal sub-Weibull assumption. \square

The results above, therefore, reveal several interesting sample complexity requirements, apart from the expected ones, under both the joint and marginal sub-Weibull assumptions. To the best of our knowledge, a unified and general set of results like this regarding the RE condition is not (easily) available/accessible within the core statistics literature. Lastly, we also point out that some of the logarithmic factors in the requirements above could possibly be improved (or removed) based on a more refined analysis which is not pursued here.

Remark 4.13 (Requirement/Relevance of Exponential Tails for RE Condition). Observe that the RE condition is only concerned with the minimum sparse eigenvalue and so, the assumption of exponential tails may not be required in its full strength; see [van de Geer and Muro \(2014\)](#) and [Oliveira \(2013\)](#) for details. In particular, for this problem, it is only required to bound (possibly exponentially), for some $\varepsilon > 0$, the probability of the event

$$\frac{1}{n} \sum_{i=1}^n \left(X_i^\top \theta \right)^2 \leq (1 - \varepsilon) \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\left(X_i^\top \theta \right)^2 \right].$$

Because this event is related to the average of non-negative random variables, it can have an exponentially small probability even under polynomial moment conditions; see Theorem 2.19 of [de la Peña et al. \(2009\)](#), for instance, for an exponential tail bound under only finite fourth moment conditions. [Oliveira \(2013\)](#) formalizes this to bound the probability of the event uniformly over all θ , and proved a general result related to the RE condition for a *normalized* covariance matrix; see Theorem 5.2 there. Some of the main differences between his result and the results in high dimensional statistics literature are listed after Theorem 5.2 therein; see also Section 6.1 of [van de Geer and Muro \(2014\)](#) for more comparisons.

Our marginal sub-Weibull assumption (a) in Theorem 4.4 is equivalent to the moment growth: $\|X_i(j)\|_r \leq C_\alpha r^{1/\alpha}$ for all $r \geq 1$. Under an additional so-called *small-ball condition*, Theorem E of [Lecué and Mendelson \(2017\)](#) shows that the same moment growth, but *only* for $1 \leq r \leq \log(wp)$, for some constant $w \geq 1$, suffices to verify the RE condition. Note that for p diverging with n , this weaker assumption of [Lecué and Mendelson \(2017\)](#) is almost equivalent to a marginal sub-Weibull requirement. It is also not clear if Theorem E of [Lecué and Mendelson \(2017\)](#), which is primarily aimed at the RE condition's verification, can be extended to prove more general and stronger RSC type bounds of the form (4.17)–(4.18), as we obtain in Theorem 4.4. Nevertheless, it must also be mentioned that their result on the RE condition requires a sample complexity of

$n \gtrsim \max\{k \log p, (\log p)^{4/\alpha-1}\}$ only. This is a weaker condition (and perhaps the weakest known) than what we can achieve here.

As we noted earlier (e.g., at the end of Remark 4.8, though in a different context), our main goal throughout Section 4 is to demonstrate ‘easy’ applications of the ready-to-use inequalities from Section 3 in handling these statistical problems and obtain unified and general, yet user-friendly, results for each of them. A more targeted problem-specific approach, possibly using different techniques (and assumptions), can perhaps lead to slightly better results or conditions for some of these problems. Given the larger focus of this paper, we do not to pursue such deeper nuanced analyses here.

Finally, we also remark that although it may be possible to prove the RE condition itself under weaker tail assumptions on the covariates and allowing for an exponential growth of p , the theoretical analysis of Lasso and other related high dimensional estimators — where this condition is perhaps most needed — usually *requires* (almost) exponential tails for the covariates anyway to ensure logarithmic dependence on p in the bounds and in the rates. \diamond

4.4 High Dimensional Linear Regression

In this section, we derive results related to the Lasso, a well-known high dimensional linear regression estimator introduced by Tibshirani (1996). Let $(X_1^\top, Y_1)^\top, \dots, (X_n^\top, Y_n)^\top$ be n independent random vectors in $\mathbb{R}^p \times \mathbb{R}$. Let $\beta_0 \in \mathbb{R}^p$ be a vector such that

$$Y_i = X_i^\top \beta_0 + \varepsilon_i \quad \text{with} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\varepsilon_i X_i] = 0 \in \mathbb{R}^p. \quad (4.28)$$

Observe that such a vector β_0 *always* exists (regardless of whether or not $\mathbb{E}(Y_i|X_i)$ is linear), as long as the population gram matrix $\sum_{i=1}^n \mathbb{E}[X_i X_i^\top]/n$ is invertible, and is given by

$$\begin{aligned} \beta_0 &= \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\left(Y_i - X_i^\top \theta \right)^2 \right] \\ &= \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[X_i X_i^\top \right] \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E} [X_i Y_i] \right). \end{aligned}$$

Differentiating the objective function above implies the second condition in (4.28). A linear model is said to be *well-specified* if $\mathbb{E}[\varepsilon_i|X_i] = 0$ in which case the second condition in (4.28) holds trivially, and $\mathbb{E}(Y_i|X_i)$ is exactly linear and equals $X_i^\top \beta_0$. Thus, the specification (4.28) is a much weaker condition and allows for a *misspecified* linear model. Note also that in (4.28), X is allowed to include 1 to account for an intercept term.

The Lasso estimator $\hat{\beta}_n(\lambda)$ of β_0 , for a regularization parameter $\lambda > 0$, is given by

$$\hat{\beta}_n(\lambda) := \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{2n} \sum_{i=1}^n \left(Y_i - X_i^\top \theta \right)^2 + \lambda \|\theta\|_1. \quad (4.29)$$

In most of the literature on Lasso, the guarantees on the estimator are usually obtained under some restrictive assumptions such as fixed or jointly sub-Gaussian covariates and/or homoscedastic Gaussian/sub-Gaussian errors, although these are not the only settings studied; see Vidaurre et al. (2013) and references therein for a detailed survey of L_1 -penalized regression methods and their computational and theoretical properties.

Outline of the Main Results. In this section, we analyze the Lasso under *much weaker* than usual tail assumptions on the covariates X_i as well as the errors ε_i . Our main results in this regard are Theorems 4.5 and 4.6 (proved in Appendix E.4) below. Both results only assume a *marginal* sub-Weibull property for the X_i 's, while for the errors, Theorem 4.5 assumes ε_i to be sub-Weibull and Theorem 4.6 only assumes ε_i to have polynomial tails (i.e. finite moments upto some order $r \geq 2$). Moreover, our analysis throughout *allows* for (a) *model misspecification*, and (b) *both fixed and random covariates* since we do not assume identical distributions of the random vectors. The main message of both the results here is that the Lasso estimator attains the rate of $\sqrt{k \log p/n}$ for a large range of k, p if β_0 is k -sparse. More details on both results and their implications, including their rates of convergence and applicability in various settings, are discussed in Remarks 4.14 and 4.15. Further extensions as well as a general *oracle inequality* for the Lasso are given in Remark 4.16. To the best of our knowledge, these results are among the very few (if not the only) results proving rates of convergence of the Lasso estimator in this generality, with one notable exception being a recent result from Han and Wellner (2019) which will be discussed later in the context of Theorem 4.6.

A very general result about Lasso is obtained by Negahban et al. (2012) that is derived based on deterministic inequalities (see Section 4.2 therein). Both our main results here are based on this general result. We present Theorem 4.5 first. Recall the definitions of $\hat{\Sigma}_n$, Σ_n from (4.2), and $\Xi_{n,s}^{(M)}$ from Theorem 4.4(a), and also that $\|\beta_0\|_0$ denotes the sparsity of β_0 .

Theorem 4.5 (Lasso with Marginally Sub-Weibull X_i 's and Sub-Weibull ε_i 's). *Consider the setting above. Suppose $\|\beta_0\|_0 \leq k$ and there exists $0 < \alpha \leq 2$, and $\vartheta, K_{n,p} > 0$ such that*

$$\max \left\{ \|X_i\|_{M, \psi_\alpha}, \|\varepsilon_i\|_{\psi_\vartheta} \right\} \leq K_{n,p} \quad \text{for all } 1 \leq i \leq n.$$

Also suppose $n \geq 2$, $k \geq 1$ and the matrix Σ_n satisfies

$$\lambda_{\min}(\Sigma_n) \geq 54 \min_{1 \leq s \leq p} \left\{ \Xi_{n,s}^{(M)} + \frac{32k\Xi_{n,s}^{(M)}}{s} \right\}, \quad (4.30)$$

with $\Xi_{n,s}^{(M)}$ as defined in Theorem 4.4(a). Then, with probability at least $1 - 3(np)^{-1} - 3n^{-1}$, the regularization parameter λ_n can be chosen to be

$$\lambda_n = 14\sqrt{2}\sigma_{n,p} \sqrt{\frac{\log(np)}{n}} + \frac{C_\gamma K_{n,p}^2 (\log(2n))^{1/\gamma} (2 \log(np))^{1/\gamma}}{n}, \quad (4.31)$$

so that the Lasso estimator $\hat{\beta}_n(\lambda_n)$ satisfies

$$\left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2 \leq \frac{84\sqrt{2}}{\lambda_{\min}(\Sigma_n)} \left[\sigma_{n,p} \sqrt{\frac{k \log(np)}{n}} + \frac{C_\gamma K_{n,p}^2 k^{1/2} (\log(np))^{2/\gamma}}{n} \right],$$

where $C_\gamma > 0$ is some constant depending only on γ and

$$\frac{1}{\gamma} := \frac{1}{\alpha} + \frac{1}{\vartheta}, \quad \text{and} \quad \sigma_{n,p}^2 := \max_{1 \leq j \leq p} \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i(j)\varepsilon_i) > 0.$$

Remark 4.14 (Rate of Convergence and a Few Other Comments on Theorem 4.5). It follows from the result that if (4.30) holds (as verified in Section 4.3.1) and if

$$(\log(np))^{4/\gamma-1} = o(n), \quad \text{as } n \rightarrow \infty,$$

then the rate of convergence of the Lasso is $\sqrt{k \log p/n}$ which is also known to be the (near) minimax optimal rate (Raskutti et al., 2011). Note further that the probability guarantee in Theorem 4.5 is converging to 1 as $n \rightarrow \infty$, and so, the bound therein has $\log(np)$ instead of the usual $\log p$. By making the probability to be $1 - O(p^{-1})$, the usual rate of $\sqrt{k \log p/n}$ can be recovered. In the special case of conditionally homoscedastic errors ε_i with $\mathbb{E}(\varepsilon_i|X_i) = 0$ and $\text{Var}(\varepsilon_i|X_i) = \sigma^2$, and with X_i 's normalized to have marginal variances of 1, we have $\sigma_{n,p} = \sigma$ and this then leads to the familiar rate of $\sigma \sqrt{k \log p/n}$ for the Lasso estimator.

Lastly, if a *joint* (instead of marginal) sub-Weibull property is assumed on the covariates in Theorem 4.5, then the same result holds with $\Xi_{n,s}^{(M)}$ in (4.30) replaced by $\Xi_{n,s}^{(J)}$, with $\Xi_{n,s}^{(J)}$ as in Theorem 4.4(b). (This is true for Theorem 4.6 as well and won't be repeated there.) With $\Xi_{n,s}^{(J)} \leq \Xi_{n,s}^{(M)}$, this version of (4.30) imposes weaker sample complexity related conditions on the growth of (n, k) , as seen from Section 4.3.1 as well. Some related results for the Lasso with jointly sub-Weibull dependent random vectors can be found in Wong et al. (2020). \diamond

Lasso under Polynomial Moments on Errors. A careful inspection of the theoretical analysis of Lasso reveals that the assumption of sub-Weibull errors in Theorem 4.5 *can be weakened* to polynomial-tailed errors. This has also been noted in the recent work of Han and Wellner (2019); see Theorem 5 and Examples 4-5 therein, where they provide a general recipe for deriving the convergence rates of Lasso allowing for much weaker tailed errors. Their results, however, are asymptotic in nature and need the restrictive assumption of ε_i 's being mean 0 and independent of X_i , $1 \leq i \leq n$, although they do allow for dependence among ε_i 's. In Theorem 4.6 below, we prove an analogue of Theorem 4.5 assuming only polynomial moments (upto some order $r \geq 2$) of ε_i . Recall Definition 2.5 and $\Xi_{n,s}^{(M)}$ from Theorem 4.4, and recall that for any random variable W , $\|W\|_r = (\mathbb{E}[|W|^r])^{1/r}$ for $r > 0$.

Theorem 4.6 (Lasso with Marginally Sub-Weibull X_i 's and Polynomial-Tailed ε_i 's). *Under the setting of Theorem 4.5, suppose $\|\beta_0\|_0 \leq k$ and there exists $0 < \alpha \leq 2, r \geq 2$ so that*

$$\max_{1 \leq i \leq n} \|X_i\|_{M, \psi_\alpha} \leq K_{n,p}, \quad \text{and} \quad \max_{1 \leq i \leq n} \|\varepsilon_i\|_r \leq K_{\varepsilon,r}.$$

Also suppose $n \geq 2, k \geq 1$ and that Σ_n satisfies (4.30). Then for $L \geq 1$, with probability at least $1 - 3(np)^{-1} - 3n^{-1} - L^{-1}$, the regularization parameter λ_n can be chosen to be

$$\lambda_n = 14\sqrt{2}\sigma_{n,p} \sqrt{\frac{\log(np)}{n}} + \frac{C_\alpha K_{n,p} K_{\varepsilon,r} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{n^{1-1/r}},$$

so that the Lasso estimator $\hat{\beta}(\lambda_n)$ satisfies

$$\begin{aligned} \left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2 &\leq \frac{84\sqrt{2}}{\lambda_{\min}(\Sigma_n)} \sigma_{n,p} \sqrt{\frac{k \log(np)}{n}} \\ &\quad + C_\alpha K_{n,p} K_{\varepsilon,r} \frac{k^{1/2} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{\lambda_{\min}(\Sigma_n) n^{1-1/r}}, \end{aligned}$$

for some constant $C_\alpha > 0$ depending only on α .

Remark 4.15 (Convergence Rates and the Special Case of Fixed Designs). Theorem 4.6 readily proves that the rate of convergence of the Lasso is $\sigma_{n,p}\sqrt{k \log p/n}$ if

$$K_{\varepsilon,r}(\log(np))^{1/\alpha-1/2}(\log(2n))^{1/\alpha} = o(n^{1/2-1/r}). \quad (4.32)$$

In comparison to Han and Wellner (2019), Theorem 4.6 provides a precise non-asymptotic extension of their (asymptotic) results under (marginally) sub-Weibull covariates, without the assumption regarding the errors being independent of the covariates. Since our result allows for (a) non-identically distributed observations, (b) both fixed and random designs, as well as (c) possibly misspecified models, it serves as a generalization (under sub-Weibull covariates) of Theorem 5 (and Example 5) of Han and Wellner (2019). Moreover, a careful inspection of their sample complexity requirement, as given in Equation (4.4) of their result, implies the condition $(\log p)^{4/\alpha+1} = O(n^{2-4/r})$ when translated into our setup and notation. This is a far stronger condition (e.g., if $\log p$ is polynomial in n) than our requirement (4.32).

Finally, note that if we are under a *fixed design*, i.e. if $X_i, 1 \leq i \leq n$ are n fixed vectors, then X_i 's simply are marginally sub-Weibull (∞) and

$$\max_{1 \leq i \leq n} \|X_i\|_{M,\psi_2} \leq \max_{1 \leq i \leq n} \|X_i\|_{M,\psi_\infty} = \max_{1 \leq i \leq n} \max_{1 \leq j \leq p} |X_i(j)|.$$

Hence, applying Theorem 4.6 with $\alpha = 2$ in this case, we observe that a rate of $\sqrt{k \log p/n}$ can be achieved under the (almost trivial) rate constraint

$$K_{\varepsilon,r}(\log(np))^{-1/2}(\log(2n))^{1/2} = o(n^{1/2-1/r}),$$

which is satisfied as long as n is large enough and $r > 2$. Similarly, for Theorem 4.5, the constraint becomes: $(\log(np))^{4/\theta-1} = o(n)$. It should be noted that for fixed designs, the RE condition is simply an explicit assumption. \diamond

Remark 4.16 (Extensions and Other Estimators). Using the probability tools from Section 3 and the method of proof in this section, it is possible to prove very general results extending Theorem 4.5 in several directions (similar extensions also apply to Theorem 4.6, although we only illustrate them for Theorem 4.5). We briefly discuss some of these below.

Theorem 4.5 is proved under the assumption of ‘hard’ sparsity in the sense that no more than k entries of β_0 are non-zero. One can actually derive a more *general oracle inequality* (without any hard sparsity condition) for $\hat{\beta}_n(\lambda_n)$ using Theorem 1 of Negahban et al. (2012).

Under the assumptions of Theorem 4.5 (except the hard sparsity), an *oracle inequality for the Lasso* is as follows. For a choice of λ_n as in (4.31), with probability converging to 1,

$$\begin{aligned} & \left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2^2 \\ & \leq \min_{S: \Xi_{n,|S|}^{(M)} = o(1)} \left[\frac{18\lambda_n^2 |S|}{\Gamma_n^2(S)} + \frac{8\lambda_n \|\beta_0(S^c)\|_1}{\Gamma_n(S)} + \frac{3456\Xi_{n,|S|}^{(M)} \|\beta_0(S^c)\|_1^2}{|S|\Gamma_n(S)} \right], \end{aligned} \quad (4.33)$$

where $\Gamma_n(S) := \lambda_{\min}(\Sigma_n) - 1755\Xi_{n,|S|}^{(M)}$. Under the condition $\Xi_{n,|S|}^{(M)} = o(1)$, for large enough n , $\Gamma_n(S) \geq \lambda_{\min}(\Sigma_n)/2$. (The constants could possibly be improved here.) This is an oracle inequality because there is *no* assumption on β_0 and the bound *adapts* to the true sparsity of β_0 . The

proof is in Appendix E.4 (Proposition E.1). As shown in Section 4.3 of Negahban et al. (2010), inequality (4.33) implies a rate of convergence if β_0 is *weakly sparse*.

Following the proof of Proposition 2 of Negahban et al. (2010), and using the proof of Theorem 4.4, it is easy to prove the restricted strong convexity property for generalized linear models when the covariates are marginally sub-Weibull. Hence, Theorem 4.5 can be easily extended to L_1 -penalized estimation methods for *generalized linear models* as well.

Finally, we mention that apart from the Lasso, there are many *other estimators* available for high dimensional linear regression, including, for instance, the Dantzig selector (Candes and Tao, 2007) and the square-root Lasso (Belloni et al., 2011), among others. The key ingredients in the analysis of all these estimators are the restricted eigenvalue condition and the gradient’s control, as shown in van de Geer (2016). Hence, the rate of convergence of these estimators can also be derived under weaker tail assumptions based on our results. \diamond

5 Conclusions and Future Work

In this paper, we proposed a new Orlicz norm that extracts a part sub-Gaussian tail behavior for sums of independent random variables. Various concentration inequalities related to sub-Weibull random variables and processes are then studied in a unified way. We hope that the exposition here amplifies the use of sub-Weibull random variables, especially the heavy-tailed ones, in the theoretical analysis of statistical methods. To illustrate this, we studied four fundamental statistical problems in high-dimensions and extended many of the by-now standard results in the literature. As mentioned earlier (e.g., in Section 4.3.1), our main goal here was to demonstrate the applications of our user-friendly concentration inequalities in handling these statistical problems under (much) weaker than usual tail assumptions, and obtaining fairly general results that still compare favorably to existing ones under stronger conditions. For some of the problems (e.g., RE condition), a more nuanced problem specific analysis can possibly lead to slightly better results or conditions than ours. But we refrain from such refined analyses given our main focus in this paper. Nevertheless, we do believe our results in Section 4.3 provide a much needed unified analysis on the RE condition that is not easily accessible in the statistics literature. Moreover, our results on the Lasso in Section 4.4 are possibly the first results in the literature that are obtained in such generality.

Throughout the paper, we have restricted the random variables/vectors to be independent to keep the presentation simple. The independence assumption, however, may not be appropriate for many econometric applications. The extensions of the results in Section 3 are available in Merlevède et al. (2011) for strong mixing random variables, and in Appendix B of Kuchibhotla et al. (2018) for functionally dependent random variables (Wu, 2005). Unfortunately, many useful processes are not strongly mixing and the results of Kuchibhotla et al. (2018) do not reduce to those in Section 3 under independence. Extensions to the case of martingales are also not fully understood. A recent progress in this direction is Fan (2017) that provides the result for martingales with $\alpha = 2$; see also Fan et al. (2017) for related results. Tail bounds for martingales matching their asymptotic normality under sub-Weibull martingale differences have important implications for concentration results related to functions of independent random variables, which in turn are useful for dependent data (Wu, 2005); see Boucheron et al. (2005) for more applications in this regard. Thus, it is worth considering possible extensions of our results in Section 3 to martingales.

In terms of further statistical applications of our results, an important problem worth considering is a complete study of the problem in Remark 3.3, including consistency of the LKAEs

in terms of the supremum norm and/or uniform-in-bandwidth consistency. These problems have been considered under an asymptotic setting by [Einmahl and Mason \(2000, 2005\)](#) using empirical process techniques. Their basic framework can indeed be adopted and combined with our results on suprema of empirical processes in [Appendix B](#) to obtain a sequence of widely applicable non-asymptotic results for LKAEs involving sub-Weibulls.

Further, it is also of interest to study the version of these problems involving the so-called “generated regressors”, wherein the kernel smoothing is only performed over (possibly) lower dimensional and/or estimated (if unknown) transformations of the original covariates. Such methods are of considerable importance in econometrics and in the sufficient dimension reduction literature. The latter can be particularly useful in high dimensional settings, where a fully non-parametric smoothing may be undesirable due to the curse of dimensionality; see [Mammen et al. \(2012, 2013\)](#) for some results and literature review on non-parametric regression over generated regressors. Using our empirical process results from [Appendix B](#) again, it would be of interest to obtain non-asymptotic tail bounds and rates of convergence for such LKAEs over generated regressors, especially in “truly” high dimensional settings where the dimension of the original covariates could be much larger than the sample size. While all these problems are interesting, a detailed analysis is far too involved for the scope of the current paper. We do hope to explore some of these problems separately in the future.

A Properties of the GBO Norm

In this section, we provide a collection of some useful basic properties of the GBO norm. Since it does not have a closed form, it is hard to directly see the part sub-Gaussian behavior captured by the GBO norm for sub-Weibulls, as shown in [\(2.4\)](#) for sub-exponentials. To resolve this issue, we first provide in [Proposition A.1](#) an equivalent norm that is based on a closed form g . (The proofs of all Propositions in this Appendix are given in [Appendix C](#).)

Proposition A.1. *Fix $\alpha, L > 0$. Define $\phi_{\alpha,L} : [0, \infty) \rightarrow [0, \infty)$ as*

$$\phi_{\alpha,L}(x) = \exp\left(\min\left\{x^2, \left(\frac{x}{L}\right)^\alpha\right\}\right) - 1.$$

Then for any random variable X , $\|X\|_{\Psi_{\alpha,L}} \leq \|X\|_{\phi_{\alpha,L}} \leq 2\|X\|_{\Psi_{\alpha,L}}$.

In the remaining part of this section, we derive various properties of $\|\cdot\|_{\Psi_{\alpha,L}}$, the proofs of which are all in [Appendix C](#). We start with simple monotonicity properties of $\|\cdot\|_{\Psi_{\alpha,L}}$.

Proposition A.2 (Monotonicity Properties). *The following monotonicity properties hold for the GBO norm:*

- (a) *If $|X| \leq |Y|$ almost surely, then $\|X\|_{\Psi_{\alpha,L}} \leq \|Y\|_{\Psi_{\alpha,L}}$ for all $\alpha, L > 0$.*
- (b) *For any random variable X , $\|X\|_{\Psi_{\alpha,L}} \leq \|X\|_{\Psi_{\alpha,K}}$ for $0 \leq L \leq K$.*

The following sequence of propositions prove the equivalence of finite $\Psi_{\alpha,L}$ -norm with a tail bound and a moment growth. The proofs are similar to those of [van de Geer and Lederer \(2013\)](#). It is worth mentioning here that although we present some of the results with explicit constants, our goal is not to provide optimal constants and they could possibly be improved.

Proposition A.3 (Equivalence of Tail and Norm Bounds). *For any random variable X with $\delta := \|X\|_{\Psi_{\alpha,L}}$, we have*

$$\mathbb{P}\left(|X| \geq \delta \left\{ \sqrt{t} + Lt^{1/\alpha} \right\}\right) \leq 2 \exp(-t), \quad \text{for all } t \geq 0. \quad (\text{A.1})$$

Conversely, if the tail bound (A.1) holds for some constants $\delta, L > 0$, then

$$\|X\|_{\Psi_{\alpha,c(\alpha)L}} \leq \sqrt{3}\delta, \quad \text{where } c(\alpha) := 3^{1/\alpha}/\sqrt{3}.$$

Proposition A.4 (Equivalence of Moment Growth and Norm Bound). *For any random variable X ,*

$$C_*(\alpha) \sup_{p \geq 1} \frac{\|X\|_p}{\sqrt{p} + Lp^{1/\alpha}} \leq \|X\|_{\Psi_{\alpha,L}} \leq C^*(\alpha) \sup_{p \geq 1} \frac{\|X\|_p}{\sqrt{p} + Lp^{1/\alpha}},$$

where $C_*(\alpha) := \frac{1}{2} \min\{1, \alpha^{1/\alpha}\}$ and $C^*(\alpha) := e \max\{2, 4^{1/\alpha}\}$.

Proposition A.5 (Quasi-Norm Property). *For any sequence of (possibly dependent) random variables $X_i, 1 \leq i \leq k$,*

$$\left\| \sum_{i=1}^k X_i \right\|_{\Psi_{\alpha,L}} \leq Q_\alpha \sum_{i=1}^k \|X_i\|_{\Psi_{\alpha,L}},$$

where

$$Q_\alpha := \begin{cases} 2e(4/\alpha)^{1/\alpha}, & \text{if } \alpha < 1, \\ 1, & \text{if } \alpha \geq 1. \end{cases}$$

One of the main advantages of Orlicz norms of the exponential type lies in their usefulness to derive maximal inequalities. The following result proves one such for the GBO norm $\|\cdot\|_{\Psi_{\alpha,L}}$.

Proposition A.6 (Maximal Inequality). *Let X_1, \dots, X_N be random variables (possibly dependent) such that $\max_{1 \leq j \leq N} \|X_j\|_{\Psi_{\alpha,L}} \leq \Delta < \infty$ for some $\alpha, L, \Delta > 0$. Set $X_N^* := \max_{1 \leq j \leq N} |X_j|$, and recall $c(\alpha)$ and Q_α from Propositions A.3 and A.5. Then for all $t \geq 0$,*

$$\mathbb{P}\left(X_N^* \geq \Delta \left\{ \sqrt{t + \log N} + L(t + \log N)^{1/\alpha} \right\}\right) \leq 2 \exp(-t),$$

and

$$\|X_N^*\|_{\Psi_{\alpha,K(\alpha)L}} \leq \Delta Q_\alpha \left\{ \sqrt{3} + \sqrt{\log N} + M(\alpha)L(\log N)^{\frac{1}{\alpha}} \right\},$$

where $K(\alpha) := c(\alpha)M(\alpha)$ with $M(\alpha) := \max\{1, 2^{(1-\alpha)/\alpha}\}$.

Remark A.1 (Bound on the Expectation of the Maximum). From Proposition A.6 it follows that

$$\|X_N^*\|_1 \leq \max_{1 \leq j \leq N} \|X_j\|_{\Psi_{\alpha,L}} C_\alpha \left\{ \sqrt{\log N} + L(\log N)^{1/\alpha} \right\},$$

for some constant C_α depending only on α . Note that if the random variables are sub-Gaussian ($\alpha = 2$), then the rate becomes $\sqrt{\log N}$. The main implication of the GBO norm is that it shows the rate can *still* be $\sqrt{\log N}$ even if $\alpha \neq 2$ as long as $L(\log N)^{1/\alpha-1/2} = o(1)$. \diamond

The next proposition provides an alternative to, and a generalization of, Proposition A.6. This is similar to Proposition 4.3.1 of [de la Peña and Giné \(1999\)](#). Note that for infinitely many random variables ($N = \infty$), Proposition A.6 does not lead to useful bounds; see the discussion following Proposition 4.3.1 of [de la Peña and Giné \(1999\)](#) for the importance of considering an alternative result as presented below.

Proposition A.7 (A Sharper Maximal Inequality). *Let X_1, X_2, \dots be any sequence of random variables (possibly dependent) such that for all $i = 1, 2, \dots$, $\|X_i\|_{\Psi_{\alpha,L}} < \infty$ for some $\alpha, L > 0$, and recall $c(\alpha)$, Q_α and $M(\alpha)$ from Propositions A.3, A.5 and A.6. Then*

$$\left\| \sup_{k \geq 1} \frac{|X_k|}{\sqrt{2} \|X_k\|_{\Psi_{\alpha,L}} \Psi_{\alpha, S(\alpha)L}^{-1}(k)} \right\|_{\Psi_{\alpha, c(\alpha)M(\alpha)L}} \leq 2.5Q_\alpha,$$

where $S(\alpha) := 2^{1/\alpha}M(\alpha)/2$.

A.1 Extensions to Tail Behaviors Involving Multiple Regimes

The GBO norm $\|\cdot\|_{\Psi_{\alpha,L}}$ introduced in Section 2 is designed to exploit two regimes in the tail of a random variable, namely, Gaussian and Weibull of order α . It is of interest to extend the theory to exploit more than two regimes in the tail of a random variable. Many examples exist where this is relevant, including in particular U -statistics based on independent variables; see, for example, [Latała \(1999\)](#), [Giné et al. \(2000\)](#) and [Boucheron et al. \(2005\)](#) for results on U -statistics and Rademacher Chaos.

For vectors $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k) \in (\mathbb{R}^+)^k$ and $\mathbf{L} = (L_1, \dots, L_k) \in (\mathbb{R}^+)^k$, for some k , define the function $\Psi_{\boldsymbol{\alpha}, \mathbf{L}}(\cdot)$ based on the inverse function

$$\Psi_{\boldsymbol{\alpha}, \mathbf{L}}^{-1}(t) := \sum_{j=1}^k L_j (\log(1+t))^{1/\alpha_j} \quad \text{for } t \geq 0.$$

The extended multiple regime GBO norm is defined by setting $g(\cdot) = \Psi_{\boldsymbol{\alpha}, \mathbf{L}}(\cdot)$ in Definition 2.1. The GBO norm $\|\cdot\|_{\Psi_{\alpha,L}}$ corresponds to $\boldsymbol{\alpha} = (1/2, \alpha)$ and $\mathbf{L} = (1, L)$. Similar to $\Psi_{\alpha,L}(\cdot)$, there is no closed form expression for $\Psi_{\boldsymbol{\alpha}, \mathbf{L}}(\cdot)$, and a function $\phi_{\boldsymbol{\alpha}, \mathbf{L}}(\cdot)$ closely related to $\Psi_{\boldsymbol{\alpha}, \mathbf{L}}(\cdot)$ is given by:

$$\phi_{\boldsymbol{\alpha}, \mathbf{L}}^{-1}(t) := \max \left\{ L_j (\log(1+t))^{1/\alpha_j} : 1 \leq j \leq k \right\}.$$

It is easy to check that $\|X\|_{\Psi_{\boldsymbol{\alpha}, \mathbf{L}}} \leq \|X\|_{\phi_{\boldsymbol{\alpha}, \mathbf{L}}} \leq k \|X\|_{\Psi_{\boldsymbol{\alpha}, \mathbf{L}}}$. All the properties stated in this section also hold for the extended GBO norm $\|\cdot\|_{\Psi_{\boldsymbol{\alpha}, \mathbf{L}}}$. Their proofs are similar and hence omitted to avoid repetition.

Supplementary Material

Supplement to “Moving Beyond Sub-Gaussianity in High Dimensional Statistics: Applications in Covariance Estimation and Linear Regression”. The supplementary material (Appendices B–F) contains additional results and technical materials that could not be accommodated in the main article. In Appendix B, we extend the study of sub-Weibulls to tail bounds for the suprema of empirical processes. In Appendices C–F, we present the proofs of all our results in the main article and the supplement.

Acknowledgements

We would like to thank the Editor, the anonymous Associate Editor and the two Reviewers for their constructive comments and useful suggestions that helped significantly improve the article. We would also like to thank Dr. Edward George for helpful initial discussions that improved the article's presentation.

References

- Adamczak, R. (2008). A tail inequality for suprema of unbounded empirical processes with applications to Markov chains. *Electron. J. Probab.*, 13:1000–1034.
- Adamczak, R., Litvak, A. E., Pajor, A., and Tomczak-Jaegermann, N. (2011). Restricted isometry property of matrices with independent columns and neighborly polytopes by random sampling. *Constr. Approx.*, 34(1):61–88.
- Alexander, K. S. (1985). Rates of growth for weighted empirical processes. In *Proceedings of the Berkeley conference in honor of Jerzy Neyman and Jack Kiefer (Berkeley, Calif., 1983)*, volume 2, pages 475–493.
- Bakhshizadeh, M., Maleki, A., and de la Peña, V. H. (2020). Sharp concentration results for heavy-tailed distributions. *arXiv preprint arXiv:2003.13819*.
- Baraniuk, R., Davenport, M., DeVore, R., and Wakin, M. (2008). A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28(3):253–263.
- Belloni, A. and Chernozhukov, V. (2013). Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547.
- Belloni, A., Chernozhukov, V., and Wang, L. (2011). Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika*, 98(4):791–806.
- Bickel, P. J. and Levina, E. (2008). Covariance regularization by thresholding. *Ann. Statist.*, 36(6):2577–2604.
- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.*, 37(4):1705–1732.
- Bogucki, R. (2015). Suprema of canonical weibull processes. *Statist. Probab. Lett.*, 107:253–263.
- Boucheron, S., Bousquet, O., Lugosi, G., and Massart, P. (2005). Moment inequalities for functions of independent random variables. *Ann. Probab.*, 33(2):514–560.
- Buja, A., Brown, L., Berk, R., George, E., Pitkin, E., Traskin, M., Zhang, K., and Zhao, L. (2019). Models as approximations I: Consequences illustrated with linear regression. *Statist. Sci.*, 34(4):523–544.
- Cai, T. T. and Liu, W. (2011). Adaptive thresholding for sparse covariance matrix estimation. *J. Amer. Statist. Assoc.*, 106(494):672–684.

- Cai, T. T., Ren, Z., and Zhou, H. H. (2016). Estimating structured high-dimensional covariance and precision matrices: optimal rates and adaptive estimation. *Electron. J. Stat.*, 10(1):1–59.
- Cai, T. T. and Yuan, M. (2012). Adaptive covariance matrix estimation through block thresholding. *Ann. Statist.*, 40(4):2014–2042.
- Candes, E. and Tao, T. (2007). The Dantzig selector: statistical estimation when p is much larger than n . *Ann. Statist.*, 35(6):2313–2351.
- Candes, E. J. and Tao, T. (2005). Decoding by linear programming. *IEEE transactions on information theory*, 51(12):4203–4215.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014). Gaussian approximation of suprema of empirical processes. *Ann. Statist.*, 42(4):1564–1597.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2017). Central limit theorems and bootstrap in high dimensions. *Ann. Probab.*, 45(4):2309–2352.
- de la Peña, V. H., Lai, T. L., and Shao, Q.-M. (2009). *Self-normalized processes: Limit theory and statistical applications*. Probability and its Applications (New York). Springer-Verlag, Berlin.
- de la Peña, V. H. and Giné, E. (1999). *Decoupling: From dependence to independence*. Probability and its Applications (New York). Springer-Verlag, New York.
- Dirksen, S. (2015). Tail bounds via generic chaining. *Electron. J. Probab.*, 20(53):1–29.
- Einmahl, U. and Mason, D. M. (2000). An empirical process approach to the uniform consistency of kernel-type function estimators. *J. Theoret. Probab.*, 13(1):1–37.
- Einmahl, U. and Mason, D. M. (2005). Uniform in bandwidth consistency of kernel-type function estimators. *Ann. Statist.*, 33(3):1380–1403.
- Fan, J., Liao, Y., and Liu, H. (2016). An overview of the estimation of large covariance and precision matrices. *Econom. J.*, 19(1):C1–C32.
- Fan, X. (2017). Gaussian martingale inequality applies to random functions and maxima of empirical processes. *ArXiv e-prints:1706.03916*.
- Fan, X., Grama, I., and Liu, Q. (2017). Deviation inequalities for martingales with applications. *J. Math. Anal. Appl.*, 448(1):538–566.
- Giné, E. and Koltchinskii, V. (2006). Concentration inequalities and asymptotic results for ratio type empirical processes. *Ann. Probab.*, 34(3):1143–1216.
- Giné, E., Latała, R., and Zinn, J. (2000). Exponential and moment inequalities for U -statistics. In *High dimensional probability II*, volume 47, pages 13–38. Springer.
- Giné, E. and Nickl, R. (2016). *Mathematical foundations of infinite-dimensional statistical models*. Cambridge University Press, New York.
- Gluskin, E. D. and Kwapien, S. (1995). Tail and moment estimates for sums of independent random variables with logarithmically concave tails. *Studia Math.*, 114(3):303–309.

- Guédon, O., Litvak, A. E., Pajor, A., and Tomczak-Jaegermann, N. (2014). Restricted isometry property for random matrices with heavy-tailed columns. *Comptes Rendus Mathématique*, 352(5):431–434.
- Guédon, O., Litvak, A. E., Pajor, A., and Tomczak-Jaegermann, N. (2015). On the interval of fluctuation of the singular values of random matrices. *arXiv preprint arXiv:1509.02322*.
- Han, Q. and Wellner, J. A. (2019). Convergence rates of least squares regression estimators with heavy-tailed errors. *Ann. Stat.*, 47(4):2286–2319.
- Hansen, B. E. (2008). Uniform convergence rates for kernel estimation with dependent data. *Econom. Theory*, 24(3):726–748.
- Hastie, T., Tibshirani, R., and Wainwright, M. (2015). *Statistical learning with sparsity: The Lasso and generalizations*. CRC Press.
- Hitczenko, P., Montgomery-Smith, S. J., and Oleszkiewicz, K. (1997). Moment inequalities for sums of certain independent symmetric random variables. *Studia Math.*, 123(1):15–42.
- Jameson, G. J. O. (2015). A simple proof of Stirling’s formula for the gamma function. *Math. Gaz.*, 99(544):68–74.
- Joly, E., Lugosi, G., and Oliveira, R. I. (2017). On the estimation of the mean of a random vector. *Electron. J. Stat.*, 11(1):440–451.
- Klein, T. and Rio, E. (2005). Concentration around the mean for maxima of empirical processes. *Ann. Probab.*, 33(3):1060–1077.
- Kolesko, K. and Latała, R. (2015). Moment estimates for chaoses generated by symmetric random variables with logarithmically convex tails. *Statist. Probab. Lett.*, 107:210–214.
- Kuchibhotla, A. K., Brown, L. D., Buja, A., George, E. I., and Zhao, L. (2018). A model free perspective for linear regression: Uniform-in-model bounds for post selection inference. *ArXiv e-prints:1802.05801*.
- Kuchibhotla, A. K. and Chakraborty, A. (2018). Moving beyond sub-gaussianity in high-dimensional statistics: Applications in covariance estimation and linear regression. *arXiv preprint arXiv:1804.02605*.
- Kuchibhotla, A. K. and Patra, R. K. (2019). On least squares estimation under heteroscedastic and heavy-tailed errors. *arXiv preprint arXiv:1909.02088*.
- Latała, R. (1997). Estimation of moments of sums of independent real random variables. *Ann. Probab.*, 25(3):1502–1513.
- Latała, R. (1999). Tail and moment estimates for some types of chaos. *Studia Math.*, 135(1):39–53.
- Lecué, G. and Mendelson, S. (2017). Sparse recovery under weak moment assumptions. *J. Eur. Math. Soc.(JEMS)*, 19(3):881–904.
- Lederer, J. and van de Geer, S. (2014). New concentration inequalities for suprema of empirical processes. *Bernoulli*, 20(4):2020–2038.

- Ledoux, M. and Talagrand, M. (1991). *Probability in Banach spaces: Isoperimetry and processes*, volume 23. Springer-Verlag, Berlin.
- Loh, P.-L. and Wainwright, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *Ann. Statist.*, 40(3):1637–1664.
- Major, P. (2005). Tail behaviour of multiple random integrals and u-statistics. *Probability Surveys*, 2:448–505.
- Mammen, E., Rothe, C., and Schienle, M. (2012). Nonparametric regression with nonparametrically generated covariates. *Ann. Statist.*, 40(2):1132–1170.
- Mammen, E., Rothe, C., and Schienle, M. (2013). Generated covariates in nonparametric estimation: A short review. In *Recent developments in modeling and applications in statistics*, pages 97–105. Springer, Heidelberg.
- McKeague, I. W. and Qian, M. (2015). An adaptive resampling test for detecting the presence of significant predictors. *J. Amer. Statist. Assoc.*, 110(512):1422–1433.
- Merlevède, F., Peligrad, M., and Rio, E. (2011). A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probab. Theory Related Fields*, 151(3-4):435–474.
- Negahban, S. N., Ravikumar, P., Wainwright, M. J., and Yu, B. (2010). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. *ArXiv e-prints:1010.2731*.
- Negahban, S. N., Ravikumar, P., Wainwright, M. J., and Yu, B. (2012). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. *Statist. Sci.*, 27(4):538–557.
- Oliveira, R. I. (2013). The lower tail of random quadratic forms, with applications to ordinary least squares and restricted eigenvalue properties. *ArXiv e-prints:1312.2903*.
- Plan, Y. and Vershynin, R. (2013). One-bit compressed sensing by linear programming. *Comm. Pure Appl. Math.*, 66(8):1275–1297.
- Pollard, D. (2002). Maximal inequalities via bracketing with adaptive truncation. *Ann. Inst. H. Poincaré Probab. Statist.*, 38(6):1039–1052.
- Portnoy, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.*, 16(1):356–366.
- Raskutti, G., Wainwright, M. J., and Yu, B. (2010). Restricted eigenvalue properties for correlated Gaussian designs. *J. Mach. Learn. Res.*, 11:2241–2259.
- Raskutti, G., Wainwright, M. J., and Yu, B. (2011). Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE Trans. Inform. Theory*, 57(10):6976–6994.
- Rudelson, M. and Vershynin, R. (2008). On sparse reconstruction from Fourier and Gaussian measurements. *Comm. Pure Appl. Math.*, 61(8):1025–1045.

- Rudelson, M. and Zhou, S. (2013). Reconstruction from anisotropic random measurements. *IEEE Trans. Inform. Theory*, 59(6):3434–3447.
- Talagrand, M. (2014). *Upper and lower bounds for stochastic processes: Modern methods and classical problems*, volume 60. Springer, Heidelberg.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.
- van de Geer, S. (2016). *Estimation and testing under sparsity*. Lecture Notes in Mathematics. Springer. Saint-Flour Probability Summer School.
- van de Geer, S. and Bühlmann, P. (2009). On the conditions used to prove oracle results for the Lasso. *Electron. J. Stat.*, 3:1360–1392.
- van de Geer, S. and Lederer, J. (2013). The Bernstein-Orlicz norm and deviation inequalities. *Probab. Theory Related Fields*, 157(1-2):225–250.
- van de Geer, S. and Muro, A. (2014). On higher order isotropy conditions and lower bounds for sparse quadratic forms. *Electron. J. Stat.*, 8(2):3031–3061.
- van de Geer, S. A. (2000). *Empirical Processes in M-Estimation*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- van der Vaart, A. and Wellner, J. A. (2011). A local maximal inequality under uniform entropy. *Electron. J. Stat.*, 5:192–203.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York.
- Vershynin, R. (2012). How close is the sample covariance matrix to the actual covariance matrix? *J. Theoret. Probab.*, 25(3):655–686.
- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press.
- Vidaurre, D., Bielza, C., and Larrañaga, P. (2013). A survey of L_1 regression. *Int. Stat. Rev.*, 81(3):361–387.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press.
- Wellner, J. A. (2017). The Bennett-Orlicz norm. *Sankhya A*, 79(2):355–383.
- Wong, K. C., Li, Z., and Tewari, A. (2020). Lasso guarantees for β -mixing heavy tailed time series. *Ann. Statist.*, 48(2):1124–1142.
- Wu, W. B. (2005). Nonlinear system theory: another look at dependence. *Proc. Natl. Acad. Sci. USA*, 102(40):14150–14154.
- Yu, G., Bien, J., and Tibshirani, R. (2019). Reluctant interaction modeling. *arXiv preprint arXiv:1907.08414*.

B Norms of Supremum of Empirical Processes

In this section, we present tail and norm bounds for the supremum of empirical processes with certain tail bounds on the envelope function. To avoid any issues about measurability, we follow the convention of [Talagrand \(2014\)](#) and define

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] := \sup \left\{ \mathbb{E} \left[\sup_{t \in S} X_t \right] : S \subseteq T \text{ is finite} \right\},$$

for any stochastic process $\{X_t\}$ indexed by $t \in T$ for some set T ; see Equation (2.2) of [Talagrand \(2014\)](#). Using this convention, we can define the g -Orlicz norm of the supremum as

$$\left\| \sup_{t \in T} X_t \right\|_g := \inf \left\{ C > 0 : \mathbb{E} \left[g \left(\left| \sup_{t \in S} \frac{X_t}{C} \right| \right) \right] \leq 1 \text{ for all } S \subseteq T \text{ finite} \right\}. \quad (\text{B.1})$$

The setting for all the results in this section is as follows. Let X_1, X_2, \dots, X_n be independent random variables with values in a measurable space $(\mathcal{X}, \mathcal{B})$ and \mathcal{F} is a class of measurable functions $f : \mathcal{X} \rightarrow \mathbb{R}$ such that $\mathbb{E}f(X_i) = 0$ for all $f \in \mathcal{F}$. Define

$$Z := \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n f(X_i) \right| \quad \text{and} \quad \Sigma_n(\mathcal{F}) := \sup_{f \in \mathcal{F}} \sum_{i=1}^n \mathbb{E} [f^2(X_i)]. \quad (\text{B.2})$$

Without loss of generality, we can assume that \mathcal{F} is finite, using (B.1). The final result will not depend on the cardinality of \mathcal{F} implying the result by (B.1). Based on the Generalized Bernstein-Orlicz norm and the generic chaining proof techniques in Section 10.2 of [Talagrand \(2014\)](#) and Section 5 of [Dirksen \(2015\)](#), one can obtain “optimal” tail bounds on the supremum of the empirical processes under a sub-Weibull envelope assumption in terms of the γ -functionals of [Talagrand \(2014\)](#). These bounds, however, require computation of the complexity of \mathcal{F} in terms of two distances and this can be hard in many examples of interest. For this reason, we first provide deviation bounds, and then bounds on the expectation (maximal inequalities), in terms of uniform covering and bracketing numbers. The proofs of all results in this section are given in Appendix F.

Before proceeding to unbounded function classes, we first state a result that provides a moment bound for the supremum of a bounded empirical process. This is essentially the Talagrand’s inequality for empirical processes. The result is based on Theorem 3.3.16 of [Giné and Nickl \(2016\)](#) and is given with explicit constants to resemble the Bernstein’s inequality for real-valued random variables; see also Theorem 1.1 and Lemma 3.4 of [Klein and Rio \(2005\)](#).

Proposition B.1. *Suppose \mathcal{F} is a class of uniformly bounded measurable functions $f : \mathcal{X} \rightarrow [-U, U]$ for some $U < \infty$. Then, under the setting above, for $p \geq 1$,*

$$\|Z\|_p \leq \mathbb{E}[Z] + p^{1/2} (2\Sigma_n(\mathcal{F}) + 4U\mathbb{E}[Z])^{1/2} + 6Up. \quad (\text{B.3})$$

Proposition B.1 can now be extended to possibly unbounded empirical processes using the proof of Theorem 4 of [Adamczak \(2008\)](#) and this is in lines with our use of the technique in the proofs of Theorems 3.2 and 3.3. Set

$$F(X_i) := \sup_{f \in \mathcal{F}} |f(X_i)| \quad \text{for } 1 \leq i \leq n \quad \text{and} \quad \rho := 8\mathbb{E} \left[\max_{1 \leq i \leq n} |F(X_i)| \right].$$

The function $F(\cdot)$ is called the envelope function of \mathcal{F} . Define the truncated part and the remaining unbounded part of Z as

$$\begin{aligned} Z_1 &:= \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \left(f(X_i) \mathbb{1}\{|f(X_i)| \leq \rho\} - \mathbb{E}[f(X_i) \mathbb{1}\{|f(X_i)| \leq \rho\}] \right) \right|, \\ Z_2 &:= \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \left(f(X_i) \mathbb{1}\{|f(X_i)| > \rho\} - \mathbb{E}[f(X_i) \mathbb{1}\{|f(X_i)| > \rho\}] \right) \right|. \end{aligned} \quad (\text{B.4})$$

Theorem B.1. *Suppose, for some $\alpha, K > 0$,*

$$\max_{1 \leq i \leq n} \left\| \sup_{f \in \mathcal{F}} |f(X_i)| \right\|_{\psi_\alpha} \leq K < \infty.$$

Then, under the notation outlined above, for $\alpha_ = \min\{\alpha, 1\}$ and $p \geq 2$,*

$$\|Z\|_p \leq 2\mathbb{E}[Z_1] + \sqrt{2}p^{1/2}\Sigma_n^{1/2}(\mathcal{F}) + C_\alpha p^{1/\alpha_*} \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_\alpha}, \quad (\text{B.5})$$

and

$$\|(Z - 2e\mathbb{E}[Z_1])_+\|_{\Psi_{\alpha_*, L_n(\alpha)}} \leq 3\sqrt{2}e\Sigma_n^{1/2}(\mathcal{F}), \quad (\text{B.6})$$

where

$$\begin{aligned} C_\alpha &:= 3\sqrt{2\pi}(1/\alpha_*)^{1/\alpha_*} K_{\alpha_*} \left[8 + (\log 2)^{1/\alpha_* - 1} \right], \\ L_n(\alpha) &:= \frac{9^{1/\alpha_*} C_\alpha}{3\sqrt{2}} \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_\alpha} \Sigma_n^{-1/2}(\mathcal{F}). \end{aligned}$$

Here the constant K_{α_} is the one used in Theorem 6.21 of [Ledoux and Talagrand \(1991\)](#).*

Remark B.1 It is clear that this result reduces to Theorems 3.2 and 3.3 if the function class \mathcal{F} contains only one function. Note that in this case, $\mathbb{E}[Z_1]$ is bounded by $\Sigma_n^{1/2}(\mathcal{F})$. There are two differences of Theorem B.1 in comparison with Theorem 4 of [Adamczak \(2008\)](#). Firstly, our result allows for the full range $\alpha \in (0, \infty)$ instead of just $\alpha \in (0, 1]$. Secondly, our result *only* involves $\mathbb{E}[Z_1]$, that is, the expectation of the supremum of *bounded* empirical processes instead of $\mathbb{E}[Z]$. This allows us to use many of the existing maximal inequalities for supremum of bounded empirical processes for the study of unbounded empirical processes as well. Also, it is interesting to note that using the bound on $\mathbb{E}[Z_1]$, and the moment bound (B.5), we can bound $\mathbb{E}[Z]$. This is similar to the results in Section 5 of [Chernozhukov et al. \(2014\)](#). \diamond

Remark B.2 The proof technique as mentioned above is truncation and using the Talagrand's inequality for the truncated part. We have taken this proof technique from [Adamczak \(2008\)](#). Even if the envelope function does not satisfy a ψ_α -norm bound, this part of the proof works. The moment bounds for the remaining unbounded part have to be obtained under whatever moment assumption the envelope function satisfies. This was done in [Lederer and van de Geer \(2014\)](#) under polynomial tails of the envelope function. The dominating term even in their bounds resemble the asymptotic Gaussian behavior as do ours. \diamond

The application of Theorem B.1 only requires bounding $\mathbb{E}[Z_1]$, the expectation of the supremum of a bounded empirical process. Most of the maximal inequalities available in the literature apply to this case. The following two results provide such inequalities based on uniform entropy and bracketing entropy (defined below). There are many classes for which uniform covering and bracketing numbers are available and these can be found in [van der Vaart and Wellner \(1996\)](#). We only give these inequalities for bounded classes and explicitly show the dependence on the bound (which in our case may increase with the sample size). In the following, we use the classical empirical processes notation. For any function f , define the linear operator

$$\mathbb{G}_n(f) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \{f(X_i) - \mathbb{E}[f(X_i)]\}.$$

Note here that we allow for non-identically distributed random variables X_1, X_2, \dots, X_n .

Given a metric or a pseudo-metric space (T, d) with metric d , for any $\epsilon > 0$, its covering number $N(\epsilon, T, d)$ is defined as the smallest number of balls of d -radius ϵ needed to cover T . More precisely, $N(\epsilon, T, d)$ is the smallest m such that there exists $t_1, t_2, \dots, t_m \in T$ satisfying

$$\sup_{t \in T} \inf_{1 \leq j \leq m} d(t, t_j) \leq \epsilon.$$

For any function class \mathcal{F} with envelope function F , the uniform entropy integral is defined for $\delta > 0$ as

$$J(\delta, \mathcal{F}, \|\cdot\|_2) := \sup_Q \int_0^\delta \sqrt{\log(2N(x \|F\|_{2,Q}, \mathcal{F}, \|\cdot\|_{2,Q}))} dx,$$

where the supremum is taken over all discrete probability measures Q and $\|h\|_{2,Q}$ denotes the $\|\cdot\|_2$ -norm of h with respect to the probability measure Q , that is, $\|h\|_{2,Q}^2 := \mathbb{E}_Q[h^2]$. To provide explicit constants we use Theorem 3.5.1 of [Giné and Nickl \(2016\)](#) along with Theorem 2.1 of [van der Vaart and Wellner \(2011\)](#).

Proposition B.2. *Suppose \mathcal{F} is a class of measurable functions with envelope function F satisfying $\|F\|_\infty \leq U < \infty$. Assume that \mathcal{F} contains the zero function. Then*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right] \leq 16\sqrt{2} \|F\|_{2,P} J(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_2) \left[1 + \frac{128\sqrt{2}U J(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_2)}{\sqrt{n}\delta_n^2(\mathcal{F}) \|F\|_{2,P}} \right],$$

where $\Sigma_n(\mathcal{F})$ is as defined in (B.2),

$$\|F\|_{2,P}^2 := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[F^2(X_i)], \quad \text{and} \quad \delta_n^2(\mathcal{F}) := \frac{\Sigma_n(\mathcal{F})}{n \|F\|_{2,P}^2}.$$

The following proposition proves an alternative to Proposition B.2 using bracketing numbers. For $\epsilon > 0$, let the set $\{[f_j^L, f_j^U] : 1 \leq j \leq N_\epsilon\}$ represents the minimal ϵ -bracketing set of \mathcal{F} with respect to $\|\cdot\|_{2,P}$ -norm if for any $f \in \mathcal{F}$, there exists an $1 \leq I \leq N_\epsilon$ such that for all x ,

$$f_I^L(x) \leq f(x) \leq f_I^U(x) \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[|f_I^U(X_i) - f_I^L(X_i)|^2] \leq \epsilon^2.$$

The number N_ϵ is the ϵ -bracketing number, usually denoted by $N_{[\cdot]}(\epsilon, \mathcal{F}, \|\cdot\|_{2,P})$. Define the bracketing entropy integral as

$$J_{[\cdot]}(\eta, \mathcal{F}, \|\cdot\|_{2,P}) := \int_0^\eta \sqrt{\log\left(2N_{[\cdot]}(x, \mathcal{F}, \|\cdot\|_{2,P})\right)} dx \quad \text{for } \eta > 0.$$

The following proposition is very similar to Proposition 3.4.2 of [van der Vaart and Wellner \(1996\)](#) and we provide it here with explicit constants allowing for non-identically distributed random variables. The proof follows that of Theorem 3.5.13 and Proposition 3.5.15 of [Giné and Nickl \(2016\)](#) and we do not repeat the proof except for necessary changes. Also, see Theorem 6 of [Pollard \(2002\)](#).

Proposition B.3. *Suppose \mathcal{F} is a class of measurable functions with envelope function F satisfying $\|F\|_\infty \leq U < \infty$. Then*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right] \leq 2J_{[\cdot]}(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_{2,P}) \left[58 + \frac{J_{[\cdot]}(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_{2,P}) U}{\sqrt{n} \delta_n^2(\mathcal{F})} \right],$$

for any $\delta_n(\mathcal{F})$ satisfying $\delta_n(\mathcal{F}) \geq \Sigma_n^{1/2}(\mathcal{F})/\sqrt{n}$ with $\Sigma_n(\mathcal{F})$ as in [\(B.2\)](#).

For the sake of completeness, we provide one last result relating the expectation of the unbounded supremum Z in terms of the expectation of the supremum Z_1 of a bounded empirical process. Theorem [B.1](#) provides such a result under a sub-Weibull envelope assumption, while the following result applies in general.

Proposition B.4. *Under the notation outlined before Theorem [B.1](#), we have*

$$\mathbb{E}[Z] \leq \mathbb{E}[Z_1] + 8\mathbb{E} \left[\max_{1 \leq i \leq n} F(X_i) \right].$$

Remark B.3 We note that only a sample of empirical process results are presented here. For many applications, the results on the statistic Z in [\(B.2\)](#) are not sufficient. The main reason for this is that these results do not allow for function dependent scaling. For example, if the variance of $\sum f(X_i)$ varies too much as f varies over \mathcal{F} , then it is desirable to obtain bounds for

$$\sup_{f \in \mathcal{F}} \left(\sum_{i=1}^n \mathbb{E}[f^2(X_i)] \right)^{-1/2} \left| \sum_{i=1}^n f(X_i) \right|.$$

This arises in uniform-in-bandwidth results related to linear kernel averages; see Theorem 1 of [Einmahl and Mason \(2005\)](#) for a precise problem. The derivation there is based on a well-known technique called the peeling device introduced by [Alexander \(1985\)](#); see [van de Geer \(2000, page 70\)](#) for more details. More general function dependent scalings in empirical processes are considered in [Giné and Koltchinskii \(2006\)](#). In both these works, the functions are taken to be uniformly bounded and extensions to sub-Weibull random variables are desirable. The problems above have a non-random function dependent scaling and there are also some interesting problems involving a random function dependent scaling. One ready example of this is related to the Nadaraya-Watson kernel smoothing estimator of the conditional expectation which is a ratio of two linear kernel averages. We hope to explore some of these in the future. \diamond

C Proofs of All Results in Section 2 and Appendix A

Proof of Proposition A.1. It is clear from the definition of $\phi_{\alpha,L}(\cdot)$ that

$$\phi_{\alpha,L}^{-1}(t) = \max \left\{ \sqrt{\log(1+t)}, L(\log(1+t))^{1/\alpha} \right\} \quad \text{for all } t \geq 0.$$

It follows that for all $t \geq 0$,

$$\phi_{\alpha,L}^{-1}(t) \leq \Psi_{\alpha,L}^{-1}(t) \leq 2\phi_{\alpha,L}^{-1}(t). \quad (\text{C.1})$$

Hence for all $x \geq 0$,

$$\phi_{\alpha,L}(x/2) \leq \Psi_{\alpha,L}(x) \leq \phi_{\alpha,L}(x). \quad (\text{C.2})$$

The result now follows by Definition 2.1. \square

Proof of Proposition A.2. (a) If $\|Y\|_{\Psi_{\alpha,L}} = \infty$, then the result is trivially true. If $\delta = \|Y\|_{\Psi_{\alpha,L}} < \infty$, then for $\eta > \delta$,

$$\mathbb{E} \left[\Psi_{\alpha,L} \left(\frac{|Y|}{\eta} \right) \right] \leq 1 \quad \Rightarrow \quad \mathbb{E} \left[\Psi_{\alpha,L} \left(\frac{|X|}{\eta} \right) \right] \leq 1.$$

Letting $\eta \downarrow \delta$ implies the result.

(b) The result trivially holds if $\|X\|_{\Psi_{\alpha,L}} = \infty$. Assume $\|X\|_{\Psi_{\alpha,L}} < \infty$. It is clear from the definition (2.5) of $\Psi_{\alpha,L}^{-1}(t)$,

$$\Psi_{\alpha,L}^{-1}(t) \leq \Psi_{\alpha,K}^{-1}(t) \quad \text{for all } t \geq 0.$$

Observe that for $\eta > \|X\|_{\Psi_{\alpha,L}}$,

$$\begin{aligned} \mathbb{E} [\Psi_{\alpha,K}(|X|/\eta)] &= \int_0^\infty \mathbb{P} \left(|X| \geq \eta \Psi_{\alpha,K}^{-1}(t) \right) dt \\ &\leq \int_0^\infty \mathbb{P} \left(|X| \geq \eta \Psi_{\alpha,L}^{-1}(t) \right) dt = \mathbb{E} \left[\Psi_{\alpha,L} \left(\frac{|X|}{\eta} \right) \right] \leq 1. \end{aligned}$$

Letting $\eta \downarrow \|X\|_{\Psi_{\alpha,L}}$ implies the result. \square

Proof of Proposition A.3. From definitions (2.1) and (2.5), for $\eta > \delta$,

$$\begin{aligned} \mathbb{P} \left(|X| \geq \eta \left[\sqrt{t} + Lt^{1/\alpha} \right] \right) &= \mathbb{P} \left(\frac{|X|}{\eta} \geq \Psi_{\alpha,L}^{-1}(e^t - 1) \right) \\ &= \mathbb{P} \left(\Psi_{\alpha,L}(|X|/\eta) + 1 \geq e^t \right) \\ &\leq \left(\mathbb{E} [\Psi_{\alpha,L}(|X|/\eta)] + 1 \right) \exp(-t) \leq 2 \exp(-t). \end{aligned}$$

Now taking limit as $\eta \downarrow \delta$ implies the first part of the result.

For the converse result, set $c(\alpha) = 3^{1/\alpha-1/2}$. Observe that

$$\begin{aligned} \mathbb{E} \left[\Psi_{\alpha, c(\alpha)L} \left(\frac{|X|}{\sqrt{3}\delta} \right) \right] &= \int_0^\infty \mathbb{P} \left(|X| \geq \sqrt{3}\delta \Psi_{\alpha, c(\alpha)L}^{-1}(t) \right) dt \\ &= \int_0^\infty \mathbb{P} \left(|X| \geq \sqrt{3}\delta \left\{ \sqrt{\log(1+t)} + c(\alpha)L (\log(1+t))^{1/\alpha} \right\} \right) dt \\ &= \int_0^\infty \mathbb{P} \left(|X| \geq \delta \left\{ \sqrt{\log(1+t)^3} + L (\log(1+t)^3)^{1/\alpha} \right\} \right) dt \\ &\leq 2 \int_0^\infty \frac{1}{(1+t)^3} dt \leq 1. \end{aligned}$$

This implies $\|X\|_{\alpha, c(\alpha)L} \leq \sqrt{3}\delta$ and completes the proof of the proposition. \square

Proof of Proposition A.4. For a proof of the first inequality in Proposition A.4, note that it holds trivially if $\|X\|_{\Psi_{\alpha,L}} = \infty$. Assume $\delta := \|X\|_{\Psi_{\alpha,L}} < \infty$. Fix $\eta > \delta$. From the hypothesis and inequality (C.2),

$$\mathbb{E} [\Psi_{\alpha,L}(|X|/\eta)] \leq 1 \quad \Rightarrow \quad \mathbb{E} \left[\exp \left(\min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right) - 1 \right] \leq 1.$$

Thus, for $p \geq 1$, (using the inequalities $x^p/p! \leq \exp(x) - 1$ and $(p!)^{1/p} \leq p$)

$$\left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_p \leq p. \quad (\text{C.3})$$

Now observe by the equivalence of inverse functions (C.1), for any $x \geq 0$

$$x \leq \Psi_{\alpha,L}^{-1}(\phi_{\alpha,L}(x)) = \left(\min \left\{ x^2, \left(\frac{x}{L} \right)^\alpha \right\} \right)^{1/2} + L \left(\min \left\{ x^2, \left(\frac{x}{L} \right)^\alpha \right\} \right)^{1/\alpha}. \quad (\text{C.4})$$

Taking $x = |X|/(2\eta)$ in (C.4) and using triangle inequality of $\|\cdot\|_p$ -norm,

$$\left\| \frac{X}{2\eta} \right\|_p \leq \left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_{\frac{p}{2}}^{\frac{1}{2}} + L \left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_{\frac{p}{\alpha}}^{\frac{1}{\alpha}}. \quad (\text{C.5})$$

If $p \geq \alpha$, then from (C.3)

$$\left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_{p/\alpha} \leq p/\alpha,$$

and for $1 \leq p < \alpha$,

$$\left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_{p/\alpha}^{1/\alpha} \leq \left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_1 \leq 1 \leq p^{1/\alpha}.$$

Combining these two inequalities, we get for $p \geq 1$,

$$\left\| \min \left\{ \left(\frac{|X|}{2\eta} \right)^2, \left(\frac{|X|}{2\eta L} \right)^\alpha \right\} \right\|_{p/\alpha}^{1/\alpha} \leq p^{1/\alpha} \max \{1, (1/\alpha)^{1/\alpha}\}. \quad (\text{C.6})$$

A similar inequality holds with $(p/\alpha, 1/\alpha)$ replaced by $(p/2, 1/2)$. Substituting inequality (C.6) in (C.5), it follows for $p \geq 1$ that

$$\|X\|_p \leq (2\eta) \left[\sqrt{p} + Lp^{1/\alpha} \max\{1, (1/\alpha)^{1/\alpha}\} \right].$$

Therefore by letting $\eta \downarrow \delta$, for $p \geq 1$,

$$\|X\|_p \leq 2 \|X\|_{\Psi_{\alpha,L}} \sqrt{p} + 2L \|X\|_{\Psi_{\alpha,L}} p^{1/\alpha} \max\{1, (1/\alpha)^{1/\alpha}\},$$

or equivalently,

$$\frac{1}{2} \min\{1, \alpha^{1/\alpha}\} \sup_{p \geq 1} \frac{\|X\|_p}{\sqrt{p} + Lp^{1/\alpha}} \leq \|X\|_{\Psi_{\alpha,L}}.$$

Converse: For a proof of the second inequality in Proposition A.4, set

$$\Delta := \sup_{p \geq 1} \frac{\|X\|_p}{\sqrt{p} + Lp^{1/\alpha}},$$

so that

$$\|X\|_p \leq \Delta \sqrt{p} + L\Delta p^{1/\alpha} \quad \text{for all } p \geq 1.$$

Note by Markov's inequality and these moment bounds that for any $t \geq 1$,

$$\mathbb{P} \left(|X| \geq e\Delta \sqrt{t} + eL\Delta t^{1/\alpha} \right) \leq \exp(-t),$$

and for $0 < t < 1$ (trivially),

$$\mathbb{P} \left(|X| \geq e\Delta \sqrt{t} + eL\Delta t^{1/\alpha} \right) \leq 1.$$

Hence, for any $t > 0$,

$$\mathbb{P} \left(|X| \geq e\Delta \sqrt{t} + eL\Delta t^{1/\alpha} \right) \leq e \exp(-t). \quad (\text{C.7})$$

Take $K = e \max\{2, 4^{1/\alpha}\}$. Observe that,

$$\begin{aligned} \mathbb{E} \left[\Psi_{\alpha,L} \left(\frac{|X|}{K\Delta} \right) \right] &= \int_0^\infty \mathbb{P} \left(|X| \geq K\Delta \Psi_{\alpha,L}^{-1}(t) \right) dt \\ &= \int_0^\infty \mathbb{P} \left(|X| \geq K\Delta \left\{ \sqrt{\log(1+t)} + L(\log(1+t))^{1/\alpha} \right\} \right) dt \\ &\leq \int_0^\infty \mathbb{P} \left(|X| \geq e\Delta \sqrt{\log(1+t)^4} + eL\Delta (\log(1+t)^4)^{1/\alpha} \right) dt \\ &\leq e \int_0^\infty \frac{1}{(1+t)^4} dt = e/3 < 1. \end{aligned}$$

Therefore, $\|X\|_{\Psi_{\alpha,L}} \leq K\Delta$. □

For the proofs of the results in Section 3, we use the following alternative result regarding inversion of moment bounds to get bounds on the GBO norm.

Proposition C.1. *If $\|X\|_p \leq C_1\sqrt{p} + C_2p^{1/\alpha}$, holds for $p \geq 1$ and some constants C_1, C_2 , then $\|X\|_{\Psi_{\alpha, K}} \leq 2eC_1$, where $K := 4^{1/\alpha}C_2/(2C_1)$.*

Proof. From the proof of Proposition A.4 (or, in particular (C.7)), we get

$$\mathbb{P}\left(|X| \geq eC_1\sqrt{t} + eC_2t^{1/\alpha}\right) \leq e \exp(-t), \quad \text{for all } t \geq 0.$$

Take $K = 4^{1/\alpha}C_2/(2C_1)$ as in the statement of the result. Observe that with $\delta := eC_1$,

$$\begin{aligned} \mathbb{E}\left[\Psi_{\alpha, K}\left(\frac{|X|}{\sqrt{4\delta}}\right)\right] &= \int_0^\infty \mathbb{P}\left(|X| \geq \sqrt{4\delta}\Psi_{\alpha, K}^{-1}(t)\right) dt \\ &= \int_0^\infty \mathbb{P}\left(|X| \geq \sqrt{4\delta}\left\{\sqrt{\log(1+t)} + K(\log(1+t))^{1/\alpha}\right\}\right) dt \\ &= \int_0^\infty \mathbb{P}\left(|X| \geq eC_1\sqrt{\log(1+t)^4} + eC_2L(\log(1+t)^4)^{1/\alpha}\right) dt \\ &\leq e \int_0^\infty \frac{1}{(1+t)^4} dt = e/3 < 1. \end{aligned}$$

Therefore, $\|X\|_{\Psi_{\alpha, K}} \leq 2eC_1$. □

Proof of Proposition A.5. Assume without loss of generality that $\|X_i\|_{\Psi_{\alpha, L}} < \infty$ for all $1 \leq i \leq n$, as otherwise the result is trivially true. If $\alpha > 1$, then $\Psi_{\alpha, L}^{-1}(\cdot)$ is a concave function and hence $\|\cdot\|_{\Psi_{\alpha, L}}$ is a proper norm proving the result. For $\alpha < 1$, the result follows trivially by noting that both sides of the inequality in Proposition A.4 are norms. This completes the proof. □

Proof of Proposition A.6. By union bound and Proposition A.3,

$$\begin{aligned} &\mathbb{P}\left(\max_{1 \leq j \leq N} |X_j| \geq \Delta \left\{\sqrt{t + \log N} + L(t + \log N)^{1/\alpha}\right\}\right) \\ &\leq \sum_{j=1}^N \mathbb{P}\left(|X_j| \geq \Delta \left\{\sqrt{t + \log N} + L(t + \log N)^{1/\alpha}\right\}\right) \\ &\leq 2N \exp(-t - \log N) \leq \frac{2N}{N} \exp(-t). \end{aligned}$$

Hence the tail bound follows. To bound the norm note that for all $\alpha > 0$,

$$(t + \log N)^{1/\alpha} \leq M(\alpha) \left(t^{1/\alpha} + (\log N)^{1/\alpha}\right),$$

and from the tail bound of the maximum,

$$\begin{aligned} &\mathbb{P}\left(Z \geq \delta \left\{\sqrt{t} + M(\alpha)Lt^{1/\alpha}\right\}\right) \\ &\leq \mathbb{P}\left(\max_{1 \leq j \leq N} |X_j| \geq \Delta \left\{\sqrt{t + \log N} + L(t + \log N)^{1/\alpha}\right\}\right) \leq 2 \exp(-t), \end{aligned}$$

where

$$Z := \left(\max_{1 \leq j \leq N} |X_j| - \Delta \left\{ \sqrt{\log N} + M(\alpha)L(\log N)^{1/\alpha} \right\} \right)_+.$$

Hence by Proposition A.3, $\|Z\|_{\Psi_{\alpha, K(\alpha)}} \leq \sqrt{3}\Delta$. The result follows by Proposition A.5 along with the fact

$$\max_{1 \leq j \leq N} |X_j| \leq Z + \Delta \left\{ \sqrt{\log N} + M(\alpha)L(\log N)^{1/\alpha} \right\},$$

and by noting that the random variables on both sides are non-negative. \square

Proof of Proposition A.7. By homogeneity, we can without loss of generality assume that

$$\|X_k\|_{\Psi_{\alpha, L}} = 1 \quad \text{for all } k \geq 1.$$

Note that by the union bound, for $t \geq 0$,

$$\begin{aligned} & \mathbb{P} \left(\sup_{k \geq 1} \left(|X_k| - \sqrt{2 \log(1+k)} - M(\alpha)L(2 \log(1+k))^{1/\alpha} \right)_+ \geq \sqrt{t} + M(\alpha)Lt^{1/\alpha} \right) \\ & \leq \sum_{k \geq 1} \mathbb{P} \left(|X_k| \geq \sqrt{t + 2 \log(1+k)} + L(t + 2 \log(1+k))^{1/\alpha} \right) \\ & \leq \sum_{k \geq 1} \frac{2}{(1+k)^2} \exp(-t) \\ & \leq \frac{2(\pi^2 - 6)}{6} \exp(-t) < 2 \exp(-t). \end{aligned}$$

Hence by Proposition A.3,

$$\left\| \sup_{k \geq 1} \left(|X_k| - \sqrt{2 \log(1+k)} - M(\alpha)L(2 \log(1+k))^{1/\alpha} \right)_+ \right\|_{\Psi_{\alpha, c(\alpha)M(\alpha)L}} \leq \sqrt{3}. \quad (\text{C.8})$$

Recall $c(\alpha) = 3^{1/\alpha}/\sqrt{3}$. Since $\sqrt{2 \log(1+k)} \geq 1$ for $k \geq 1$, using (C.8), it follows that

$$\left\| \sup_{k \geq 1} \left(\frac{|X_k|}{\sqrt{2 \log(1+k)} + M(\alpha)L(2 \log(1+k))^{1/\alpha}} - 1 \right)_+ \right\|_{\Psi_{\alpha, c(\alpha)M(\alpha)L}} \leq 1.5.$$

The result now follows by an application of Proposition A.5. \square

D Proofs of All Results in Section 3

Proof of Theorem 3.1. Since $a_i X_i = (a_i \|X_i\|_{\psi_\alpha})(X_i / \|X_i\|_{\psi_\alpha})$, we can without loss of generality assume $\|X_i\|_{\psi_\alpha} = 1$. Define $Y_i = (|X_i| - \eta)_+$ with $\eta = (\log 2)^{1/\alpha}$. This implies that

$$\mathbb{P}(|X_i| \geq t) \leq 2 \exp(-t^\alpha) \quad \Rightarrow \quad \mathbb{P}(Y_i \geq t) \leq \exp(-t^\alpha). \quad (\text{D.1})$$

By symmetrization inequality (Proposition 6.3 of Ledoux and Talagrand (1991)),

$$\left\| \sum_{i=1}^n a_i X_i \right\|_p \leq 2 \left\| \sum_{i=1}^n \varepsilon_i a_i X_i \right\|_p,$$

for independent Rademacher random variables $\varepsilon_i, 1 \leq i \leq n$ independent of $X_i, 1 \leq i \leq n$. Using the fact that $\varepsilon_i X_i$ is identically distributed as $\varepsilon_i |X_i|$ and by Theorem 1.3.1 of [de la Peña and Giné \(1999\)](#), it follows that

$$\begin{aligned} \left\| \sum_{i=1}^n a_i X_i \right\|_p &\leq 2 \left\| \sum_{i=1}^n \varepsilon_i a_i |X_i| \right\|_p \leq 2 \left\| \sum_{i=1}^n \varepsilon_i a_i (\eta + Y_i) \right\|_p \\ &\leq 2 \left\| \sum_{i=1}^n \varepsilon_i a_i Y_i \right\|_p + 2\eta \left\| \sum_{i=1}^n \varepsilon_i a_i \right\|_p \\ &\leq 2 \left\| \sum_{i=1}^n \varepsilon_i a_i Y_i \right\|_p + 2\eta \sqrt{p} \|a\|_2. \end{aligned} \quad (\text{D.2})$$

By inequality [\(D.1\)](#),

$$\left\| \sum_{i=1}^n a_i \varepsilon_i Y_i \right\|_p \leq \left\| \sum_{i=1}^n a_i Z_i \right\|_p,$$

for symmetric independent random variables $Z_i, 1 \leq i \leq n$ satisfying $\mathbb{P}(|Z_i| \geq t) = \exp(-t^\alpha)$ for all $t \geq 0$. Now we apply the bound in Examples 3.2 and 3.3 of [Latała \(1997\)](#) in combination with Theorem 2 there.

Case $\alpha \leq 1$: Example 3.3 of [Latała \(1997\)](#) shows that for $p \geq 2$,

$$\left\| \sum_{i=1}^n a_i Z_i \right\|_p \leq \max \left\{ p^{1/2} \sqrt{2} \|a\|_2 2^{1/\alpha}, \frac{p^{1/\alpha} \|a\|_p}{\exp(1/(2e))} \right\} \frac{e^3 (2\pi)^{1/4} e^{1/24}}{\alpha^{1/\alpha}}.$$

Using the proof of Corollary 1.2 of [Bogucki \(2015\)](#) and substituting the resulting bound in [\(D.2\)](#), it follows that for $p \geq 2$,

$$\left\| \sum_{i=1}^n a_i X_i \right\|_p \leq \sqrt{8} e^3 (2\pi)^{1/4} e^{1/24} (e^{2/e}/\alpha)^{1/\alpha} \left[\sqrt{p} \|a\|_2 + p^{1/\alpha} \|a\|_\infty \right].$$

Corollary 1.2 of [Bogucki \(2015\)](#) uses the inequality $p^{1/p} \leq e$ but using $p^{1/p} \leq e^{1/e}$ gives the bound above; see also Remark 3, Equation (3) of [Kolesko and Latała \(2015\)](#). Set $C'(\alpha) := \sqrt{8} e^3 (2\pi)^{1/4} e^{1/24} (e^{2/e}/\alpha)^{1/\alpha}$. For $p = 1$ note that

$$\left\| \sum_{i=1}^n a_i X_i \right\|_1 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_2 \leq C'(\alpha) \left[\sqrt{2} \|a\|_1 + 2^{1/\alpha} \|a\|_\infty \right]$$

Thus for $p \geq 1$,

$$\left\| \sum_{i=1}^n a_i X_i \right\|_p \leq C'(\alpha) \max\{\sqrt{2}, 2^{1/\alpha}\} \left[\sqrt{p} \|a\|_2 + p^{1/\alpha} \|a\|_\infty \right].$$

Hence the result follows by Proposition [C.1](#).

Case $\alpha \geq 1$: it follows from (13) of Example 3.2 of [Latała \(1997\)](#) that for $p \geq 2$,

$$\left\| \sum_{i=1}^n a_i Z_i \right\|_p \leq 4e \left[\sqrt{p} \left(\sum_{i=1}^n a_i^2 \right)^{1/2} + p^{1/\alpha} \left(\sum_{i=1}^n |a_i|^\beta \right)^{1/\beta} \right],$$

with $\beta = \beta(\alpha)$ as mentioned in the statement. Therefore, for $p \geq 2$,

$$\left\| \sum_{i=1}^n a_i X_i \right\|_p \leq (4e + 2\eta) \sqrt{p} \|a\|_2 + 4ep^{1/\alpha} \|a\|_\beta.$$

For $p = 1$, note that

$$\left\| \sum_{i=1}^n a_i X_i \right\|_1 \leq \left\| \sum_{i=1}^n a_i X_i \right\|_2 \leq \max\{\sqrt{2}, 2^{1/\alpha}\} \left[(4e + 2\eta) \|a\|_2 + 4e \|a\|_\beta \right],$$

and so, for $p \geq 1$,

$$\left\| \sum_{i=1}^n a_i X_i \right\|_p \leq \max\{\sqrt{2}, 2^{1/\alpha}\} \left[(4e + 2\eta) \sqrt{p} \|a\|_2 + 4ep^{1/\alpha} \|a\|_\beta \right].$$

The result now follows by an application of Proposition [C.1](#). The tail bound follows from Proposition [A.3](#). \square

Before proving moment inequalities with unbounded variables, we first provide the Bernstein moment bounds for bounded random variables since this forms an integral part of our proofs.

Proposition D.1. (*Bernstein's Inequality for Bounded Random Variables*) Suppose Z_1, Z_2, \dots, Z_n are independent random variables with mean zero and uniformly bounded by U in absolute value. Then for $p \geq 1$,

$$\left\| \sum_{i=1}^n Z_i \right\|_p \leq \sqrt{6p} \left(\sum_{i=1}^n \mathbb{E}[Z_i^2] \right)^{1/2} + 10pU.$$

Proof of Proposition D.1. By Theorem 3.1.7 of [Giné and Nickl \(2016\)](#),

$$\mathbb{P}(|S_n| \geq t) \leq 2 \exp\left(-\frac{t^2}{2\sigma_n^2 + 2Ut/3}\right), \quad \text{for all } t \geq 0.$$

where

$$S_n := \sum_{i=1}^n Z_i, \quad \text{and} \quad \sigma_n^2 := \sum_{i=1}^n \mathbb{E}[Z_i^2].$$

To bound the moments of $|S_n|$, note that

(a) If $2Ut/3 \leq 2\delta\sigma_n^2$ (or equivalently, $t \leq 3\delta\sigma_n^2/U$), then

$$\exp\left(-\frac{t^2}{2\sigma_n^2 + 2Ut/3}\right) \leq \exp\left(-\frac{t^2}{2(1+\delta)\sigma_n^2}\right).$$

(b) If $2Ut/3 \geq 2\delta\sigma_n^2$, then

$$\exp\left(-\frac{t^2}{2\sigma_n^2 + 2Ut/3}\right) \leq \exp\left(-t\frac{3\delta}{2U(1+\delta)}\right).$$

Set $t_0 := 3\delta\sigma_n^2/U$. Now observe that for $p \geq 2$,

$$\begin{aligned} & \mathbb{E}[|S_n|^p] \\ &= \int_0^\infty pt^{p-1}\mathbb{P}(|S_n| \geq t) dt \\ &\leq 2 \int_0^\infty pt^{p-1} \exp\left(-\frac{t^2}{2\sigma_n^2 + 2Ut/3}\right) dt \\ &= 2 \int_0^{t_0} pt^{p-1} \exp\left(-\frac{t^2}{2(1+\delta)\sigma_n^2}\right) dt + 2 \int_{t_0}^\infty pt^{p-1} \exp\left(-\frac{3\delta t}{2U(1+\delta)}\right) dt \\ &\leq 2 \int_0^\infty pt^{p-1} \exp\left(-\frac{t^2}{2(1+\delta)\sigma_n^2}\right) dt + 2 \int_0^\infty pt^{p-1} \exp\left(-\frac{3\delta t}{2U(1+\delta)}\right) dt \\ &=: \mathbf{I} + \mathbf{II}. \end{aligned}$$

By a change of variable for **I**, we have

$$\begin{aligned} \mathbf{I} &= 2 \int_0^\infty pt^{p-1} \exp\left(-\frac{t^2}{2(1+\delta)\sigma_n^2}\right) dt \\ &= 2 \left(\sqrt{(1+\delta)\sigma_n^2}\right)^p \int_0^\infty pz^{p-1} \exp\left(-\frac{z^2}{2}\right) dz \\ &\stackrel{(1)}{=} 2 \left(\sqrt{(1+\delta)\sigma_n^2}\right)^p 2^{p/2} \Gamma\left(1 + \frac{p}{2}\right) \\ &\stackrel{(2)}{\leq} 2 \left(\sqrt{2(1+\delta)\sigma_n^2}\right)^p \sqrt{2\pi} \exp\left(-1 - \frac{p}{2}\right) \left(1 + \frac{p}{2}\right)^{(p+1)/2} \exp\left(\frac{1}{12(1+p/2)}\right). \end{aligned}$$

Inequality (1) above can be found in Exercise 3.3.4(a) of [Giné and Nickl \(2016\)](#) and inequality (2) follows from Theorem 1.1 of [Jameson \(2015\)](#). Simplifying the above bound for $p \geq 2$, we get

$$\begin{aligned} \mathbf{I}^{1/p} &\leq \sqrt{\frac{2(1+\delta)\sigma_n^2}{e}} \left(\frac{2\sqrt{2\pi}}{e}\right)^{1/p} \left(1 + \frac{p}{2}\right)^{\frac{1}{2} + \frac{1}{2p}} \exp\left(\frac{1}{12p(1+p/2)}\right) \\ &\leq \sqrt{\frac{2(1+\delta)}{e}} \left(\frac{2\sqrt{2\pi}}{e}\right)^{1/2} 2^{1/4} \exp\left(\frac{1}{48}\right) \sigma_n \sqrt{p} \leq 2.1\sigma_n \sqrt{p}, \quad (\text{taking } \delta = 1). \end{aligned}$$

To bound **II**, note that by change of variable

$$\begin{aligned} \mathbf{II} &= 2 \int_0^\infty pt^{p-1} \exp\left(-\frac{3\delta t}{2U(1+\delta)}\right) dt = 2 \left(\frac{2U(1+\delta)}{3\delta}\right)^p \int_0^\infty pz^{p-1} \exp(-z) dz \\ &\stackrel{(1')}{=} 2 \left(\frac{2U(1+\delta)}{3\delta}\right)^p \Gamma(1+p) \\ &\stackrel{(2')}{\leq} 2 \left(\frac{2U(1+\delta)}{3\delta}\right)^p \sqrt{2\pi}(1+p)^{p+\frac{1}{2}} \exp(-p-1) \exp\left(\frac{1}{12(p+1)}\right). \end{aligned}$$

Here again inequalities (1') and (2') follows from Exercise 3.3.4 (a) of [Giné and Nickl \(2016\)](#) and Theorem 1.1 of [Jameson \(2015\)](#), respectively. This implies for $p \geq 2$, that

$$\begin{aligned} \mathbf{II}^{1/p} &\leq \left(\frac{2eU(1+\delta)}{3\delta} \right) \left(\frac{2\sqrt{2\pi}}{e} \right)^{1/p} (1+p)^{1+\frac{1}{2p}} \exp\left(\frac{1}{12p(p+1)} \right) \\ &\leq \frac{2eU(1+\delta)}{3\delta} \left(\frac{2\sqrt{2\pi}}{e} \right)^{1/2} \left(\frac{3p}{2} \right) (1+p)^{1/(2p)} \exp\left(\frac{1}{72} \right) \\ &\leq \frac{2e(1+\delta)}{3\delta} \left(\frac{2\sqrt{2\pi}}{e} \right)^{1/2} \left(\frac{3}{2} \right) 3^{1/4} \exp\left(\frac{1}{72} \right) Up \leq 10Up, \end{aligned}$$

also for $\delta = 1$. Therefore, for $p \geq 1$,

$$(\mathbb{E}[|S_n|^p])^{1/p} \leq 2.1\sigma_n\sqrt{p} + 10pU \leq \sqrt{6p\sigma_n^2} + 10pU.$$

□

Proof of Theorem 3.2. The method of proof is a combination of truncation and Hoffmann-Jorgensen's inequality. Define

$$Z = \max_{1 \leq i \leq n} |X_i|, \quad \rho = 8\mathbb{E}[Z], \quad K = \max_{1 \leq i \leq n} \|X_i\|_{\psi_\alpha},$$

$$X_{i,1} = X_i \mathbb{1}\{|X_i| \leq \rho\} - \mathbb{E}[X_i \mathbb{1}\{|X_i| \leq \rho\}], \quad \text{and} \quad X_{i,2} = X_i - X_{i,1}.$$

It is clear that $X_i = X_{i,1} + X_{i,2}$ and $|X_{i,1}| \leq 2\rho$ for $1 \leq i \leq n$. Also by triangle inequality, for $p \geq 1$,

$$\left\| \sum_{i=1}^n X_i \right\|_p \leq \left\| \sum_{i=1}^n X_{i,1} \right\|_p + \left\| \sum_{i=1}^n X_{i,2} \right\|_p.$$

Now note that for $1 \leq i \leq n$,

$$\mathbb{E}[X_{i,1}^2] = \text{Var}(X_{i,1}) = \text{Var}(X_i \mathbb{1}\{|X_i| \leq \rho\}) \leq \mathbb{E}[X_i^2].$$

Thus, applying Bernstein's inequality (Proposition D.1), for $p \geq 1$,

$$\left\| \sum_{i=1}^n X_{i,1} \right\|_p \leq \sqrt{6p} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{1/2} + 20p\rho.$$

By Hoffmann-Jorgensen's inequality (Proposition 6.8 of [Ledoux and Talagrand \(1991\)](#)) and by the choice of ρ ,

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_1 \leq 2 \left\| \sum_{i=1}^n |X_i| \mathbb{1}\{|X_i| \geq \rho\} \right\|_1 \leq 16 \|Z\|_1,$$

since

$$\mathbb{P} \left(\max_{1 \leq k \leq n} \sum_{i=1}^k |X_i| \mathbb{1}\{|X_i| \geq \rho\} > 0 \right) \leq \mathbb{P}(Z \geq \rho) \leq 1/8.$$

Therefore, by Theorem 6.21 of [Ledoux and Talagrand \(1991\)](#),

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_{\psi_\alpha} \leq 17K_\alpha \|Z\|_{\psi_\alpha},$$

where the constant K_α is given in Theorem 6.21 of [Ledoux and Talagrand \(1991\)](#). Hence, for $p \geq 1$,

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_p \leq C_\alpha K_\alpha K (\log(n+1))^{1/\alpha} p^{1/\alpha},$$

for some constant C_α depending on α . Therefore, for $p \geq 1$,

$$\left\| \sum_{i=1}^n X_i \right\|_p \leq \sqrt{6p} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{1/2} + C_\alpha K_\alpha K (\log(n+1))^{1/\alpha} p^{1/\alpha}, \quad (\text{D.3})$$

for some constant $C_\alpha > 0$ (possibly different from the previous line). Hence the result follows by [Proposition C.1](#). \square

Proof of Theorem 3.3. The proof follows the same technique as that of [Theorem 3.2](#). Define $Z = \max_{1 \leq i \leq n} |X_i|$, $\rho = 8\mathbb{E}[Z]$, $K = \max_{1 \leq i \leq n} \|X_i\|_{\psi_\alpha}$,

$$X_{i,1} = X_i \mathbb{1}\{|X_i| \leq \rho\} - \mathbb{E}[X_i \mathbb{1}\{|X_i| \leq \rho\}], \quad \text{and} \quad X_{i,2} = X_i - X_{i,1}.$$

Following the same argument as in the proof of [Proposition 3.2](#), for $p \geq 1$,

$$\left\| \sum_{i=1}^n X_{i,1} \right\|_p \leq \sqrt{6p} \left(\sum_{i=1}^n \mathbb{E}[X_i^2] \right)^{1/2} + 20p\rho. \quad (\text{D.4})$$

By Hoffmann-Jorgensen's inequality ([Proposition 6.8](#) of [Ledoux and Talagrand \(1991\)](#)) and by the choice of ρ ,

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_1 \leq 2 \left\| \sum_{i=1}^n |X_i| \mathbb{1}\{|X_i| \geq \rho\} \right\|_1 \leq 16 \|Z\|_1.$$

Since $\|X_i\|_{\psi_\alpha} < \infty$ for $\alpha > 1$, $\|X_i\|_{\psi_1} < \infty$. Hence applying [Theorem 6.21](#) of [Ledoux and Talagrand \(1991\)](#), with $\alpha = 1$,

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_{\psi_1} \leq K_1 \left[16 \|Z\|_1 + \|Z\|_{\psi_1} \right] \leq 17K_1 \|Z\|_{\psi_1}.$$

By [Problem 5](#) of [Chapter 2.2](#) of [van der Vaart and Wellner \(1996\)](#),

$$\|Z\|_{\psi_1} \leq \|Z\|_{\psi_\alpha} (\log 2)^{1/\alpha-1} \quad \text{for } \alpha \geq 1,$$

and so,

$$\left\| \sum_{i=1}^n X_{i,2} \right\|_{\psi_1} \leq 17K_1 (\log 2)^{1/\alpha-1} \|Z\|_{\psi_\alpha} \leq C_\alpha (\log(n+1))^{1/\alpha} \max_{1 \leq i \leq n} \|X_i\|_{\psi_\alpha}, \quad (\text{D.5})$$

for some constant $C_\alpha > 0$ depending only on α . Therefore, combining inequalities (D.4) and (D.5) with $\rho \leq 8C_\alpha(\log(n+1))^{1/\alpha}$, for $p \geq 1$

$$\left\| \sum_{i=1}^n X_i \right\|_p \leq \sqrt{6p} \left(\sum_{i=1}^n \mathbb{E} [X_i^2] \right)^{1/2} + C_\alpha p (\log(n+1))^{1/\alpha}, \quad (\text{D.6})$$

for some constant $C_\alpha > 0$ (possibly different from that in (D.5)) depending only on α . Now the result follows by Proposition C.1 with $\alpha = 1$. \square

Proof of Theorem 3.4. Case $\alpha \leq 1$: Using the moment bound (D.3) in the proof of Theorem 3.2, it follows that for all $1 \leq j \leq q$ and $t \geq 0$,

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_{i(j)} \right| \geq e \sqrt{\frac{6\Gamma_{n,q}t}{n}} + \frac{C_\alpha K_{n,q} (\log(2n))^{1/\alpha} t^{1/\alpha}}{n} \right) \leq ee^{-t},$$

for some constant C_α depending only on α (see, for example, the proof of (C.7) for inversion of moment bounds to tail bounds). Hence by the union bound,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_\infty \geq 7 \sqrt{\frac{\Gamma_{n,q}(t + \log q)}{n}} + \frac{C_\alpha K_{n,q} (\log(2n))^{1/\alpha} (t + \log q)^{1/\alpha}}{n} \right) \\ \leq \sum_{j=1}^q \frac{ee^{-t}}{q} \leq 3e^{-t}. \end{aligned}$$

Case $\alpha \geq 1$: Using the moment bound (D.6) in the proof of Theorem 3.3, it follows that for all $1 \leq j \leq q$ and $t \geq 0$,

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_{i(j)} \right| \geq e \sqrt{\frac{6\Gamma_{n,q}t}{n}} + \frac{C_\alpha K_{n,q} (\log(2n))^{1/\alpha} t}{n} \right) \leq ee^{-t},$$

for some constant C_α depending only on α . Hence by the union bound,

$$\mathbb{P} \left(\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_\infty \geq 7 \sqrt{\frac{\Gamma_{n,q}(t + \log q)}{n}} + \frac{C_\alpha K_{n,q} (\log(2n))^{1/\alpha} (t + \log q)}{n} \right) \leq 3e^{-t}.$$

This completes the proof. \square

Proof of the Claim (3.5) Regarding the Orlicz Norm of Products. The following result establishes the bound (3.5) on the Orlicz norm of a product of random variables.

Proposition D.2. *If W_i , $1 \leq i \leq k$ are (possibly dependent) random variables satisfying $\|W_i\|_{\psi_{\alpha_i}} < \infty$ for some $\alpha_i > 0$, then*

$$\left\| \prod_{i=1}^k W_i \right\|_{\psi_\beta} \leq \prod_{i=1}^k \|W_i\|_{\psi_{\alpha_i}} \quad \text{where} \quad \frac{1}{\beta} := \sum_{i=1}^k \frac{1}{\alpha_i}.$$

Proof. The bound is trivial for $k = 1$ and it holds for $k > 2$ if it holds for $k = 2$ by recursion. For $k = 2$, set $\delta_i = \|W_i\|_{\psi_{\alpha_i}}$ for $i = 1, 2$. Fix $\eta_i > \delta_i$ for $i = 1, 2$. By definition of δ_i , this implies

$$\mathbb{E} \left[\exp \left(\left| \frac{W_i}{\eta_i} \right|^{\alpha_i} \right) \right] \leq 2 \quad \text{for } i = 1, 2. \quad (\text{D.7})$$

Observe that

$$\begin{aligned} \mathbb{E} \left[\exp \left(\left| \frac{W_1}{\eta_1} \cdot \frac{W_2}{\eta_2} \right|^\beta \right) \right] &\stackrel{(1)}{\leq} \mathbb{E} \left[\exp \left(\left| \frac{W_1}{\alpha_1} \right|^{\alpha_1} \frac{\beta}{\alpha_1} + \left| \frac{W_2}{\eta_2} \right|^{\alpha_2} \frac{\beta}{\alpha_2} \right) \right] \\ &\stackrel{(2)}{\leq} \left(\mathbb{E} \left[\exp \left(\left| \frac{W_1}{\eta_1} \right|^{\alpha_1} \right) \right] \right)^{\beta/\alpha_1} \left(\mathbb{E} \left[\exp \left(\left| \frac{W_2}{\eta_2} \right|^{\alpha_2} \right) \right] \right)^{\beta/\alpha_2} \\ &\stackrel{(3)}{\leq} 2. \end{aligned}$$

Here, (1) and (2) are applications of Young's inequality and Hölder's inequality respectively, while (3) follows by (D.7) and the definition of β . By taking limit as $\eta_i \downarrow \delta_i$, the result (3.5) now follows for $k = 2$. For $k > 2$, the result then follows by recursion, as noted earlier. \square

Proof of the Bound (3.4) in Remark 3.3. For each fixed $x \in \mathbb{R}^p$, let $T_h(Z; x) := h^{-p} Y K((X - x)/h)$, where $Z := (Y, X)$. Then under our assumed conditions, using the quasi-norm property and moment bounds for the $\|\cdot\|_{\psi_\alpha}$ norm (see, for instance, Chapter 2.2 of [van der Vaart and Wellner \(1996\)](#)) along with Proposition D.2, we have: for all $x \in \mathbb{R}^p$,

$$\begin{aligned} \|T_h(Z; x) - \mathbb{E}\{T_h(Z; x)\}\|_{\psi_\alpha} &\leq A_\alpha \left[\|T_h(Z; x)\|_{\psi_\alpha} + \mathbb{E}\{|T_h(Z; x)|\} \right] \\ &\leq A_\alpha \left[\|T_h(Z; x)\|_{\psi_\alpha} + B_\alpha \|T_h(Z; x)\|_{\psi_\alpha} \right] \\ &= D_\alpha \|T_h(Z; x)\|_{\psi_\alpha} \leq D_\alpha h^{-p} C_Y C_K, \end{aligned} \quad (\text{D.8})$$

where $A_\alpha, B_\alpha > 0$ are some constants depending only on α , and $D_\alpha := A_\alpha(1 + B_\alpha) > 0$.

Further, $\text{Var}\{T_h(Z; x)\} \leq \mathbb{E}\{T_h^2(Z; x)\}$ and $\mathbb{E}\{T_h^2(Z; x)\}$ satisfies: for all $x \in \mathbb{R}^p$,

$$\begin{aligned} \mathbb{E}\{T_h^2(Z; x)\} &= \mathbb{E} \left[\mathbb{E} \{T_h^2(Z; x) | X\} \right] = \frac{1}{h^{2p}} \int_{\mathbb{R}^p} \{ \mathbb{E} (Y^2 | X = u) \} K^2 \left(\frac{u - x}{h} \right) f(u) du \\ &= \frac{1}{h^{2p}} \int_{\mathbb{R}^p} \{ \mathbb{E} (Y^2 | X = x + h\varphi) \} K^2(\varphi) f(\varphi) h^p d\varphi \leq \frac{R_K M_Y}{h^p}, \end{aligned} \quad (\text{D.9})$$

where the final bound is due to our assumptions. The result (3.4) now follows by simply applying Theorem 3.4 to the random variables $T_h(Z_i; x) - \mathbb{E}\{T_h(Z_i; x)\}$, $1 \leq i \leq n$, and by using the bounds (D.8) and (D.9). This completes the proof. \square

E Proofs of All Results in Section 4

E.1 Proofs of the Results in Section 4.1

Proof of Theorem 4.1. Under assumption (4.3), it follows from Proposition D.2 that

$$\max_{1 \leq i \leq n} \max_{1 \leq j \leq k \leq p} \|X_i(j)X_i(k)\|_{\psi_{\alpha/2}} \leq K_{n,p}^2,$$

and so, Theorem 3.4 with $q = p^2$ implies the result. \square

Proof of Theorem 4.2. It is easy to verify that

$$\begin{aligned}\hat{\Sigma}_n^* &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{\mu}_n) (X_i - \bar{\mu}_n)^\top - (\bar{X}_n - \bar{\mu}_n) (\bar{X}_n - \bar{\mu}_n)^\top \\ &=: \tilde{\Sigma}_n^* - (\bar{X}_n - \bar{\mu}_n) (\bar{X}_n - \bar{\mu}_n)^\top.\end{aligned}$$

Clearly,

$$\Delta_n^* \leq \|\tilde{\Sigma}_n^* - \Sigma_n^*\|_\infty + \|\bar{X}_n - \bar{\mu}_n\|_\infty^2, \quad (\text{E.1})$$

where $\|x\|_\infty$ represents the maximum absolute element of x . Since $\tilde{\Sigma}_n^*$ is the gram matrix corresponding to the random vectors $X_i - \bar{\mu}_n$, Theorem 4.1 applies for the first term on the right hand side of (E.1). For the second term, Theorem 3.4 applies. Combining these two bounds, we get that for any $t \geq 0$, with probability at least $1 - 6e^{-t}$,

$$\begin{aligned}\Delta_n^* &\leq 7A_{n,p}^* \sqrt{\frac{t + 2 \log p}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + 2 \log p)^{2/\alpha}}{n} \\ &\quad + 98 \left(\frac{B_{n,p}(t + \log p)}{n} \right) + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + \log p)^{2/\alpha}}{n^2},\end{aligned} \quad (\text{E.2})$$

where $B_{n,p} := \max_{1 \leq j \leq p} \sum_{i=1}^n \text{Var}(X_i(j)) / n$. As before, it is easy to show that $B_{n,p} \leq C_\alpha K_{n,p}^2$ and so, the last two terms of inequality (E.2) are of lower order than the second term and hence, we obtain that with probability at least $1 - 6e^{-t}$,

$$\Delta_n^* \leq 7A_{n,p}^* \sqrt{\frac{t + 2 \log p}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + 2 \log p)^{2/\alpha}}{n},$$

with a possibly increased constant $C_\alpha > 0$. \square

E.2 Proofs of the Results in Section 4.2

Proof of Theorem 4.3. To prove the result, note that

$$\text{RIP}_n(k) = \sup_{\substack{\theta \in \mathbb{R}^p, \\ \|\theta\|_0 \leq k, \|\theta\|_2 \leq 1}} \left| \frac{1}{n} \sum_{i=1}^n \left\{ (X_i^\top \theta)^2 - \mathbb{E} \left[(X_i^\top \theta)^2 \right] \right\} \right|,$$

and define the set

$$\Theta_k := \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq k, \|\theta\|_2 = 1\} \subseteq \mathbb{R}^p.$$

For every $\varepsilon > 0$, let \mathcal{N}_ε denote the ε -net of Θ_k , that is, every $\theta \in \Theta_k$ can be written as $\theta = x_\theta + z_\theta$ where $\|x_\theta\|_2 \leq 1$, $x_\theta \in \mathcal{N}_\varepsilon$ and $\|z_\theta\|_2 \leq \varepsilon$. In this representation x_θ and z_θ can be taken to have the same support as that of θ . By Lemma 3.3 of Plan and Vershynin (2013), it follows that

$$|\mathcal{N}_{1/4}| \leq \left(\frac{36p}{k} \right)^k.$$

By Proposition 2.2 of [Vershynin \(2012\)](#), it is easy to see that $\text{RIP}_n(k)$ can be bounded by a finite maximum as

$$\text{RIP}_n(k) \leq 2 \sup_{\theta \in \mathcal{N}_{1/4}} \left| \frac{1}{n} \sum_{i=1}^n \left\{ \left(X_i^\top \theta \right)^2 - \mathbb{E} \left[\left(X_i^\top \theta \right)^2 \right] \right\} \right|.$$

This implies that $\text{RIP}_n(k)$ can be controlled by controlling a finite maximum of averages. Set

$$\Lambda_n(k) := \sup_{\theta \in \mathcal{N}_{1/4}} \left| \frac{1}{n} \sum_{i=1}^n \left\{ \left(X_i^\top \theta \right)^2 - \mathbb{E} \left[\left(X_i^\top \theta \right)^2 \right] \right\} \right|.$$

- (a) Under the marginal ψ_α -bound, it is easy to see that for $\theta \in \Theta_k$ with support $S \subseteq \{1, \dots, p\}$ of size k ,

$$\left\| \left(X_i^\top \theta \right)^2 \right\|_{\psi_{\alpha/2}} \leq \left\| \sum_{j \in S} X_i^2(j) \right\|_{\psi_{\alpha/2}} \leq C_\alpha \sum_{j \in S} \|X_i(j)\|_{\psi_\alpha}^2 \leq C_\alpha K_{n,p}^2 k,$$

for some constant C_α depending only on α . Hence by [Theorem 3.4](#), it follows that for any $t > 0$, with probability at least $1 - 3e^{-t}$,

$$\Lambda_n(k) \leq 7 \sqrt{\frac{\Upsilon_{n,k}(t + k \log(36p/k))}{n}} + \frac{C_\alpha K_{n,p}^2 k (\log(2n))^{2/\alpha} (t + k \log(36p/k))^{2/\alpha}}{n}.$$

- (b) Under the joint ψ_α -bound, it readily follows that

$$\sup_{\theta \in \Theta_k} \left\| \left(X_i^\top \theta \right)^2 \right\|_{\psi_{\alpha/2}} \leq K_{n,p}^2.$$

Hence by [Theorem 3.4](#), we get that for any $t > 0$, with probability at least $1 - 3e^{-t}$,

$$\Lambda_n(k) \leq 7 \sqrt{\frac{\Upsilon_{n,k}(t + k \log(36p/k))}{n}} + \frac{C_\alpha K_{n,p}^2 (\log(2n))^{2/\alpha} (t + k \log(36p/k))^{2/\alpha}}{n}.$$

The result now follows since $\text{RIP}_n(k) \leq 2\Lambda_n(k)$. □

E.3 Proofs of the Results in [Section 4.3](#)

Proof of [Theorem 4.4](#). The proof follows using [Theorem 4.3](#) and Lemma 12 of [Loh and Wainwright \(2012\)](#).

- (a) From part (a) of [Theorem 4.3](#), we have with probability at least $1 - 3s(np)^{-1}$,

$$\text{RIP}_n(s) \leq \Xi_{n,s}^{(M)}.$$

On the event where this inequality holds, applying Lemma 12 of [Loh and Wainwright \(2012\)](#) with $\Gamma = \hat{\Sigma}_n - \Sigma_n$ and $\delta = \Xi_{n,s}^{(M)}$ proves that with probability at least $1 - 3s(np)^{-1}$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_{n,s}^{(M)} \right) \|\theta\|_2^2 - \frac{54\Xi_{n,s}^{(M)}}{s} \|\theta\|_1^2 \quad \text{for all } \theta \in \mathbb{R}^p.$$

(b) From part (b) of Theorem 4.3, we get with probability at least $1 - 3s(np)^{-1}$,

$$\text{RIP}_n(s) \leq \Xi_{n,s}^{(J)}.$$

By a similar argument as above, the result follows.

This completes the proof of Theorem 4.4. \square

E.4 Proofs of the Results in Section 4.4

The following is a general result of [Negahban et al. \(2012\)](#) and this particular form is taken from [Hastie et al. \(2015\)](#).

Lemma E.1 (Theorem 11.1 of [Hastie et al. \(2015\)](#)). *Assume that the matrix $\hat{\Sigma}_n$ satisfies the restricted eigenvalue bound (4.15) with $\delta = 3$. Fix any vector $\beta \in \mathbb{R}^p$ with $\|\beta\|_0 \leq k$. Given a regularization parameter λ_n satisfying*

$$\lambda_n \geq 2 \left\| \frac{1}{n} \sum_{i=1}^n X_i (Y_i - X_i^\top \beta) \right\|_\infty > 0,$$

any estimator $\hat{\beta}_n(\lambda_n)$ from the Lasso (4.29) satisfies the bound

$$\left\| \hat{\beta}_n(\lambda_n) - \beta \right\|_2 \leq \frac{3}{\gamma_n} \sqrt{k} \lambda_n.$$

Lemma E.1 holds for any of the minimizers $\hat{\beta}_n(\lambda_n)$ in case of non-uniqueness.

Proof of Theorem 4.5. Using Proposition D.2, it follows that

$$\max_{1 \leq i \leq n} \|X_i(j)\varepsilon_j\|_{\psi_\gamma} \leq K_{n,p}^2.$$

By Theorem 3.4, it follows that with probability at least $1 - 3(np)^{-1}$,

$$\left\| \frac{1}{n} \sum_{i=1}^n X_i \varepsilon_i \right\|_\infty \leq 7\sqrt{2}\sigma_{n,p} \sqrt{\frac{\log(np)}{n}} + \frac{C_\gamma K_{n,p}^2 (\log(2n))^{1/\gamma} (2\log(np))^{1/\gamma}}{n}. \quad (\text{E.3})$$

Next, under assumption (4.30), and using similar arguments as those used to prove (4.21) in the analysis in Section 4.3.1, we have: if \bar{s}_M minimizes $\Xi_{n,s}^{(M)} + 32k\Xi_{n,s}^{(M)}/s$ over $1 \leq s \leq p$, then (4.20) implies that with probability at least $1 - 3\bar{s}_M/(np) \geq 1 - 3/n$,

$$\theta^\top \hat{\Sigma}_n \theta \geq \left(\lambda_{\min}(\Sigma_n) - 27 \min_{1 \leq s \leq p} \left\{ \Xi_s + 32 \frac{k\Xi_s}{s} \right\} \right) \|\theta\|_2^2 \geq \frac{\lambda_{\min}(\Sigma_n)}{2} \|\theta\|_2^2 \quad \forall \theta \in \bigcup_{|S| \leq k} \mathcal{C}(S; 3),$$

so that the $\text{RE}(k)$ condition (4.15) holds with probability at least $1 - 3n^{-1}$ and with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$ and $\delta = 3$. Therefore, applying Lemma E.1, along with the bound (E.3), it follows that with probability at least $1 - 3(np)^{-1} - 3n^{-1}$, there exists a λ_n (given by twice the upper bound in (E.3)) and hence, a Lasso estimator $\hat{\beta}_n(\lambda_n)$, satisfying:

$$\left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2 \leq \frac{84\sqrt{2}}{\lambda_{\min}(\Sigma_n)} \left[\sigma_{n,p} \sqrt{\frac{k \log(np)}{n}} + 2^{1/\gamma} C_\gamma K_{n,p}^2 \frac{k^{1/2} (\log(np))^{2/\gamma}}{n} \right].$$

Note that the choice of λ_n above is as claimed in the result. This completes the proof. \square

Proof of Theorem 4.6. We already showed in the proof of Theorem 4.5 above that under assumption (4.30), the $\text{RE}(k)$ condition (4.15) holds with probability at least $1 - 3n^{-1}$ with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$ and $\delta = 3$. Hence, to apply Lemma E.1, it is enough to show that the λ_n in the statement of Theorem 4.6 is a valid choice for Lemma E.1. For this, we prove that with probability at least $1 - 3(np)^{-1} - L^{-1}$, for any $L \geq 1$,

$$\begin{aligned} \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i X_i(j) \right| &\leq 7\sqrt{2} \sigma_{n,p} \sqrt{\frac{\log(np)}{n}} \\ &\quad + \frac{C_\alpha K_{n,p} K_{\varepsilon,r} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{n^{1-1/r}}. \end{aligned}$$

We follow the proof technique of Theorem 3.2 to reduce the assumption on ε_i to polynomial moments, as follows. Define

$$C_{n,\varepsilon} := 8\mathbb{E} \left[\max_{1 \leq i \leq n} |\varepsilon_i| \right] \leq 8n^{1/r} \max_{1 \leq i \leq n} \|\varepsilon_i\|_r \leq 8n^{1/r} K_{\varepsilon,r}.$$

Note that under the setting of Theorem 4.5, for $1 \leq j \leq p$,

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i X_i(j) = \frac{1}{n} \sum_{i=1}^n \{ \varepsilon_i X_i(j) - \mathbb{E}[\varepsilon_i X_i(j)] \}.$$

Set for $1 \leq i \leq n$, $S_i := \varepsilon_i X_i - \mathbb{E}[\varepsilon_i X_i] \in \mathbb{R}^p$, and for $1 \leq j \leq p$,

$$\begin{aligned} S_i^{(1)}(j) &:= S_i(j) \mathbb{1}_{\{|\varepsilon_i| \leq C_{n,\varepsilon}\}} - \mathbb{E} [S_i(j) \mathbb{1}_{\{|\varepsilon_i| \leq C_{n,\varepsilon}\}}], \\ S_i^{(2)}(j) &:= S_i(j) \mathbb{1}_{\{|\varepsilon_i| > C_{n,\varepsilon}\}} - \mathbb{E} [S_i(j) \mathbb{1}_{\{|\varepsilon_i| > C_{n,\varepsilon}\}}]. \end{aligned}$$

Therefore, by triangle inequality,

$$\begin{aligned} \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i X_i(j) \right| &= \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i(j) \right| \\ &\leq \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(1)}(j) \right| + \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(2)}(j) \right|. \end{aligned} \tag{E.4}$$

For the summands of the first term, note that

$$\begin{aligned} \text{Var}(S_i^{(1)}(j)) &\leq \mathbb{E} [S_i^2(j) \mathbb{1}_{\{|\varepsilon_i| \leq C_{n,\varepsilon}\}}] \\ &\leq \mathbb{E} [S_i^2(j)] = \text{Var}(S_i(j)) = \text{Var}(\varepsilon_i X_i(j)), \end{aligned}$$

and for some constant B_α (depending only on α),

$$\begin{aligned}
\left\| S_i^{(1)}(j) \right\|_{\psi_\alpha} &\leq 2 \left\| S_i(j) \mathbb{1}_{\{|\varepsilon_i| \leq C_{n,\varepsilon}\}} \right\|_{\psi_\alpha} \\
&\leq 2B_\alpha \left\| \varepsilon_i X_i(j) \mathbb{1}_{\{|\varepsilon_i| \leq C_{n,\varepsilon}\}} \right\|_{\psi_\alpha} + 2B_\alpha |\mathbb{E}[\varepsilon_i X_i(j)]| \\
&\leq 2B_\alpha C_{n,\varepsilon} K_{n,p} + 2B_\alpha \|\varepsilon_i\|_2 \|X_i(j)\|_2 \\
&\leq 2B_\alpha C_{n,\varepsilon} K_{n,p} + 2B_\alpha \|\varepsilon_i\|_2 K_{n,p} = 2B_\alpha K_{n,p} [C_{n,\varepsilon} + \|\varepsilon_i\|_2] \\
&\leq 4B_\alpha K_{n,p} C_{n,\varepsilon} \leq 32n^{1/r} B_\alpha K_{n,p} K_{\varepsilon,r}.
\end{aligned}$$

Therefore, by Theorem 3.4, it follows that with probability at least $1 - 3(np)^{-1}$,

$$\begin{aligned}
\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(1)}(j) \right| &\leq 7\sqrt{2}\sigma_{n,p} \sqrt{\frac{\log(np)}{n}} \\
&\quad + \frac{C_\alpha K_{n,p} K_{\varepsilon,r} (\log(2n))^{1/\alpha} (\log(np))^{1/\alpha}}{n^{1-1/r}}.
\end{aligned} \tag{E.5}$$

For the second term in (E.4), note that

$$\left\| \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(2)}(j) \right| \right\|_1 \leq 2 \left\| \max_{1 \leq j \leq p} \frac{1}{n} \sum_{i=1}^n |\varepsilon_i X_i(j)| \mathbb{1}_{\{|\varepsilon_i| > C_{n,\varepsilon}\}} \right\|_1.$$

By the definition of $C_{n,\varepsilon}$, we have

$$\mathbb{P} \left(\max_{1 \leq j \leq p} \frac{1}{n} \sum_{i=1}^n |\varepsilon_i X_i(j)| \mathbb{1}_{\{|\varepsilon_i| > C_{n,\varepsilon}\}} > 0 \right) \leq \mathbb{P} \left(\max_{1 \leq i \leq n} |\varepsilon_i| > C_{n,\varepsilon} \right) \leq 1/8.$$

Thus by Hoffmann-Jorgensen's inequality, we have

$$\begin{aligned}
\left\| \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(2)}(j) \right| \right\|_1 &\leq \frac{2}{n} \left\| \max_{1 \leq j \leq p} \max_{1 \leq i \leq n} |\varepsilon_i X_i(j)| \right\|_1 \\
&\leq \frac{2}{n} \left\| \max_{1 \leq i \leq n} |\varepsilon_i| \right\|_2 \left\| \max_{\substack{1 \leq i \leq n, \\ 1 \leq j \leq p}} |X_i(j)| \right\|_2 \\
&\leq \frac{2C_\alpha (\log(np))^{1/\alpha}}{n^{1-1/r}} K_{\varepsilon,r} K_{n,p},
\end{aligned}$$

for some constant $C_\alpha > 0$. So, for any $L \geq 1$, with probability at least $1 - L^{-1}$,

$$\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i^{(2)}(j) \right| \leq \frac{2LC_\alpha (\log(np))^{1/\alpha} K_{\varepsilon,r} K_{n,p}}{n^{1-1/r}}. \tag{E.6}$$

From inequalities (E.5) and (E.6), we get with probability at least $1 - 3(np)^{-1} - L^{-1}$,

$$\begin{aligned}
\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n S_i(j) \right| &\leq 7\sqrt{2}\sigma_{n,p} \sqrt{\frac{\log(np)}{n}} \\
&\quad + \frac{C_\alpha K_{n,p} K_{\varepsilon,r} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{n^{1-1/r}}.
\end{aligned} \tag{E.7}$$

Taking together the events on which the RE condition and the inequality (E.7) hold, we have: with probability at least $1 - 3(np)^{-1} - 3n^{-1} - L^{-1}$, the RE(k) condition (4.15) is satisfied with $\gamma_n = \lambda_{\min}(\Sigma_n)/2$ and $\delta = 3$, and λ_n can be chosen as

$$\lambda_n = 14\sqrt{2}\sigma_{n,p}\sqrt{\frac{\log(np)}{n}} + \frac{C_\alpha K_{n,p} K_{\varepsilon,r} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{n^{1-1/r}},$$

so that the lasso estimator $\hat{\beta}_n(\lambda_n)$ satisfies (by Lemma E.1),

$$\begin{aligned} \left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2 &\leq \frac{84\sqrt{2}}{\lambda_{\min}(\Sigma_n)} \sigma_{n,p} \sqrt{\frac{k \log(np)}{n}} \\ &\quad + C_\alpha K_{n,p} K_{\varepsilon,r} \frac{k^{1/2} (\log(np))^{1/\alpha} [(\log(2n))^{1/\alpha} + L]}{\lambda_{\min}(\Sigma_n) n^{1-1/r}}. \end{aligned}$$

This completes the proof of Theorem 4.6. \square

Proof of Remark 4.16. The following result proves the oracle inequality stated in Remark 4.16.

Proposition E.1 (Oracle Inequality for Lasso). *Consider the setting of Theorem 4.5 (except the hard sparsity on β_0). For the choice of λ_n as in (4.31), with probability converging to 1,*

$$\begin{aligned} &\left\| \hat{\beta}_n(\lambda_n) - \beta_0 \right\|_2^2 \\ &\leq \min_{S: \Xi_{n,|S|}^{(M)} = o(1)} \left[\frac{18\lambda_n^2 |S|}{\Gamma_n^2(S)} + \left(\frac{8\lambda_n \|\beta_0(S^c)\|_1}{\Gamma_n(S)} + \frac{3456\Xi_{n,|S|}^{(M)} \|\beta^*(S^c)\|_1^2}{|S|\Gamma_n(S)} \right) \right], \end{aligned}$$

where

$$\Gamma_n(S) := \lambda_{\min}(\Sigma_n) - 1755\Xi_{n,|S|}^{(M)}.$$

Proof. The proof closely follows the arguments of Theorem 11.1 of Hastie et al. (2015) and Section 4.3 of Negahban et al. (2010). Set for $\nu \in \mathbb{R}^p$,

$$G(\nu) := \frac{1}{2n} \sum_{i=1}^n \left(Y_i - X_i^\top (\beta_0 + \nu) \right)^2 + \lambda_n \|\beta_0 + \nu\|_1,$$

and $\hat{\nu} := \hat{\beta}_n(\lambda_n) - \beta_0$. Also, fix any subset $S \subseteq \{1, 2, \dots, p\}$ with $\Xi_{n,|S|}^{(M)} = o(1)$. Note that with probability at least $1 - 3(np)^{-1}$,

$$\lambda_n \geq 2 \left\| \frac{1}{n} \sum_{i=1}^n X_i \varepsilon_i \right\|_\infty,$$

as shown in the proof of Theorem 4.5. On this event the following calculations hold true. By definition $G(\hat{\nu}) \leq G(0)$ and so,

$$\frac{\hat{\nu}^\top \hat{\Sigma}_n \hat{\nu}}{2} \leq \hat{\nu}^\top \left(\frac{1}{n} \sum_{i=1}^n X_i \varepsilon_i \right) + \lambda_n [\|\beta_0\|_1 - \|\beta_0 + \hat{\nu}\|_1]. \quad (\text{E.8})$$

Now observe that

$$\begin{aligned}\|\beta_0 + \hat{\nu}\|_1 &\geq \|\beta_0(S) + \hat{\nu}(S)\|_1 - \|\beta_0(S^c)\|_1 + \|\hat{\nu}(S^c)\|_1 \\ &\geq \|\beta_0(S)\|_1 - \|\hat{\nu}(S)\|_1 - \|\beta_0(S^c)\|_1 + \|\hat{\nu}(S^c)\|_1.\end{aligned}$$

Since $\|\beta_0\|_1 = \|\beta_0(S)\|_1 + \|\beta_0(S^c)\|_1$, the above inequality substituted in (E.8) implies

$$\begin{aligned}\frac{\hat{\nu}^\top \hat{\Sigma}_n \hat{\nu}}{2} &\leq \hat{\nu}^\top \left(\frac{1}{n} \sum_{i=1}^n X_i \varepsilon_i \right) + \lambda_n [2 \|\beta_0(S^c)\|_1 + \|\hat{\nu}(S)\|_1 - \|\hat{\nu}(S^c)\|_1] \\ &\leq \|\hat{\nu}\|_1 \left\| \frac{1}{n} \sum_{i=1}^n X_i \varepsilon_i \right\|_\infty + \lambda_n [2 \|\beta_0(S^c)\|_1 + \|\hat{\nu}(S)\|_1 - \|\hat{\nu}(S^c)\|_1] \\ &\leq \frac{\lambda_n}{2} \|\hat{\nu}(S)\|_1 + \frac{\lambda_n}{2} \|\hat{\nu}(S^c)\|_1 + \lambda_n [2 \|\beta_0(S^c)\|_1 + \|\hat{\nu}(S)\|_1 - \|\hat{\nu}(S^c)\|_1] \\ &\leq \frac{3\lambda_n}{2} \|\hat{\nu}(S)\|_1 - \frac{\lambda_n}{2} \|\hat{\nu}(S^c)\|_1 + 2\lambda_n \|\beta_0(S^c)\|_1.\end{aligned}\tag{E.9}$$

This inequality has two implications that prove the result. Firstly, the left hand side of (E.9) is non-negative and so,

$$\|\hat{\nu}(S^c)\|_1 \leq 3 \|\hat{\nu}(S)\|_1 + 4 \|\beta_0(S^c)\|_1.\tag{E.10}$$

For the second implication, note that inequality (E.10) implies that

$$\begin{aligned}\|\hat{\nu}\|_1 &= \|\hat{\nu}(S)\|_1 + \|\hat{\nu}(S^c)\|_1 \\ &\leq 4 \|\hat{\nu}(S)\|_1 + 4 \|\beta_0(S^c)\|_1 \\ &\leq 4\sqrt{|S|} \|\hat{\nu}(S)\|_2 + 4 \|\beta_0(S^c)\|_1 \leq 4\sqrt{|S|} \|\hat{\nu}\|_2 + 4 \|\beta_0(S^c)\|_1.\end{aligned}$$

Therefore, applying Theorem 4.4 with $s = |S|$, we get that with probability at least $1 - |S|(np)^{-1}$,

$$\begin{aligned}\hat{\nu}^\top \hat{\Sigma}_n \hat{\nu} &\geq \left(\lambda_{\min}(\Sigma_n) - 27\Xi_{n,|S|}^{(M)} \right) \|\hat{\nu}\|_2^2 - \frac{54\Xi_{n,|S|}^{(M)}}{|S|} \left(32|S| \|\hat{\nu}\|_2^2 + 32 \|\beta_0(S^c)\|_2^2 \right) \\ &= \left(\lambda_{\min}(\Sigma_n) - 1755\Xi_{n,|S|}^{(M)} \right) \|\hat{\nu}\|_2^2 - \frac{1728\Xi_{n,|S|}^{(M)}}{|S|} \|\beta_0(S^c)\|_1^2 \\ &= \Gamma_n(S) \|\hat{\nu}\|_2^2 - \frac{1728\Xi_{n,|S|}^{(M)}}{|S|} \|\beta_0(S^c)\|_1^2.\end{aligned}\tag{E.11}$$

Combining inequality (E.11) with inequality (E.9), we obtain

$$\frac{\Gamma_n(S)}{2} \|\hat{\nu}\|_2^2 \leq \frac{3\lambda_n \sqrt{|S|}}{2} \|\hat{\nu}\|_2 + 2\lambda_n \|\beta_0(S^c)\|_1 + \frac{864\Xi_{n,|S|}^{(M)}}{|S|} \|\beta_0(S^c)\|_1^2$$

Hence,

$$\|\hat{\nu}\|_2 \leq \frac{3\lambda_n \sqrt{|S|}}{\Gamma_n(S)} + \sqrt{\frac{2}{\Gamma_n(S)}} \left(2\lambda_n \|\beta_0(S^c)\|_1 + \frac{864\Xi_{n,|S|}^{(M)}}{|S|} \|\beta_0(S^c)\|_1^2 \right)^{1/2},$$

and so the result follows. \square

F Proofs of All Results in Appendix B

Proof of Proposition B.1. By Theorem 3 of Adamczak (2008), we get for all $t \geq 0$ that

$$\mathbb{P}((Z - \mathbb{E}[Z])_+ \geq t) \leq \exp\left(-\frac{t^2}{2(\Sigma_n(\mathcal{F}) + 2U\mathbb{E}[Z]) + 3Ut}\right).$$

Set $A := 2(\Sigma_n(\mathcal{F}) + 2U\mathbb{E}[Z])$ and $B := 3U$. Then using the arguments in Proposition D.1, we get for $p \geq 2$, (and any $\delta > 0$)

$$\begin{aligned} \mathbb{E}[(Z - \mathbb{E}[Z])_+^p] &\leq \int_0^\infty pt^{p-1} \exp\left(-\frac{t^2}{A(1+\delta)}\right) dt + \int_0^\infty pt^{p-1} \exp\left(-\frac{t\delta}{B(\delta+1)}\right) dt \\ &\leq \left(\sqrt{A(1+\delta)/e}\right)^p \frac{\sqrt{2\pi}}{e} \left(1 + \frac{p}{2}\right)^{(p+1)/2} \exp\left(\frac{1}{12(1+p/2)}\right) \\ &\quad + \left(\frac{B(1+\delta)}{e\delta}\right)^p \frac{\sqrt{2\pi}}{e} (1+p)^{p+\frac{1}{2}} \exp\left(\frac{1}{12(p+1)}\right) =: \mathbf{I} + \mathbf{II}. \end{aligned}$$

So, for $p \geq 2$,

$$\begin{aligned} \mathbf{I}^{1/p} &\leq p^{1/2} \sqrt{A(1+\delta)/e} \left(\frac{\sqrt{2\pi}}{e}\right)^{1/p} \left(1 + \frac{p}{2}\right)^{\frac{1}{2p}} \exp\left(\frac{1}{12p(1+p/2)}\right) \\ &\leq p^{1/2} \sqrt{A} \quad \text{by taking } \delta = 1/2. \end{aligned}$$

Also, regarding \mathbf{II} , for $p \geq 2$,

$$\mathbf{II}^{1/p} \leq p \left(\frac{3B(1+\delta)}{2e\delta}\right) \left(\frac{\sqrt{2\pi}}{e}\right)^{1/p} (1+p)^{1/(2p)} \exp\left(\frac{1}{12p(p+1)}\right) \leq 2Bp.$$

Therefore, for $p \geq 2$,

$$\|(Z - \mathbb{E}[Z])_+\|_p \leq (2\Sigma_n(\mathcal{F}) + 4U\mathbb{E}[Z])^{1/2} \sqrt{p} + 6pU,$$

and since $\|Z\|_p \leq \|(Z - \mathbb{E}[Z])_+\|_p + \mathbb{E}[Z]$, for $p \geq 1$,

$$\|Z\|_p \leq \mathbb{E}[Z] + (2\Sigma_n(\mathcal{F}) + 4U\mathbb{E}[Z])^{1/2} \sqrt{p} + 6Up,$$

proving (B.3). □

Proof of Theorem B.1. By triangle inequality, $Z \leq Z_1 + Z_2$. Note that Z_1 is a supremum of bounded empirical process and so, by Proposition B.1 for $p \geq 2$

$$\begin{aligned} \|(Z_1 - \mathbb{E}[Z_1])_+\|_p &\leq p^{1/2} (2\Sigma_n(\mathcal{F}) + 8\rho\mathbb{E}[Z_1])^{1/2} + 12\rho p \\ &\leq \sqrt{2}p^{1/2}\Sigma_n^{1/2}(\mathcal{F}) + 2\sqrt{2}p^{1/2}\rho^{1/2} (\mathbb{E}[Z_1])^{1/2} + 12\rho p \\ &\leq \sqrt{2}p^{1/2}\Sigma_n^{1/2}(\mathcal{F}) + (2\rho p + \mathbb{E}[Z_1]) + 12\rho p \\ &= \mathbb{E}[Z_1] + \sqrt{2}p^{1/2}\Sigma_n^{1/2}(\mathcal{F}) + 14\rho p, \end{aligned} \tag{F.1}$$

where we used the arithmetic-geometric mean inequality and the fact that

$$\text{Var}(f(X_i)\mathbb{1}\{|f(X_i)| \leq \rho\}) \leq \mathbb{E}[f^2(X_i)\mathbb{1}\{|f(X_i)| \leq \rho\}] \leq \mathbb{E}[f^2(X_i)] \leq \text{Var}(f(X_i)).$$

To deal with Z_2 , observe that

$$\|Z_2\|_{\psi_{\alpha_*}} \leq 2 \left\| \sum_{i=1}^n F(X_i)\mathbb{1}\{F(X_i) \geq \rho\} \right\|_{\psi_{\alpha_*}}$$

Since $\alpha_* \leq 1$ for all $\alpha > 0$ and $\|F(X_i)\|_{\psi_{\alpha_*}} < \infty$, it follows from Theorem 6.21 of [Ledoux and Talagrand \(1991\)](#) that

$$\|Z_2\|_{\psi_{\alpha_*}} \leq 2K_{\alpha_*} \left\{ \mathbb{E} \left[\sum_{i=1}^n F(X_i)\mathbb{1}\{F(X_i) \geq \rho\} \right] + \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha_*}} \right\},$$

with the constant K_{α_*} as in the cited theorem. Additionally by Hoffmann-Jorgensen inequality combined with the definition of ρ , we have

$$\mathbb{E} \left[\sum_{i=1}^n F(X_i)\mathbb{1}\{F(X_i) \geq \rho\} \right] \leq 8\mathbb{E} \left[\max_{1 \leq i \leq n} F(X_i) \right] \leq 8 \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}.$$

And by Problem 5 of Chapter 2.2 of [van der Vaart and Wellner \(1996\)](#),

$$\left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha_*}} \leq (\log 2)^{1/\alpha - 1/\alpha_*} \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}.$$

Therefore,

$$\|Z_2\|_{\psi_{\alpha_*}} \leq 2K_{\alpha_*} \left[8 + (\log 2)^{1/\alpha - 1/\alpha_*} \right] \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}.$$

This implies that for $p \geq 1$,

$$\|Z_2\|_p \leq 2\sqrt{2\pi}(p/\alpha_*)^{1/\alpha_*} K_{\alpha_*} \left[8 + (\log 2)^{1/\alpha - 1/\alpha_*} \right] \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}. \quad (\text{F.2})$$

Note that for all $\alpha > 0$, $(p/\alpha_*)^{1/\alpha_*} \geq p$ for all $p \geq 1$ and so, for all $\alpha > 0$,

$$14p\rho \leq \sqrt{2\pi}(p/\alpha_*)^{1/\alpha_*} K_{\alpha_*} \left[8 + (\log 2)^{1/\alpha - 1/\alpha_*} \right] \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}.$$

Therefore, combining bounds (F.1) and (F.2), we obtain for $p \geq 2$,

$$\begin{aligned} \|Z\|_p &\leq \mathbb{E}[Z_1] + \|(Z_1 - \mathbb{E}[Z_1])_+\|_p + \|Z_2\|_p \\ &\leq 2\mathbb{E}[Z_1] + \sqrt{2}p^{1/2}\Sigma_n^{1/2}(\mathcal{F}) \\ &\quad + 3\sqrt{2\pi}(p/\alpha_*)^{1/\alpha_*} K_{\alpha_*} \left[8 + (\log 2)^{1/\alpha - 1/\alpha_*} \right] \left\| \max_{1 \leq i \leq n} F(X_i) \right\|_{\psi_{\alpha}}. \end{aligned}$$

This proves (B.5). Using the reasoning as in Proposition B.1,

$$\mathbb{E} \left[\Psi_{\alpha_*, L_n(\alpha)} \left(\frac{(Z - 2e\mathbb{E}[Z_1])_+}{3\sqrt{2}e\Sigma_n^{1/2}(\mathcal{F})} \right) \right] \leq 1,$$

with $L_n(\alpha)$ is as defined in the statement. This proves (B.6). \square

Proof of Proposition B.2. By Theorem 3.5.1 and inequality (3.167) of [Giné and Nickl \(2016\)](#),

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right] \leq 8\sqrt{2} \mathbb{E} \left[\int_0^{\eta_n(\mathcal{F})} \sqrt{\log(2N(x, \mathcal{F}, \|\cdot\|_{2, P_n}))} dx \right],$$

where P_n represents the empirical measure of X_1, X_2, \dots, X_n , that is, $P_n(\{X_i\}) = 1/n$. Here

$$\eta_n^2(\mathcal{F}) := \sup_{f \in \mathcal{F}} \|f\|_{2, P_n}^2 = \sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n f^2(X_i) \right)^{1/2}.$$

Using a change-of-variable formula,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right] \leq 8\sqrt{2} \mathbb{E} \left[\|F\|_{2, P_n} J(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_2) \right],$$

where

$$\delta_n^2(\mathcal{F}) := \frac{1}{\|F\|_{2, P_n}^2} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f^2(X_i), \quad \text{and} \quad \|F\|_{2, P_n}^2 := \frac{1}{n} \sum_{i=1}^n F^2(X_i).$$

Now an application of Lemma 3.5.3 (c) of [Giné and Nickl \(2016\)](#) implies that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right] \leq 8\sqrt{2} \|F\|_{2, P} J \left(\frac{\Delta}{\|F\|_{2, P}}, \mathcal{F}, \|\cdot\|_2 \right), \quad (\text{F.3})$$

where

$$\Delta^2 := \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f^2(X_i) \right] \quad \text{and} \quad \|F\|_{2, P}^2 := \frac{1}{n} \sum_{i=1}^n \mathbb{E} [F^2(X_i)].$$

Note that by symmetrization and contraction principle

$$\begin{aligned} \Delta^2 &\leq \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [f^2(X_i)] + \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \{f^2(X_i) - \mathbb{E} [f^2(X_i)]\} \right| \right] \\ &\leq n^{-1} \Sigma_n(\mathcal{F}) + \frac{16U}{\sqrt{n}} \mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \right]. \end{aligned}$$

See Lemma 6.3 and Theorem 4.12 of [Ledoux and Talagrand \(1991\)](#). Substitute (F.3) in this inequality, we obtain

$$\frac{\Delta^2}{\|F\|_{2, P}^2} \leq \frac{n^{-1} \Sigma_n(\mathcal{F})}{\|F\|_{2, P}^2} + \frac{128\sqrt{2}U}{\sqrt{n} \|F\|_{2, P}} J \left(\frac{\Delta}{\|F\|_{2, P}}, \mathcal{F}, \|\cdot\|_2 \right).$$

For notation convenience, let

$$H(\tau) := J(\tau, \mathcal{F}, \|\cdot\|_2), \quad A^2 := \frac{n^{-1} \Sigma_n(\mathcal{F})}{\|F\|_{2, P}^2}, \quad \text{and} \quad B^2 := \frac{128\sqrt{2}U}{\sqrt{n} \|F\|_{2, P}}.$$

Following the proof of Lemma 2.1 of [van der Vaart and Wellner \(2011\)](#) with $r = 1$, it follows that

$$H\left(\frac{\Delta}{\|F\|_{2,P}}\right) \leq H(A) + \frac{B}{A}H(A)H^{1/2}\left(\frac{\Delta}{\|F\|_{2,P}}\right).$$

Solving the quadratic inequality, we get

$$H\left(\frac{\Delta}{\|F\|_{2,P}}\right) \leq 2\frac{B^2}{A^2}H^2(A) + 2H(A).$$

Substituting this bound in [\(F.3\)](#), it follows that

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)|\right] \leq 16\sqrt{2}\|F\|_{2,P}J(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_2)\left[1 + \frac{128\sqrt{2}UJ(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_2)}{\sqrt{n}\delta_n^2(\mathcal{F})\|F\|_{2,P}}\right].$$

This proves the result. \square

Proof of Proposition B.3. In the proof of Theorem 3.5.13 of [Giné and Nickl \(2016\)](#), the decomposition (3.206) holds as it is and the calculations that follow have to be done for averages of non-identically distributed random variables. For example, the display after (3.206) should be replaced by Lemma 4 of [Pollard \(2002\)](#). (The inequality in Lemma 4 of [Pollard \(2002\)](#) is written for $\sqrt{n}\mathbb{G}_n(f)$ not $n^{-1/2}\mathbb{G}_n(f)$). The variance calculations after (3.209) of [Giné and Nickl \(2016\)](#) should be done as

$$\text{Var}(\mathbb{G}_n(\Delta_k f I_{\{\tau f = k\}})) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\Delta_k f)^2 I(\Delta_k f \leq \alpha_{n,k-1}, \Delta_k(f) > \alpha_{n,k})].$$

See, for example, Lemma 5 of [Pollard \(2002\)](#). There is a typo in Proposition 3.5.15 in the statement; it should be $Pf^2 \leq \delta^2$ for all $f \in \mathcal{F}$. In our case this δ would be the one defined in the statement. The final result follows by noting the concavity of $J_{[\cdot]}(\cdot, \mathcal{F}, \|\cdot\|_{2,P})$,

$$J_{[\cdot]}(2\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_{2,P}) \leq 2J_{[\cdot]}(\delta_n(\mathcal{F}), \mathcal{F}, \|\cdot\|_{2,P}).$$

This completes the proof. \square

Proof of Proposition B.4. It is clear by the triangle inequality that $Z \leq Z_1 + Z_2$ and so,

$$\mathbb{E}[Z] \leq \mathbb{E}[Z_1] + \mathbb{E}[Z_2].$$

From the definition [\(B.4\)](#) of Z_2 , we get

$$\mathbb{E}[Z_2] \leq 2\mathbb{E}\left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n |f(X_i)| \mathbb{1}_{\{|f(X_i)| \geq \rho\}}\right] \leq 2\mathbb{E}\left[\sum_{i=1}^n F(X_i) \mathbb{1}_{\{F(X_i) \geq \rho\}}\right].$$

Using Hoffmann-Jorgensen's inequality along with the definition of ρ , we have

$$\mathbb{E}\left[\sum_{i=1}^n F(X_i) \mathbb{1}_{\{F(X_i) \geq \rho\}}\right] \leq 8\mathbb{E}\left[\max_{1 \leq i \leq n} F(X_i)\right].$$

Therefore,

$$\mathbb{E}[Z] \leq \mathbb{E}[Z_1] + 8\mathbb{E}\left[\max_{1 \leq i \leq n} F(X_i)\right].$$

This completes the proof. \square