

Quantizing Heavy-tailed Data in Statistical Estimation: (Near) Minimax Rates, Covariate Quantization, and Uniform Recovery

Junren Chen, Michael K. Ng, *Senior Member, IEEE*, and Di Wang, *Member, IEEE*

Abstract—Modern datasets often exhibit heavy-tailed behavior, while quantization is inevitable in digital signal processing and many machine learning problems. This paper studies the quantization of heavy-tailed data in several fundamental statistical estimation problems where the underlying distributions have bounded moments of some order (no greater than 4). We propose to truncate and properly dither the data prior to a uniform quantization. Our major standpoint is that (near) minimax rates of estimation error could be achieved by computationally tractable estimators based on the quantized data produced by the proposed scheme. In particular, concrete results are worked out for covariance estimation, compressed sensing (also interpreted as sparse linear regression), and matrix completion, all agreeing that the quantization only slightly worsens the multiplicative factor. Additionally, while prior results focused on the quantization of responses (i.e., measurements), we study compressed sensing where the covariates (i.e., sensing vectors) are also quantized; in this case, though our recovery program is non-convex (since our covariance matrix estimator lacks positive semi-definiteness), we prove that all local minimizers enjoy near-optimal estimation error. Moreover, by the concentration inequality of the product process and a covering argument, we establish a near minimax uniform recovery guarantee for quantized compressed sensing with heavy-tailed noise. Finally, numerical simulations are provided to corroborate our theoretical results.

I. INTRODUCTION

Heavy-tailed distributions are ubiquitous in modern datasets, especially those arising in economy, finance, imaging, biology, see [5], [45], [57], [84], [88], [93] for instance. In the recent literature, heavy-tailed distribution is often captured by bounded l -th moment, where l is some fixed small scalar; this is essentially weaker than the sub-Gaussian assumption. As a result, outliers and extreme values appear much more frequently in data from heavy-tailed distributions (referred to as heavy-tailed data), which poses challenges for statistical

analysis. In fact, many standard statistical procedures developed for sub-Gaussian data suffer from performance degradation in the heavy-tailed regime. Fortunately, the past decade has witnessed considerable progress in statistical estimation methods that are robust to heavy-tailedness, see [9], [17], [35], [36], [44], [65], [68], [73], [97] for instance.

Departing momentarily from heavy-tailed data, quantization, which maps signals to bitstreams so that they can be stored, processed, and transmitted, is an inevitable process in the era of digital signal processing. In particular, the resolution of quantization should be selected to achieve a trade-off between accuracy and various data processing costs. In some applications, relatively low resolution would be preferable. For instance, in a distributed learning setting or a MIMO system, the frequent information transmission among multiple parties often results in prohibitive communication cost [56], [69], and quantizing signals or data to fairly low resolution is an effective approach to reduce the cost [43], [96]. Under such a big picture, in recent years there has been rapidly growing literature on high-dimensional signal recovery from quantized data (see, e.g., [7], [22], [26], [32], [33], [50], [89] for 1-bit quantization, [46], [50], [89], [94] for multi-bit uniform quantization), trying to understand the interplay between quantization and signal reconstruction in some fundamental estimation problems.

Independently, a set of robustifying techniques has been developed to overcome the challenge posed by heavy-tailed data, and uniform data quantization under uniform dither was shown to cost very little in some recovery problems. Considering the ubiquitousness of heavy-tailed behavior and data quantization, a natural question is to design a quantization scheme for heavy-tailed data that only incurs minor information loss. For instance, when applied to statistical estimation problems with heavy-tailed data, an appropriate quantization scheme should enable at least one faithful estimator from the quantized data, and ideally the estimator that could nearly achieve the optimal error rate. Despite the vast literature in this field, prior results that simultaneously take heavy-tailed data and quantization into account are surprisingly rare — only the ones presented in [33] and our earlier work [22] regarding the dithered 1-bit quantizer, to the best of our knowledge. These results remain incomplete and exhibit some downsides. Specifically, [33] considered a computationally intractable program for quantized compressed sensing and used techniques hard to generalize to other problems, while the error rates in [22] are inferior to the corresponding minimax ones (under unquantized sub-Gaussian

Junren Chen is with Department of Mathematics, The University of Hong Kong. Michael K. Ng is with Department of Mathematics, Hong Kong Baptist University. Di Wang is with Division of CEMSE, King Abdullah University of Science and Technology (KAUST). (e-mail: chenjr58@connect.hku.hk; michael-ng@hkbu.edu.hk; di.wang@kaust.edu.sa)

Junren Chen was supported by the Hong Kong Ph.D. Fellowship from the Hong Kong Research Grants Council (HKRGC). Michael K. Ng was partially supported by the HKRGC GRF 17201020, 17300021, CRF C7004-21GF and Joint NSFC-RGC N-HKU76921. Di Wang was supported in part by the baseline funding BAS/1/1689-01-01, funding from the CRG grand URF/1/4663-01-01, REI/1/5232-01-01, REI/1/5332-01-01, FCC/1/1976-49-01 from CBRC of King Abdullah University of Science and Technology (KAUST), and also supported by the funding RGC/3/4816-09-01 of the SDAIA-KAUST Center of Excellence in Data Science and Artificial Intelligence (SDAIA-KAUST AI). (*Corresponding author: Junren Chen.*)

data), as will be discussed in Section I-A3. In a nutshell, a quantization scheme for heavy-tailed data arising in statistical estimation problems that allows for computationally tractable near-minimax estimators is still lacking.

This paper aims to provide a solution to the above question and narrow the gap between heavy-tailed data and data quantization in the literature. In particular, we propose a unified quantization scheme for heavy-tailed data which, when applied to the canonical estimation problems of (sparse) covariance matrix estimation, compressed sensing (or sparse linear regression) and matrix completion, allows for (near) minimax estimators that are either in closed-form or can be solved from convex programs. Additionally, we present novel developments concerning covariate (or sensing vector) quantization and uniform signal recovery in quantized compressed sensing with heavy-tailed data.

A. Related Works

This section is devoted to a review of the most relevant works. Before that we note that a heavy-tailed random variable in this work is formulated by the moment constraint $\mathbb{E}|X|^l \leq M$, where M is an absolute constant and l is some fixed small scalar (specifically, $l \leq 4$ in the present paper). In Sections I-A1 to I-A3 we focus on estimation problems from possibly heavy-tailed and/or quantized data; then, in Sections I-A4 to I-A5 we touch on two specific aspects of quantized compressed sensing, namely the lesser-known covariate quantization problem and the highly sought-after notion of uniform recovery.

1) *Statistical Estimation under Heavy-Tailed Data:* Developing estimation methods that are robust to heavy-tailedness has become a recent focus in the statistics literature, where heavy-tailed distributions are often only assumed to have bounded moments of some small order. In particular, significant efforts have been devoted to the fundamental problem of mean estimation for heavy-tailed distribution. For instance, effective techniques available in the literature include Catoni's mean estimator [17], [35], median of means [68], [73], and trimmed mean [28], [66]. Indeed, these methods share the same core spirit of making the outliers less influential. To this end, the trimmed method (also referred to as truncation or shrinkage) may be the most intuitive — it truncates overlarge data to some threshold so that they are more benign for the estimation procedure. For more in-depth discussions we refer to the recent survey [65]. Furthermore, these robust methods for estimating the mean have been applied to empirical risk minimization [9], [44] and various high-dimensional estimation problems [36], [97], achieving near optimal guarantees. For instance, by invoking M-estimators with truncated data, (near) minimax rates can be achieved in high-dimensional sparse linear regression, matrix completion and covariance estimation [36].

While we capture heavy-tailedness by bounded moment of some small order, there has been a line of works considering sub-exponential or more generally sub-Weibull distributions [39], [58], [80], [85], which have heavier tail than sub-Gaussian ones but still possess finite moment up to arbitrary

order. Specifically, without truncation and quantization, sparse linear regression was studied under sub-exponential data in [85] and under sub-Weibull data in [58], and the obtained error rates match the ones in the sub-Gaussian case up to logarithmic factors. Additionally, under sub-exponential measurement matrix and noise, [80] established a uniform guarantee for 1-bit generative compressed sensing, while [39] analyzed generalized Lasso for a general nonlinear model. Because the tail assumptions in these works are substantially stronger than ours, there is not a common fair stage for further comparison.

2) *Statistical Estimation from Quantized Data: Quantized Compressed Sensing.* While there have been other quantization methods, we only review the most relevant memoryless quantization schemes¹ that allow for simple hardware design, with an emphasis on the benefit of dithering. An important model is 1-bit compressed sensing where only the sign of the measurement is retained [7], [48], [75], [76]; more precisely, this model concerns the recovery of sparse $\theta^* \in \mathbb{R}^d$ from $\text{sign}(\mathbf{X}\theta^*)$ with the sensing matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$. However, 1-bit compressed sensing associated with the direct $\text{sign}(\cdot)$ quantization suffers from some frustrating limitations, e.g., the loss of signal norm information, and the identifiability issue under some regular sensing matrix (e.g., under Bernoulli sensing matrix, see [33]).² Fortunately, these limitations can be overcome by random dithering prior to the quantization, under which the 1-bit measurements read as $\text{sign}(\mathbf{X}\theta^* + \tau)$, with $\tau \in \mathbb{R}^n$ being some suitably chosen random dither. Specifically, under Gaussian dither $\tau \sim \mathcal{N}(0, \mathbf{I}_n)$ and standard Gaussian sensing matrix \mathbf{X} , full reconstruction with norm information could be achieved [54]. More surprisingly, under a uniform random dither, recovery with norm can be achieved under a rather general sub-Gaussian sensing matrix [22], [33], [50], [89] even with near optimal error rate.

Besides the 1-bit quantizer that retains the sign, the uniform quantizer maps $a \in \mathbb{R}$ to $\mathcal{Q}_\Delta(a) = \Delta(\lfloor \frac{a}{\Delta} \rfloor + \frac{1}{2})$ for some pre-specified $\Delta > 0$; here and hereafter, we refer to Δ as the quantization level, and note that smaller Δ represents higher resolution. While recovering θ^* from $\mathcal{Q}_\Delta(\mathbf{X}\theta^*)$ encounters identifiability issue,³ it is again beneficial to use random dithering to obtain the measurements $\mathcal{Q}_\Delta(\mathbf{X}\theta^* + \tau)$. More specifically, by using uniform dither the Lasso estimator [87], [89] and projected back projection (PBP) method [94] achieve minimax rate in certain cases, and the derived error bounds for these estimators demonstrate that the dithered uniform quantization does not affect the scaling law but only slightly worsens the multiplicative factor. Although the aforementioned progress was recently made, the technique of dithering in

¹This means that the quantization methods for different measurements are independent. For other quantization schemes, we refer to the recent survey [29].

²In fact, almost all existing guarantees using the 1-bit observations $\text{sign}(\mathbf{X}\theta^*)$ are restricted to standard Gaussian sensing matrix consisting of i.i.d. $\mathcal{N}(0, 1)$ entries, with the exceptions of [1] for sub-Gaussian sensing matrix and [30], [86] for partial Gaussian circulant matrix.

³For instance, if $\mathbf{X} \in \{-1, 1\}^{n \times d}$ (typical example is the Bernoulli design where entries of \mathbf{X} are i.i.d. zero-mean) and $\Delta = 1$, then $\theta_1 := 1.1e_1$ and $\theta_2 := 1.2e_1 + 0.1e_2$ can never be distinguished because $\mathcal{Q}_1(\mathbf{X}\theta_1) = \mathcal{Q}_1(\mathbf{X}\theta_2)$ always holds.

quantization indeed has a long history and (at least) dates back to some early engineering work (e.g., [83]), see [41] for a brief introduction.

Other Estimation Problems with Quantized Data. Some other statistical estimation problems were also investigated under dithered 1-bit quantization. Specifically, [24] studied a general signal estimation problem in a traditional setting where sample size tends to infinity, showing that dithered 1-bit quantization incurs merely logarithmic rate loss. Inspired by potential application in the reduction of power consumption in a large scale massive MIMO system, [32] proposed to collect 2 bits per entry from each sub-Gaussian sample and developed an estimator that is in general near minimax optimal. Their estimator was extended to the high-dimensional sparse case in [22], and its tuning-free version was devised in [31]. More recently, a parameter-free covariance estimator with improved operator norm error rate was proposed in [19]. Next, considering the ubiquitousness of binary observations in many recommendation systems, the authors of [26] first approached the 1-bit matrix completion problem by maximum likelihood estimation with a nuclear norm constraint. Their method was developed in a series of follow-up works by using different regularizers/constraints to encourage low-rankness, or considering multi-bit quantization on a finite alphabet [3], [14], [16], [52], [59]. Using a uniformly dithered 1-bit quantizer, the 1-bit matrix completion result in [22] essentially departs from the standard likelihood approach and can tolerate pre-quantization noise with unknown distribution.

3) *Quantization of Heavy-Tailed Data in Statistical Estimation:* Note that the results we just reviewed are for estimation problems from either unquantized heavy-tailed data (Section I-A1) or quantized sub-Gaussian data (Section I-A2). In this part, we turn to existing results under quantized heavy-tailed data that are more closely related to this work. Despite being a natural question with practical relevance, quantization of heavy-tailed data was rarely studied in prior work; to our best knowledge, the only results were presented in [22], [33] concerning dithered 1-bit quantization. Specifically, [33, Thm. 1.11] considered heavy-tailed noise and possibly heavy-tailed covariate, implying that a sharp uniform error rate is achievable (see their Example 1.10). However, their result is for a computationally intractable program (Hamming distance minimization) and hence of limited practical value. Another limitation is that their techniques (based on random hyperplane tessellations) are specialized to 1-bit compressed sensing but do not generalize to other estimation problems. In contrast, [22] proposed a unified quantization scheme that first truncates the data and then invokes a dithered 1-bit quantizer. Although this quantization scheme could (at least) be applied to sparse covariance matrix estimation, compressed sensing, and matrix completion while still enabling practical estimators, the main drawback is that the convergence rates of estimation errors are essentially slower than the corresponding minimax optimal ones (e.g., $\tilde{O}\left(\frac{\sqrt{s}}{n^{1/3}}\right)$ for 1-bit compressed sensing under heavy-tailed noise [22, Thm. 10]), and in certain cases the rates cannot be improved without changing the quantization process (e.g., [22, Thm. 11] complements [22, Thm. 10] with a

nearly matching lower bound). In a nutshell, [33] proved a sharp rate for 1-bit compressed sensing but used highly intractable program and its techniques are not extendable to other estimation regimes, while the more widely applicable scheme and practical estimators in [22] suffer from slow error rates.

4) *Covariate Quantization in Compressed Sensing:* In the rest of the review, we will concentrate on quantized compressed sensing, i.e., the recovery of a sparse signal $\theta^* \in \mathbb{R}^d$ from the quantized version of $(x_k, y_k := x_k^\top \theta^* + \epsilon_k)_{k=1}^n$ where x_k, y_k, ϵ_k are the sensing vector, measurement and noise, respectively. Note that this formulation also models the sparse linear regression problem (e.g., [72], [81]) where one wants to learn a sparse parameter $\theta^* \in \mathbb{R}^d$ from the given data $(x_k, y_k)_{k=1}^n$, which are believed to follow the linear model $y_k = x_k^\top \theta^* + \epsilon_k$. In this regression setting, x_k, y_k are commonly referred to as covariate and response, respectively. We are interested in both settings in this work (as further explained in Section III-B), but for clearer presentation, we simply refer to the problem as *quantized compressed sensing*, and term x_k, y_k as *covariate* and *response*, respectively.

Note that studying “*how quantization of covariate affects the recovery/learning*” is meaningful especially when the problem is interpreted as sparse linear regression — working with low-precision data in some (distributed) learning systems could significantly reduce communication cost and power consumption [43], [96], which we will further demonstrate in Section IV-A. However, almost all of the prior works are restricted to response quantization. To the best of our knowledge, the only existing rigorous guarantees involving covariate quantization were obtained in [22, Thms. 7-8]. Nevertheless, these results require $\mathbb{E}(x_k x_k^\top)$ to be sparse [22, Assumption 3] (so that their sparse covariance matrix estimator is applicable). Note that this assumption is non-standard and rarely assumed in sparse linear regression and compressed sensing.⁴ It is thus desired to develop theoretical guarantees without resorting to the sparsity on $\mathbb{E}(x_k x_k^\top)$.

5) *Uniform Signal Recovery in Compressed Sensing:* It is standard in compressed sensing to leverage a random sensing matrix, so a recovery guarantee can be uniform or non-uniform. More precisely, a uniform guarantee ensures the recovery of all structured signals of interest with a single draw of the sensing ensemble, while a non-uniform guarantee is only valid for a structured signal fixed before drawing the random ensemble, with the implication that a new realization of the sensing matrix is required for sensing a new signal. Uniformity is a highly desired property in compressed sensing, since in applications the measurement ensemble is typically fixed and is expected to work for all signals [21], [40]. Besides, the derivation of a uniform guarantee is often significantly harder than a non-uniform one, making uniformity an interesting theoretical problem in its own right. (As a result, in the literature of compressed sensing, theoretical results for

⁴In fact, although isotropic sensing vector (i.e., $\mathbb{E}(x_k x_k^\top) = I_d$) has been conventional in compressed sensing, many results in the literature can be extended to sensing vector with general unknown covariance matrix and hence do not really rely on the sparsity of $\mathbb{E}(x_k x_k^\top)$.

the nonuniform setting are often a precursor to the uniform recovery guarantees.)

A classical fact in linear compressed sensing is that the restricted isometry property (RIP) of the sensing matrix implies uniform recovery of all sparse signals (e.g., [38]), but this is unfortunately not the case when it comes to nonlinear compressed sensing models. For instance, in the specific quantization model, the more general single index model $y_k = f(\mathbf{x}_k^\top \boldsymbol{\theta}^*)$ with possibly unknown $f(\cdot)$, and the lesser-known phase-only model $y_k = \text{sign}(\mathbf{x}_k^\top \boldsymbol{\theta}^*)$ with $\mathbf{x}_k \in \mathbb{C}^d$, most representative results are non-uniform (e.g., [47], [75], [78], [79], [87], [89], [94]). We refer to [20], [33], [50], [75], [94] for concrete uniform guarantees, some of which remain (near) optimal (e.g., [94, Sect. 7.2A]), while others suffer from essential degradation compared to the non-uniform ones (e.g., [75, Thm. 1.3]). It is worth noting that the interesting recent work [40] provided a unified approach to uniform guarantee for a series of non-linear models; however, their uniform guarantees typically exhibit a decaying rate of $O(n^{-1/4})$ that is slower than the non-uniform one of $O(n^{-1/2})$ (Section 4 therein). As a follow-up work, [21] obtained a near optimal uniform $O(n^{-1/2})$ error rate for various nonlinear generative compressed sensing models.

The above-reviewed uniform guarantees are restricted to sub-Gaussian data. Regarding dithered 1-bit quantization of heavy-tailed data, results in [22, Sect. III] are non-uniform, while [33, Thm. 1.11] presents a sharp uniform guarantee for the intractable program of hamming distance minimization.

B. Our Contributions

We summarize our contributions in this part. Our primary contribution is a unified quantization scheme (for heavy-tailed data) that allows for near-minimax estimators. Besides, we present new developments regarding covariate quantization and uniform recovery in a heavy-tailed quantized setting. This work also provides some notable innovations that prove useful in related studies.

A Unified Quantization Scheme and Estimators with (Near) Minimax Rates. We propose a unified quantization scheme for heavy-tailed data which allows for *practical* and (near) *optimal* estimators in *multiple* estimation problems. The proposed scheme consists of three steps: 1) *truncation* that shrinks data to some threshold, 2) *dithering* that adds suitable random noise to the truncated data, and 3) *uniform quantization*. Note that the proposed scheme replaces the 1-bit quantizer in [22] with the less extreme (multi-bit) uniform quantizer $\mathcal{Q}_\Delta(\cdot)$, but the gain turns out to be significant — we are now able to derive (near) optimal rates that are essentially faster than the ones in [22], see Theorems 2-8. As a concrete example, for quantized compressed sensing with sub-Gaussian sensing vector \mathbf{x}_k but heavy-tailed measurement y_k satisfying $\mathbb{E}|y_k|^{2+\nu} \leq M$ for some $\nu > 0$, we derive the ℓ_2 -norm error rate $O((M^{1/(2l)} + \Delta)\sqrt{\frac{s \log d}{n}})$ (Theorem 5, s, d, n are respectively the sparsity, signal dimension, measurement number).

Compared to [33], our major advantages are that our estimators are computationally feasible and that our method applies

to estimation problems beyond quantized compressed sensing. Concerning the effect of quantization, our results (in all the considered models) suggest a unified conclusion — *dithered uniform quantization does not affect the scaling law but only slightly worsens the multiplicative factor*, which generalizes similar findings for quantized compressed sensing in [87], [89], [94] towards two directions, i.e., to the case where heavy-tailed data present and to some other estimation problems (i.e., covariance matrix estimation, matrix completion).

New Results on Covariate Quantization. Besides the above main contributions, we establish the estimation guarantees for quantized compressed sensing under covariate quantization that are free of the non-standard assumption on the sparsity of $\mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$. This provides important relaxations to [22, Thms. 7-8]. More specifically, unlike [22] that relies on the sparsity of $\mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$ to ensure convexity, we deal with the non-convex program directly and manage to prove that all local minimizers deliver near minimax estimation errors (Theorems 9-10). Our analysis is motivated by a line of works on non-convex M-estimator [62]–[64] but also exhibits some essential differences (Remark 5). Further, we extract our techniques as a deterministic framework (Proposition 1) and then use it to establish guarantees for dithered 1-bit quantization as byproducts (Theorems 11-12), which are comparable to [22, Thms. 7-8] but free of sparsity on $\mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$.

Uniform Recovery Guarantee. We additionally contribute to the literature a uniform guarantee for constrained Lasso under the dithered uniform quantization of heavy-tailed response. Specifically, we upgrade our non-uniform Theorem 5 to its uniform version Theorem 13, which states that using a single realization of the sub-Gaussian sensing matrix, heavy-tailed noise and uniform dither, all s -sparse signals within an ℓ_2 -ball can be uniformly recovered up to an ℓ_2 -norm error of $\tilde{O}(\sqrt{\frac{s}{n}})$, thus matching the near minimax non-uniform rate in Theorem 5 up to logarithmic factors. The proof relies on a concentration inequality for product process [67] and a careful covering argument inspired by [94]. Due to the heavy-tailed noise, a new treatment is needed before invoking the concentration result from [67].

Some Notable Innovations. An important innovation of our work is to use triangular dither when covariance estimation is necessary, which departs from the uniform dither commonly adopted in prior works (e.g., [22], [32], [89], [94]) and is novel to the literature. From a technical side, many of our analyses on the dithered quantizer are much cleaner than prior works because we make full use of the nice statistical properties of the quantization error and quantization noise (Theorem 1),⁵ see Section II-B. Based on these two innovations (and thus after the first appearance of the present paper), a clean analysis on quantized low-rank multivariate regression *with possibly quantized covariates* is provided in [23], and more recently a 2-bit covariance estimator based on triangular dithering is devised in [19].

⁵Prior work that did not fully leverage these properties may incur extra technical complication, e.g., the symmetrization and contraction in [89, Lem. A.2].

C. Outline

The remainder of this paper is structured as follows. We provide the notation and preliminaries in Section II. We present the first set of main results (concerning the (near) optimal guarantees for three estimation problems under quantized heavy-tailed data) in Section III. Our second set of results (concerning covariate quantization and uniform recovery in quantized compressed sensing) is then presented in Section IV. To corroborate our theory, numerical results on synthetic data are reported in Section V. We give some remarks to conclude the paper in Section VI. All the proofs are postponed to the Appendices.

II. PRELIMINARIES

We adopt the following conventions throughout the paper:

1) We use boldface symbols (e.g., \mathbf{A} , \mathbf{x}) to denote matrices and vectors, and regular letters (e.g., a , x) for scalars. We write $[m] = \{1, \dots, m\}$ for positive integer m . We denote the complex unit by i . The i -th entry for a vector \mathbf{x} (likewise, \mathbf{y} , $\boldsymbol{\tau}$) is denoted by x_i (likewise, y_i , τ_i).

2) Notation with “ \star ” as superscript denotes the desired underlying parameter or signal, e.g., $\boldsymbol{\Sigma}^\star$, $\boldsymbol{\theta}^\star$. Moreover, notation marked by a tilde (e.g., $\tilde{\mathbf{x}}$) and a dot (e.g., $\dot{\mathbf{x}}$) stands for the truncated data and quantized data, respectively.

3) We reserve d and n for the problem dimension and sample size, respectively. In many cases $\hat{\mathbf{Y}}$ denotes the estimation error, e.g., $\hat{\mathbf{Y}} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star$ if $\hat{\boldsymbol{\theta}}$ is the estimator for the desired signal $\boldsymbol{\theta}^\star$. We use Σ_s to denote the set of d -dimensional s -sparse signals.

4) For vector $\mathbf{x} \in \mathbb{R}^d$, we work with its transpose \mathbf{x}^\top , ℓ_p -norm $\|\mathbf{x}\|_p = (\sum_{i \in [d]} |x_i|^p)^{1/p}$ ($p \geq 1$), max norm $\|\mathbf{x}\|_\infty = \max_{i \in [d]} |x_i|$. We define the standard Euclidean sphere as $\mathbb{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$.

5) For matrix $\mathbf{A} = [a_{ij}] \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}}$, recall the operator norm $\|\mathbf{A}\|_{op} = \sup_{\mathbf{v} \in \mathbb{S}^{n-1}} \|\mathbf{A}\mathbf{v}\|_2 = \sigma_1$, Frobenius norm $\|\mathbf{A}\|_F = (\sum_{i,j} a_{ij}^2)^{1/2}$, nuclear norm $\|\mathbf{A}\|_{nu} = \sum_{k=1}^{\min\{m,n\}} \sigma_k$, and max norm $\|\mathbf{A}\|_\infty = \max_{i,j} |a_{ij}|$. $\lambda_{\min}(\mathbf{A})$ (resp. $\lambda_{\max}(\mathbf{A})$) stands for the minimum eigenvalue (resp. maximum eigenvalue) of a symmetric \mathbf{A} .

6) We denote universal constants by C , c , C_i and c_i , whose value may vary from line to line. We write $T_1 \lesssim T_2$ or $T_1 = O(T_2)$ if $T_1 \leq CT_2$. Conversely, if $T_1 \geq CT_2$ we write $T_1 \gtrsim T_2$ or $T_1 = \Omega(T_2)$. Also, we write $T_1 \asymp T_2$ if $T_1 = O(T_2)$ and $T_2 = \Omega(T_1)$ simultaneously hold.

7) We use $\mathcal{U}(\Omega)$ to denote the uniform distribution over $\Omega \subset \mathbb{R}^N$, $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to denote Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, $t(\nu)$ to denote student's t distribution with degrees of freedom ν .

8) Our technique to handle heavy-tailedness is a data truncation step, for which we introduce the operator $\mathcal{T}_\zeta(\cdot)$ for some threshold $\zeta > 0$. It is defined as $\mathcal{T}_\zeta(a) = \text{sign}(a) \min\{|a|, \zeta\}$ for some $a \in \mathbb{R}$. To truncate vectors we apply $\mathcal{T}_\zeta(\cdot)$ entry-wisely in most cases, with the exception of covariance matrix estimation under operator norm error (Theorem 3).

9) $\mathcal{Q}_\Delta(\cdot)$ is the uniform quantizer with quantization level $\Delta > 0$. It applies to scalar a by $\mathcal{Q}_\Delta(a) = \Delta(\lfloor \frac{a}{\Delta} \rfloor + \frac{1}{2})$, and

we set $\mathcal{Q}_0(a) = a$. Given a threshold μ , the hard thresholding of scalar a is $\mathcal{T}_\mu(a) = a \cdot \mathbb{1}(|a| \geq \mu)$. Both functions element-wisely apply to vectors or matrices.

A. High-Dimensional Statistics

Let X be a real random variable, we present some basic knowledge of the sub-Gaussian and sub-exponential random variables. Then we also precisely formulate the heavy-tailed distribution.

1) The sub-Gaussian norm is defined as $\|X\|_{\psi_2} = \inf\{t > 0 : \mathbb{E} \exp(\frac{X^2}{t^2}) \leq 2\}$. A random variable X with finite $\|X\|_{\psi_2}$ is said to be sub-Gaussian. Analogously to the Gaussian variable, a sub-Gaussian random variable exhibits an exponentially decaying probability tail and satisfies a moment constraint:

$$\mathbb{P}(|X| \geq t) \leq 2 \exp\left(-\frac{ct^2}{\|X\|_{\psi_2}^2}\right); \quad (1)$$

$$(\mathbb{E}|X|^p)^{1/p} \leq C\|X\|_{\psi_2} \sqrt{p}, \quad \forall p \geq 1. \quad (2)$$

Note that these two properties can also define $\|\cdot\|_{\psi_2}$ up to multiplicative constant, e.g., $\|X\|_{\psi_2} \asymp \sup_{p \geq 1} \frac{(\mathbb{E}|X|^p)^{1/p}}{\sqrt{p}}$ (see [92, Prop. 2.5.2]). For a d -dimensional random vector \mathbf{x} we define its sub-Gaussian norm as $\|\mathbf{x}\|_{\psi_2} = \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \|\mathbf{v}^\top \mathbf{x}\|_{\psi_2}$.

2) The sub-exponential norm is defined as $\|X\|_{\psi_1} = \inf\{t > 0 : \mathbb{E} \exp(\frac{|X|}{t}) \leq 2\}$, and X is sub-exponential if $\|X\|_{\psi_1} < \infty$. The sub-exponential X satisfies the following properties:

$$\mathbb{P}(|X| \geq t) \leq 2 \exp\left(-\frac{ct}{\|X\|_{\psi_1}}\right);$$

$$(\mathbb{E}|X|^p)^{1/p} \leq C\|X\|_{\psi_1} p, \quad \forall p \geq 1. \quad (3)$$

To relate $\|\cdot\|_{\psi_1}$ and $\|\cdot\|_{\psi_2}$ one has $\|XY\|_{\psi_2} \leq \|X\|_{\psi_1} \|Y\|_{\psi_1}$ [92, Lem. 2.7.7].

3) In contrast to the moment constraints in (2) and (3), heavy-tailed distributions in this work are only assumed to satisfy bounded moments of some small order no greater than 4, formulated for a random variable X as $\mathbb{E}|X|^l \leq M$ for some $M > 0$ and $l \in (0, 4]$. Following [58, Def. 2.4, 2.5], we consider the following two moment assumptions for a heavy-tailed random vector $\mathbf{x} \in \mathbb{R}^d$ (again, $M > 0$, $l \in (0, 4]$):

- **Marginal Moment Constraint.** The weaker assumption that constrains the moment of each coordinate is formulated by $\sup_{i \in [d]} \mathbb{E}|x_i|^l \leq M$.
- **Joint Moment Constraint.** The stronger assumption that constrains the moments “toward all directions $\mathbf{v} \in \mathbb{S}^{d-1}$,” is formulated by $\sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}|\mathbf{v}^\top \mathbf{x}|^l \leq M$.

B. Dithered Uniform Quantization

In this part, we describe the dithered uniform quantizer and its properties in detail. We also specify the choices of random dither in this work.

1) We first provide the detailed procedure of dithered quantization and its general property. Let $\mathbf{x} \in \mathbb{R}^N$ be the input signal with dimension $N \geq 1$ whose entries may be random and dependent. Independent of \mathbf{x} , we generate the random

dither $\tau \in \mathbb{R}^N$ with i.i.d. entries from some distribution,⁶ and then quantize x to $\hat{x} = \mathcal{Q}_\Delta(x + \tau)$. Following [41], we refer to $w := \hat{x} - (x + \tau)$ as the quantization error, and $\xi := \hat{x} - x$ as the quantization noise. The principal properties of dithered quantization are provided in Theorem 1.

Theorem 1. (Adapted from [41, Thms. 1-2]). *Consider the dithered uniform quantization described above for the input signal x , with random dither $\tau = [\tau_i]$, quantization error w and quantization noise $\xi = [\xi_i]$. Use i to denote the imaginary unit, and let Y be the random variable having the same distribution as the random dither τ_i .*

(a) (Quantization Error). *If $f(u) := \mathbb{E}(\exp(iuY))$ satisfies $f(\frac{2\pi l}{\Delta}) = 0$ for all non-zero integer l , then $w \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^N)$ is independent of x .*⁷

(b) (Quantization Noise). *Assume that $Z \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ is independent of Y . Let $g(u) := \mathbb{E}(\exp(iuY))\mathbb{E}(\exp(iuZ))$. Given positive integer p , if the p -th order derivative $g^{(p)}(u)$ satisfies $g^{(p)}(\frac{2\pi l}{\Delta}) = 0$ for all non-zero integer l , then the p -th conditional moment of ξ_i does not depend on x : $\mathbb{E}[\xi_i^p | x] = \mathbb{E}(Y + Z)^p$.*

We note that Theorem 1 serves as the cornerstone for our analysis on the dithered uniform quantizer; for instance, (a) allows for applications of concentration inequalities in our analyses, and (b) inspires us to develop a covariance matrix estimator from quantized samples. The take-home message is that, adding appropriate dither before quantization can make the quantization error and quantization noise behave in a statistically nice manner. For example, the elementary form of Theorem 1(a) is that under a dither τ_i satisfying the condition there, the quantization noise $\mathcal{Q}_\Delta(x_i + \tau_i) - (x_i + \tau_i)$ follows $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ under any given scalar x_i [41, Lem. 1].

2) We use *uniform dither* for quantization of the response in compressed sensing and matrix completion. More specifically, under $\Delta > 0$, we adopt the uniform dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ for the response $y_k \in \mathbb{R}$, which is also a common choice in previous works (e.g., [33], [50], [89], [94]). For $Y \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$, it can be calculated that

$$\begin{aligned} \mathbb{E}(\exp(iuY)) &= \int_{-\Delta/2}^{\Delta/2} \frac{1}{\Delta} (\cos(ux) + i \sin(ux)) dx \\ &= \frac{2}{\Delta u} \sin\left(\frac{\Delta u}{2}\right), \end{aligned} \quad (4)$$

and hence $\mathbb{E}(\exp(i\frac{2\pi l}{\Delta}Y)) = 0$ holds for all non-zero integer l . Therefore, the benefit of using $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ is that the quantization errors $w_k = \mathcal{Q}_\Delta(y_k + \tau_k) - (y_k + \tau_k)$ i.i.d. follow $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$, and are independent of $\{y_k\}$.

3) We use *triangular dither* for quantization of the covariate, i.e., the sample in covariance estimation or the covariate in compressed sensing. Particularly, when considering the uniform quantizer $\mathcal{Q}_\Delta(\cdot)$ for the covariate $x_k \in \mathbb{R}^d$, we adopt

the dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d) + \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$,⁸ which is the sum of two independent $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$ and referred to as a triangular dither [41]. Simple calculations verify that the triangular dither respects not only the condition in Theorem 1(a), but also the one in Theorem 1(b) with $p = 2$; specifically, let $Y = Y_1 + Y_2$ where Y_1 and Y_2 are independent and follow $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$, and let $Z \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ be independent of Y , then based on (4), we know that

$$f(u) = \mathbb{E}(\exp(iuY)) = \left[\frac{2}{\Delta u} \sin \frac{\Delta u}{2} \right]^2$$

satisfies $f(\frac{2\pi l}{\Delta}) = 0$, and that

$$g(u) = \mathbb{E}(\exp(iuY))\mathbb{E}(\exp(iuZ)) = \left[\frac{2}{\Delta u} \sin \frac{\Delta u}{2} \right]^3$$

satisfies $g''(\frac{2\pi l}{\Delta}) = 0$, where l is any non-zero integer. Thus, at the cost of a dithering variance larger than uniform dither, the triangular dither brings the additional nice property of signal-independent variance for the quantization noise — $\mathbb{E}(\xi_{ki}^2) = \frac{1}{4}\Delta^2$, where ξ_{ki} is the i -th entry of $\xi_k = \mathcal{Q}_\Delta(x_k + \tau_k) - (x_k + \tau_k)$.

To the best of our knowledge, the triangular dither is new to the literature of quantized compressed sensing. We will explain its necessity if covariance estimation is involved. This is also complemented by numerical simulation (see Figure 5(a)).

III. (NEAR) MINIMAX ERROR RATES

In this section, we derive (near) optimal error rates for several canonical statistical estimation problems. Our novelty is that by using the proposed quantization scheme for heavy-tailed data, (near) optimal error rates could be achieved by computationally feasible estimators.

A. Quantized Covariance Matrix Estimation

Given $\mathcal{X} := \{x_1, \dots, x_n\}$ as i.i.d. copies of a zero-mean random vector $x \in \mathbb{R}^d$, one often encounters the covariance matrix estimation problem, i.e., to estimate $\Sigma^* = \mathbb{E}(xx^\top)$. This estimation problem is of fundamental importance in multivariate analysis and has attracted much research interest (e.g., [4], [10], [12], [13]). However, the practically useful setting (e.g., in a massive MIMO system [95]) where the samples undergo certain quantization process remains underdeveloped, for which we are only aware of the 1-bit quantization results in [22], [32]. This setting poses the problem of quantized covariance matrix estimation (QCME), in which *one aims to design a quantization scheme for x_k that allows for accurate estimation of Σ^* only based on the quantized samples*. We consider heavy-tailed x_k that possesses bounded fourth moments either marginally or jointly, but note that our estimation methods and theoretical results appear to be new even for sub-Gaussian x_k (Remark 1).

As introduced before, we overcome the heavy-tailedness of x_k by a data truncation step, i.e., we first truncate x_k to \tilde{x}_k in order to make the outliers less influential. Here, we defer the

⁶Throughout this work, we suppose that a random dither is drawn independent of anything else (particularly, the signal to be quantized and other dithers), and the dither has i.i.d. entries if it is a vector.

⁷Although the statement is a bit different, it can be implied by [41, Thm. 1] and the proof therein.

⁸An equivalent statement is that entries of τ_k are i.i.d. distributed as $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]) + \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$. The equivalence can be clearly seen by comparing the joint probability density functions.

precise definition of $\tilde{\mathbf{x}}_k$ to concrete results because it should be well suited to the error metric. After the truncation, we dither and quantize $\tilde{\mathbf{x}}_k$ to $\hat{\mathbf{x}}_k = \mathcal{Q}_\Delta(\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k)$ with the triangular dither $\boldsymbol{\tau}_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d) + \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$. Different from the uniform dither adopted in the literature (e.g., [22], [33], [50], [89], [94]), first let us explain our choice of triangular dither. Recall that the quantization noise and quantization error are respectively defined as $\boldsymbol{\xi}_k := \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k$ and $\mathbf{w}_k := \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k - \boldsymbol{\tau}_k$, thus giving $\boldsymbol{\xi}_k = \boldsymbol{\tau}_k + \mathbf{w}_k$. Under uniform dither or triangular dither, \mathbf{w}_k is independent of $\tilde{\mathbf{x}}_k$ and follows $\mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$ (see Section 2.2), thus allowing us to calculate that

$$\begin{aligned} \mathbb{E}(\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top) &= \mathbb{E}((\tilde{\mathbf{x}}_k + \boldsymbol{\xi}_k)(\tilde{\mathbf{x}}_k + \boldsymbol{\xi}_k)^\top) \\ &= \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top) + \mathbb{E}(\tilde{\mathbf{x}}_k \boldsymbol{\xi}_k^\top) + \mathbb{E}(\boldsymbol{\xi}_k \tilde{\mathbf{x}}_k^\top) + \mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top) \\ &\stackrel{(i)}{=} \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top) + \mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top). \end{aligned} \quad (5)$$

Note that (i) is because

$$\begin{aligned} \mathbb{E}(\boldsymbol{\xi}_k \tilde{\mathbf{x}}_k^\top) &= \mathbb{E}(\boldsymbol{\tau}_k \tilde{\mathbf{x}}_k^\top) + \mathbb{E}(\mathbf{w}_k \tilde{\mathbf{x}}_k^\top) \\ &= \mathbb{E}(\boldsymbol{\tau}_k) \mathbb{E}(\tilde{\mathbf{x}}_k^\top) + \mathbb{E}(\mathbf{w}_k) \mathbb{E}(\tilde{\mathbf{x}}_k^\top) = 0 \end{aligned}$$

due to the previously noted fact that $\boldsymbol{\tau}_k$ and \mathbf{w}_k are independent of $\tilde{\mathbf{x}}_k$ and zero-mean. While with suitable choice of the truncation threshold $\mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top)$ is expected to well approximate $\boldsymbol{\Sigma}^*$, the remaining $\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top)$ gives rise to constant bias. To address the issue, a straightforward idea is to remove the bias, which requires the full knowledge of $\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top)$, i.e., the covariance matrix of the quantization noise. For $i \neq j$, because $\boldsymbol{\tau}_k, \mathbf{w}_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$ and

$$\mathbb{E}(w_{ki} \tau_{kj}) = \mathbb{E}_{\tilde{\mathbf{x}}_{ki}}(\mathbb{E}[w_{ki} \tau_{kj} | \tilde{x}_{ki}]) = 0$$

(note that conditionally on \tilde{x}_{ki} , $w_{ki} = \mathcal{Q}_\Delta(\tilde{x}_{ki} + \tau_{ki}) - (\tilde{x}_{ki} + \tau_{ki})$ and τ_{kj} are independent), we have

$$\begin{aligned} \mathbb{E}(\xi_{ki} \xi_{kj}) &= \mathbb{E}((w_{ki} + \tau_{ki})(w_{kj} + \tau_{kj})) \\ &= \mathbb{E}(w_{ki} w_{kj}) + \mathbb{E}(w_{ki} \tau_{kj}) + \mathbb{E}(\tau_{ki} w_{kj}) + \mathbb{E}(\tau_{ki} \tau_{kj}) \\ &= 0, \end{aligned}$$

showing that $\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top)$ is diagonal. Moreover, under triangular dither the i -th diagonal entry is also known as $\mathbb{E}|\xi_{ki}|^2 = \frac{\Delta^2}{4}$, see Section II-B. Taken collectively, we arrive at

$$\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top) = \frac{\Delta^2}{4} \mathbf{I}_d; \quad (6)$$

Based on (5) we thus propose the following estimator

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\Delta^2}{4} \mathbf{I}_d, \quad (7)$$

which is the sample covariance of the quantized sample $\hat{\mathcal{X}} := \{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_n\}$ followed by a correction step. On the other hand, the reason why the standard uniform dither is not suitable for QCME becomes self-evident — the diagonal of $\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top)$ remains unknown⁹ and hence there is no hope to precisely remove the bias.

We are now ready to present error bounds for $\hat{\boldsymbol{\Sigma}}$ under max-norm, operator norm. We will also investigate the high-dimensional setting by assuming some sparse structure of

⁹It depends on the input signal, see [41, Page 3].

$\boldsymbol{\Sigma}^*$, for which we propose a thresholding estimator. More concretely, our first result provides the error rate under $\|\cdot\|_\infty$, in which we assume \mathbf{x}_k satisfies the marginal fourth-moment constraint and utilize an element-wise truncation $\tilde{\mathbf{x}}_k = \mathcal{T}_\zeta(\mathbf{x}_k)$.

Theorem 2. (Element-Wise Error). *Given $\Delta > 0$ and $\delta > 4$, we consider the problem of QCME described above. We suppose that $\mathbf{x}_{k,s}$ are i.i.d. zero-mean and satisfy the marginal moment constraint $\mathbb{E}|x_{ki}|^4 \leq M$ for any $i \in [d]$, where x_{ki} is the i -th entry of \mathbf{x}_k . We truncate \mathbf{x}_k to $\tilde{\mathbf{x}}_k = [\tilde{x}_{ki}] = \mathcal{T}_\zeta(\mathbf{x}_k)$ with threshold $\zeta \asymp (\frac{nM}{\delta \log d})^{1/4}$, then quantize $\tilde{\mathbf{x}}_k$ to $\hat{\mathbf{x}}_k = \mathcal{Q}_\Delta(\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k)$ with triangular dither $\boldsymbol{\tau}_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d) + \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$. If $n \gtrsim \delta \log d$, then the estimator in (7) satisfies*

$$\mathbb{P} \left(\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_\infty \geq C \mathcal{L} \sqrt{\frac{\delta \log d}{n}} \right) \leq 2d^{2-\delta},$$

where $\mathcal{L} := \sqrt{M} + \Delta^2$.

Notably, despite the heavy-tailedness and quantization, the estimator achieves an element-wise rate $O(\sqrt{\frac{\log d}{n}})$ coincident with the one for the sub-Gaussian case. One can clearly position quantization level Δ in the multiplicative factor $\mathcal{L} = \sqrt{M} + \Delta^2$. Thus, the information loss incurred by quantization is inessential in that it does not affect the key scaling law but only slightly worsens the leading factor. These remarks on the (near) optimality and the information loss incurred by quantization remain valid in our subsequent theorems.

Our next result concerns the operator norm estimation error, under which we impose a stronger joint moment constraint on \mathbf{x}_k and truncate \mathbf{x}_k regarding ℓ_4 -norm, i.e., $\tilde{\mathbf{x}}_k = \frac{\mathbf{x}_k}{\|\mathbf{x}_k\|_4} \min\{\|\mathbf{x}_k\|_4, \zeta\}$ for some threshold ζ . After the dithered uniform quantization, we still define the estimator as (7).

Theorem 3. (Operator Norm Error). *Given $\Delta > 0$ and $\delta > 0$, we consider the problem of QCME described above. Suppose that the i.i.d. zero-mean $\mathbf{x}_{k,s}$ satisfy $\mathbb{E}|\mathbf{v}^\top \mathbf{x}_k|^4 \leq M$ for any $\mathbf{v} \in \mathbb{S}^{d-1}$. We truncate \mathbf{x}_k to $\tilde{\mathbf{x}}_k = \frac{\mathbf{x}_k}{\|\mathbf{x}_k\|_4} \min\{\|\mathbf{x}_k\|_4, \zeta\}$ with threshold $\zeta \asymp (M^{1/4} + \Delta) (\frac{n}{\delta \log d})^{1/4}$, then quantize $\tilde{\mathbf{x}}_k$ to $\hat{\mathbf{x}}_k = \mathcal{Q}_\Delta(\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k)$ with triangular dither $\boldsymbol{\tau}_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d) + \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}]^d)$. If $n \gtrsim \delta d \log d$, then the estimator in (7) satisfies*

$$\mathbb{P} \left(\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{op} \geq C \mathcal{L} \sqrt{\frac{\delta d \log d}{n}} \right) \leq 2d^{-\delta},$$

with $\mathcal{L} := \sqrt{M} + \Delta^2$.

The operator norm error rate in Theorem 3 is near minimax optimal, e.g., compared to the lower bound in [36, Thm. 7], which states that for any estimator $\hat{\boldsymbol{\Sigma}}$ of the positive semi-definite matrix $\boldsymbol{\Sigma}^*$ based on i.i.d. zero-mean $\{\mathbf{x}_k\}_{k=1}^n$ with covariance matrix $\boldsymbol{\Sigma}^*$, there exists some $\mathbf{v}_0 \in \mathbb{S}^{d-1}$ such that $\mathbb{P}(\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{op} \geq \frac{1}{48} \sqrt{\frac{6d}{n}}) \geq \frac{1}{3}$, where $\boldsymbol{\Sigma}^* = \mathbf{I}_d + \mathbf{v}_0 \mathbf{v}_0^\top$. Again, the quantization only affects the multiplicative factor \mathcal{L} . Nevertheless, one still needs (at least) $n \gtrsim d$ to achieve a small operator norm error. In fact, in a high-dimensional

setting where d may exceed n , even the sample covariance $\frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^\top$ for sub-Gaussian zero-mean \mathbf{x}_k may have extremely bad performance. To achieve small operator norm error in a high-dimensional regime, we resort to additional structure on Σ^* , and specifically we use column-wise sparsity as an example, which corresponds to the situations where dependencies among different coordinates are weak. Based on the estimator in Theorem 2, we further invoke a thresholding regularization [4], [12] to promote sparsity.

Theorem 4. (Sparse QCME). *Under conditions and estimator $\widehat{\Sigma}$ in Theorem 2, we additionally assume that all columns of $\Sigma^* = [\sigma_{ij}^*]$ are s -sparse and consider the thresholding estimator $\widehat{\Sigma}_s := \mathcal{T}_\mu(\widehat{\Sigma})$ for some μ (recall that $\mathcal{T}_\mu(a) = a \cdot \mathbb{1}(|a| \geq \mu)$ for $a \in \mathbb{R}$). If $\mu = C_1(\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$ with sufficiently large C_1 , then $\widehat{\Sigma}_s$ satisfies*

$$\mathbb{P} \left(\|\widehat{\Sigma}_s - \Sigma^*\|_{op} \leq C \mathcal{L} s \sqrt{\frac{\delta \log d}{n}} \right) \geq 1 - \exp(-0.25\delta),$$

where $\mathcal{L} := \sqrt{M} + \Delta^2$.

Notably, our estimator $\widehat{\Sigma}_s$ achieves minimax rates $O\left(s \sqrt{\frac{\log d}{n}}\right)$ under operator norm, e.g., compared to the minimax lower bound derived in [12, Thm. 2], which states that (under some regular scaling) for any covariance estimator Σ_{es} based on n i.i.d. samples of $\mathcal{N}(\boldsymbol{\mu}, \Sigma^*)$ where Σ^* is the true covariance matrix, there exists some covariance matrix Σ^* with s -sparse columns such that $\mathbb{E} \|\Sigma_{es} - \Sigma^*\|_{op}^2 \gtrsim s^2 \frac{\log d}{n}$. Note that this lower bound is for a general sparse covariance matrix with column-wise sparsity, and it is possible to achieve faster rate over covariance matrices with more specific sparse structures, see [11] for instance.

To analyze the thresholding estimator, our proof resembles the ones developed in prior works (e.g., [12]) but requires more effort like bounding the additional bias terms arising from the data truncation and quantization. We also point out that the results for the full-data unquantized regime immediately follow by setting $\Delta = 0$, thus Theorems 2-3 represent the strict extension of [36, Sect. 4], and Theorem 4 complements [36] with a high-dimensional sparse setting.

As the parameter choices in Theorems 2-4 rely on M that somehow connects to $\|\Sigma^*\|_\infty$, our methods require certain prior estimate of the unknown Σ^* , or otherwise a careful tuning of the parameter. It would be interesting to investigate how to address this in practice.

Remark 1. (Sub-Gaussian Case). *While we concentrate on the quantization of heavy-tailed data in this work, our results can be readily adjusted to sub-Gaussian \mathbf{x}_k , for which the truncation step is inessential and can be removed (i.e., $\zeta = \infty$). These results are also new to the literature but will not be presented here.*

B. Quantized Compressed Sensing

We consider the linear model

$$y_k = \mathbf{x}_k^\top \boldsymbol{\theta}^* + \epsilon_k, \quad k = 1, \dots, n, \quad (8)$$

where \mathbf{x}_k s are the covariates, y_k s are responses, $\boldsymbol{\theta}^*$ is the sparse signal in compressed sensing or sparse parameter vector in high-dimensional linear regression that we want to estimate. In the quantized compressed sensing (QCS) problem, we are interested in *developing quantization scheme for (\mathbf{x}_k, y_k) s (mainly for y_k in prior works) that enables accurate recovery of $\boldsymbol{\theta}^*$ based on the quantized data.*

In spite of the same mathematical formulation, there are some important differences between compressed sensing and sparse linear regression that we should clarify first. Specifically, different from sensing vectors in compressed sensing that are generated by some analog measuring device and can oftentimes be designed, \mathbf{x}_k s in sparse linear regression represent the sample data from certain datasets that are believed to affect the responses y_k s through (8). While the sparsity of $\boldsymbol{\theta}^*$ is arguably the most classical signal structure for compressed sensing, due to good interpretability it is also commonly adopted to achieve dimension reduction in high-dimensional statistics. In this work, we are interested in both problem settings. Thus, we do not adopt the isotropic convention (i.e., $\mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top) = \mathbf{I}_d$) from compressed sensing but instead deal with \mathbf{x}_k having general unknown covariance matrix. While the study of quantization and heavy-tailed noise is meaningful in both settings, we note that some of our subsequent results are mainly of interest to the specific sensing or regression problem. For instance, the heavy-tailed covariate considered in Theorem 6 is primarily motivated by the regression setting, in which \mathbf{x}_k may come from a dataset that exhibits a much heavier tail than sub-Gaussian data. Moreover, as will be elaborated in Section IV when appropriate, our subsequent results on covariate quantization (resp., uniform signal recovery guarantee) may prove more useful to the regression problem (resp., compressed sensing problem).

To fix idea, we assume that \mathbf{x}_k s are i.i.d. drawn from some multi-variate distribution, ϵ_k s are i.i.d. statistical noise independent of the \mathbf{x}_k s, and we truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ and then quantize it to $\hat{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}\left(-\frac{\Delta}{2}, \frac{\Delta}{2}\right)$. Under these statistical assumptions and dithered quantization, near-optimal recovery guarantees have been established in [89], [94] for the regime where both \mathbf{x}_k and ϵ_k are drawn from sub-Gaussian distributions (hence the truncation is not needed). In contrast, our focus is on the quantization of heavy-tailed data. Particularly, we always assume that the noise ϵ_k s are i.i.d. drawn from some heavy-tailed distribution, resulting in heavy-tailed responses. We will separately deal with the case of sub-Gaussian covariate and a more challenging situation where \mathbf{x}_k s are also heavy-tailed.

To estimate the sparse $\boldsymbol{\theta}^*$, a classical approach is via the regularized M-estimator known as Lasso [70], [72], [90]

$$\arg \min_{\boldsymbol{\theta}} \frac{1}{2n} \sum_{k=1}^n (y_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2 + \lambda \|\boldsymbol{\theta}\|_1,$$

whose objective combines the ℓ_2 -loss for data fidelity and ℓ_1 -norm that encourages sparsity. Because we can only access the quantized data $(\mathbf{x}_k, \hat{y}_k)$ (or even $(\tilde{\mathbf{x}}_k, \hat{y}_k)$ if covariate quantization is involved, see Section IV), the main issue lies in the ℓ_2 -loss $\frac{1}{2n} \sum_{k=1}^n (y_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2$ that requires the unquantized

data (\mathbf{x}_k, y_k) . To resolve the issue, we calculate the expected ℓ_2 -loss:

$$\begin{aligned} \mathbb{E}(y_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2 &\stackrel{(i)}{=} \boldsymbol{\theta}^\top \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top) \boldsymbol{\theta} - 2\mathbb{E}(y_k \mathbf{x}_k)^\top \boldsymbol{\theta} : \\ &\stackrel{(ii)}{=} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^* \boldsymbol{\theta} - 2\boldsymbol{\Sigma}_{yx}^\top \boldsymbol{\theta}, \end{aligned} \quad (9)$$

where (i) holds up to an inessential constant $\mathbb{E}|y_k|^2$, and in (ii) we let $\boldsymbol{\Sigma}^* := \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$, $\boldsymbol{\Sigma}_{yx} = \mathbb{E}(y_k \mathbf{x}_k)$. This inspires us to generalize the ℓ_2 loss to $\frac{1}{2} \boldsymbol{\theta}^\top \mathbf{Q} \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta}$ and consider the following program

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathcal{S}} \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{Q} \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta} + \lambda \|\boldsymbol{\theta}\|_1. \quad (10)$$

Compared to (9) we will use (\mathbf{Q}, \mathbf{b}) that well approximates $(\boldsymbol{\Sigma}^*, \boldsymbol{\Sigma}_{yx})$, and we also introduce the constraint $\boldsymbol{\theta} \in \mathcal{S}$ to allow more flexibility. It is important to note that this is the general strategy in this work to design estimators in different QCS settings, see more discussions in Remark 3.

The next theorem is concerned with QCS under sub-Gaussian covariate but heavy-tailed response. Note that the heavy-tailedness of y_k stems from the noise distribution assumed to have bounded $2 + \nu$ moment ($\nu = 2(l - 1) > 0$ in the theorem statement), but following [22], [36], [97] we directly impose the moment constraint on the response.

Theorem 5. (Sub-Gaussian Covariate, Heavy-Tailed Response). *Given some $\delta > 0, \Delta > 0$, in (8) we suppose that \mathbf{x}_k s are i.i.d., zero-mean sub-Gaussian with $\|\mathbf{x}_k\|_{\psi_2} \leq \sigma$, $\kappa_0 \leq \lambda_{\min}(\boldsymbol{\Sigma}^*) \leq \lambda_{\max}(\boldsymbol{\Sigma}^*) \leq \kappa_1$ for some $\kappa_1 > \kappa_0 > 0$ where $\boldsymbol{\Sigma}^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$, $\boldsymbol{\theta}^* \in \mathbb{R}^d$ is s -sparse, the noise ϵ_k s are i.i.d. heavy-tailed and independent of \mathbf{x}_k s, and we assume $\mathbb{E}|y_k|^{2l} \leq M$ for some fixed $l > 1$. In the quantization, we truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with threshold $\zeta_y \asymp \left(\frac{nM^{1/l}}{\delta \log d}\right)^{1/2}$, then quantize \tilde{y}_k to $\dot{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}\left[-\frac{\Delta}{2}, \frac{\Delta}{2}\right]$. For recovery, we define the estimator $\hat{\boldsymbol{\theta}}$ as (10) with*

$$\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^\top, \quad \mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \mathbf{x}_k, \quad \mathcal{S} = \mathbb{R}^d.$$

We set

$$\lambda = C_1 \frac{\sigma^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 . If $n \gtrsim \delta s \log d$ for some hidden constant only depending on (κ_0, σ) , then with probability at least $1 - 9d^{1-\delta}$, the estimation error $\hat{\boldsymbol{\Upsilon}} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ satisfies

$$\|\hat{\boldsymbol{\Upsilon}}\|_2 \leq C_3 \mathcal{L} \sqrt{\frac{\delta s \log d}{n}} \quad \text{and} \quad \|\hat{\boldsymbol{\Upsilon}}\|_1 \leq C_4 \mathcal{L} s \sqrt{\frac{\delta \log d}{n}}$$

where $\mathcal{L} := \frac{\sigma^2 (\Delta + M^{1/(2l)})}{\kappa_0^{3/2}}$.

The rate $O\left(\sqrt{\frac{s \log d}{n}}\right)$ for ℓ_2 -norm error is minimax optimal up to logarithmic factor (e.g., compared to [81]). Note that a random noise bounded by Δ roughly contributes Δ to $(\mathbb{E}|y_k|^{2l})^{1/(2l)}$, and the latter is bounded by $M^{1/(2l)}$; because in the error bound Δ and $M^{1/(2l)}$ almost play the same role, the effect of uniform quantization can be readily interpreted

as an additional bounded noise, analogously to the error rate in [87].

Next, we switch to the more challenging situation where both \mathbf{x}_k and y_k are heavy-tailed, assuming that they both possess bounded fourth moments (a marginal moment constraint for \mathbf{x}_k). The consideration of this setting is motivated by the setting of sparse linear regression, where the covariates \mathbf{x}_k s may oftentimes exhibit heavy-tailed behavior. Specifically, we element-wisely truncate \mathbf{x}_k to $\tilde{\mathbf{x}}_k$ and set $\mathbf{Q} := \frac{1}{n} \sum_{k=1}^n \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top$ as a robust covariance matrix estimator, whose estimation performance under $\|\cdot\|_\infty$ follows immediately from Theorem 2 by setting $\Delta = 0$.

Theorem 6. (Heavy-Tailed Covariate, Heavy-Tailed Response). *Given some $\delta > 0, \Delta > 0$, in (8) we suppose that \mathbf{x}_k s are i.i.d. zero-mean satisfying a marginal fourth moment constraint $\sup_{i \in [d]} \mathbb{E}|x_{ki}|^4 \leq M$, $\kappa_0 \leq \lambda_{\min}(\boldsymbol{\Sigma}^*) \leq \lambda_{\max}(\boldsymbol{\Sigma}^*) \leq \kappa_1$ for some $\kappa_1 > \kappa_0 > 0$ where $\boldsymbol{\Sigma}^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$, $\boldsymbol{\theta}^* \in \Sigma_s$ satisfies $\|\boldsymbol{\theta}^*\|_1 \leq R$, the noise ϵ_k s are i.i.d. heavy-tailed and independent of \mathbf{x}_k s, and we assume $\mathbb{E}|y_k|^4 \leq M$. In the quantization, we truncate \mathbf{x}_k, y_k respectively to $\tilde{\mathbf{x}}_k = [\tilde{x}_{ki}] = \mathcal{T}_{\zeta_x}(\mathbf{x}_k)$, $\tilde{y}_k := \mathcal{T}_{\zeta_y}(y_k)$ with $\zeta_x, \zeta_y \asymp \left(\frac{nM}{\delta \log d}\right)^{1/4}$, then we quantize \tilde{y}_k to $\dot{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}\left[-\frac{\Delta}{2}, \frac{\Delta}{2}\right]$. For recovery, we define the estimator $\hat{\boldsymbol{\theta}}$ as (10) with*

$$\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top, \quad \mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \tilde{\mathbf{x}}_k, \quad \mathcal{S} = \mathbb{R}^d.$$

We set

$$\lambda = C_1 (R\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 . If $n \gtrsim \delta s^2 \log d$ for some hidden constant only depending on (κ_0, M) , then with probability at least $1 - 4d^{2-\delta}$, the estimation error $\hat{\boldsymbol{\Upsilon}} := \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ satisfies

$$\|\hat{\boldsymbol{\Upsilon}}\|_2 \leq C_2 \mathcal{L} \sqrt{\frac{\delta s \log d}{n}} \quad \text{and} \quad \|\hat{\boldsymbol{\Upsilon}}\|_1 \leq C_3 \mathcal{L} s \sqrt{\frac{\delta \log d}{n}}$$

where $\mathcal{L} := \frac{R\sqrt{M} + \Delta^2}{\kappa_0}$.

Theorem 6 generalizes [36, Thm. 2(b)] to the uniform quantization setting. Clearly, the obtained rate remains near minimax optimal if R is of minor scaling (e.g., bounded or logarithmic factors). Nevertheless, such near optimality in Theorem 6 comes at the cost of more restricted conditions and stronger scaling, as remarked in the following.

Remark 2. (Comparing Theorems 5-6). *Compared with $n \gtrsim s \log d$ in Theorem 5, the first downside of Theorem 6 is the sub-optimal sample complexity $n \gtrsim s^2 \log d$, and note that $n \gtrsim s^2 \log d$ is also required in [36, Thm. 2(b)]. But indeed, it can be improved to $n \gtrsim s \log d$ by explicitly adding the constraint $\|\boldsymbol{\theta}\|_1 \leq R$ to the recovery program, as will be noted as an interesting side finding in Remark 6. Secondly, following [36] we impose an ℓ_1 -norm constraint $\|\boldsymbol{\theta}^*\|_1 \leq R$ that is stronger than $\|\boldsymbol{\theta}^*\|_2 \lesssim \frac{M^{1/(2l)}}{\sigma}$ used in the proof of Theorem 5. In fact, when replacing the ℓ_1 constraint in Theorem 6 with an ℓ_2 -norm bound $\|\boldsymbol{\theta}^*\|_2 \leq R$, then our proof technique leads*

to an error rate $\|\hat{\mathbf{Y}}\|_2 = O(\sqrt{\frac{s^2 \log d}{n}})$ that exhibits worse dependence on s .

Remark 3. (Modification of ℓ_2 -loss). Recall that we generalize the regular ℓ_2 -loss $\frac{1}{2n} \sum_{k=1}^n (y_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2$ to $\frac{1}{2} \boldsymbol{\theta}^\top \mathbf{Q} \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta}$ as loss function in (10). Note that the choice of (\mathbf{Q}, \mathbf{b}) in Theorem 5 is tantamount to using the loss function $\frac{1}{2n} \sum_{k=1}^n (\hat{y}_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2$ that replaces y_k with the quantized response \hat{y}_k ; this idea is analogous to the generalized Lasso investigated for single index model [78] and dithered quantized model [89], and will be used again in quantized matrix completion, see (12) below. However, our generalized ℓ_2 -loss provides more flexibility to deal with heavy-tailedness or quantization of \mathbf{x}_k , e.g., (\mathbf{Q}, \mathbf{b}) in Theorem 6 amounts to adopting $\frac{1}{2n} \sum_{k=1}^n (\hat{y}_k - \tilde{\mathbf{x}}_k^\top \boldsymbol{\theta})^2$ as loss function, and under quantized covariate more delicate modifications are required in Theorems 9-12, which is beyond the range of prior works on generalized Lasso.

C. Quantized Matrix Completion

Completing a low-rank matrix from only a partial observation of its entries is known as the matrix completion problem, which has found many applications including recommendation systems, image inpainting, quantum state tomography [2], [18], [27], [42], [74], to name just a few. Mathematically, let $\boldsymbol{\Theta}^* \in \mathbb{R}^{d \times d}$ be the underlying matrix satisfying $\text{rank}(\boldsymbol{\Theta}^*) \leq r$, the matrix completion problem can be formulated as

$$y_k = \langle \mathbf{X}_k, \boldsymbol{\Theta}^* \rangle + \epsilon_k, \quad k = 1, 2, \dots, n, \quad (11)$$

where \mathbf{X}_k s are distributed on $\mathcal{X} := \{\mathbf{e}_i \mathbf{e}_j^\top : i, j \in [d]\}$ (\mathbf{e}_i is the i -th column of \mathbf{I}_d), ϵ_k is observation noise. Note that for $\mathbf{X}_k = \mathbf{e}_{i(k)} \mathbf{e}_{j(k)}^\top$ one has $\langle \mathbf{X}_k, \boldsymbol{\Theta}^* \rangle = \theta_{i(k), j(k)}^*$, so each observation is a noisy entry. Our main interest is in quantized matrix completion (QMC), where our goal is to design a quantizer for the observation y_k that allows for accurate estimation of $\boldsymbol{\Theta}^*$ from the quantized observations.

Unlike in compressed sensing, additional condition (besides the low-rankness) on $\boldsymbol{\Theta}^*$ is needed to ensure the well-posedness of the matrix completion problem. More specifically, certain incoherence conditions are required if we pursue exact recovery (e.g., [15], [82]), whereas a faithful estimation can be achieved as long as the underlying matrix is not overly spiky and sufficiently diffuse (e.g., [51], [71]). The latter condition is also known as ‘‘low spikiness’’ and is formulated by $\frac{d \|\boldsymbol{\Theta}^*\|_\infty}{\|\boldsymbol{\Theta}^*\|_F} \leq \alpha$ [36], [71], which has been noted to be necessary for the well-posedness of matrix completion problem [27], [71]. In subsequent works, the low-spikiness condition is often formulated as the simpler max-norm constraint $\|\boldsymbol{\Theta}^*\|_\infty \leq \alpha$ [18], [26], [37], [51], [53].

In this work, we consider the uniform sampling scheme $\mathbf{X}_k \sim \mathcal{U}(\mathcal{X})$, but with a little bit more work it generalizes to a more general sampling scheme [51]. We apply the proposed quantization scheme to possibly heavy-tailed y_k — we truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with some threshold ζ_y , and then quantize \tilde{y}_k to $\hat{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$. Because we do not pursue exact recovery (which is impossible under quantization), we do not assume any incoherence condition like [82]. Instead, we only hope to

accurately estimate $\boldsymbol{\Theta}^*$, and following [18], [26], [37], [51], [53] we impose a max-norm constraint

$$\|\boldsymbol{\Theta}^*\|_\infty \leq \alpha.$$

Overall, we estimate $\boldsymbol{\Theta}^*$ from $(\mathbf{X}_k, \hat{y}_k)$ by the regularized M-estimator [70], [72]

$$\hat{\boldsymbol{\Theta}} = \arg \min_{\|\boldsymbol{\Theta}\|_\infty \leq \alpha} \frac{1}{2n} \sum_{k=1}^n (\hat{y}_k - \langle \mathbf{X}_k, \boldsymbol{\Theta} \rangle)^2 + \lambda \|\boldsymbol{\Theta}\|_{nu} \quad (12)$$

that combines an ℓ_2 -loss and nuclear norm regularizer.

In the literature, there has been a line of works on 1-bit or multi-bit matrix completion related to our results to be presented [3], [14], [16], [52], [59]. While the referenced works commonly adopted a likelihood approach, our method is an essential departure and embraces some advantage, see a precise comparison in Remark 4. Considering such novelty, we include the result for sub-exponential ϵ_k in Theorem 7, for which the truncation of y_k becomes unnecessary and we simply set $\zeta_y = \infty$.

Theorem 7. (QMC under Sub-Exponential Noise). Given some $\Delta > 0, \delta > 0$, in (11) we suppose that \mathbf{X}_k s are i.i.d. uniformly distributed over $\mathcal{X} = \{\mathbf{e}_i \mathbf{e}_j^\top : i, j \in [d]\}$, $\boldsymbol{\Theta}^* \in \mathbb{R}^{d \times d}$ satisfies $\text{rank}(\boldsymbol{\Theta}^*) \leq r$ and $\|\boldsymbol{\Theta}^*\|_\infty \leq \alpha$, the noise ϵ_k s are i.i.d. zero-mean sub-exponential satisfying $\|\epsilon_k\|_{\psi_1} \leq \sigma$, and are independent of \mathbf{X}_k s. In the quantization, we do not truncate y_k but directly quantize it to $\hat{y}_k = \mathcal{Q}_\Delta(y_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$. We choose $\lambda = C_1(\sigma + \Delta) \sqrt{\frac{\delta \log d}{nd}}$ with sufficiently large C_1 , and define $\hat{\boldsymbol{\Theta}}$ as (12). If $\delta d \log^3 d \lesssim n \lesssim \delta r^2 d^2 \log d$, then with probability at least $1 - 4d^{-\delta}$, the estimation error $\hat{\mathbf{Y}} := \hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*$ satisfies

$$\frac{\|\hat{\mathbf{Y}}\|_F}{d} \leq C_2 \mathcal{L} \sqrt{\frac{\delta r d \log d}{n}} \quad \text{and} \quad \frac{\|\hat{\mathbf{Y}}\|_{nu}}{d} \leq C_3 \mathcal{L} r \sqrt{\frac{\delta d \log d}{n}}$$

where $\mathcal{L} := \alpha + \sigma + \Delta$.

By contrast, under heavy-tailed noise only assumed to have bounded variance, we truncate y_k with a suitable threshold before the dithered quantization to achieve an optimal trade-off between bias and variance.

Theorem 8. (QMC under Heavy-tailed Noise). Given some $\Delta > 0, \delta > 0$, we consider (11) in the setting of Theorem 7 but with the assumption $\|\epsilon_k\|_{\psi_1} \leq \sigma$ replaced by $\mathbb{E}|\epsilon_k|^2 \leq M$. In the quantization, we truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with $\zeta_y \asymp (\sqrt{M} + \alpha) \sqrt{\frac{n}{\delta d \log d}}$, and then quantize \tilde{y}_k to $\hat{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$. We choose $\lambda = C_1(\alpha + \sqrt{M} + \Delta) \sqrt{\frac{\delta \log d}{nd}}$ with sufficiently large C_1 , and define $\hat{\boldsymbol{\Theta}}$ as (12). If $\delta d \log d \lesssim n \lesssim \delta r^2 d^2 \log d$, then with probability at least $1 - 6d^{-\delta}$, the estimation error $\hat{\mathbf{Y}} := \hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*$ satisfies

$$\frac{\|\hat{\mathbf{Y}}\|_F}{d} \leq C_2 \mathcal{L} \sqrt{\frac{\delta r d \log d}{n}} \quad \text{and} \quad \frac{\|\hat{\mathbf{Y}}\|_{nu}}{d} \leq C_3 \mathcal{L} r \sqrt{\frac{\delta d \log d}{n}}$$

where $\mathcal{L} := \alpha + \sqrt{M} + \Delta$.

Compared to the information-theoretic lower bounds in [55], [71], the error rates obtained in Theorems 7-8 are minimax

optimal up to logarithmic factors. Specifically, Theorem 8 derives a near optimal guarantee for QMC with heavy-tailed observations, as the key standpoint of this paper. Note that, the 1-bit quantization counterpart of these two Theorems was derived in our previous work [22]; in sharp contrast to Theorem 8, for 1-bit QMC under heavy-tailed noise, the error rate under $\frac{\|\hat{\mathbf{x}}\|_F}{d}$ in [22, Thm. 13] reads as $O\left(\left(\frac{r^2 d \log d}{n}\right)^{1/4}\right)$ and is essentially slower; using the 1-bit observations therein, this slow error rate is indeed nearly tight due to the lower bound in [22, Thm. 14].

To close this section, we give a remark to illustrate the novelty and advantage of our QMC method by a careful comparison with prior works.

Remark 4. *QMC with 1-bit or multi-bit quantized observations has received considerable research interest [3], [14], [16], [26], [52], [59]. Adapted to our notation, these works studied the model $\dot{y}_k = \mathcal{Q}(\langle \mathbf{X}_k, \Theta^* \rangle + \tau_k)$ under general random dither τ_k and quantizer $\mathcal{Q}(\cdot)$, and they commonly adopted regularized (or constrained) maximum likelihood estimation for estimating Θ^* . By contrast, with the random dither and quantizer specialized to $\tau_k \sim \mathcal{U}([- \frac{\Delta}{2}, \frac{\Delta}{2}])$ and $\mathcal{Q}_\Delta(\cdot)$, our model is formulated as $\dot{y}_k = \mathcal{Q}_\Delta(\mathcal{T}_{\zeta_y}(\langle \mathbf{X}_k, \Theta^* \rangle + \epsilon_k) + \tau_k)$. Thus, while suffering from less generality in (τ_k, \mathcal{Q}) , our method embraces the advantage of robustness to pre-quantization noise ϵ_k , whose distribution is unknown and can even be heavy-tailed. Note that such unknown ϵ_k evidently forbids the likelihood approach.*

IV. COVARIATE QUANTIZATION AND UNIFORM SIGNAL RECOVERY IN QUANTIZED COMPRESSED SENSING

By now we have presented near optimal results in the contexts of QCME, QCS, and QMC under heavy-tailed data that further undergo the proposed quantization scheme, which we position as the primary contribution of this work. In this section, we further provide two additional developments to enhance our results on heavy-tailed QCS.

A. Covariate Quantization

In the area of QCS, almost all prior works merely focused on the quantization of response y_k , see the recent survey [29]; here, we consider a setting of “complete quantization” — meaning that the covariate \mathbf{x}_k is also quantized. To motivate our study of “complete quantization”, we interpret compressed sensing as sparse linear regression. Indeed, to reduce the power consumption and computational cost, it is sometimes preferable to work with low-precision data in a machine learning system, e.g., the sample quantization scheme developed in [96] led to experimental success in training linear model. Also, it was shown that direct gradient quantization may not be efficient in certain distributed learning systems where the terminal nodes are connected to the server only through a very weak communication fabric and the number of parameters is extremely huge; rather, quantizing and transmitting some important samples could provably reduce communication cost [43]. In fact, the process of data collection may already appeal to quantization due to certain limits of the data acquisition

device (e.g., a low-resolution analog-to-digital module used in distributed signal processing [25]). Our main goal is to understand how quantization of (\mathbf{x}_k, y_k) s affects the subsequent recovery/learning process, particularly showing that the simple dithered uniform quantization scheme still allows for an accurate estimator that may even provide near minimax error rate. To our best knowledge, the only prior rigorous estimation guarantees for QCS with covariate quantization are [22, Thms. 7-8]; these two results require a restricted and unnatural assumption, which we will also relax later.

1) *Multi-bit QCS with Quantized Covariate:* Since we will also consider the 1-bit quantization, we more precisely refer to the QCS under uniform quantizer as multi-bit QCS. We will generalize Theorems 5-6 to covariate quantization in the next two theorems.

Let $(\hat{\mathbf{x}}_k, \dot{y}_k)$ be the quantized covariate-response pair, we first quickly sketch the idea of our approach. Specifically, we stick to the framework of M-estimator in (10), which appeals to accurate surrogates for $\Sigma^* = \mathbb{E}(x_k x_k^\top)$ and $\Sigma_{yx} = \mathbb{E}(y_k x_k)$ based on $(\hat{\mathbf{x}}_k, \dot{y}_k)$, where $\hat{\mathbf{x}}_k$ represents the quantized covariate. Fortunately, the surrogates can be constructed analogously to our QCME estimator when triangular dither is used for quantizing \mathbf{x}_k . Let us first state our quantization scheme as follows:

- **Response Quantization.** This is the same as Theorems 5-6. We truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with threshold ζ_y , and then quantize \tilde{y}_k to $\dot{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \phi_k)$ with uniform dither $\phi_k \sim \mathcal{U}([- \frac{\Delta}{2}, \frac{\Delta}{2}])$ and quantization level $\Delta \geq 0$.
- **Covariate Quantization.** This is the same as Theorem 2. We truncate \mathbf{x}_k to $\tilde{\mathbf{x}}_k = \mathcal{T}_{\zeta_x}(\mathbf{x}_k)$ with threshold ζ_x , and then quantize $\tilde{\mathbf{x}}_k$ to $\hat{\mathbf{x}}_k = \mathcal{Q}_{\bar{\Delta}}(\tilde{\mathbf{x}}_k + \tau_k)$ with triangular dither $\tau_k \sim \mathcal{U}([- \frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d) + \mathcal{U}([- \frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d)$ and quantization level $\bar{\Delta} \geq 0$.
- **Notation.** We write the quantization noise as $\varphi_k = \dot{y}_k - \tilde{y}_k$ and $\xi_k = \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k$, the quantization error as $\vartheta_k = \dot{y}_k - (\tilde{y}_k + \phi_k)$ and $\mathbf{w}_k = \hat{\mathbf{x}}_k - (\tilde{\mathbf{x}}_k + \tau_k)$.

We will adopt the above notation in subsequent developments. Based on the quantized covariate-response pairs $(\hat{\mathbf{x}}_k, \dot{y}_k)$ s, we specify our estimator by setting (\mathbf{Q}, \mathbf{b}) in (10) as

$$\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d \quad \text{and} \quad \mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \hat{\mathbf{x}}_k. \quad (13)$$

Note that the choice of \mathbf{Q} is due to the estimator in Theorem 2, while \mathbf{b} is inspired by the calculation

$$\begin{aligned} \mathbb{E}(\dot{y}_k \hat{\mathbf{x}}_k) &= \mathbb{E}((\tilde{y}_k + \varphi_k)(\tilde{\mathbf{x}}_k + \xi_k)) \\ &= \mathbb{E}(\tilde{y}_k \tilde{\mathbf{x}}_k) + \mathbb{E}(\tilde{y}_k \xi_k) + \mathbb{E}(\varphi_k \tilde{\mathbf{x}}_k) + \mathbb{E}(\varphi_k \xi_k) \\ &= \mathbb{E}(\tilde{y}_k \tilde{\mathbf{x}}_k), \end{aligned}$$

where the last equality can be seen by conditioning on $\tilde{\mathbf{x}}_k$ or \tilde{y}_k . However, the issue is that \mathbf{Q} is not positive semi-definite, hence the resulting program is non-convex. To explain this, note that the rank of $\frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top$ does not exceed n , so when $d > n$ at least $d - n$ eigenvalues of \mathbf{Q} are $-\frac{\bar{\Delta}^2}{4}$. Alternatively, the non-convexity can also be seen from the

observation that setting (\mathbf{Q}, \mathbf{b}) as in (13) is tantamount to replacing the regular ℓ_2 -loss $\frac{1}{2n} \sum_{k=1}^n (y_k - \mathbf{x}_k^\top \boldsymbol{\theta})^2$ with

$$\frac{1}{2n} \sum_{k=1}^n (\hat{y}_k - \hat{\mathbf{x}}_k^\top \boldsymbol{\theta})^2 - \frac{\bar{\Delta}}{8} \|\boldsymbol{\theta}\|_2^2.$$

We mention that the lack of positive semi-definiteness of \mathbf{Q} is problematic in both statistics and optimization aspects: 1) Statistically, Lemma 4 used to derive the error rates in Theorems 5-6 requires \mathbf{Q} to be positive semi-definite, and is hence no longer applicable here; 2) From the optimization side, it is, in general, unknown how to globally optimize a non-convex program.

Motivated by a line of previous works on non-convex M-estimator [62]–[64], we add an ℓ_1 -norm constraint to (10) by setting $\mathcal{S} = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_1 \leq R\}$, where R represents the prior estimation on $\|\boldsymbol{\theta}^*\|_1$. Let $\partial\|\boldsymbol{\theta}\|_1$ be a subdifferential of $\|\boldsymbol{\theta}\|_1$ at $\boldsymbol{\theta} = \boldsymbol{\theta}_1$,¹⁰ we consider the local minimizer of the proposed recovery program,¹¹ or more generally put, $\tilde{\boldsymbol{\theta}} \in \mathcal{S}$ that satisfies¹²

$$\langle \mathbf{Q}\tilde{\boldsymbol{\theta}} - \mathbf{b} + \lambda \cdot \partial\|\tilde{\boldsymbol{\theta}}\|_1, \boldsymbol{\theta} - \tilde{\boldsymbol{\theta}} \rangle \geq 0, \quad \forall \boldsymbol{\theta} \in \mathcal{S}. \quad (14)$$

We will prove a fairly strong guarantee stating that all $\tilde{\boldsymbol{\theta}} \in \mathcal{S}$ satisfying (14) (of course including all local minimizers) enjoy near minimax error rate. While this guarantee bears resemblance to the ones in [64], we point out that, [64] only derived concrete results for the sub-Gaussian regime; because of the heavy-tailed data and quantization in our setting, some essentially different ingredients are required for the technical analysis (see Remark 5). As before, our results for sub-Gaussian \mathbf{x}_k and heavy-tailed \mathbf{x}_k are presented separately.

Theorem 9. (Quantized Sub-Gaussian Covariate). *Given $\Delta \geq 0$, $\bar{\Delta} \geq 0$, $\delta > 0$, we consider (8) with the same assumptions on $(\mathbf{x}_k, y_k, \boldsymbol{\theta}^*)$ as Theorem 5, and additionally assume that $\|\boldsymbol{\theta}^*\|_2 \leq R$. The quantization of (\mathbf{x}_k, y_k) is described above, and we set $\zeta_x = \infty$, $\zeta_y \asymp \sqrt{\frac{nM^{1/l}}{\delta \log d}}$. For recovery, we let*

$$\begin{aligned} \mathbf{Q} &= \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d, \\ \mathbf{b} &= \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k, \quad \mathcal{S} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_1 \leq R\sqrt{s}\}, \end{aligned}$$

and set

$$\lambda = C_1 \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 . If $n \gtrsim \delta s \log d$ for some hidden constant only depending on $(\kappa_0, \sigma, \Delta, \bar{\Delta}, M, R)$, with probability at least $1 - 8d^{1-\delta} - C_2 \exp(-C_3 n)$, all $\tilde{\boldsymbol{\theta}} \in \mathcal{S}$ satisfying (14) have estimation error $\tilde{\boldsymbol{\Upsilon}} := \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ bounded by

$$\|\tilde{\boldsymbol{\Upsilon}}\|_2 \leq C \mathcal{L} \sqrt{\frac{\delta s \log d}{n}} \quad \text{and} \quad \|\tilde{\boldsymbol{\Upsilon}}\|_1 \leq C' \mathcal{L} s \sqrt{\frac{\delta \log d}{n}}$$

¹⁰Thus, $\partial\|\tilde{\boldsymbol{\theta}}\|_1$ in (14) below should be understood as ‘‘there exists one element in $\partial\|\boldsymbol{\theta}\|_1$ such that (14) holds.’’

¹¹The existence of local minimizer is guaranteed because of the additional ℓ_1 -constraint.

¹²To distinguish the global minimizer in (10), we denote by $\tilde{\boldsymbol{\theta}}$ the estimator in QCS with quantized covariate.

where $\mathcal{L} := \frac{(\sigma + \bar{\Delta})^2 (\Delta + M^{1/(2l)})}{\kappa_0^{3/2}}$.

Similarly, the next result extends Theorem 6 to a setting involving covariate quantization.

Theorem 10. (Quantized Heavy-Tailed Covariate). *Given $\Delta \geq 0$, $\bar{\Delta} \geq 0$, $\delta > 0$, we consider (8) with the same assumptions on $(\mathbf{x}_k, y_k, \boldsymbol{\theta}^*)$ as Theorem 6. The quantization of (\mathbf{x}_k, y_k) is described above, and we set $\zeta_x, \zeta_y \asymp \left(\frac{nM}{\delta \log d}\right)^{1/4}$. For recovery, we let*

$$\begin{aligned} \mathbf{Q} &= \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d, \\ \mathbf{b} &= \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k, \quad \mathcal{S} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_1 \leq R\}, \end{aligned}$$

and set

$$\lambda = C_1 (R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 . If $n \gtrsim \delta s \log d$ for some hidden constant only depending on (κ_0, M) , then with probability at least $1 - 8d^{1-\delta}$, all $\tilde{\boldsymbol{\theta}} \in \mathcal{S}$ satisfying (14) have estimation error $\tilde{\boldsymbol{\Upsilon}} := \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ bounded by

$$\|\tilde{\boldsymbol{\Upsilon}}\|_2 \leq C_3 \mathcal{L} \sqrt{\frac{\delta s \log d}{n}} \quad \text{and} \quad \|\tilde{\boldsymbol{\Upsilon}}\|_1 \leq C_4 \mathcal{L} s \sqrt{\frac{\delta \log d}{n}}$$

where $\mathcal{L} := \frac{R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2}{\kappa_0}$.

Remark 5. (Comparing Our Analyses with [64]). *The above results are motivated by a line of works on nonconvex M-estimator [62]–[64], and our guarantee for the whole set of stationary points (14) resembles [64] most. While the main strategy for proving Theorem 9 is adjusted from [64], the proof of Theorem 10 does involve an essentially different RSC condition, see our (42). In particular, compared with [64, equation (4)], the leading factor of $\|\tilde{\boldsymbol{\Upsilon}}\|_1^2$ in (42) degrades from $O(\frac{\log d}{n})$ to $O(\sqrt{\frac{\log d}{n}})$. To retain near optimal rate we need to impose a stronger scaling $\|\boldsymbol{\theta}^*\|_1 \leq R$ with proper changes in the proof. Although Theorem 10 is presented for a concrete setting, it sheds light on an extension of [64] to a weaker RSC condition that could accommodate covariate with a heavier tail. Such extension is formally presented as a deterministic framework in Proposition 1.*

Proposition 1. *Suppose that the s -sparse $\boldsymbol{\theta}^* \in \mathbb{R}^d$ satisfies $\|\boldsymbol{\theta}^*\|_1 \leq R$, and the positive definite matrix $\boldsymbol{\Sigma}^* \in \mathbb{R}^{d \times d}$ satisfies $\lambda_{\min}(\boldsymbol{\Sigma}^*) \geq \kappa_0$. If for some $\mathbf{Q} \in \mathbb{R}^{d \times d}$, $\mathbf{b} \in \mathbb{R}^d$ we have*

$$\lambda \geq C_1 \max \{ \|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty, R \cdot \|\mathbf{Q} - \boldsymbol{\Sigma}^*\|_\infty \} \quad (15)$$

holds for sufficiently large C_1 , then all $\tilde{\boldsymbol{\theta}} \in \mathcal{S}$ satisfying (14) with $\mathcal{S} = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_1 \leq R\}$ have estimation error $\tilde{\boldsymbol{\Upsilon}} := \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ bounded by

$$\|\tilde{\boldsymbol{\Upsilon}}\|_2 \leq C_2 \frac{\sqrt{s}\lambda}{\kappa_0} \quad \text{and} \quad \|\tilde{\boldsymbol{\Upsilon}}\|_1 \leq C_3 \frac{s\lambda}{\kappa_0}.$$

By extracting the ingredients that guarantee (14) to be accurate, interestingly, Proposition 1 is now independent of

the model assumption (8). Particularly, we could set $\Sigma^* = \mathbb{E}[\mathbf{x}_k \mathbf{x}_k^\top]$ when we apply Proposition 1 to (8). Compared with the framework [64, Thm. 1], the key strength of Proposition 1 is that it does not explicitly assume the RSC condition on the loss function that is hard to verify without assuming a sub-Gaussian covariate. Instead, the role of the RSC assumption in [64] is now played by $\lambda \gtrsim R \|\mathbf{Q} - \Sigma^*\|_\infty$, which immediately yields a kind of RSC condition by simple argument as (43). Although this RSC condition is often essentially weaker than the conventional one in terms of the leading factor of $\|\tilde{\mathbf{Y}}\|_1^2$ (see Remark 5), along this line one can still derive non-trivial (or even near optimal) error rate. The gain of replacing RSC assumption with $\lambda \gtrsim R \|\mathbf{Q} - \Sigma^*\|_\infty$ is that the latter amounts to constructing element-wise estimator for Σ^* , which is often much easier for heavy-tailed covariate (e.g., due to many existing robust covariance estimators).

We conclude this part with a side interesting observation.

Remark 6. By setting $\bar{\Delta} = 0$, Theorem 10 produces a result (with a convex program) for the setting of Theorem 6. Interestingly, with the additional ℓ_1 -constraint, a notable improvement is that the sub-optimal $n \gtrsim s^2 \log d$ in Theorem 6 is sharpened to the near optimal one in Theorem 10. More concretely, this is because (ii) in (30) can be tightened to (ii) of (43). Going back to the full-data unquantized regime, Theorem 10 with $\Delta = \bar{\Delta} = 0$ recovers [36, Theorem 2(b)] with improved sample complexity requirement.

2) *1-bit QCS with Quantized Covariate:* Our consideration of covariate quantization in QCS seems fairly new to the literature. To the best of our knowledge, the only related results are [22, Thms. 7-8] for QCS with 1-bit quantized covariate and response. The assumption there, however, is quite restrictive. Specifically, it is assumed that $\Sigma^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$ has sparse columns (see [22, Assumption 3]), which is non-standard in both compressed sensing and sparse linear regression. Departing momentarily from our focus on dithered uniform quantization, we consider QCS under dithered 1-bit quantization and will apply Proposition 1 to derive results comparable to [22, Thms. 7-8] without resorting to the sparsity of Σ^* .

We first review the 1-bit quantization scheme developed in [22]:

- **Response Quantization.** We truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with some threshold ζ_y , and then quantize \tilde{y}_k to $\dot{y}_k = \text{sign}(\tilde{y}_k + \phi_k)$ with uniform dither $\phi_k \sim \mathcal{U}([-\gamma_y, \gamma_y])$.
- **Covariate Quantization.** We truncate \mathbf{x}_k to $\tilde{\mathbf{x}}_k = \mathcal{T}_{\zeta_x}(\mathbf{x}_k)$ with some threshold ζ_x , and then quantize $\tilde{\mathbf{x}}_k$ to $\dot{\mathbf{x}}_{k1} = \text{sign}(\tilde{\mathbf{x}}_k + \tau_{k1})$ and $\dot{\mathbf{x}}_{k2} = \text{sign}(\tilde{\mathbf{x}}_k + \tau_{k2})$, where $\tau_{k1}, \tau_{k2} \sim \mathcal{U}([-\gamma_x, \gamma_x]^d)$ are independent uniform dithers. (Note that we collect 2 bits per entry).

The following two results refine [22, Thms. 7-8] by deriving comparable error rates without using sparsity of Σ^* .

Theorem 11. (1-bit Quantized Sub-Gaussian Covariate). *Given $\delta > 0$, we consider (8) where the s -sparse θ^* satisfies $\|\theta^*\|_1 \leq R$, \mathbf{x}_{kS} are i.i.d. zero-mean sub-Gaussian with $\|\mathbf{x}_k\|_{\psi_2} \leq \sigma$, and $\Sigma^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$ satisfies $\lambda_{\min}(\Sigma^*) \geq \kappa_0$ for some $\kappa_0 > 0$, the noise ϵ_{kS} are independent of \mathbf{x}_{kS} and*

i.i.d. sub-Gaussian, while for simplicity we assume $\|y_k\|_{\psi_2} \leq \sigma$. In the quantization of (\mathbf{x}_k, y_k) described above, we set $\zeta_x = \zeta_y = \infty$ and $\gamma_x, \gamma_y \asymp \sigma \sqrt{\log(\frac{n}{2\delta \log d})}$. For recovery we let

$$\mathbf{Q} := \frac{\gamma_x^2}{2n} \sum_{k=1}^n (\dot{\mathbf{x}}_{k1} \dot{\mathbf{x}}_{k2}^\top + \dot{\mathbf{x}}_{k2} \dot{\mathbf{x}}_{k1}^\top),$$

$$\mathbf{b} := \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_{k1}, \quad \mathcal{S} = \{\theta : \|\theta\|_1 \leq R\}$$

and set

$$\lambda = C_1 \sigma^2 R \sqrt{\frac{\delta \log d(\log n)^2}{n}}$$

with sufficiently large C_1 . If $n \gtrsim \delta s \log d(\log n)^2$, then with probability at least $1 - 4d^{2-\delta}$, all $\tilde{\theta} \in \mathcal{S}$ satisfying (14) have estimation error $\tilde{\mathbf{Y}} := \tilde{\theta} - \theta^*$ bounded by

$$\|\tilde{\mathbf{Y}}\|_2 \leq C_2 \frac{\sigma^2}{\kappa_0} \cdot R \sqrt{\frac{\delta s \log d(\log n)^2}{n}},$$

$$\|\tilde{\mathbf{Y}}\|_1 \leq C_3 \frac{\sigma^2}{\kappa_0} \cdot R s \sqrt{\frac{\delta \log d(\log n)^2}{n}}.$$

Theorem 12. (1-bit Quantized Heavy-Tailed Covariate). *Given $\delta > 0$, we consider (8) where the s -sparse θ^* satisfies $\|\theta^*\|_1 \leq R$, \mathbf{x}_{kS} are i.i.d. zero-mean heavy-tailed satisfying the joint fourth-moment constraint $\sup_{v \in \mathbb{S}^{d-1}} \mathbb{E}|v^\top \mathbf{x}_k|^4 \leq M$, and $\Sigma^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$ satisfies $\lambda_{\min}(\Sigma^*) \geq \kappa_0$ for some $\kappa_0 > 0$, the noise ϵ_{kS} are independent of \mathbf{x}_{kS} and i.i.d. heavy-tailed with a bounded fourth-moment, while for simplicity we assume $\mathbb{E}|y_k|^4 \leq M$. In the quantization of (\mathbf{x}_k, y_k) described above, we set $\zeta_x, \zeta_y, \gamma_x, \gamma_y \asymp (\frac{nM^2}{\delta \log d})^{1/8}$ and enforce $\zeta_x < \gamma_x$, $\zeta_y < \gamma_y$. For recovery we let*

$$\mathbf{Q} := \frac{\gamma_x^2}{2n} \sum_{k=1}^n (\dot{\mathbf{x}}_{k1} \dot{\mathbf{x}}_{k2}^\top + \dot{\mathbf{x}}_{k2} \dot{\mathbf{x}}_{k1}^\top),$$

$$\mathbf{b} := \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_{k1}, \quad \mathcal{S} = \{\theta : \|\theta\|_1 \leq R\}$$

and set

$$\lambda = C_1 \sqrt{MR} \left(\frac{\delta \log d}{n}\right)^{1/4}$$

with sufficiently large C_1 . If $n \gtrsim \delta s^2 \log d$, then with probability at least $1 - 4d^{2-\delta}$, all $\tilde{\theta} \in \mathcal{S}$ satisfying (14) have estimation error $\tilde{\mathbf{Y}} := \tilde{\theta} - \theta^*$ bounded by

$$\|\tilde{\mathbf{Y}}\|_2 \leq C_2 \frac{\sqrt{M}}{\kappa_0} \cdot R \left(\frac{\delta s^2 \log d}{n}\right)^{1/4},$$

$$\|\tilde{\mathbf{Y}}\|_1 \leq C_3 \frac{\sqrt{M}}{\kappa_0} \cdot R s \left(\frac{\delta \log d}{n}\right)^{1/4}.$$

B. Uniform Recovery Guarantee

Uniformity is a highly desired property for a compressed sensing guarantee because it allows one to use a fixed (possibly randomly drawn) measurement ensemble for all sparse signals. Unfortunately, as with many other results for nonlinear compressed sensing in the literature, our earlier recovery guarantees are non-uniform and only ensure the accurate

recovery of a sparse signal fixed before drawing the random measurement ensemble.

We provide another additional development to QCS in this part. Specifically, we establish a uniform recovery guarantee which, despite the heavy-tailed noise and nonlinear quantization scheme, notably retains a near minimax error rate. This is done by upgrading Theorem 5 to be uniform over all sparse θ^* by more in-depth technical tools and a careful covering argument. Part of the techniques is inspired by prior works [40], [94], but certain technical innovations are required:

1) Like the recent work [40], one crucial technical tool in our proof is a powerful concentration inequality for product process due to Mendelson [67], as adapted in the present Lemma 9. However, [40] only studied sub-Gaussian distribution, and the results produced by their unified approach typically exhibit a decaying rate of $O(n^{-1/4})$ [40, Sect. 4]. By contrast, our problem involves heavy-tailed noise only having bounded $(2 + \nu)$ -th moment ($\nu > 0$), and we aim to establish a near minimax uniform error bound — cautiousness and new treatment are thus needed in the application of Lemma 9. More specifically, in the proof we need to bound

$$I_1 = \sup_{\theta \in \Sigma_{s,R_0}} \sup_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n (\tilde{y}_k \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\tilde{y}_k \mathbf{x}_k^\top \mathbf{v}]),$$

where $\mathcal{V} = \{\mathbf{v} : \|\mathbf{v}\|_2 = 1, \|\mathbf{v}\|_1 \leq 2\sqrt{s}\}$, and $\Sigma_{s,R_0} = \Sigma_s \cap \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq R_0\}$ is the signal space of interest, and recall that $\tilde{y}_k = \mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \theta + \epsilon_k)$ with sub-Gaussian \mathbf{x}_k . It is natural to invoke Lemma 9 to bound I_1 straightforwardly, but the issue is on lack of good bound for $\|\tilde{y}_k\|_{\psi_2}$ due to the heavy-tailedness of ϵ_k ; indeed, one only has the trivial estimate as $\|\tilde{y}_k\|_{\psi_2} = O(\zeta_y)$, which is much worse than an $O(1)$ bound since $\zeta_y \asymp \sqrt{\frac{n}{\delta \log d}}$, and using Lemma 9 with this estimate leads to a loose bound for I_1 and finally a sub-optimal error rate. To address the issue, our main idea is to introduce the truncated heavy-tailed noise $\mathcal{T}_{\zeta_y}(\epsilon_k)$ and define $\tilde{z}_k = \tilde{y}_k - \mathcal{T}_{\zeta_y}(\epsilon_k)$, which enables us to decompose I_1 as

$$\begin{aligned} I_1 &\leq \underbrace{\sup_{\theta \in \Sigma_{s,R_0}} \sup_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n (\tilde{z}_k \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\tilde{z}_k \mathbf{x}_k^\top \mathbf{v}])}_{:=I_{11}} \\ &\quad + \underbrace{\sup_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n (\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k^\top \mathbf{v}])}_{:=I_{12}}. \end{aligned}$$

Now, the benefits of working with I_{11}, I_{12} are that: i) We can directly invoke Lemma 9 to bound I_{11} since we have a good sub-Gaussian norm estimate $\|\tilde{z}_k\|_{\psi_2} \leq \|\mathbf{x}_k^\top \theta\|_{\psi_2} \lesssim \|\mathbf{x}_k\|_{\psi_2} R_0$, see Step 2.1.1 in the proof; ii) I_{12} becomes the supremum of a process that is independent of θ and only indexed by \mathbf{v} , hence Bernstein's inequality suffices for bounding I_{12} (Step 2.1.2 in the proof), analogously to the proof of the non-uniform guarantee (Theorem 5).

2) Like [94, Prop. 6.1], we invoke a covering argument with similar techniques to bound

$$I_0 = \sup_{\theta \in \Sigma_{s,R_0}} \sup_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n \xi_k \mathbf{x}_k^\top \mathbf{v},$$

where $\xi_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k) - \tilde{y}_k$ is the quantization noise. Nevertheless, our Lasso estimator is different from their projected back projection estimator, and it turns out that we need to directly handle “ $\sup_{\mathbf{v} \in \mathcal{V}}$ ” by Lemma 10, unlike [94, Prop. 6.2] that again used a covering argument for this purpose. See more discussions in Step 2.4 of the proof.

We are in a position to present our uniform recovery guarantee. We follow most assumptions in Theorem 5 but specify the signal space as $\theta^* \in \Sigma_{s,R_0}$ and impose the $(2l)$ -th moment constraint on ϵ_k . Following prior works on QCS (e.g., [40], [89]), we consider constrained Lasso that utilizes an ℓ_1 -constraint $\|\theta\|_1 \leq \|\theta^*\|_1$ (rather than (10)) to pursue uniform recovery.

Theorem 13. (Uniform Version of Theorem 5). *Given some $\delta > 0, \Delta > 0$, in (8) we suppose that \mathbf{x}_k s are i.i.d., zero-mean sub-Gaussian with $\|\mathbf{x}_k\|_{\psi_2} \leq \sigma, \kappa_0 \leq \lambda_{\min}(\Sigma^*) \leq \lambda_{\max}(\Sigma^*) \leq \kappa_1$ for some $\kappa_1 \geq \kappa_0 > 0$ where $\Sigma^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$, $\theta^* \in \Sigma_{s,R_0} := \Sigma_s \cap \{\theta : \|\theta\|_2 \leq R_0\}$ for some absolute constant R_0 , ϵ_k s are i.i.d. noise that are independent of \mathbf{x}_k s and satisfy $\mathbb{E}|\epsilon_k|^{2l} \leq M$ for some fixed $l > 1$. In quantization, we truncate y_k to $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k)$ with threshold $\zeta_y \asymp \left(\frac{n(M^{1/l} + \sigma^2)}{\delta \log d}\right)^{1/2}$, then quantize \tilde{y}_k to $\hat{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$ with uniform dither $\tau_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$. For recovery, we define the estimator $\hat{\theta}$ as the solution to constrained Lasso*

$$\hat{\theta} = \arg \min_{\|\theta\|_1 \leq \|\theta^*\|_1} \frac{1}{2n} \sum_{k=1}^n (\hat{y}_k - \mathbf{x}_k^\top \theta)^2$$

If $n \gtrsim \delta s \log \mathcal{W}$ for $\mathcal{W} = \frac{\kappa_1 d^2 n^3}{\Delta^2 s^5 \delta^3}$ and some hidden constant depending on (κ_0, σ) , then with probability at least $1 - Cd^{1-\delta}$ on a single random draw of $(\mathbf{x}_k, \epsilon_k, \tau_k)_{k=1}^n$, it holds uniformly for all $\theta^* \in \Sigma_{s,R_0}$ that the estimation error $\hat{\mathbf{Y}} := \hat{\theta} - \theta^*$ satisfy

$$\begin{aligned} \|\hat{\mathbf{Y}}\|_2 &\leq \frac{C_3 \sigma (\sigma + M^{\frac{1}{2l}})}{\kappa_0} \sqrt{\frac{\delta s \log d}{n}} + \frac{C_3 \sigma \Delta}{\kappa_0} \sqrt{\frac{\delta s \log \mathcal{W}}{n}}, \\ \|\hat{\mathbf{Y}}\|_1 &\leq \frac{C_4 \sigma (\sigma + M^{\frac{1}{2l}})}{\kappa_0} s \sqrt{\frac{\delta \log d}{n}} + \frac{C_4 \sigma \Delta}{\kappa_0} s \sqrt{\frac{\delta \log \mathcal{W}}{n}}. \end{aligned}$$

Notably, our uniform guarantee is still minimax optimal up to some additional logarithmic factors (i.e., $\sqrt{\log \mathcal{W}}$) arising from the covering argument (Step 2.4 of the proof), whose main aim is to show that one uniform dither $\tau = [\tau_k]$ suffices for all signals. Thus naturally, $\sqrt{\log \mathcal{W}}$ is associated with a leading factor of the quantization level Δ , meaning that the logarithmic gap between uniform recovery and non-uniform recovery closes when $\Delta \rightarrow 0$. In particular, Theorem 13 implies a uniform error rate matching the non-uniform one in Theorem 5 (up to some multiplicative factors) when Δ is small enough or in an unquantized case.

To the best of our knowledge, the only existing uniform guarantee for heavy-tailed QCS is [33, Thm. 1.11], but the following distinctions make it impossible to closely compare their result with our Theorem 13: 1) [33, Thm. 1.11] is for dithered 1-bit quantization, but ours is for dithered uniform quantizer; 2) We handle heavy-tailedness by truncation, while [33, Thm. 1.11] does not involve this kind of special treatment; 3) [33, Thm. 1.11] considers a highly intractable program with

hamming distance as objective and $\theta \in \Sigma_s$ as constraint (when specialized to sparse signal), while our Theorem 13 is for the convex program Lasso; 4) Their analysis is based on an in-depth result on random hyperplane tessellations (see also [34], [77]), while our proof follows the more standard strategy (i.e., to upgrade each piece in a non-uniform proof to be uniform) and requires certain technical innovations (e.g., the treatment to deal with the truncation step).¹³ Note that [33, Thm. 1.11] is robust to adversarial bit flips, and we leave it future work to extend our results to a setting where adversarial noise presents.

V. NUMERICAL SIMULATIONS

In this section, we provide two sets of experimental results to support and demonstrate our theoretical developments. The first set of our simulations is devoted to validating our major standpoint that near minimax rates are achievable in quantized heavy-tailed settings. Then, the second set of results is presented to illustrate the crucial role played by the appropriate dither (i.e., triangular dither for covariate, uniform dither for response) before uniform quantization. For the importance of data truncation we refer to in [36, Sect. 5], which includes three estimation problems in this work and contrasts the estimations with or without the data truncation.

A. (Near) Minimax Error Rates

Each data point in our results is set to be the mean value of 50 or 100 independent trials.

1) *Quantized Covariance Matrix Estimation:* We start from covariance matrix estimation, specifically, we verify the element-wise rate $\mathcal{B}_1 := O(\mathcal{L} \sqrt{\frac{\log d}{n}})$ and operator norm rate $\mathcal{B}_2 := O(\mathcal{L} \sqrt{\frac{d \log d}{n}})$ in Theorems 2-3.

For estimator in Theorem 2, we draw $\mathbf{x}_k = (x_{ki})$ such that the first two coordinates are independently drawn from $t(4.5)$, $(x_{ki})_{i=3,4}$ are from $t(6)$ with covariance $\mathbb{E}(x_{k3}x_{k4}) = 1.2$, and the remaining $d - 4$ coordinates are i.i.d. following $t(6)$. We test different choices of (d, Δ) under $n = 80 : 20 : 220$, and the log-log plots are shown in Figure 1(a). Clearly, for each (d, Δ) the experimental points roughly exhibit a straight line that is well aligned with the dashed line representing the $n^{-1/2}$ rate. As predicted by the factor $\mathcal{L} = \sqrt{M} + \Delta^2$, the curves with larger Δ are higher, but note that the error decreasing rates remain unchanged. In addition, the curves of $(d, \Delta) = (100, 1), (120, 1)$ are extremely close, which is consistent with the logarithmic dependence of \mathcal{B}_1 on d .

For the error bound \mathcal{B}_2 , the coordinates of \mathbf{x}_k are independently drawn from a scaled version of $t(4.5)$ such that $\Sigma^* = \text{diag}(2, 2, 1, \dots, 1)$, and we test different settings of (d, Δ) under $n = 200 : 100 : 1000$. As shown in Figure 1(b), the operator norm error decreases with n in the optimal rate $n^{-1/2}$, and using a coarser dithered quantizer (i.e., larger Δ) only slightly lifts the curves. Indeed, the effect seems consistent with \mathcal{L} 's quadratic dependence on Δ . To validate the relative scaling of n and d , in addition to the setting

$(d, \Delta) = (100, 1)$ under $n = 200 : 100 : 1000$, we try $(d, \Delta) = (150, 1)$ under 1.5 times the original sample size $n = 1.5 \times (200 : 100 : 1000)$ (but in Figure 1(b) we still plot the curve according to the sample size of $200 : 100 : 1000$ without the multiplicative factor of 1.5), and surprisingly the obtained curve coincides with the one for $(d, \Delta) = (100, 1)$. Thus, ignoring the logarithmic factor $\log d$, the operator norm error can be characterized by \mathcal{B}_2 fairly well.

Additionally, we want to compare \mathcal{B}_1 and \mathcal{B}_2 regarding the dependence on d more clearly. Specifically, we generate the samples \mathbf{x}_k s as in Figure 1(a) and test the fixed sample size $n = 180$ and varying dimension $d = 80 : 20 : 260$. The max norm estimation errors of $\widehat{\Sigma}$ in Theorem 2 and the operator norm errors (under $d = 80 : 20 : 180$ to ensure $n \geq d$) of the estimator in Theorem 3 are reported in Figure 1(c). It is clear that the max norm error increases with d rather slowly, while the operator norm error increases much more significantly under larger d . This is consistent with the logarithmic dependence of \mathcal{B}_1 on d and the more essential dependence of \mathcal{B}_2 on d .

2) *Quantized Compressed Sensing:* We now switch to QCS with unquantized covariate and aim to verify the ℓ_2 -norm error rate $\mathcal{B}_3 = O(\mathcal{L} \sqrt{\frac{s \log d}{n}})$ obtained in Theorems 5-6. We let the support of the s -sparse $\theta^* \in \mathbb{R}^d$ be $[s]$, and then draw the non-zero entries from a uniform distribution over \mathbb{S}^{s-1} (hence $\|\theta^*\|_2 = 1$). For the setting of Theorem 5 we adopt $\mathbf{x}_k \sim \mathcal{N}(0, \mathbf{I}_d)$ and $\epsilon_k \sim \frac{1}{\sqrt{6}}t(3)$, while $x_{ki} \stackrel{iid}{\sim} \frac{\sqrt{5}}{3}t(4.5)$ and $\epsilon_k \sim \frac{1}{\sqrt{3}}t(4.5)$ for Theorem 6. We simulate different choices of (d, s, Δ) under $n = 100 : 100 : 1000$, and the proposed convex program (10) is solved with the framework of ADMM (we refer to the review [8]). Experimental results are shown as log-log plots in Figure 2. Consistent with the theoretical bound \mathcal{B}_3 , the errors in both cases decrease in a rate of $n^{-1/2}$, whereas the effect of uniform quantization is merely on the multiplicative factor \mathcal{L} . Interestingly, it seems that the gaps between $\Delta = 0, 0.5$ and $\Delta = 0.5, 1$ are in agreement with the explicit form of \mathcal{L} , i.e., $\mathcal{L} \asymp M^{1/(2l)} + \Delta$ for Theorem 5, and $\mathcal{L} \asymp \sqrt{M} + \Delta^2$ for Theorem 6. In addition, note that the curves of $(d, s) = (150, 5), (180, 5)$ are close, whereas increasing $s = 8$ suffers from a significantly larger error. This is consistent with the scaling law of (n, d, s) in \mathcal{B}_3 .

Then, we simulate the complete quantization setting where both covariate and response are quantized (Theorems 9-10). The simulation details are the same as before except that \mathbf{x}_k is also quantized with a quantization level the same as y_k . We provide the best ℓ_1 -norm constraint for recovery, i.e., $\mathcal{S} := \{\theta : \|\theta\|_1 \leq \|\theta^*\|_1\}$. Then, composite gradient descent [63], [64] is invoked to handle the non-convex estimation program. We show the log-log plots in Figure 3. Note that these results have implications similar to Figure 2, in terms of the $n^{-1/2}$ rate, the effect of quantization, and the relative scaling of (n, d, s) .

3) *Quantized Matrix Completion:* Finally, we simulate QMC and demonstrate the error bound $\mathcal{B}_4 = O(\mathcal{L} \sqrt{\frac{rd \log d}{n}})$ for $\|\widehat{\mathbf{Y}}\|_F/d$ in Theorems 7-8. We generate the rank- r $\Theta^* \in \mathbb{R}^{d \times d}$ as follows: we first generate $\Theta_0 \in \mathbb{R}^{d \times r}$ with i.i.d. standard Gaussian entries to obtain the rank- r $\Theta_1 := \Theta_0 \Theta_0^\top$,

¹³It is possible to use such a standard strategy to upgrade Theorem 6 to a uniform result; we suspect that the error rate may exhibit worse dependence on s due to covering argument.

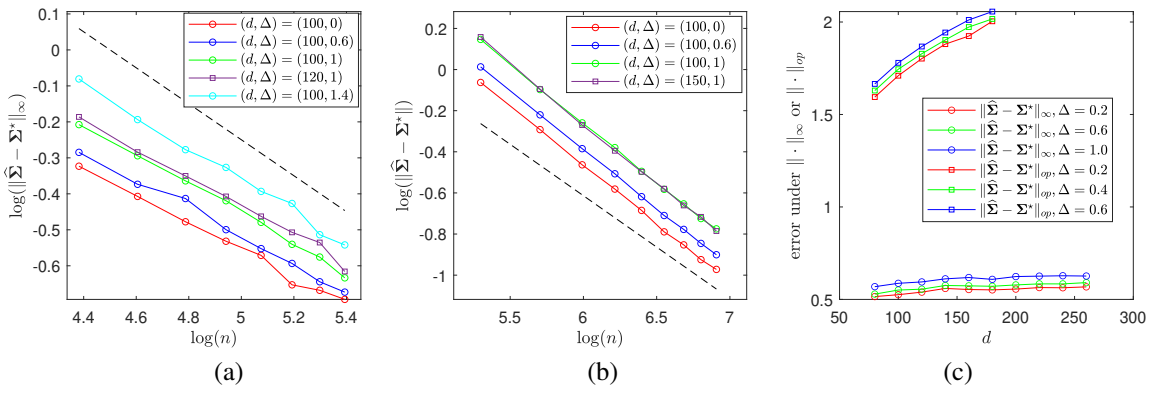


Fig. 1. (a): Element-wise error (Theorem 2); (b): operator norm error (Theorem 3); (c): the dependence on d of both error metrics.

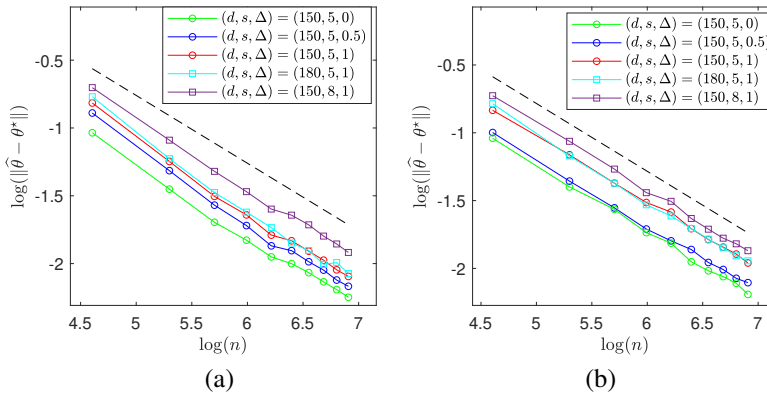


Fig. 2. (a): QCS in Theorem 5; (b): QCS in Theorem 6.

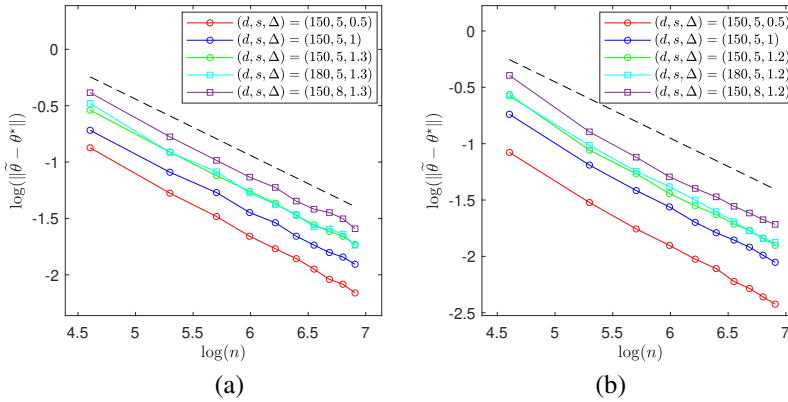


Fig. 3. (a): QCS in Theorem 9; (b): QCS in Theorem 10.

then we rescale it to $\Theta^* := k_1 \Theta_1$ such that $\|\Theta^*\|_F = d$. We use $\epsilon_k \sim \mathcal{N}(0, \frac{1}{4})$ to simulate the sub-exponential noise in Theorem 7, while $\epsilon_k \sim \frac{1}{\sqrt{6}} \mathbf{t}(3)$ for Theorem 8. The convex program (12) is fed with $\alpha = \|\Theta^*\|_\infty$ and optimized by the ADMM algorithm. We test different choices of (d, r, Δ) under $n = 2000 : 1000 : 8000$, with the log-log error plots displayed in Figure 4. Firstly, the experimental curves are well aligned with the dashed line that represents the optimal $n^{-1/2}$ rate. Then, comparing the results for $\Delta = 0, 0.5, 1$, we conclude that quantization only affects the multiplicative factor \mathcal{L} in the estimation error. It should also be noted that, increasing either d or r leads to a significantly larger error, which is

consistent with the \mathcal{B}_4 's essential dependence on d and r .

B. Importance of Appropriate Dithering

To demonstrate the crucial role played by the suitable dither, we provide the second set of simulations. In order to observe more significant phenomena and then conclude evidently, we may test a huge sample size but a rather simple estimation problem under coarse quantization (i.e., large Δ).

Specifically, for covariance matrix estimation we set $d = 1$ and i.i.d. draw X_1, \dots, X_n from $\mathcal{N}(0, 1)$. Thus, the problem boils down to estimating $\mathbb{E}|X_k|^2$, for which the estimators in Theorems 2-3 coincide. Since X_k is sub-Gaussian, we

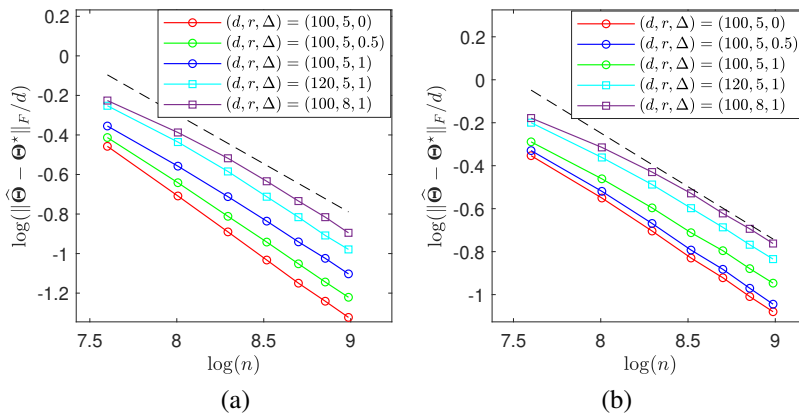


Fig. 4. (a): QMC in Theorem 7; (b): QMC in Theorem 8.

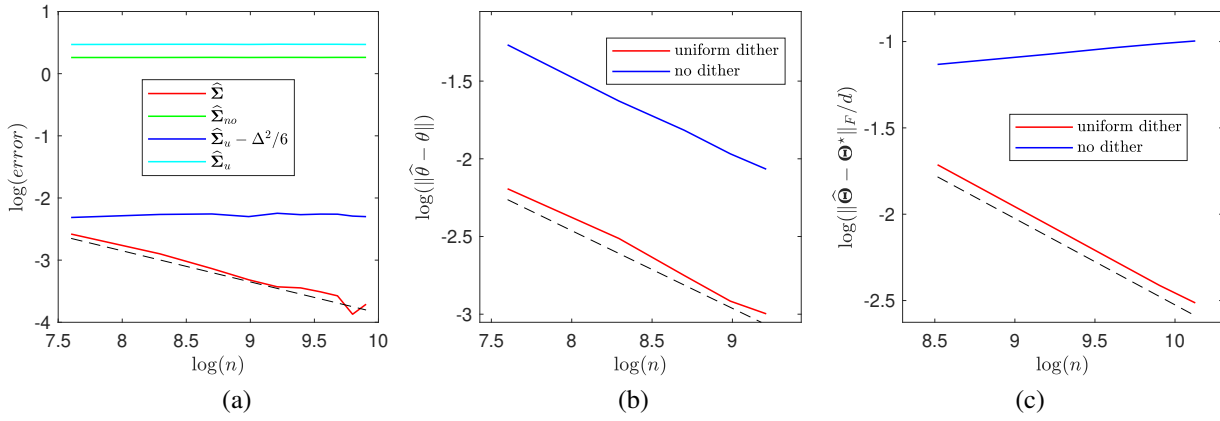


Fig. 5. (a): covariance matrix estimation; (b): QCS in Theorem 5; (c): QMC in Theorem 7.

do not perform data truncation before dithered quantization. Besides our estimator $\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n \hat{X}_k^2 - \frac{\Delta^2}{4}$ where $\hat{X}_k = \mathcal{Q}_\Delta(X_k + \tau_k)$ and τ_k is triangular dither, we invite the following competitors:

- $\hat{\Sigma}_{no} = \frac{1}{n} \sum_{k=1}^n (\hat{X}'_k)^2$, where $\hat{X}'_k = \mathcal{Q}_\Delta(X_k)$ is the direct quantization without dithering;
- $\hat{\Sigma}_u - \frac{\Delta^2}{6}$ and $\hat{\Sigma}_u$, where $\hat{\Sigma}_u = \frac{1}{n} \sum_{k=1}^n (\hat{X}''_k)^2$, and $\hat{X}''_k = \mathcal{Q}_\Delta(X_k + \tau''_k)$ is quantized under uniform dither $\tau''_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$.

To illustrate the choice of $\hat{\Sigma}_u - \frac{\Delta^2}{6}$ and $\hat{\Sigma}_u$, we write

$$\hat{X}''_k = X_k + \tau''_k + w_k = X_k + \xi_k$$

with quantization error $w_k \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ (due to Theorem 1(a)) and quantization noise $\xi_k = \tau''_k + w_k$, then (5) gives

$$\mathbb{E}(\hat{X}''_k)^2 = \mathbb{E}|X_k|^2 + \mathbb{E}|\xi_k|^2,$$

while $\mathbb{E}|\xi_k|^2$ remains unknown. Thus, we consider $\hat{\Sigma}_u - \frac{\Delta^2}{6}$ because of an *unjustified* guess

$$\mathbb{E}|\xi_k|^2 \approx \mathbb{E}|\tau''_k|^2 + \mathbb{E}|w_k|^2 = \frac{\Delta^2}{6},$$

while $\hat{\Sigma}_u$ simply gives up the correction of $\mathbb{E}|\xi_k|^2$. We test $\Delta = 3$ under $n = (2 : 2 : 20) \cdot 10^3$. From the results shown in Figure 5(a), the proposed estimator based on quantized data under triangular dither embraces the lowest estimation errors

and the optimal rate of $n^{-1/2}$, whereas other competitors are not consistent, i.e., they all reach some error floors under a large sample size.

For the two remaining signal recovery problems, we simply focus on the quantization of the response y_k . In particular, we simulate QCS in the setting of Theorem 5, with $(d, s, \Delta) = (50, 3, 2)$ under $n := (2 : 2 : 20) \cdot 10^3$. Other experimental details are as previously stated. We compare our estimator $\hat{\theta}$ with its counterpart $\hat{\theta}'$ defined by (10) with the same \mathcal{Q}, \mathcal{S} but $\mathbf{b}' = \frac{1}{n} \sum_{k=1}^n \hat{y}'_k \mathbf{x}_k$, where $\hat{y}'_k = \mathcal{Q}_\Delta(\tilde{y}_k)$ is a direct uniform quantization with no dither. Evidently, the simulation results in Figure 5(b) confirm that the application of a uniform dither significantly lessens the recovery errors. Without dithering, although our results under the Gaussian covariate still exhibit $n^{-1/2}$ decreasing rate, the identifiability issue unavoidably arises under the Bernoulli covariate. In that case, the simulation without dithering will evidently deviate from the $n^{-1/2}$ rate, see [87, Figure 1] for instance.

In analogy, we simulate QMC (Theorem 7) with data generated as previous experiments, and specifically we try $(d, r, \Delta) = (30, 5, 1.5)$ under $n = (5 : 5 : 25) \cdot 10^3$. While our estimator $\hat{\Theta}$ is defined in (12) involving \hat{y}_k from a dithered quantizer, we simulate the performance of its counterpart without dithering, i.e., $\hat{\Theta}'$ defined in (12) with \hat{y}_k substituted by $\hat{y}'_k = \mathcal{Q}_\Delta(y_k)$. From the experimental results displayed in Figure 5(c), one shall clearly see that $\hat{\Theta}$ performs much

better in terms of the decreasing rate of $n^{-1/2}$ and the estimation error; while the curve without dithering even does not decrease.

VI. CONCLUDING REMARKS

In digital signal processing and many distributed machine learning systems, data quantization is an indispensable process. On the other hand, many modern datasets exhibit heavy-tailedness, and the past decade has witnessed an increasing interest in statistical estimation methods robust to heavy-tailed data. In this work, we try to bridge these two developments by studying the quantization of heavy-tailed data. We propose to truncate the heavy-tailed data prior to a uniform quantizer with random dither well suited to the problem at hand. Applying our quantization scheme to covariance matrix estimation, compressed sensing, and matrix completion, we have proposed (near) optimal estimators based on quantized data, and they are computationally feasible. These results suggest a unified conclusion that the dithered quantization does not affect the key scaling law in the error rate but only slightly worsens the multiplicative factor, which was complemented by numerical simulations. Further, in two respects, we presented additional developments for quantized compressed sensing. Firstly, we study a novel setting that involves covariate quantization. Because our quantized covariance matrix estimator is not positive semi-definite, the proposed recovery program is non-convex, but we proved that all local minimizers enjoy near minimax rates. At a higher level, this development extends a line of works on non-convex M-estimator [62]–[64] to accommodate heavy-tailed covariate, see the deterministic framework Proposition 1. As an application, we derive results for (dithered) 1-bit compressed sensing as byproducts. Secondly, we established a near minimax uniform recovery guarantee for QCS under heavy-tailed noise, which states that all sparse signals within an ℓ_2 -ball can be uniformly recovered up to near optimal ℓ_2 -norm error, using a single realization of the measurement ensemble. We believe the developments presented in this work will prove useful in many other estimation problems, for instance, the triangular dither and the quantization scheme apply to multi-task learning, as shown by subsequent works [23], [60].

REFERENCES

- [1] Albert Ai, Alex Lapanowski, Yaniv Plan, and Roman Vershynin. One-bit compressed sensing with non-gaussian measurements. *Linear Algebra and its Applications*, 441:222–239, 2014.
- [2] James Bennett, Stan Lanning, et al. The netflix prize. In *Proceedings of KDD cup and workshop*, volume 2007, page 35. New York, 2007.
- [3] Sonia A Bhaskar. Probabilistic low-rank matrix completion from quantized measurements. *The Journal of Machine Learning Research*, 17(1):2131–2164, 2016.
- [4] Peter J Bickel and Elizaveta Levina. Covariance regularization by thresholding. *The Annals of statistics*, 36(6):2577–2604, 2008.
- [5] Atanu Biswas, Sujay Datta, Jason P Fine, and Mark R Segal. *Statistical advances in the biomedical sciences: clinical trials, epidemiology, survival analysis, and bioinformatics*. John Wiley & Sons, 2007.
- [6] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [7] Petros T Boufounos and Richard G Baraniuk. 1-bit compressive sensing. In *2008 42nd Annual Conference on Information Sciences and Systems*, pages 16–21. IEEE, 2008.
- [8] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- [9] Christian Brownlees, Emilien Joly, and Gábor Lugosi. Empirical risk minimization for heavy-tailed losses. *The Annals of Statistics*, 43(6):2507–2536, 2015.
- [10] T Tony Cai, Cun-Hui Zhang, and Harrison H Zhou. Optimal rates of convergence for covariance matrix estimation. *The Annals of Statistics*, 38(4):2118–2144, 2010.
- [11] T Tony Cai and Harrison H Zhou. Minimax estimation of large covariance matrices under ℓ_1 -norm. *Statistica Sinica*, pages 1319–1349, 2012.
- [12] T Tony Cai and Harrison H Zhou. Optimal rates of convergence for sparse covariance matrix estimation. *The Annals of Statistics*, pages 2389–2420, 2012.
- [13] Tony Cai and Weidong Liu. Adaptive thresholding for sparse covariance matrix estimation. *Journal of the American Statistical Association*, 106(494):672–684, 2011.
- [14] Tony Cai and Wen-Xin Zhou. A max-norm constrained minimization approach to 1-bit matrix completion. *J. Mach. Learn. Res.*, 14(1):3619–3647, 2013.
- [15] Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- [16] Yang Cao and Yao Xie. Categorical matrix completion. In *2015 IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, pages 369–372. IEEE, 2015.
- [17] Olivier Catoni. Challenging the empirical mean and empirical variance: a deviation study. In *Annales de l’IHP Probabilités et statistiques*, volume 48, pages 1148–1185, 2012.
- [18] Junren Chen and Michael K Ng. Color image inpainting via robust pure quaternion matrix completion: Error bound and weighted loss. *SIAM Journal on Imaging Sciences*, 15(3):1469–1498, 2022.
- [19] Junren Chen and Michael K Ng. A parameter-free two-bit covariance estimator with improved operator norm error rate. *arXiv preprint arXiv:2308.16059*, 2023.
- [20] Junren Chen and Michael K. Ng. Uniform exact reconstruction of sparse signals and low-rank matrices from phase-only measurements. *IEEE Transactions on Information Theory*, 69(10):6739–6764, 2023.
- [21] Junren Chen, Jonathan Scarlett, Michael K. Ng, and Zhaoqiang Liu. A unified framework for uniform signal recovery in nonlinear generative compressed sensing. *Advances in Neural Information Processing Systems*, 2023.
- [22] Junren Chen, Cheng-Long Wang, Michael K. Ng, and Di Wang. High dimensional statistical estimation under uniformly dithered one-bit quantization. *IEEE Transactions on Information Theory*, 69(8):5151–5187, 2023.
- [23] Junren Chen, Yueqi Wang, and Michael K. Ng. Quantized low-rank multivariate regression with random dithering. *IEEE Transactions on Signal Processing*, pages 1–16, 2023.
- [24] Onkar Dabeer and Aditya Karnik. Signal parameter estimation using 1-bit dithered quantization. *IEEE Transactions on Information Theory*, 52(12):5389–5405, 2006.
- [25] Alireza Danaee, Rodrigo C de Lamare, and Vitor Heloiz Nascimento. Distributed quantization-aware rls learning with bias compensation and coarsely quantized signals. *IEEE Transactions on Signal Processing*, 70:3441–3455, 2022.
- [26] Mark A Davenport, Yaniv Plan, Ewout Van Den Berg, and Mary Wootters. 1-bit matrix completion. *Information and Inference: A Journal of the IMA*, 3(3):189–223, 2014.
- [27] Mark A Davenport and Justin Romberg. An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):608–622, 2016.
- [28] Luc Devroye, Matthieu Lerasle, Gabor Lugosi, and Roberto I Oliveira. Sub-gaussian mean estimators. *The Annals of Statistics*, 44(6):2695–2725, 2016.
- [29] Sjoerd Dirksen. Quantized compressed sensing: a survey. In *Compressed Sensing and Its Applications*, pages 67–95. Springer, 2019.
- [30] Sjoerd Dirksen, Hans Christian Jung, and Holger Rauhut. One-bit compressed sensing with partial gaussian circulant matrices. *Information and Inference: A Journal of the IMA*, 9(3):601–626, 2020.
- [31] Sjoerd Dirksen and Johannes Maly. Tuning-free one-bit covariance estimation using data-driven dithering. *arXiv preprint arXiv:2307.12613*, 2023.

- [32] Sjoerd Dirksen, Johannes Maly, and Holger Rauhut. Covariance estimation under one-bit quantization. *The Annals of Statistics*, 50(6):3538–3562, 2022.
- [33] Sjoerd Dirksen and Shahar Mendelson. Non-gaussian hyperplane tessellations and robust one-bit compressed sensing. *Journal of the European Mathematical Society*, 23(9):2913–2947, 2021.
- [34] Sjoerd Dirksen, Shahar Mendelson, and Alexander Stollenwerk. Sharp estimates on random hyperplane tessellations. *SIAM Journal on Mathematics of Data Science*, 4(4):1396–1419, 2022.
- [35] Jianqing Fan, Qufeng Li, and Yuyan Wang. Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(1):247–265, 2017.
- [36] Jianqing Fan, Weichen Wang, and Ziwei Zhu. A shrinkage principle for heavy-tailed data: High-dimensional robust low-rank matrix recovery. *Annals of statistics*, 49(3):1239, 2021.
- [37] Simon Foucart, Deanna Needell, Reese Pathak, Yaniv Plan, and Mary Wootters. Weighted matrix completion from non-random, non-uniform sampling patterns. *IEEE Transactions on Information Theory*, 67(2):1264–1290, 2020.
- [38] Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*. Springer New York, New York, NY, 2013.
- [39] Martin Genzel and Christian Kipp. Generic error bounds for the generalized lasso with sub-exponential data. *Sampling Theory, Signal Processing, and Data Analysis*, 20(2):15, 2022.
- [40] Martin Genzel and Alexander Stollenwerk. A unified approach to uniform signal recovery from nonlinear observations. *Foundations of Computational Mathematics*, pages 1–74, 2022.
- [41] Robert M Gray and Thomas G Stockham. Dithered quantizers. *IEEE Transactions on Information Theory*, 39(3):805–812, 1993.
- [42] David Gross, Yi-Kai Liu, Steven T Flammia, Stephen Becker, and Jens Eisert. Quantum state tomography via compressed sensing. *Physical review letters*, 105(15):150401, 2010.
- [43] Osama A Hanna, Yahya H Ezzeldin, Christina Fragouli, and Suhas Diggavi. Quantization of distributed data for learning. *IEEE Journal on Selected Areas in Information Theory*, 2(3):987–1001, 2021.
- [44] Daniel Hsu and Sivan Sabato. Loss minimization and parameter estimation with heavy tails. *The Journal of Machine Learning Research*, 17(1):543–582, 2016.
- [45] Marat Ibragimov, Rustam Ibragimov, and Johan Walden. *Heavy-tailed distributions and robustness in economics and finance*, volume 214. Springer, 2015.
- [46] Laurent Jacques and Valerio Cambareri. Time for dithering: fast and quantized random embeddings via the restricted isometry property. *Information and Inference: A Journal of the IMA*, 6(4):441–476, 2017.
- [47] Laurent Jacques and Thomas Feullen. The importance of phase in complex compressive sensing. *IEEE Transactions on Information Theory*, 67(6):4150–4161, 2021.
- [48] Laurent Jacques, Jason N Laska, Petros T Boufounos, and Richard G Baraniuk. Robust 1-bit compressive sensing via binary stable embeddings of sparse vectors. *IEEE transactions on information theory*, 59(4):2082–2102, 2013.
- [49] Halyun Jeong, Xiaowei Li, Yaniv Plan, and Ozgur Yilmaz. Sub-gaussian matrices on sets: Optimal tail dependence and applications. *Communications on Pure and Applied Mathematics*, 75(8):1713–1754, 2022.
- [50] Hans Christian Jung, Johannes Maly, Lars Palzer, and Alexander Stollenwerk. Quantized compressed sensing by rectified linear units. *IEEE transactions on information theory*, 67(6):4125–4149, 2021.
- [51] Olga Klopp. Noisy low-rank matrix completion with general sampling distribution. *Bernoulli*, 20(1):282–303, 2014.
- [52] Olga Klopp, Jean Lafond, Éric Moulines, and Joseph Salmon. Adaptive multinomial matrix completion. *Electronic Journal of Statistics*, 9(2):2950–2975, 2015.
- [53] Olga Klopp, Karim Lounici, and Alexandre B Tsybakov. Robust matrix completion. *Probability Theory and Related Fields*, 169(1):523–564, 2017.
- [54] Karin Knudson, Rayan Saab, and Rachel Ward. One-bit compressive sensing with norm estimation. *IEEE Transactions on Information Theory*, 62(5):2748–2758, 2016.
- [55] Vladimir Koltchinskii, Karim Lounici, and Alexandre B Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics*, 39(5):2302–2329, 2011.
- [56] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [57] Piotr Kruczek, Radosław Zimroz, and Agnieszka Wyłomańska. How to detect the cyclostationarity in heavy-tailed distributed signals. *Signal Processing*, 172:107514, 2020.
- [58] Arun Kumar Kuchibhotla and Abhishek Chakraborty. Moving beyond sub-gaussianity in high-dimensional statistics: Applications in covariance estimation and linear regression. *Information and Inference: A Journal of the IMA*, 11(4):1389–1456, 2022.
- [59] Jean Lafond, Olga Klopp, Eric Moulines, and Joseph Salmon. Probabilistic low-rank matrix completion on finite alphabets. *Advances in Neural Information Processing Systems*, 27, 2014.
- [60] Kangqiang Li and Yuxuan Wang. Two results on low-rank heavy-tailed multiresponse regressions. *arXiv preprint arXiv:2305.13897*, 2023.
- [61] Christopher Liaw, Abbas Mehrabian, Yaniv Plan, and Roman Vershynin. A simple tool for bounding the deviation of random matrices on geometric sets. In *Geometric aspects of functional analysis*, pages 277–299. Springer, 2017.
- [62] Po-Ling Loh. Statistical consistency and asymptotic normality for high-dimensional robust m -estimators. *The Annals of Statistics*, 45(2):866–896, 2017.
- [63] Po-Ling Loh and Martin J. Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *The Annals of statistics*, 40(3):1637–1664, 2012.
- [64] Po-Ling Loh and Martin J. Wainwright. Regularized m -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16(19):559–616, 2015.
- [65] Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
- [66] Gabor Lugosi and Shahar Mendelson. Robust multivariate mean estimation: the optimality of trimmed mean. *The Annals of Statistics*, 49(1):393–410, 2021.
- [67] Shahar Mendelson. Upper bounds on product and multiplier empirical processes. *Stochastic Processes and their Applications*, 126(12):3652–3680, 2016.
- [68] Stanislav Minsker. Geometric median and robust estimation in banach spaces. *Bernoulli*, 21(4):2308–2335, 2015.
- [69] Jianhua Mo and Robert W Heath. Limited feedback in single and multi-user mimo systems with finite-bit adcs. *IEEE Transactions on Wireless Communications*, 17(5):3284–3297, 2018.
- [70] Sahand Negahban and Martin J Wainwright. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069–1097, 2011.
- [71] Sahand Negahban and Martin J Wainwright. Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *The Journal of Machine Learning Research*, 13(1):1665–1697, 2012.
- [72] Sahand N Negahban, Pradeep Ravikumar, Martin J Wainwright, and Bin Yu. A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Statistical science*, 27(4):538–557, 2012.
- [73] Arkadij Semenovič Nemirovskij and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.
- [74] Luong Trung Nguyen, Junhan Kim, and Byonghyo Shim. Low-rank matrix completion: A contemporary survey. *IEEE Access*, 7:94215–94237, 2019.
- [75] Yaniv Plan and Roman Vershynin. Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach. *IEEE Transactions on Information Theory*, 59(1):482–494, 2012.
- [76] Yaniv Plan and Roman Vershynin. One-bit compressed sensing by linear programming. *Communications on Pure and Applied Mathematics*, 66(8):1275–1297, 2013.
- [77] Yaniv Plan and Roman Vershynin. Dimension reduction by random hyperplane tessellations. *Discrete & Computational Geometry*, 51(2):438–461, 2014.
- [78] Yaniv Plan and Roman Vershynin. The generalized lasso with non-linear observations. *IEEE Transactions on information theory*, 62(3):1528–1537, 2016.
- [79] Yaniv Plan, Roman Vershynin, and Elena Yudovina. High-dimensional estimation with geometric constraints. *Information and Inference: A Journal of the IMA*, 6(1):1–40, 2017.
- [80] Shuang Qiu, Xiaohan Wei, and Zhuoran Yang. Robust one-bit recovery via relu generative networks: Near-optimal statistical rate and global landscape analysis. In *International Conference on Machine Learning*, pages 7857–7866. PMLR, 2020.
- [81] Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE transactions on information theory*, 57(10):6976–6994, 2011.

- [82] Benjamin Recht. A simpler approach to matrix completion. *Journal of Machine Learning Research*, 12(12), 2011.
- [83] Lawrence Roberts. Picture coding using pseudo-random noise. *IRE Transactions on Information Theory*, 8(2):145–154, 1962.
- [84] Sima Sahu, Harsh Vikram Singh, Basant Kumar, and Amit Kumar Singh. De-noising of ultrasound image using bayesian approached heavy-tailed cauchy distribution. *Multimedia Tools and Applications*, 78(4):4089–4106, 2019.
- [85] Vidyashankar Sivakumar, Arindam Banerjee, and Pradeep K Ravikumar. Beyond sub-gaussian measurements: High-dimensional structured estimation with sub-exponential designs. *Advances in neural information processing systems*, 28, 2015.
- [86] Alexander Stollenwerk. *One-bit compressed sensing and fast binary embeddings*. PhD thesis, Dissertation, RWTH Aachen University, 2019, 2019.
- [87] Zhongxing Sun, Wei Cui, and Yulong Liu. Quantized corrupted sensing with random dithering. *IEEE Transactions on Signal Processing*, 70:600–615, 2022.
- [88] Ananthram Swami and Brian M Sadler. On some detection and estimation problems in heavy-tailed noise. *Signal Processing*, 82(12):1829–1846, 2002.
- [89] Christos Thrampoulidis and Ankit Singh Rawat. The generalized lasso for sub-gaussian measurements with dithered quantization. *IEEE Transactions on Information Theory*, 66(4):2487–2500, 2020.
- [90] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [91] Joel A Tropp. An introduction to matrix concentration inequalities. *arXiv preprint arXiv:1501.01571*, 2015.
- [92] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [93] Robert F Woolson and William R Clarke. *Statistical methods for the analysis of biomedical data*. John Wiley & Sons, 2011.
- [94] Chunlei Xu and Laurent Jacques. Quantized compressive sensing with rip matrices: The benefit of dithering. *Information and Inference: A Journal of the IMA*, 9(3):543–586, 2020.
- [95] Tianyu Yang, Johannes Maly, Sjoerd Dirksen, and Giuseppe Caire. Plug-in channel estimation with dithered quantized signals in spatially non-stationary massive mimo systems. *arXiv preprint arXiv:2301.04641*, 2023.
- [96] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. Zipml: Training linear models with end-to-end low precision, and a little bit of deep learning. In *International Conference on Machine Learning*, pages 4035–4043. PMLR, 2017.
- [97] Ziwei Zhu and Wenjing Zhou. Taming heavy-tailed features by shrinkage. In *International Conference on Artificial Intelligence and Statistics*, pages 3268–3276. PMLR, 2021.

APPENDIX

A. Proofs in Section III

1) **Quantized Covariance Matrix Estimation:** We first provide Bernstein's inequality that is recurring in our proofs. In application, we will choose the more convenient one from (16) and (17).

Lemma 1. (Bernstein's inequality, [6, Thm. 2.10, Coro. 2.11]). *Let X_1, \dots, X_n be independent random variables, and assume that there exist positive numbers v and c such that $\sum_{i=1}^n \mathbb{E}[X_i^2] \leq v$ and*

$$\sum_{i=1}^n \mathbb{E}|X_i|^q \leq \frac{q!}{2} v c^{q-2} \text{ for all integers } q \geq 3,$$

then for any $t > 0$ we have

$$\mathbb{P} \left(\left| \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \geq \sqrt{2vt} + ct \right) \leq 2 \exp(-t) \quad (16)$$

$$\mathbb{P} \left(\left| \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \geq t \right) \leq 2 \exp \left(-\frac{t^2}{2(v+ct)} \right) \quad (17)$$

We will also use the Matrix Bernstein's inequality.

Lemma 2. (Matrix Bernstein, [91, Thm. 6.1.1]). *Let $\mathbf{S}_1, \dots, \mathbf{S}_n$ be independent zero-mean random variables with common dimension $d_1 \times d_2$. We assume that $\|\mathbf{S}_k\|_{op} \leq L$ for $k \in [n]$ and introduce the matrix variance statistic*

$$\nu = \max \left\{ \left\| \sum_{k=1}^n \mathbb{E}(\mathbf{S}_k \mathbf{S}_k^\top) \right\|_{op}, \left\| \sum_{k=1}^n \mathbb{E}(\mathbf{S}_k^\top \mathbf{S}_k) \right\|_{op} \right\}.$$

Then for any $t \geq 0$, we have

$$\mathbb{P} \left(\left\| \sum_{k=1}^n \mathbf{S}_k \right\|_{op} \geq t \right) \leq (d_1 + d_2) \exp \left(\frac{-\frac{1}{2}t^2}{\nu + \frac{Lt}{3}} \right).$$

Proof of Theorem 2:

Proof. Recall that $\boldsymbol{\xi}_k = \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k$ is the quantization noise, and $\mathbb{E}(\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top) = \frac{\Delta^2}{4} \mathbf{I}_d$, which implies $\mathbb{E}\hat{\boldsymbol{\Sigma}} = \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top)$. Thus, by using triangle inequality we obtain

$$\begin{aligned} \|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \|_\infty &\leq \|\hat{\boldsymbol{\Sigma}} - \mathbb{E}\hat{\boldsymbol{\Sigma}}\|_\infty + \|\mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top)\|_\infty \\ &:= I_1 + I_2. \end{aligned}$$

Step 1. Bounding I_1 .

Note that

$$\|\hat{\boldsymbol{\Sigma}} - \mathbb{E}\hat{\boldsymbol{\Sigma}}\|_\infty = \left\| \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \mathbb{E}(\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top) \right\|_\infty,$$

so for any $(i, j) \in [d] \times [d]$ we aim to bound the (i, j) -th entry error

$$|\hat{\sigma}_{ij} - \mathbb{E}\hat{\sigma}_{ij}| = \left| \sum_{k=1}^n \frac{1}{n} \hat{x}_{ki} \hat{x}_{kj} - \mathbb{E}(\hat{x}_{ki} \hat{x}_{kj}) \right|$$

Observe that the quantization noise is bounded as follows

$$\|\boldsymbol{\xi}_k\|_\infty \leq \|\mathcal{Q}_\Delta(\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k) - (\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k)\|_\infty + \|\boldsymbol{\tau}_k\|_\infty \leq \frac{3\Delta}{2},$$

which implies $\mathbb{E}|\xi_{ki}|^4 \leq (\frac{3\Delta}{2})^4$ and

$$\|\hat{\mathbf{x}}_k\|_\infty \leq \|\tilde{\mathbf{x}}_k\|_\infty + \|\boldsymbol{\xi}_k\|_\infty \leq \zeta + \frac{3\Delta}{2}.$$

By the moment constraint on x_{ki} we have $\mathbb{E}|\tilde{x}_{ki}|^4 \leq \mathbb{E}|x_{ki}|^4 \leq M$. Thus, for any positive integer $p \geq 2$ we have the following bound

$$\begin{aligned} &\sum_{k=1}^n \mathbb{E} \left| \frac{\hat{x}_{ki} \hat{x}_{kj}}{n} \right|^q \\ &\leq \frac{(\zeta + \frac{3}{2}\Delta)^{2(q-2)}}{n^q} \sum_{k=1}^n \mathbb{E}(\hat{x}_{ki} \hat{x}_{kj})^2 \\ &\leq \frac{(\zeta + \frac{3}{2}\Delta)^{2(q-2)}}{2n^q} \sum_{k=1}^n (\mathbb{E}|\hat{x}_{ki}|^4 + \mathbb{E}|\hat{x}_{kj}|^4) \\ &\leq \frac{4(\zeta + \frac{3}{2}\Delta)^{2(q-2)}}{n^q} \\ &\quad \cdot \sum_{k=1}^n (\mathbb{E}|\tilde{x}_{ki}|^4 + \mathbb{E}|\tilde{x}_{kj}|^4 + \mathbb{E}|\xi_{ki}|^4 + \mathbb{E}|\xi_{kj}|^4) \\ &\leq \frac{q!}{2} v_0 c_0^{q-2}, \end{aligned} \quad (18)$$

for some $v_0 = O(\frac{M+\Delta^4}{n})$, $c_0 = O(\frac{(\zeta+\Delta^2)}{n})$. With these preparations, we can invoke Bernstein's inequality (Lemma 1) to obtain that, for any $t \geq 0$, with probability at least $1 - 2\exp(-t)$,

$$|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| \leq C_1 \left(\sqrt{\frac{(M + \Delta^4)t}{n}} + \frac{(\zeta^2 + \Delta^2)t}{n} \right).$$

Taking $t = \delta \log d$ and using the choice $\zeta \asymp (\frac{nM}{\delta \log d})^{1/4}$, then applying a union bound over $(i, j) \in [d] \times [d]$, under the scaling $n \gtrsim \delta \log d$, we obtain that

$$I_1 \lesssim (\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 2d^{2-\delta}$.

Step 2. Bounding I_2 .

We aim to bound $|\mathbb{E}(\widetilde{x}_{ki}\widetilde{x}_{kj} - x_{ki}x_{kj})|$ for any $(i, j) \in [d] \times [d]$. First by the definition of truncation we have

$$\begin{aligned} & |\mathbb{E}(\widetilde{x}_{ki}\widetilde{x}_{kj} - x_{ki}x_{kj})| \\ & \leq \mathbb{E}[|x_{ki}x_{kj}|(\mathbb{1}(|x_{ki}| \geq \zeta) + \mathbb{1}(|x_{kj}| \geq \zeta))]; \end{aligned}$$

then applying Cauchy-Schwarz to $\mathbb{E}[|x_{ki}x_{kj}| \mathbb{1}(|x_{ki}| \geq \zeta)]$, we obtain

$$\begin{aligned} & \mathbb{E}[|x_{ki}x_{kj}| \mathbb{1}(|x_{ki}| \geq \zeta)] \\ & \leq [\mathbb{E}|x_{ki}x_{kj}|^2]^{1/2} [\mathbb{P}(|x_{ki}| \geq \zeta)]^{1/2} \\ & \leq \sqrt{M} \sqrt{\frac{M}{\zeta^4}} = \frac{M}{\zeta^2}, \end{aligned}$$

where the second inequality is due to Markov's inequality. Note that this bound remains valid for $\mathbb{E}[|x_{ki}x_{kj}| \mathbb{1}(|x_{kj}| \geq \zeta)]$. Since this holds for any $(i, j) \in [d] \times [d]$, combining with $\zeta \asymp (\frac{nM}{\delta \log d})^{1/4}$, we obtain

$$\|\mathbb{E}(\widetilde{\mathbf{x}}_k \widetilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top)\|_\infty \leq \frac{2M}{\zeta^2} \lesssim \sqrt{M} \sqrt{\frac{\delta \log d}{n}}.$$

By putting pieces together, we have

$$\|\widehat{\Sigma} - \Sigma^*\|_\infty \lesssim (\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$$

with probability at least $1 - 2d^{2-\delta}$, as claimed. \square

Proof of Theorem 3:

Proof. Note that the calculations in (5) and (6) remain valid (but the truncated samples are denoted by $\check{\mathbf{x}}_k$ rather than $\widetilde{\mathbf{x}}_k$), so we have $\mathbb{E}\widehat{\Sigma} = \mathbb{E}(\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top)$. Using triangle inequality we first decompose the error as

$$\begin{aligned} \|\widehat{\Sigma} - \Sigma^*\|_{op} & \leq \|\widehat{\Sigma} - \mathbb{E}\widehat{\Sigma}\|_{op} + \|\mathbb{E}(\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top)\|_{op} \\ & =: I_1 + I_2. \end{aligned}$$

Step 1. Bounding I_1 .

We first write that

$$\widehat{\Sigma} - \mathbb{E}\widehat{\Sigma} = \frac{1}{n} \sum_{k=1}^n \mathbf{S}_k \text{ where } \mathbf{S}_k = \check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top - \mathbb{E}(\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top).$$

Recall that we define quantization error as $\mathbf{w}_k = \check{\mathbf{x}}_k - \widetilde{\mathbf{x}}_k - \boldsymbol{\tau}_k$ and quantization noise as $\boldsymbol{\xi}_k = \check{\mathbf{x}}_k - \widetilde{\mathbf{x}}_k$, and observe that

the quantization noise is bounded $\|\boldsymbol{\xi}_k\|_\infty = \|\check{\mathbf{x}}_k - \widetilde{\mathbf{x}}_k\|_\infty = \|\boldsymbol{\tau}_k + \mathbf{w}_k\|_\infty \leq \frac{3}{2}\Delta$. Thus, by $\|\mathbf{a}\|_2^2 \leq \sqrt{d}\|\mathbf{a}\|_4^2$ that holds for any $\mathbf{a} \in \mathbb{R}^d$, we obtain

$$\begin{aligned} \|\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top\|_{op} & = \|\check{\mathbf{x}}_k\|_2^2 = \|\widetilde{\mathbf{x}}_k + \boldsymbol{\xi}_k\|_2^2 \leq 2\|\widetilde{\mathbf{x}}_k\|_2^2 + 2\|\boldsymbol{\xi}_k\|_2^2 \\ & \leq 2\sqrt{d} \cdot \|\widetilde{\mathbf{x}}_k\|_4^2 + 2d \cdot \left(\frac{3\Delta}{2}\right)^2 \leq 2\sqrt{d}\zeta^2 + \frac{9}{2}d\Delta^2, \end{aligned}$$

which implies

$$\|\mathbf{S}_k\|_{op} \leq \|\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top\|_{op} + \mathbb{E}\|\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top\|_{op} \leq 4\sqrt{d}\zeta^2 + 9d\Delta^2.$$

Moreover, we estimate the matrix variance statistic. Since \mathbf{S}_k is symmetric, we simply deal with $\|\mathbb{E}\mathbf{S}_k^2\|_{op}$ and some algebra gives $\mathbb{E}\mathbf{S}_k^2 = \mathbb{E}[\|\check{\mathbf{x}}_k\|_2^2 \check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] - (\mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top])^2$. First let us note that

$$\begin{aligned} \left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op}^2 & = \left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op}^2 \\ & = \left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] + \frac{\Delta^2}{4} \mathbf{I}_d \right\|_{op}^2 \\ & \leq \left(\left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op} + \frac{\Delta^2}{4} \right)^2 \\ & \leq 2 \left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op}^2 + \frac{\Delta^4}{8}. \end{aligned}$$

Combining with the observation that

$$\begin{aligned} \left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op} & = \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}(\mathbf{v}^\top \check{\mathbf{x}}_k)^2 \\ & \leq \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \sqrt{\mathbb{E}(\mathbf{v}^\top \mathbf{x}_k)^4} \leq \sqrt{M}, \end{aligned}$$

we obtain $\left\| \mathbb{E}[\check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top] \right\|_{op}^2 = O(M + \Delta^4)$. Then we turn to the operator norm of $\mathbb{E}[\|\check{\mathbf{x}}_k\|_2^2 \check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top]$. We apply Cauchy-Schwarz to estimate

$$\begin{aligned} & \left\| \mathbb{E}(\|\check{\mathbf{x}}_k\|_2^2 \check{\mathbf{x}}_k \check{\mathbf{x}}_k^\top) \right\|_{op} \\ & = \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}(\|\check{\mathbf{x}}_k\|_2^2 (\mathbf{v}^\top \check{\mathbf{x}}_k)^2) \\ & \leq \sqrt{\mathbb{E}\|\check{\mathbf{x}}_k\|_2^4} \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \sqrt{\mathbb{E}(\mathbf{v}^\top \check{\mathbf{x}}_k)^4}. \end{aligned} \tag{19}$$

By $\|\mathbf{a}\|_2^2 \leq \sqrt{d}\|\mathbf{a}\|_4^2$ that holds for any $\mathbf{a} \in \mathbb{R}^d$, $\mathbb{E}|\check{x}_{ki}|^4 \leq \mathbb{E}|x_{ki}|^4 \leq M$, $\check{\mathbf{x}}_k = \widetilde{\mathbf{x}}_k + \boldsymbol{\xi}_k$ and $\|\boldsymbol{\xi}_k\|_\infty \leq \frac{3\Delta}{2}$, we obtain

$$\begin{aligned} \mathbb{E}\|\check{\mathbf{x}}_k\|_2^4 & \leq \mathbb{E}(\|\widetilde{\mathbf{x}}_k\|_2 + \|\boldsymbol{\xi}_k\|_2)^4 \lesssim \mathbb{E}(\|\widetilde{\mathbf{x}}_k\|_4^4 + \|\boldsymbol{\xi}_k\|_4^4) \\ & \leq d\mathbb{E}(\|\widetilde{\mathbf{x}}_k\|_4^4 + \|\boldsymbol{\xi}_k\|_4^4) \lesssim d^2(M + \Delta^4). \end{aligned} \tag{20}$$

For any $\mathbf{v} \in \mathbb{S}^{d-1}$, we write $\check{\mathbf{x}}_k = \widetilde{\mathbf{x}}_k + \boldsymbol{\tau}_k + \mathbf{w}_k$ and then have the bound

$$\begin{aligned} \mathbb{E}(\mathbf{v}^\top \check{\mathbf{x}}_k)^4 & \leq \mathbb{E}(\mathbf{v}^\top \widetilde{\mathbf{x}}_k)^4 + \mathbb{E}(\mathbf{v}^\top \boldsymbol{\tau}_k)^4 + \mathbb{E}(\mathbf{v}^\top \mathbf{w}_k)^4 \\ & \stackrel{(i)}{\lesssim} M + \Delta^4, \end{aligned} \tag{21}$$

where (i) is because

$$\mathbb{E}(\mathbf{v}^\top \widetilde{\mathbf{x}}_k)^4 \leq \mathbb{E}(\mathbf{v}^\top \mathbf{x}_k)^4 \leq M$$

and

$$\boldsymbol{\tau}_k \sim \mathcal{U}([- \frac{\Delta}{2}, \frac{\Delta}{2}]^d) + \mathcal{U}([- \frac{\Delta}{2}, \frac{\Delta}{2}]^d),$$

and the quantization error \mathbf{w}_k follows $\mathcal{U}([- \frac{\Delta}{2}, \frac{\Delta}{2}]^d)$; in more detail, $\|\mathbf{v}^\top \boldsymbol{\tau}_k\|_{\psi_2}, \|\mathbf{v}^\top \mathbf{w}_k\|_{\psi_2} = O(\Delta)$ and then the

moment constraint of sub-Gaussian random variable implies $\mathbb{E}(\mathbf{v}^\top \boldsymbol{\tau}_k)^4 = O(\Delta^4)$ and $\mathbb{E}(\mathbf{v}^\top \mathbf{w}_k) = O(\Delta^4)$. From (19), (20) and (21), we obtain

$$\|\mathbb{E}(\|\dot{\mathbf{x}}_k\|_2^2 \dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top)\|_{op} = O(d(\Delta^4 + M)).$$

Further combining with

$$\mathbb{E}S_k^2 = \mathbb{E}[\|\dot{\mathbf{x}}_k\|_2^2 \dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top] - (\mathbb{E}[\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top])^2$$

and $\|(\mathbb{E}[\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top])^2\|_{op} = O(M + \Delta^4)$, we arrive at $\|\mathbb{E}S_k^2\|_{op} \lesssim d(\Delta^4 + M)$ and hence $\|\sum_{k=1}^n \mathbb{E}S_k^2\|_{op} \lesssim nd(\Delta^4 + M)$. With these preparations, Matrix Bernstein's inequality (Lemma 2) yields the following inequality that holds for any $t \geq 0$

$$\begin{aligned} & \mathbb{P}(\|\widehat{\boldsymbol{\Sigma}} - \mathbb{E}\widehat{\boldsymbol{\Sigma}}\|_{op} \geq t) \\ & \leq 2d \exp\left(-\frac{C_1 n t^2}{(M + \Delta^4)d + (\sqrt{d}\zeta^2 + d\Delta^2)t}\right). \end{aligned}$$

Setting $t = C_2(\sqrt{M} + \Delta^2)\sqrt{\frac{\delta d \log d}{n}}$ with sufficiently large C_2 , under the scaling of $n \gtrsim \delta d \log d$ and the threshold $\zeta \asymp (M^{1/4} + \Delta)\left(\frac{n}{\delta \log d}\right)^{1/4}$, we obtain that

$$I_1 = \|\widehat{\boldsymbol{\Sigma}} - \mathbb{E}\widehat{\boldsymbol{\Sigma}}\|_{op} \leq C_2(\sqrt{M} + \Delta^2)\sqrt{\frac{\delta d \log d}{n}}$$

holds with probability at least $1 - 2d^{1-\delta}$.

Step 2. Bounding I_2 .

Having bounded the concentration term I_1 , we now switch to the bias term

$$I_2 = \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} |\mathbf{v}^\top \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top) \mathbf{v}|.$$

For any $\mathbf{v} \in \mathbb{S}^{d-1}$, because $\tilde{\mathbf{x}}_k$ is obtained from truncating \mathbf{x}_4 regarding ℓ_4 -norm, we have

$$\begin{aligned} & |\mathbf{v}^\top \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top) \mathbf{v}| \\ & = \left| \mathbb{E}\left[\left((\mathbf{v}^\top \tilde{\mathbf{x}}_k)^2 - (\mathbf{v}^\top \mathbf{x}_k)^2\right) \mathbb{1}(\|\mathbf{x}_k\|_4 \geq \zeta)\right] \right| \\ & \leq \mathbb{E}\left[(\mathbf{v}^\top \mathbf{x}_k)^2 \mathbb{1}(\|\mathbf{x}_k\|_4 \geq \zeta)\right] \\ & \stackrel{(i)}{\leq} \sqrt{\mathbb{E}(\mathbf{v}^\top \mathbf{x}_k)^4} \sqrt{\mathbb{P}(\|\mathbf{x}_k\|_4 \geq \zeta^4)} \\ & \stackrel{(ii)}{\leq} \sqrt{M \frac{\mathbb{E}\|\mathbf{x}_k\|_4^4}{\zeta^4}} \stackrel{(iii)}{\lesssim} \sqrt{\frac{M \delta d \log d}{n}}, \end{aligned}$$

where (i) and (ii) are respectively by Cauchy-Schwarz and Markov's, and in (iii) we use $\zeta \asymp (M^{1/4} + \Delta)\left(\frac{\delta d \log d}{n}\right)^{1/4}$. This leads to the bound $I_2 \lesssim \sqrt{\frac{M \delta d \log d}{n}}$. Combining the bounds of I_1, I_2 completes the proof. \square

Proof of Theorem 4: This small appendix is devoted to the proof of Theorem 4, for which we need a Lemma concerning the element-wise error rate of $\widehat{\boldsymbol{\Sigma}}_s$, i.e., $|\check{\sigma}_{ij} - \sigma_{ij}^*|$ where we write $\widehat{\boldsymbol{\Sigma}}_s = [\check{\sigma}_{ij}]$, $\boldsymbol{\Sigma}^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top) = [\sigma_{ij}^*]$. Recalling that $\widehat{\boldsymbol{\Sigma}}_s = \mathcal{T}_\mu(\widehat{\boldsymbol{\Sigma}})$, the key message from Lemma 3 is that due to the thresholding operator $\mathcal{T}_\mu(\cdot)$, $\widehat{\boldsymbol{\Sigma}}_s$ respects an element-wise bound tighter than $O\left(\sqrt{\frac{\delta \log d}{n}}\right)$ in Theorem 2, as can be seen from the additional branch $|\sigma_{ij}^*|$ in (22).

Lemma 3. (Element-wise Error Rate of $\widehat{\boldsymbol{\Sigma}}_s$). *For any $i, j \in [d]$, the thresholding estimator $\widehat{\boldsymbol{\Sigma}}_s = [\check{\sigma}_{ij}]$ in Theorem 4 satisfies for some C that*

$$\mathbb{P}\left(|\check{\sigma}_{ij} - \sigma_{ij}^*| \leq C \min\left\{|\sigma_{ij}^*|, \mathcal{L} \sqrt{\frac{\delta \log d}{n}}\right\}\right) \geq 1 - 2d^{-\delta} \quad (22)$$

where $\mathcal{L} := \sqrt{M} + \Delta^2$.

Proof. Recall that $\widehat{\boldsymbol{\Sigma}}_s = [\check{\sigma}_{ij}] = \mathcal{T}_\mu(\widehat{\boldsymbol{\Sigma}}) = \mathcal{T}_\mu([\widehat{\sigma}_{ij}])$ and hence $\check{\sigma}_{ij} = \mathcal{T}_\mu(\widehat{\sigma}_{ij})$. Given (i, j) , the proof of Theorem 2 delivers

$$|\widehat{\sigma}_{ij} - \sigma_{ij}^*| \leq C_1 \mathcal{L} \sqrt{\frac{\delta \log d}{n}}$$

with probability at least $1 - 2d^{-\delta}$. Assume that we are on this event in the following analyses. As stated in Theorem 4, we set $\mu = C_2 \mathcal{L} \sqrt{\frac{\delta \log d}{n}}$ with $C_2 > C_1$, $\mathcal{L} = \sqrt{M} + \Delta^2$. Since $\check{\sigma}_{ij} = \mathcal{T}_\mu(\widehat{\sigma}_{ij})$, we discuss whether $|\widehat{\sigma}| \geq \mu$ holds.

Case 1. when $|\widehat{\sigma}_{ij}| < \mu$ holds.

In this case we have $\check{\sigma}_{ij} = 0$, thus $|\check{\sigma}_{ij} - \sigma_{ij}^*| \leq |\sigma_{ij}^*|$. Further note that

$$\begin{aligned} |\sigma_{ij}^*| & \leq |\sigma_{ij}^* - \widehat{\sigma}_{ij}| + |\widehat{\sigma}_{ij}| \\ & \leq C_1 \mathcal{L} \sqrt{\frac{\delta \log d}{n}} + \mu \lesssim \mathcal{L} \sqrt{\frac{\delta \log d}{n}}, \end{aligned}$$

so we also have $|\check{\sigma}_{ij} - \sigma_{ij}^*| \lesssim \mathcal{L} \sqrt{\frac{\delta \log d}{n}}$.

Case 2. when $|\widehat{\sigma}_{ij}| \geq \mu$ holds.

We consider $|\widehat{\sigma}_{ij}| \geq \mu$ that implies $\check{\sigma}_{ij} = \widehat{\sigma}_{ij}$, then we have

$$|\check{\sigma}_{ij} - \sigma_{ij}^*| = |\widehat{\sigma}_{ij} - \sigma_{ij}^*| \leq C_1 \mathcal{L} \sqrt{\frac{\delta \log d}{n}}.$$

Moreover, note that

$$\begin{aligned} |\sigma_{ij}^*| & \geq |\widehat{\sigma}_{ij}| - |\widehat{\sigma}_{ij} - \sigma_{ij}^*| \geq \mu - |\widehat{\sigma}_{ij} - \sigma_{ij}^*| \\ & \geq (C_2 - C_1) \mathcal{L} \sqrt{\frac{\delta \log d}{n}} \end{aligned}$$

so we also have $|\check{\sigma}_{ij} - \sigma_{ij}^*| = O(|\sigma_{ij}^*|)$.

Therefore, in both cases we have proved that $|\check{\sigma}_{ij} - \sigma_{ij}^*| \lesssim \min\left\{|\sigma_{ij}^*|, \mathcal{L} \sqrt{\frac{\delta \log d}{n}}\right\}$, which completes the proof. \square

We are now in a position to present the proof.

Proof of Theorem 4. We let $p = \frac{\delta}{4} \geq 1$ (just assume $\delta \geq 4$) and use $B_0 := \mathcal{L} \sqrt{\frac{\delta \log d}{n}}$ as shorthand. For $(i, j) \in [d] \times [d]$ we define the event \mathcal{A}_{ij} as

$$\mathcal{A}_{ij} = \left\{|\check{\sigma}_{ij} - \sigma_{ij}^*| \leq C_1 \min\left\{|\sigma_{ij}^*|, B_0\right\}\right\}.$$

By Lemma 3 we can choose C_1 to be sufficiently large such that $C_1 B_0 > 3\mu$ and $\mathbb{P}(\mathcal{A}_{ij}^c) \leq 2d^{-\delta}$; here, by convention we let \mathcal{A}_{ij}^c be the complement of \mathcal{A}_{ij} . Our proof strategy is to first bound the p -th order moment $\mathbb{E}\|\widehat{\boldsymbol{\Sigma}}_s - \boldsymbol{\Sigma}^*\|_{op}^p$, and then invoke Markov's inequality to derive a high probability bound. We start with the simple estimate displayed in (23), where (i) and (ii) are due to $\|\mathbf{A}\|_{op} \leq \sup_{j \in [d]} \sum_{i \in [d]} |a_{ij}|$ for symmetric \mathbf{A} and $(a+b)^p \leq (2a)^p + (2b)^p$. In this proof,

$$\begin{aligned} \mathbb{E}\|\widehat{\Sigma}_s - \Sigma^*\|_{op}^p &\stackrel{(i)}{\leq} \mathbb{E}\left(\sup_{j \in [d]} \sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}) + \sup_{j \in [d]} \sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}^c)\right)^p \\ &\stackrel{(ii)}{\leq} 2^p \mathbb{E} \sup_{j \in [d]} \left(\sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij})\right)^p + 2^p \mathbb{E} \sup_{j \in [d]} \left(\sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}^c)\right)^p := I_1 + I_2 \end{aligned} \quad (23)$$

$$\begin{aligned} W_j &\stackrel{(i)}{\leq} \left(\sum_{i=1}^d |\sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}^c) \mathbb{1}(|\widehat{\sigma}_{ij}| < \mu) + \sum_{i=1}^d |\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| \mathbb{1}(\mathcal{A}_{ij}^c) + \sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}^c)\right)^p \\ &\leq (3d)^{p-1} \left(\sum_{i=1}^d |\sigma_{ij}^*|^p \mathbb{1}(\mathcal{A}_{ij}^c) \mathbb{1}(|\widehat{\sigma}_{ij}| < \mu) + \sum_{i=1}^d |\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}|^p \mathbb{1}(\mathcal{A}_{ij}^c) + \sum_{i=1}^d |\check{\sigma}_{ij} - \sigma_{ij}^*|^p \mathbb{1}(\mathcal{A}_{ij}^c)\right), \end{aligned} \quad (25)$$

the ranges of indices in summation or supremum, if omitted, are $[d]$.

Step 1. Bounding I_1 .

By the definition of \mathcal{A}_{ij} , $|\check{\sigma}_{ij} - \sigma_{ij}^*| = 0$ if $|\sigma_{ij}^*| = 0$. Because the columns of Σ^* are s -sparse, we can straightforwardly bound I_1 as follows:

$$\begin{aligned} I_1 &= 2^p \mathbb{E} \sup_j \left(\sum_{i: |\sigma_{ij}^*| > 0} |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij})\right)^p \\ &\leq (2C_1 s B_0)^p. \end{aligned} \quad (24)$$

Step 2. Bounding I_2 .

We first write $I_2 = 2^p \mathbb{E} \sup_j W_j$ with

$$W_j := \left(\sum_i |\check{\sigma}_{ij} - \sigma_{ij}^*| \mathbb{1}(\mathcal{A}_{ij}^c)\right)^p,$$

then start from the display (25), where in (i) we define

$$\mathbb{E}\widehat{\sigma}_{ij} = \mathbb{E}(\tilde{x}_{ki} \tilde{x}_{kj}) := \tilde{\sigma}_{ij}.$$

By replacing \sup_j with \sum_j , this further gives

$$\begin{aligned} I_2 &\leq 6^p d^{p-1} \left(\sum_{i,j} |\sigma_{ij}^*|^p \mathbb{E}[\mathbb{1}(\mathcal{A}_{ij}^c) \mathbb{1}(|\widehat{\sigma}_{ij}| < \mu)] \right. \\ &\quad \left. + \sum_{i,j} \mathbb{E}[|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}|^p \mathbb{1}(\mathcal{A}_{ij}^c)] \right. \\ &\quad \left. + \sum_{i,j} |\check{\sigma}_{ij} - \sigma_{ij}^*|^p \mathbb{P}(\mathcal{A}_{ij}^c)\right) \\ &:= 6^p d^{p-1} (I_{21} + I_{22} + I_{23}). \end{aligned} \quad (26)$$

Step 2.1. Bounding I_{21} .

Note that \mathcal{A}_{ij}^c means $|\check{\sigma}_{ij} - \sigma_{ij}^*| > C_1 \min\{|\sigma_{ij}^*|, B_0\}$, and $|\widehat{\sigma}_{ij}| < \mu$ implies $\check{\sigma}_{ij} = 0$, their combination thus allows us to proceed as the following (i) and (iii):

$$|\sigma_{ij}^*| \stackrel{(i)}{>} C_1 B_0 \stackrel{(ii)}{>} 3\mu \stackrel{(iii)}{>} 3|\widehat{\sigma}_{ij}| \geq 3|\sigma_{ij}^*| - 3|\widehat{\sigma}_{ij} - \sigma_{ij}^*|,$$

where (ii) is due to our choice of C_1 . Thus, $\mathcal{A}_{ij}^c \cap \{|\widehat{\sigma}_{ij}| < \mu\}$ implies $|\widehat{\sigma}_{ij} - \sigma_{ij}^*| > \frac{2}{3}|\sigma_{ij}^*|$ and $|\sigma_{ij}^*| > 3\mu$. Note that Step 2 in the proof of Theorem 2 gives $|\widehat{\sigma}_{ij} - \sigma_{ij}^*| = O(B_0)$, and hence

we can assume $\mu > |\widehat{\sigma}_{ij} - \sigma_{ij}^*|$ and so $|\sigma_{ij}^*| > 3|\widehat{\sigma}_{ij} - \sigma_{ij}^*|$. Using these relations and triangle inequality, we obtain

$$\begin{aligned} \frac{2}{3}|\sigma_{ij}^*| &< |\widehat{\sigma}_{ij} - \sigma_{ij}^*| \leq |\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| + |\tilde{\sigma}_{ij} - \sigma_{ij}^*| \\ &< |\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| + \frac{1}{3}|\sigma_{ij}^*|, \end{aligned}$$

which implies $|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| > \frac{1}{3}|\sigma_{ij}^*|$. Now we conclude that, $\mathcal{A}_{ij}^c \cap \{|\widehat{\sigma}_{ij}| < \mu\}$ implies $|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| > \frac{1}{3}|\sigma_{ij}^*|$ and $|\sigma_{ij}^*| > 3\mu$, which allows us to bound I_{21} as

$$I_{21} = \sum_{i,j} |\sigma_{ij}^*|^p \mathbb{1}(|\sigma_{ij}^*| > 3\mu) \mathbb{P}\left(|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| > \frac{1}{3}|\sigma_{ij}^*|\right). \quad (27)$$

Analogously to the proof of Theorem 2, we can apply Bernstein's inequality to $\mathbb{P}\left(|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| > \frac{1}{3}|\sigma_{ij}^*|\right)$. More specifically, by preparations as in (18), we can use (17) in Lemma 1 with

$$v = O\left(\frac{M + \Delta^4}{n}\right), \quad c = O\left(\frac{\zeta^2 + \Delta^2}{n}\right) = O\left(\frac{\Delta^2}{n} + \sqrt{\frac{M}{n\delta \log d}}\right)$$

(recall that $\zeta \asymp \left(\frac{nM}{\Delta \log d}\right)^{1/4}$). For some absolute constants C_2, C_3 , it gives the display (28), where in (i) we use

$$\min\left\{\frac{|\sigma_{ij}^*|}{M + \Delta^4}, \Delta^{-2}, \sqrt{\frac{\delta \log d}{nM}}\right\} \gtrsim \frac{1}{\sqrt{M + \Delta^2}} \sqrt{\frac{\delta \log d}{n}}$$

that holds because $|\sigma_{ij}^*| > 3\mu$ and $n \gtrsim \delta \log d$. We substitute (28) into (27) and perform some estimates as in the display (29), where in (i) we substitute $|\sigma_{ij}^*| > 3\mu$ from the indicator function into the exponent, (ii) is because $\sup_{t \geq 0} t^p \exp\left(-\frac{t}{2}\right) \leq p^p$, $p = \frac{\delta}{4}$, and we consider $\mu = C_4(\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$ with C_4 large enough.

Step 2.2. Bounding I_{22} .

Then, we deal with I_{22} by Cauchy-Schwarz

$$I_{22} \leq \sum_{i,j} \sqrt{\mathbb{E}|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}|^{2p}} \sqrt{\mathbb{P}(\mathcal{A}_{ij}^c)}.$$

As in (28), we can use (17) in Lemma 1 with $v = O\left(\frac{M + \Delta^4}{n}\right)$ and $c = O\left(\frac{\Delta^2}{n} + \sqrt{\frac{M}{n\delta \log d}}\right)$, yielding that for any $t \geq 0$, $\mathbb{P}\left(|\widehat{\sigma}_{ij} - \mathbb{E}\widehat{\sigma}_{ij}| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(v+ct)}\right) \leq 2 \exp\left(-\frac{t^2}{4v}\right) +$

$$\begin{aligned} \mathbb{P}\left(|\hat{\sigma}_{ij} - \mathbb{E}\hat{\sigma}_{ij}| > \frac{1}{3}|\sigma_{ij}^*|\right) &\leq 2 \exp\left(-\frac{|\sigma_{ij}^*|^2}{C_2\left\{\frac{M+\Delta^4}{n} + \frac{\Delta^2|\sigma_{ij}^*|}{n} + \sqrt{\frac{M}{n\delta\log d}}|\sigma_{ij}^*|\right\}}\right) \\ &\leq 2 \exp\left(-\frac{3n|\sigma_{ij}^*|}{C_2} \min\left\{\frac{|\sigma_{ij}^*|}{M+\Delta^4}, \frac{1}{\Delta^2}, \sqrt{\frac{\delta\log d}{nM}}\right\}\right) \stackrel{(i)}{\leq} 2 \exp\left(-\frac{C_3|\sigma_{ij}^*|\sqrt{n\delta\log d}}{\sqrt{M+\Delta^2}}\right), \end{aligned} \quad (28)$$

$$\begin{aligned} I_{21} &\leq 2 \sum_{i,j} |\sigma_{ij}^*|^p \mathbb{1}(|\sigma_{ij}^*| > 3\mu) \exp\left(-\frac{C_3|\sigma_{ij}^*|\sqrt{n\delta\log d}}{\sqrt{M+\Delta^2}}\right) \\ &= 2 \sum_{i,j} \left(\frac{\sqrt{M+\Delta^2}}{C_3\sqrt{n\delta\log d}}\right)^p \cdot \left(\frac{C_3|\sigma_{ij}^*|\sqrt{n\delta\log d}}{\sqrt{M+\Delta^2}}\right)^p \exp\left(-\frac{0.5C_3|\sigma_{ij}^*|\sqrt{n\delta\log d}}{\sqrt{M+\Delta^2}}\right) \\ &\quad \cdot \exp\left(-\frac{0.5C_3|\sigma_{ij}^*|\sqrt{n\delta\log d}}{\sqrt{M+\Delta^2}}\right) \mathbb{1}(|\sigma_{ij}^*| > 3\mu) \\ &\stackrel{(i)}{\leq} 2 \sum_{i,j} \left(\frac{\sqrt{M+\Delta^2}}{C_3\sqrt{n\delta\log d}}\right)^p \cdot \left(\sup_{t \geq 0} t^p \exp\left(-\frac{t}{2}\right)\right) \cdot \exp\left(-\frac{3C_3\sqrt{n\delta\log d} \cdot \mu}{2\sqrt{M+\Delta^2}}\right) \\ &\stackrel{(ii)}{\leq} 2d^{2-10\delta} \left(\frac{\sqrt{M+\Delta^2}}{C_3} \sqrt{\frac{\delta}{n\log d}}\right)^p \leq 2d^{2-10\delta} (C_3^{-1}B_0)^p, \end{aligned} \quad (29)$$

$2 \exp\left(-\frac{t}{4c}\right)$. Based on this probability tail bound, we can bound the moment via integral as follows

$$\begin{aligned} &\mathbb{E}|\hat{\sigma}_{ij} - \mathbb{E}\hat{\sigma}_{ij}|^{2p} \\ &= 2p \int_0^\infty t^{2p-1} \mathbb{P}(|\hat{\sigma}_{ij} - \mathbb{E}\hat{\sigma}_{ij}| > t) dt \\ &\leq 4p \int_0^\infty t^{2p-1} \left(\exp\left(-\frac{t^2}{4v}\right) + \exp\left(-\frac{t}{4c}\right)\right) dt \\ &= 2[(4v)^p \Gamma(p+1) + (4c)^{2p} \Gamma(2p+1)] \\ &\stackrel{(i)}{\leq} 2[(4vp)^p + (8cp)^{2p}], \end{aligned}$$

where we use $\Gamma(p+1) \leq p^p$, $\Gamma(2p+1) \leq (2p)^{2p}$ in (i) under suitably large p . Thus, it follows that

$$\begin{aligned} I_{22} &\leq \sum_{i,j} 2d^{-\frac{\delta}{2}} \sqrt{(4vp)^p + (8cp)^{2p}} \\ &\leq 2d^{2-\frac{\delta}{2}} [(2\sqrt{pv})^p + (8cp)^p] \stackrel{(i)}{\leq} 2d^{2-\frac{\delta}{2}} (C_4B_0)^p, \end{aligned}$$

where (i) is due to $2\sqrt{pv} \leq (\sqrt{M} + \Delta^2)\sqrt{\frac{\delta}{n}}$ and $8cp = \frac{2\Delta^2\delta}{n} + 2\sqrt{\frac{\delta M}{n\log d}}$ (recall that $p = \frac{\delta}{4}$).

Step 2.3. Bounding I_{23} .

From Step 2 in the proof of Theorem 2 we have $|\tilde{\sigma}_{ij} - \sigma_{ij}^*| \leq C_5B_0$. This directly leads to

$$I_{23} \leq d^2 \cdot 2d^{-\delta} \cdot (C_5B_0)^p = 2d^{2-\delta} (C_5B_0)^p.$$

We are in a position to combine everything and conclude the proof. Putting all pieces into (26), it follows that $I_2 \leq d^{1-\frac{\delta}{4}} (C_6B_0)^p$. Assuming $\delta \geq 4$, such upper bound is dominated by (24) for I_1 , we can hence conclude that

$\mathbb{E}\|\hat{\Sigma}_s - \Sigma^*\|_{op}^p \leq (C_6sB_0)^p$. Therefore, by Markov's inequality,

$$\begin{aligned} &\mathbb{P}(\|\hat{\Sigma}_s - \Sigma^*\|_{op} \geq C_6esB_0) \\ &\leq \frac{\mathbb{E}\|\hat{\Sigma}_s - \Sigma^*\|_{op}^p}{(C_6esB_0)^p} \leq \exp(-p) = \exp\left(-\frac{\delta}{4}\right), \end{aligned}$$

which completes the proof. \square

2) **Quantized Compressed Sensing:** Note that our estimation procedure in QCS, QMC falls in the framework of regularized M-estimator, see [22], [36], [70] for instance. Particularly, we introduce the following deterministic result for analysing the estimator (10).

Lemma 4. (Adapted from [22, Coro. 2]). *Consider (8) and the estimator $\hat{\theta}$ defined in (10), let $\hat{\Upsilon} := \hat{\theta} - \theta^*$ be the estimation error. If \mathbf{Q} is positive semi-definite, and $\lambda \geq 2\|\mathbf{Q}\theta^* - \mathbf{b}\|_\infty$, then it holds that $\|\hat{\Upsilon}\|_1 \leq 10\sqrt{s}\|\hat{\Upsilon}\|_2$. Moreover, if for some $\kappa > 0$ we have the restricted strong convexity (RSC) $\hat{\Upsilon}^\top \mathbf{Q} \hat{\Upsilon} \geq \kappa\|\hat{\Upsilon}\|_2^2$, then we have the error bounds $\|\hat{\Upsilon}\|_2 \leq 30\sqrt{s}\left(\frac{\lambda}{\kappa}\right)$ and $\|\hat{\Upsilon}\|_1 \leq 300s\left(\frac{\lambda}{\kappa}\right)$.¹⁴*

To establish the RSC condition, a convenient way is to use the matrix deviation inequality. The following Lemma is adapted from [61], by combining Theorem 3 and Remark 1 therein.¹⁵

Lemma 5. (Adapted from [61, Thm. 3]). *Assume $\mathbf{A} \in \mathbb{R}^{n \times d}$ has independent zero-mean sub-Gaussian rows α_k^\top s satisfying $\|\alpha_k\|_{\psi_2} \leq K$, and the eigenvalues of $\Sigma := \mathbb{E}(\alpha_k \alpha_k^\top)$ are between $[\kappa_0, \kappa_1]$ for some $\kappa_1 \geq \kappa_0 > 0$. For $\mathcal{T} \subset \mathbb{R}^d$ we let*

¹⁴We do not optimize the constants in Lemmas 4, 6 for easy reference.

¹⁵The dependence on K can be further refined [49], while this is not pursued in the present paper.

$\text{rad}(\mathcal{T}) = \sup_{\mathbf{x} \in \mathcal{T}} \|\mathbf{x}\|_2$ be its radius. Then with probability at least $1 - \exp(-u^2)$, it holds that

$$\sup_{\mathbf{x} \in \mathcal{T}} \left| \|\mathbf{A}\mathbf{x}\|_2 - \sqrt{n} \|\sqrt{\Sigma}\mathbf{x}\|_2 \right| \leq \frac{C\sqrt{\kappa_1}K^2}{\kappa_0} \left(\omega(\mathcal{T}) + u \cdot \text{rad}(\mathcal{T}) \right),$$

where $\omega(\mathcal{T}) = \mathbb{E} \sup_{\mathbf{v} \in \mathcal{T}} [\mathbf{g}^\top \mathbf{v}]$ with $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ is the Gaussian width of \mathcal{T} .

Based on Lemma 4, the proofs of Theorems 5-6 are divided into two steps, i.e., showing $\lambda \geq 2\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$ and verifying the RSC. While we still have full \mathbf{x}_k in Theorems 5-6, we will study the more challenging settings where the covariates \mathbf{x}_k s are also quantized via $\mathcal{Q}_{\bar{\Delta}}(\cdot)$ in Theorems 9-10, in which we can take $\bar{\Delta} = 0$ to return the settings of Theorems 5-6. Using such perspective, for most technical ingredients (e.g., the verification of $\lambda \geq 2\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$) in the proofs of Theorems 5-6 we can simply refer to the counterparts established in the proofs of Theorems 9-10. This avoids repetition and will be explained in the proofs more clearly.

Proof of Theorem 5: *Proof.* We divide the proofs into two steps.

Step 1. Proving $\lambda \geq 2\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$

Recall that we choose $\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^\top$ and $\mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \mathbf{x}_k$. In the setting of Theorem 9, the process of obtaining \dot{y}_k remains the same, while the covariates \mathbf{x}_k s are further quantized to $\hat{\mathbf{x}}_k = \mathcal{Q}_{\bar{\Delta}}(\mathbf{x}_k + \boldsymbol{\tau}_k)$ for some $\bar{\Delta} > 0$ under triangular dither $\boldsymbol{\tau}_k \sim \mathcal{U}([-\frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d) + \mathcal{U}([-\frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d)$, and we choose

$$\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d, \quad \mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \hat{\mathbf{x}}_k$$

there. As a result, by considering $\bar{\Delta} = 0$, it can be implied by Step 1 in the proof of Theorem 9 that under the choice

$$\lambda = C_1 \frac{\sigma^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 ,

$$\lambda \geq 2 \left\| \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^\top \boldsymbol{\theta}^* - \frac{1}{n} \sum_{k=1}^n \dot{y}_k \mathbf{x}_k \right\|_\infty$$

holds with probability at least $1 - 8d^{1-\delta}$. Then, by using Lemma 4 we obtain $\|\hat{\mathbf{Y}}\|_1 \leq 10\sqrt{s}\|\hat{\mathbf{Y}}\|_2$.

Step 2. Verifying the RSC $\hat{\mathbf{Y}}^\top \mathbf{Q} \hat{\mathbf{Y}} \geq \kappa \|\hat{\mathbf{Y}}\|_2^2$

We refer to Step 2 in the proof of Theorem 9. In particular, with the choices $\bar{\Delta} = 0$ and $\mathbf{v} = \hat{\mathbf{Y}}$ in (38), combined with $\|\hat{\mathbf{Y}}\|_1 \leq 10\sqrt{s}\|\hat{\mathbf{Y}}\|_2$, we obtain

$$\begin{aligned} \frac{1}{\sqrt{n}} \|\mathbf{X}\hat{\Delta}\|_2 &\geq \sqrt{\kappa_0} \|\hat{\mathbf{Y}}\|_2 - \frac{C_2 \sqrt{\kappa_1} \sigma^2}{\kappa_0} \sqrt{\frac{\delta s \log d}{n}} \|\hat{\mathbf{Y}}\|_2 \\ &\geq \frac{1}{2} \sqrt{\kappa_0} \|\hat{\mathbf{Y}}\|_2, \end{aligned}$$

where the last inequality is due to the assumed scaling $n \gtrsim \delta s \log d$. With these preparations, a direct application of Lemma 4 completes the proof. \square

Proof of Theorem 6: *Proof.* The proof is similarly based on Lemma 4.

Step 1. Proving $\lambda \geq 2\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$

Recall that we choose $\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top$ and $\mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \tilde{\mathbf{x}}_k$. In the setting of Theorem 10, the process of obtaining \dot{y}_k remains the same, while the truncated covariates $\tilde{\mathbf{x}}_k$ s are further quantized to $\hat{\mathbf{x}}_k = \mathcal{Q}_{\bar{\Delta}}(\tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k)$ for some $\bar{\Delta} \geq 0$ under triangular dither $\boldsymbol{\tau}_k \sim \mathcal{U}([-\frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d) + \mathcal{U}([-\frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d)$, and we choose

$$\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d, \quad \mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \hat{\mathbf{x}}_k$$

there. As a result, by considering $\bar{\Delta} = 0$, it can be implied by step 1 in the proof of Theorem 10 that, our choice

$$\lambda = C_1 (R\sqrt{M} + \Delta^2) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 ensures $\lambda \geq 2\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$ with the promised probability. By Lemma 4 we obtain $\|\hat{\mathbf{Y}}\|_1 \leq 10\sqrt{s}\|\hat{\mathbf{Y}}\|_2$.

Step 2. Verifying the RSC $\hat{\mathbf{Y}}^\top \mathbf{Q} \hat{\mathbf{Y}} \geq \kappa \|\hat{\mathbf{Y}}\|_2^2$

Unlike the case of sub-Gaussian covariate that is based on matrix deviation inequality (Lemma 5), here we establish a lower bound for $\hat{\mathbf{Y}}^\top \mathbf{Q} \hat{\mathbf{Y}}$ using the bound on $\|\mathbf{Q} - \Sigma^*\|_\infty$ (Theorem 2). Specifically, setting $\Delta = 0$ in Theorem 2 yields that, $\|\mathbf{Q} - \Sigma^*\|_\infty \lesssim \sqrt{\frac{\delta M \log d}{n}}$ holds with probability at least $1 - 2d^{2-\delta}$, which allows us to proceed as follows:

$$\begin{aligned} \hat{\mathbf{Y}}^\top \mathbf{Q} \hat{\mathbf{Y}} &= \hat{\mathbf{Y}}^\top \Sigma^* \hat{\mathbf{Y}} - \hat{\mathbf{Y}}^\top (\Sigma^* - \mathbf{Q}) \hat{\mathbf{Y}} \\ &\stackrel{(i)}{\geq} \kappa_0 \|\hat{\mathbf{Y}}\|_2^2 - \sqrt{\frac{\delta M \log d}{n}} \|\hat{\mathbf{Y}}\|_1^2 \\ &\stackrel{(ii)}{\geq} \left(\kappa_0 - C_6 s \sqrt{\frac{\delta M \log d}{n}} \right) \|\hat{\mathbf{Y}}\|_2^2 \\ &\stackrel{(iii)}{\geq} \frac{\kappa_0}{2} \|\hat{\mathbf{Y}}\|_2^2, \end{aligned} \tag{30}$$

where (i) is because $\hat{\mathbf{Y}}^\top (\Sigma^* - \mathbf{Q}) \hat{\mathbf{Y}} \leq \|\hat{\mathbf{Y}}\|_1^2 \|\mathbf{Q} - \Sigma^*\|_\infty$, (ii) is due to $\|\hat{\mathbf{Y}}\|_1 \leq 10\sqrt{s}\|\hat{\mathbf{Y}}\|_2$, (iii) is due to the the assumed scaling $n \gtrsim \delta s^2 \log d$. Now the desired results follow immediately from Lemma 4. \square

3) Quantized Matrix Completion: Under the observation model (11), we first provide a deterministic framework for analysing the estimator (12).

Lemma 6. (Adapted from [22, Coro. 3]). *Let $\hat{\mathbf{Y}} := \hat{\Theta} - \Theta^*$. If*

$$\lambda \geq 2 \left\| \frac{1}{n} \sum_{k=1}^n \langle \mathbf{X}_k, \Theta^* \rangle - \dot{y}_k \mathbf{X}_k \right\|_{op}, \tag{31}$$

then it holds that $\|\hat{\mathbf{Y}}\|_{nu} \leq 10\sqrt{r}\|\hat{\mathbf{Y}}\|_F$. Moreover, if for some $\kappa > 0$ we have the restricted strong convexity (RSC) $\frac{1}{n} \sum_{k=1}^n \langle \mathbf{X}_k, \hat{\mathbf{Y}} \rangle^2 \geq \kappa \|\hat{\mathbf{Y}}\|_F^2$, then we have the error bounds $\|\hat{\mathbf{Y}}\|_F \leq 30\sqrt{r}(\frac{\lambda}{\kappa})$ and $\|\hat{\mathbf{Y}}\|_{nu} \leq 300r(\frac{\lambda}{\kappa})$.

Clearly, to derive statistical error rate of $\hat{\Theta}$ from Lemma 6, the key ingredients are (31) and the RSC. Specialized to the

covariate $\mathbf{X}_k \sim \mathcal{U}(\{e_i e_j^\top : i, j \in [d]\})$ in matrix completion, we will use the following lemma to establish RSC.

Lemma 7. (Adapted from [22, Lem. 4] with $q = 0$). *Given some $\alpha > 0, \delta > 0$, we define the constraint set ψ with sufficiently large ψ as*

$$\mathcal{C}(\psi) = \left\{ \Theta \in \mathbb{R}^{d \times d} : \begin{aligned} &\|\Theta\|_\infty \leq 2\alpha, \\ &\|\Theta\|_{nu} \leq 10\sqrt{r}\|\Theta\|_F, \\ &\|\Theta\|_F^2 \geq (\alpha d)^2 \sqrt{\frac{\psi \delta \log d}{n}} \end{aligned} \right\}. \quad (32)$$

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. uniformly distributed on $\{e_i e_j^\top : i, j \in [d]\}$, then there exist absolute constants $\kappa \in (0, 1)$ and C , such that with probability at least $1 - d^{-\delta}$ we have

$$\frac{1}{n} \sum_{k=1}^n |\langle \mathbf{X}_k, \Theta \rangle|^2 \geq \frac{\kappa \|\Theta\|_F^2}{d^2} - \frac{C \alpha^2 r d \log d}{n}, \quad \forall \Theta \in \mathcal{C}(\psi). \quad (33)$$

Matrix completion with sub-exponential noise was studied in [51], and we make use of the following Lemma in the sub-exponential case.

Lemma 8. (Adapted from [51, Lem. 5]). *Given some $\delta > 0$. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. uniformly distributed on $\{e_i e_j^\top : i, j \in [d]\}$, independent of \mathbf{X}_k s, $\epsilon_1, \dots, \epsilon_n$ are i.i.d. zero-mean and satisfy $\|\epsilon_k\|_{\psi_1} \leq \sigma$. If $n \gtrsim \delta d \log^3 d$, with probability at least $1 - d^{-\delta}$ we have*

$$\left\| \frac{1}{n} \sum_{k=1}^n \epsilon_k \mathbf{X}_k \right\|_{op} \leq \sigma \sqrt{\frac{\delta \log d}{nd}}.$$

Proof of Theorem 7: *Proof.* We divide the proof into two steps.

Step 1. Proving (31)

Defining $w_k := \dot{y}_k - y_k - \tau_k$ as the quantization error, from Theorem 1(a) we know that w_k s are independent of \mathbf{X}_k and i.i.d. uniformly distributed on $[-\frac{\Delta}{2}, \frac{\Delta}{2}]$. Thus, we can further write that

$$\dot{y}_k = y_k + \tau_k + w_k = \langle \mathbf{X}_k, \Theta^* \rangle + \epsilon_k + \tau_k + w_k,$$

which allows us to decompose I into

$$\begin{aligned} I &\leq \left\| \frac{1}{n} \sum_{k=1}^n \epsilon_k \mathbf{X}_k \right\|_{op} + \left\| \frac{1}{n} \sum_{k=1}^n \tau_k \mathbf{X}_k \right\|_{op} + \left\| \frac{1}{n} \sum_{k=1}^n w_k \mathbf{X}_k \right\|_{op} \\ &:= I_1 + I_2 + I_3. \end{aligned}$$

Because ϵ_k s are independent of \mathbf{X}_k s and i.i.d. sub-exponential noise satisfying $\|\epsilon_k\|_{\psi_1} \leq \sigma$, under the scaling $n \gtrsim \delta d \log^3 d$, Lemma 8 implies that $I_1 \lesssim \sigma \sqrt{\frac{\delta \log d}{nd}}$ holds with probability at least $1 - d^{-\delta}$. Analogously, τ_k s and w_k s are independent of $\{\mathbf{X}_k : k \in [n]\}$ and are i.i.d. uniformly distributed on $[-\frac{\Delta}{2}, \frac{\Delta}{2}]$, Lemma 8 also applies to I_2 and I_3 , yielding that with the promised probability $I_2 + I_3 \lesssim \Delta \sqrt{\frac{\delta \log d}{nd}}$. Taken collectively, $I \lesssim (\sigma + \Delta) \sqrt{\frac{\delta \log d}{nd}}$, so setting $\lambda = C_1(\sigma + \Delta) \sqrt{\frac{\delta \log d}{nd}}$ with sufficiently large C_1 ensures $\lambda \geq 2I$, with

probability at least $1 - 3d^{-\delta}$. Further, Lemma 6 gives $\|\hat{\mathbf{Y}}\|_{nu} \leq 10\sqrt{r}\|\hat{\mathbf{Y}}\|_F$.

Step 2. Verifying RSC

First note that $\|\hat{\mathbf{Y}}\|_\infty \leq \|\hat{\Theta}\|_\infty + \|\Theta^*\|_\infty \leq 2\alpha$; and as proved before, $\|\hat{\mathbf{Y}}\|_{nu} \leq 10\sqrt{r}\|\hat{\mathbf{Y}}\|_F$. To proceed we define the constraint set $\mathcal{C}(\psi)$ as in (32) with some properly chosen constant ψ . Then using Lemma 7, for some absolute constants κ, C , (33) holds with probability at least $1 - d^{-\delta}$. Then we discuss several cases.

1) If $\hat{\mathbf{Y}} \notin \mathcal{C}(\psi)$, because $\hat{\mathbf{Y}}$ satisfies the first two constraints in the definition of $\mathcal{C}(\psi)$, it must violate the third constraint and satisfy $\|\hat{\mathbf{Y}}\|_F^2 \leq (\alpha d)^2 \sqrt{\frac{\psi \delta \log d}{n}}$, which gives

$$\|\hat{\mathbf{Y}}\|_F \lesssim \alpha d \left(\frac{\delta \log d}{n}\right)^{1/4} \stackrel{(i)}{\lesssim} \alpha d \sqrt{\frac{\delta r d \log d}{n}},$$

as desired. Note that (i) is due to the scaling $n \lesssim \delta r^2 d^2 \log d$.

2) If $\hat{\mathbf{Y}} \in \mathcal{C}(\psi)$, (33) implies that $\frac{1}{n} \sum_{k=1}^n |\langle \mathbf{X}_k, \hat{\mathbf{Y}} \rangle|^2 \geq \frac{\kappa \|\hat{\mathbf{Y}}\|_F^2}{d^2} - C \frac{\alpha^2 r d \log d}{n}$, and we further consider the following two cases.

2.1) If $C \frac{\alpha^2 r d \log d}{n} \geq \frac{\kappa \|\hat{\mathbf{Y}}\|_F^2}{2d^2}$, we have $\|\hat{\mathbf{Y}}\|_F \lesssim \alpha d \sqrt{\frac{r d \log d}{n}}$, as desired.

2.2) If $C \frac{\alpha^2 r d \log d}{n} < \frac{\kappa \|\hat{\mathbf{Y}}\|_F^2}{2d^2}$, then the RSC condition holds: $\frac{1}{n} \sum_{k=1}^n |\langle \mathbf{X}_k, \hat{\mathbf{Y}} \rangle|^2 \geq \frac{\kappa \|\hat{\mathbf{Y}}\|_F^2}{2d^2}$. This allows us to apply Lemma 6 to obtain $\|\hat{\mathbf{Y}}\|_F \lesssim (\sigma + \Delta) d \sqrt{\frac{\delta r d \log d}{n}}$.

Thus, in any case, we have shown $\|\hat{\mathbf{Y}}\|_F = O((\alpha + \sigma + \Delta) d \sqrt{\frac{\delta r d \log d}{n}})$. The bound on $\|\hat{\mathbf{Y}}\|_{nu}$ follows immediately from $\|\hat{\mathbf{Y}}\|_{nu} \leq 10\sqrt{r}\|\hat{\mathbf{Y}}\|_F$. The proof is complete. \square

Proof of Theorem 8: *Proof.* The proof is based on Lemma 6 and divided into two steps.

Step 1. Proving (31)

Recall that the quantization error $w_k := \dot{y}_k - \tilde{y}_k - \tau_k$ is zero-mean and independent of \mathbf{X}_k (Theorem 1(a)), thus we have $\mathbb{E}(\dot{y}_k \mathbf{X}_k) = \mathbb{E}(\tilde{y}_k \mathbf{X}_k) + \mathbb{E}(\tau_k \mathbf{X}_k) + \mathbb{E}(w_k \mathbf{X}_k) = \mathbb{E}(\tilde{y}_k \mathbf{X}_k)$. Combining with $\mathbb{E}(\langle \mathbf{X}_k, \Theta^* \rangle \mathbf{X}_k) = \mathbb{E}(y_k \mathbf{X}_k)$, triangle inequality can first decompose the target term into

$$\begin{aligned} &\left\| \frac{1}{n} \sum_{k=1}^n (\dot{y}_k - \langle \mathbf{X}_k, \Theta^* \rangle) \mathbf{X}_k \right\|_{op} \\ &\leq \left\| \frac{1}{n} \sum_{k=1}^n \dot{y}_k \mathbf{X}_k - \mathbb{E}(\dot{y}_k \mathbf{X}_k) \right\|_{op} + \left\| \mathbb{E}(y_k \mathbf{X}_k - \tilde{y}_k \mathbf{X}_k) \right\|_{op} \\ &\quad + \left\| \frac{1}{n} \sum_{k=1}^n \langle \mathbf{X}_k, \Theta^* \rangle \mathbf{X}_k - \mathbb{E}(\langle \mathbf{X}_k, \Theta^* \rangle \mathbf{X}_k) \right\|_{op} \\ &:= I_1 + I_2 + I_3. \end{aligned}$$

Step 1.1. Bounding I_1 and I_3

We write $I_1 = \left\| \sum_{k=1}^n \mathbf{S}_k \right\|_{op}$ and $I_3 = \left\| \sum_{k=1}^n \mathbf{W}_k \right\|_{op}$ by defining

$$\begin{aligned} \mathbf{S}_k &= \frac{1}{n} (\dot{y}_k \mathbf{X}_k - \mathbb{E}(\dot{y}_k \mathbf{X}_k)), \\ \mathbf{W}_k &= \frac{1}{n} (\langle \mathbf{X}_k, \Theta^* \rangle \mathbf{X}_k - \mathbb{E}[\langle \mathbf{X}_k, \Theta^* \rangle \mathbf{X}_k]). \end{aligned}$$

By $|\dot{y}_k| \leq |\tilde{y}_k| + |\tau_k| + |w_k| \leq \zeta_y + \Delta$ we have

$$\begin{aligned} \|\mathbf{S}_k\|_{op} &\leq \frac{1}{n} \|\dot{y}_k \mathbf{X}_k\|_{op} + \frac{1}{n} \|\mathbb{E}(\dot{y}_k \mathbf{X}_k)\|_{op} \\ &\leq \frac{1}{n} \|\dot{y}_k \mathbf{X}_k\|_{op} + \frac{1}{n} \mathbb{E} \|\dot{y}_k \mathbf{X}_k\|_{op} \leq \frac{2(\zeta_y + \Delta)}{n}. \end{aligned}$$

Analogously, we have $\|\mathbf{W}_k\|_{op} \leq \frac{2\alpha}{n}$ since $|\langle \mathbf{X}_k, \Theta^\star \rangle| \leq \|\Theta^\star\|_\infty \leq \alpha$. In addition, by

$$\|\mathbb{E}\{(\mathbf{A} - \mathbb{E}\mathbf{A})^\top (\mathbf{A} - \mathbb{E}\mathbf{A})\}\|_{op} \leq \|\mathbb{E}(\mathbf{A}^\top \mathbf{A})\|_{op} \quad (\forall \mathbf{A})$$

and the simple fact $\mathbb{E}(\mathbf{X}_k \mathbf{X}_k^\top) = \mathbb{E}(\mathbf{X}_k^\top \mathbf{X}_k) = \mathbf{I}_d/d$, we estimate the matrix variance statistic as follows

$$\begin{aligned} &\left\| \sum_{k=1}^n \mathbb{E}(\mathbf{S}_k \mathbf{S}_k^\top) \right\|_{op} = n \|\mathbb{E}(\mathbf{S}_k \mathbf{S}_k^\top)\|_{op} \\ &\leq \frac{1}{n} \|\mathbb{E}(\dot{y}_k^2 \mathbf{X}_k \mathbf{X}_k^\top)\|_{op} = \frac{1}{n} \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}(\dot{y}_k^2 \cdot \|\mathbf{X}_k^\top \mathbf{v}\|_2^2) \\ &= \frac{1}{n} \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}_{\mathbf{X}_k} \left(\left[\mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\dot{y}_k^2) \right] \|\mathbf{X}_k^\top \mathbf{v}\|_2^2 \right) \\ &\stackrel{(i)}{\leq} \frac{4}{n} (\alpha^2 + M + \Delta^2) \sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbb{E}_{\mathbf{X}_k} \|\mathbf{X}_k^\top \mathbf{v}\|_2^2 \\ &\leq \frac{4(\alpha^2 + M + \Delta^2)}{nd}, \end{aligned}$$

where (i) is because given \mathbf{X}_k we can estimate $\mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\dot{y}_k^2) \leq 2(\mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\dot{y}_k) + \Delta^2)$ since $|\dot{y}_k - \tilde{y}_k| \leq \Delta$, and moreover we have

$$\begin{aligned} \mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\dot{y}_k^2) &\leq \mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\tilde{y}_k^2) \\ &\leq 2(\mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\langle \mathbf{X}_k, \Theta^\star \rangle)^2) + \mathbb{E}_{\dot{y}_k | \mathbf{X}_k}(\epsilon_k^2) \\ &\leq 2(\alpha^2 + M). \end{aligned}$$

It is not hard to see that this bound remains valid for $\|\sum_{k=1}^n \mathbb{E}(\mathbf{S}_k^\top \mathbf{S}_k)\|_{op}$. Also, by similar arguments one can prove

$$\max \left\{ \left\| \sum_{k=1}^n \mathbb{E}(\mathbf{W}_k^\top \mathbf{W}_k) \right\|_{op}, \left\| \sum_{k=1}^n \mathbb{E}(\mathbf{W}_k \mathbf{W}_k^\top) \right\|_{op} \right\} \leq \frac{\alpha^2}{nd}.$$

Thus, Matrix Bernstein's inequality (Lemma 2) gives

$$\begin{aligned} \mathbb{P}(I_1 \geq t) &\leq 2d \cdot \exp \left(- \frac{C_4 n d t^2}{(\alpha^2 + M + \Delta^2) + (\zeta_y + \Delta) d t} \right) \\ \mathbb{P}(I_3 \geq t) &\leq 2d \cdot \exp \left(- \frac{C_5 n d t^2}{\alpha^2 + \alpha d t} \right) \end{aligned}$$

Thus, setting $t = C_6(\alpha + \sqrt{M} + \Delta) \sqrt{\frac{\delta \log d}{nd}}$ in the two inequalities above with sufficiently large C_6 , combining with the scaling that $\sqrt{\frac{\delta d \log d}{n}} = O(1)$, we obtain

$$I_1 + I_3 \lesssim (\alpha + \sqrt{M} + \Delta) \sqrt{\frac{\delta \log d}{nd}}$$

with probability at least $1 - 4d^{1-\delta}$.

Step 1.2. Bounding I_2

Let us bound $\|\mathbb{E}((y_k - \tilde{y}_k) \mathbf{X}_k)\|_\infty$ first. Write (i, j) -th entry of \mathbf{X}_k as $x_{k,ij}$, then for given (i, j) , $\mathbb{P}(x_{k,ij} = 1) = d^{-2}$,

$x_{k,ij} = 0$ otherwise. We can thus proceed by the following estimations:

$$\begin{aligned} &|\mathbb{E}((y_k - \tilde{y}_k) x_{k,ij})| \\ &= |\mathbb{E}((y_k - \tilde{y}_k) x_{k,ij} \mathbb{1}(|y_k| \geq \zeta_y))| \\ &\leq \mathbb{E}(|y_k| x_{k,ij} \mathbb{1}(|y_k| \geq \zeta_y)) \\ &= \mathbb{E}_{x_{k,ij}} \left(\left\{ \mathbb{E}_{y_k | x_{k,ij}} |y_k| \mathbb{1}(|y_k| \geq \zeta_y) \right\} x_{k,ij} \right) \\ &= d^{-2} \mathbb{E}_{y_k | x_{k,ij}=1} (|y_k| \mathbb{1}(|y_k| \geq \zeta_y)) \\ &\stackrel{(i)}{\leq} d^{-2} \sqrt{\mathbb{E}_{y_k \sim \theta_{ij}^* + \epsilon_k} (y_k^2)} \sqrt{\mathbb{P}_{y_k \sim \theta_{ij}^* + \epsilon_k} (y_k^2 \geq \zeta_y^2)} \\ &\stackrel{(ii)}{\leq} d^{-2} \frac{\alpha^2 + M}{\zeta_y} \lesssim \frac{\alpha + \sqrt{M}}{d^2} \sqrt{\frac{\delta d \log d}{n}}, \end{aligned}$$

where (i), (ii) is by Cauchy-Schwarz and Markov's, respectively. Since this holds for any (i, j) , we obtain

$$\|\mathbb{E}((y_k - \tilde{y}_k) \mathbf{X}_k)\|_\infty = O((\alpha + \sqrt{M}) d^{-2} \sqrt{\frac{\delta d \log d}{n}}).$$

which further gives

$$I_2 = O((\alpha + \sqrt{M}) \sqrt{\frac{\delta \log d}{nd}})$$

by using $\|\mathbf{A}\|_{op} \leq d \|\mathbf{A}\|_\infty$ ($\forall \mathbf{A} \in \mathbb{R}^{d \times d}$). Putting pieces together, with probability at least $1 - 4d^{1-\delta}$ we have

$$\left\| \frac{1}{n} \sum_k (\dot{y}_k - \langle \mathbf{X}_k, \Theta^\star \rangle) \mathbf{X}_k \right\|_{op} \lesssim (\alpha + \sqrt{M} + \Delta) \sqrt{\frac{\delta \log d}{nd}},$$

hence $\lambda = C_1(\alpha + \sqrt{M} + \Delta) \sqrt{\frac{\delta \log d}{nd}}$ ensures (31) under the same probability. Further, Lemma 6 gives $\|\hat{\mathbf{Y}}\|_{nu} \leq 10\sqrt{r} \|\hat{\mathbf{Y}}\|_F$.

Step 2. Verifying RSC

The remaining part is almost the same as Step 2 in the proof of Theorem 7 — defining the constraint set $\mathcal{C}(\psi)$ as (32) and then discussing several cases based on whether $\hat{\mathbf{Y}} \in \mathcal{C}(\psi)$ holds. Thus, we conclude the proof without providing the details. \square

B. Proofs in Section IV

This appendix collects the proofs in Section IV concerning covariate quantization and uniform signal recovery in QCS.

1) Covariate Quantization: Because of the non-convexity, the proofs in this part can no longer be based on Lemma 4. Indeed, bounding the estimation errors of $\tilde{\theta}$ s satisfying (14) require more tedious manipulations essentially due to the additional ℓ_1 constraint (induced by the constraint \mathcal{S} in (14)).

Proof of Theorem 9: *Proof.* The proof is divided into three steps — the first two steps resemble the previous proofs that are based on Lemma 4, while we bound the estimation errors in the last step.

Step 1. Proving $\lambda \geq \beta \|\mathbf{Q} \Theta^\star - \mathbf{b}\|_\infty$ for some pre-specified $\beta > 2$

Recall that (\mathbf{Q}, \mathbf{b}) are constructed from the quantized data as $\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d$ and $\mathbf{b} = \frac{1}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_k$. We will show that,

$$\lambda = C_1 \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\log d}{n}}$$

guarantees

$$\lambda \geq \beta \left\| \frac{1}{n} \sum_{k=1}^n \left(\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d \right) \boldsymbol{\theta}^\star - \frac{1}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_k \right\|_\infty$$

holds with the promised probability, where $\beta > 2$ is any pre-specified constant. Recall the notation we introduced: $\dot{y}_k = \tilde{y}_k + \phi_k + \vartheta_k$ with the quantization error $\vartheta_k \sim \mathcal{U}([- \frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}])$ being independent of \tilde{y}_k , $\dot{\mathbf{x}}_k = \mathbf{x}_k + \boldsymbol{\tau}_k + \mathbf{w}_k$ with the quantization error $\mathbf{w}_k \sim \mathcal{U}([- \frac{\bar{\Delta}}{2}, \frac{\bar{\Delta}}{2}]^d)$ being independent of \mathbf{x}_k . Combining with the assumptions that the dithers are independent of (\mathbf{x}_k, y_k) and that ϕ_k s and $\boldsymbol{\tau}_k$ s are independent, we have

$$\begin{aligned} \mathbb{E}(\dot{y}_k \dot{\mathbf{x}}_k) &= \mathbb{E}((\tilde{y}_k + \phi_k + \vartheta_k)(\mathbf{x}_k + \boldsymbol{\tau}_k + \mathbf{w}_k)) = \mathbb{E}(\tilde{y}_k \mathbf{x}_k), \\ \mathbb{E}\left(\left[\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d\right] \boldsymbol{\theta}^\star\right) &= \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top \boldsymbol{\theta}^\star) = \mathbb{E}(y_k \mathbf{x}_k), \end{aligned} \quad (34)$$

which allows us to decompose the target term as two concentration terms (I_1, I_3) and a bias term (I_2)

$$\begin{aligned} &\left\| \frac{1}{n} \sum_{k=1}^n \left[\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d \right] \boldsymbol{\theta}^\star - \frac{1}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_k \right\|_\infty \\ &\leq \left\| \frac{1}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_k - \mathbb{E}(\dot{y}_k \dot{\mathbf{x}}_k) \right\|_\infty + \left\| \mathbb{E}(y_k \mathbf{x}_k - \tilde{y}_k \mathbf{x}_k) \right\|_\infty \\ &+ \left\| \frac{1}{n} \sum_{k=1}^n \dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top \boldsymbol{\theta}^\star - \mathbb{E}(\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top \boldsymbol{\theta}^\star) \right\|_\infty := I_1 + I_2 + I_3. \end{aligned}$$

Step 1.1. Bounding I_1

Denote the i -th entry of $\mathbf{x}_k, \dot{\mathbf{x}}_k, \boldsymbol{\tau}_k, \mathbf{w}_k$ by $x_{ki}, \dot{x}_{ki}, \tau_{ki}, w_{ki}$, respectively. For I_1 , the i -th entry reads $\frac{1}{n} \sum_{k=1}^n \dot{y}_k \dot{x}_{ki} - \mathbb{E}(\dot{y}_k \dot{x}_{ki})$. By using the relations

$$|\dot{y}_k| \leq |\tilde{y}_k| + |\phi_k| + |\vartheta_k| \leq \zeta_y + \Delta,$$

and

$$\|\dot{\mathbf{x}}_k\|_{\psi_2} \leq \|\mathbf{x}_k\|_{\psi_2} + \|\boldsymbol{\tau}_k\|_{\psi_2} + \|\mathbf{w}_k\|_{\psi_2} \lesssim \sigma + \bar{\Delta}$$

and

$$\mathbb{E}|\dot{y}_k|^{2l} \lesssim \mathbb{E}|\tilde{y}_k|^{2l} + \mathbb{E}|\phi_k + \vartheta_k|^{2l} \lesssim M + \Delta^{2l},$$

for any integer $q \geq 2$ we can bound that

$$\begin{aligned} &\sum_{k=1}^n \mathbb{E} \left| \frac{\dot{y}_k \dot{x}_{ki}}{n} \right|^q \\ &\leq \frac{(\zeta_y + \Delta)^{q-2}}{n^q} \sum_{k=1}^n \mathbb{E} |\dot{y}_k^2 \dot{x}_{ki}^q| \\ &\stackrel{(i)}{\leq} \frac{(\zeta_y + \Delta)^{q-2}}{n^q} \sum_{k=1}^n \left\{ \mathbb{E} |\dot{y}_k|^{2l} \right\}^{\frac{1}{l}} \left\{ \mathbb{E} |\dot{x}_{ki}|^{\frac{lq}{l-1}} \right\}^{1-\frac{1}{l}} \quad (35) \\ &\stackrel{(ii)}{\lesssim} \left(\frac{(\sigma + \bar{\Delta})(\zeta_y + \Delta)}{n} \right)^{q-2} \\ &\quad \cdot \left(\frac{(\sigma + \bar{\Delta})^2 (M^{\frac{1}{l}} + \Delta^2)}{n} \right) \left(\sqrt{\frac{lq}{l-1}} \right)^q; \end{aligned}$$

combining with Stirling's approximation and treating l as absolute constant, this provides

$$\begin{aligned} \sum_{k=1}^n \mathbb{E} \left| \frac{\dot{y}_k \dot{x}_{ki}}{n} \right|^q &\leq \frac{q!}{2} v_0 c_0^{q-2} \\ \text{where } v_0 &= O\left(\frac{(\sigma + \bar{\Delta})^2 (M^{1/l} + \Delta^2)}{n}\right), \\ c_0 &= O\left(\frac{(\sigma + \bar{\Delta})(\zeta_y + \Delta)}{n}\right). \end{aligned}$$

In (35), (i) is due to Holder's, and in (ii) we use the moment constraint of sub-Gaussian variable (2). With these preparations, we invoke Bernstein's inequality (Lemma 1) and then a union bound over $i \in [d]$ to obtain

$$\begin{aligned} \mathbb{P}\left(I_1 \lesssim (\sigma + \bar{\Delta})(M^{\frac{1}{2l}} + \Delta) \sqrt{\frac{t}{n}} + \frac{(\sigma + \bar{\Delta})(\zeta_y + \Delta)t}{n}\right) \\ \geq 1 - 2d \cdot \exp(-t), \end{aligned}$$

Thus, taking $t = \delta \log d$ and plug in $\zeta_y \asymp \sqrt{\frac{nM^{1/l}}{\delta \log d}}$, we obtain

$$\mathbb{P}\left(I_1 \lesssim (\sigma + \bar{\Delta})(M^{1/(2l)} + \Delta) \sqrt{\frac{\delta \log d}{n}}\right) \geq 1 - 2d^{1-\delta}.$$

Step 1.2. Bounding I_2

Moreover, we estimate the i -th entry of I_2 by

$$\begin{aligned} &|\mathbb{E}((y_k - \tilde{y}_k)x_{ki})| \leq \mathbb{E}|y_k x_{ki}| \mathbb{1}(|y_k| \geq \zeta_y) \\ &\stackrel{(i)}{\leq} (\mathbb{E}|y_k|^{\frac{2l}{2l-1}} |x_{ki}|^{\frac{2l}{2l-1}})^{1-\frac{1}{2l}} (\mathbb{P}(|y_k| \geq \zeta_y))^{\frac{1}{2l}} \\ &\stackrel{(ii)}{\leq} \left([\mathbb{E}|y_k|^{2l}]^{\frac{1}{2l-1}} [\mathbb{E}|x_{ki}|^{\frac{l}{l-1}}]^{\frac{2l-2}{2l-1}} \right)^{1-\frac{1}{2l}} \left(\mathbb{P}(|y_k|^{2l} \geq \zeta_y^{2l}) \right)^{\frac{1}{2l}} \\ &\stackrel{(iii)}{\leq} \frac{\sigma M^{1/l}}{\zeta_y} \lesssim \sigma M^{\frac{1}{2l}} \sqrt{\frac{\delta \log d}{n}}, \end{aligned} \quad (36)$$

where (i), (ii) are due to Holder's, (iii) is due to Markov's. Since this holds for $i \in [d]$, it gives $I_2 \lesssim \sigma M^{1/(2l)} \sqrt{\frac{\delta \log d}{n}}$.

Step 1.3. Bounding I_3

We first derive a bound for $\|\boldsymbol{\theta}^\star\|_2$ that is implicitly implied by other conditions:

$$\begin{aligned} M^{1/l} &\geq \mathbb{E}|y_k|^2 \geq \mathbb{E}(\mathbf{x}_k^\top \boldsymbol{\theta}^\star)^2 = (\boldsymbol{\theta}^\star)^\top \boldsymbol{\Sigma}^\star \boldsymbol{\theta}^\star \geq \kappa_0 \|\boldsymbol{\theta}^\star\|_2^2 \\ \implies \|\boldsymbol{\theta}^\star\|_2 &= O\left(\frac{M^{1/(2l)}}{\sqrt{\kappa_0}}\right). \end{aligned}$$

Hence, we can estimate

$$\begin{aligned} \|(\dot{\mathbf{x}}_k^\top \boldsymbol{\theta}^\star) \dot{x}_{ki}\|_{\psi_1} &\leq \|\dot{\mathbf{x}}_k^\top \boldsymbol{\theta}^\star\|_{\psi_2} \|\dot{x}_{ki}\|_{\psi_2} \\ &\leq \|\dot{\mathbf{x}}_k\|_{\psi_2}^2 \|\boldsymbol{\theta}^\star\|_2 \lesssim (\sigma + \bar{\Delta})^2 \frac{M^{1/(2l)}}{\sqrt{\kappa_0}}. \end{aligned}$$

Then, we invoke Bernstein's inequality regarding the independent sum of sub-exponential random variables (e.g., [92, Thm.

2.8.1)]¹⁶ to obtain that for any $t \geq 0$ we have

$$\begin{aligned} & \mathbb{P} \left(\left| \frac{1}{n} \sum_{k=1}^n (\hat{\mathbf{x}}_k^\top \boldsymbol{\theta}^*) \dot{x}_{ki} - \mathbb{E}\{(\hat{\mathbf{x}}_k^\top \boldsymbol{\theta}^*) \dot{x}_{ki}\} \right| \geq t \right) \\ & \leq 2 \exp \left(-C_5 n \min \left\{ \frac{\sqrt{\kappa_0} t}{(\sigma + \bar{\Delta})^2 M^{\frac{1}{2l}}}, \left(\frac{\sqrt{\kappa_0} t}{(\sigma + \bar{\Delta})^2 M^{\frac{1}{2l}}} \right)^2 \right\} \right) \end{aligned}$$

Hence, we can set $t = C_6 \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} M^{1/(2l)} \sqrt{\frac{\delta \log d}{n}}$ with sufficiently large C_6 , and further apply union bound over $i \in [d]$, under the scaling that $\frac{\delta \log d}{n}$ is small enough, we obtain

$$I_3 \lesssim \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} M^{1/(2l)} \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 2d^{1-\delta}$.

Putting pieces together, since $\kappa_0 \lesssim \sigma^2$, it is immediate that

$$I \lesssim \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 8d^{1-\delta}$. Since

$$\lambda = C_1 \frac{(\sigma + \bar{\Delta})^2}{\sqrt{\kappa_0}} (\Delta + M^{1/(2l)}) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_1 , $\lambda \geq \beta \cdot \|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$ holds w.h.p.

Step 2. Verifying RSC

We provide a lower bound for $\mathbf{v}^\top \mathbf{Q}\mathbf{v} = \frac{1}{n} \|\dot{\mathbf{X}}\mathbf{v}\|_2^2 - \frac{\bar{\Delta}^2}{4} \|\mathbf{v}\|_2^2$ by using the matrix deviation inequality (Lemma 5). First note that the rows of $\dot{\mathbf{X}}$ are sub-Gaussian $\|\dot{\mathbf{x}}_k\|_{\psi_2} \lesssim \sigma + \bar{\Delta}$. Since $\mathbb{E}(\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top) = \boldsymbol{\Sigma}^* + \frac{\bar{\Delta}^2}{4} \mathbf{I}_d$, all eigenvalues of $\dot{\boldsymbol{\Sigma}} := \mathbb{E}(\dot{\mathbf{x}}_k \dot{\mathbf{x}}_k^\top)$ are between $[\kappa_0 + \frac{1}{4}\bar{\Delta}^2, \kappa_1 + \frac{1}{4}\bar{\Delta}^2]$. Thus, we invoke Lemma 5 for $\mathcal{T} := \{\mathbf{v} \in \mathbb{R}^d : \|\mathbf{v}\|_1 = 1\}$ with $u = \sqrt{\delta \log d}$; due to the well-known Gaussian width estimate $\omega(\mathcal{T}) \lesssim \sqrt{\log d}$ [92, Example 7.5.9], with probability at least $1 - d^{-\delta}$ the following event holds

$$\begin{aligned} & \sup_{\|\mathbf{v}\|_1=1} \left| \|\dot{\mathbf{X}}\mathbf{v}\|_2 - \sqrt{n} \|\dot{\boldsymbol{\Sigma}}^{1/2}\mathbf{v}\|_2 \right| \\ & \leq \frac{c_1 \sqrt{\kappa_1 + \frac{1}{4}\bar{\Delta}^2} (\sigma + \bar{\Delta})^2}{\kappa_0 + \frac{1}{4}\bar{\Delta}^2} \sqrt{\delta \log d} := c_1 \mathcal{L}_1 \sqrt{\delta \log d}. \end{aligned}$$

Under the same probability, a simple rescaling then provides

$$\left| \frac{1}{\sqrt{n}} \|\dot{\mathbf{X}}\mathbf{v}\|_2 - \|\dot{\boldsymbol{\Sigma}}^{1/2}\mathbf{v}\|_2 \right| \leq c_1 \mathcal{L}_1 \sqrt{\frac{\delta \log d}{n}} \|\mathbf{v}\|_1, \quad \forall \mathbf{v} \in \mathbb{R}^d, \quad (37)$$

which implies

$$\begin{aligned} & \frac{1}{\sqrt{n}} \|\dot{\mathbf{X}}\mathbf{v}\|_2 \geq \|\dot{\boldsymbol{\Sigma}}^{1/2}\mathbf{v}\|_2 - c_1 \mathcal{L}_1 \left(\frac{\delta \log d}{n} \right)^{1/2} \|\mathbf{v}\|_1 \\ & \geq \left(\kappa_0 + \frac{1}{4}\bar{\Delta}^2 \right)^{1/2} \|\mathbf{v}\|_2 - c_1 \mathcal{L}_1 \left(\frac{\delta \log d}{n} \right)^{1/2} \|\mathbf{v}\|_1, \quad \forall \mathbf{v} \in \mathbb{R}^d. \end{aligned} \quad (38)$$

¹⁶The application of Bernstein's inequality leads to the σ^2 (σ is the upper bound on $\|\dot{\mathbf{x}}_k\|_{\psi_2}$) dependence in the multiplicative factor \mathcal{L} . It is possible to refine this quadratic dependence by using a new Bernstein's inequality developed in [49, Thm. 1.3], but we do not pursue this in the present paper.

Based on (38), we let $\hat{c} := \frac{2\kappa_0 + \bar{\Delta}^2}{4\kappa_0 + \bar{\Delta}^2}$ and use the inequality $(a-b)^2 \geq \hat{c}a^2 - \frac{\hat{c}}{1-\hat{c}}b^2$ to obtain

$$\begin{aligned} \mathbf{v}^\top \mathbf{Q}\mathbf{v} &= \frac{1}{n} \|\dot{\mathbf{X}}\mathbf{v}\|_2^2 - \frac{\bar{\Delta}^2}{4} \|\mathbf{v}\|_2^2 \\ &\geq \hat{c}(\kappa_0 + \frac{1}{4}\bar{\Delta}^2) \|\mathbf{v}\|_2^2 - c_1^2 \mathcal{L}_1^2 \frac{\hat{c}}{1-\hat{c}} \frac{\delta \log d}{n} \|\mathbf{v}\|_1^2 - \frac{\bar{\Delta}^2}{4} \|\mathbf{v}\|_2^2 \\ &\geq \frac{\kappa_0}{2} \|\mathbf{v}\|_2^2 - c_1^2 \mathcal{L}_1^2 \left(1 + \frac{\bar{\Delta}^2}{2\kappa_0}\right) \frac{\delta \log d}{n} \|\mathbf{v}\|_1^2 \\ &:= \frac{\kappa_0}{2} \|\mathbf{v}\|_2^2 - c_2(\kappa_0, \sigma, \bar{\Delta}) \cdot \frac{\delta \log d}{n} \|\mathbf{v}\|_1^2, \end{aligned}$$

which holds for all $\mathbf{v} \in \mathbb{R}^d$ and $\hat{c}_2 := c_2(\kappa_0, \sigma, \bar{\Delta})$ is a constant depending on $\kappa_0, \sigma, \bar{\Delta}$ (we remove the dependence on κ_1 by $\kappa_1 \lesssim \sigma^2$).

Step 3. Bounding the Estimation Error

We are in a position to bound the estimation error of any $\tilde{\boldsymbol{\theta}}$ satisfying (14). Note that definition of $\partial\|\tilde{\boldsymbol{\theta}}\|$ gives $\lambda\|\boldsymbol{\theta}^*\|_1 - \lambda\|\tilde{\boldsymbol{\theta}}\|_1 \geq \langle \lambda \cdot \partial\|\tilde{\boldsymbol{\theta}}\|_1, -\tilde{\mathbf{Y}} \rangle$. Thus, we set $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ in (14) and proceed as follows

$$\begin{aligned} 0 &\geq \langle \mathbf{Q}\tilde{\boldsymbol{\theta}} - \mathbf{b} + \lambda \cdot \partial\|\tilde{\boldsymbol{\theta}}\|_1, \tilde{\mathbf{Y}} \rangle \\ &= \tilde{\mathbf{Y}}^\top \mathbf{Q}\tilde{\mathbf{Y}} + \langle \mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}, \tilde{\mathbf{Y}} \rangle + \langle \lambda \cdot \partial\|\tilde{\boldsymbol{\theta}}\|_1, \tilde{\mathbf{Y}} \rangle \\ &\stackrel{(i)}{\geq} \frac{\kappa_0}{2} \|\tilde{\mathbf{Y}}\|_2^2 - \frac{2\hat{c}_2 R \sqrt{s} \cdot \delta \log d}{n} \|\tilde{\mathbf{Y}}\|_1 \\ &\quad - \frac{\lambda}{\beta} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\tilde{\boldsymbol{\theta}}\|_1 - \|\boldsymbol{\theta}^*\|_1) \\ &\stackrel{(ii)}{\geq} \frac{\kappa_0}{2} \|\tilde{\mathbf{Y}}\|_2^2 - \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\tilde{\boldsymbol{\theta}}\|_1 - \|\boldsymbol{\theta}^*\|_1), \end{aligned} \quad (39)$$

where we use $\lambda \geq \beta \|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$ ($\beta > 2$) and $\|\tilde{\mathbf{Y}}\|_1 \leq \|\tilde{\boldsymbol{\theta}}\|_1 + \|\boldsymbol{\theta}^*\|_1 \leq 2R\sqrt{s}$ in (i), and (ii) is due to the scaling

$$2\hat{c}_2 R \delta \sqrt{s} \frac{\log d}{n} \leq \left(\frac{1}{2} - \frac{1}{\beta} \right) \lambda$$

that holds under the assumed $n \gtrsim \delta \log d$ for some hidden constant depending on $(\kappa_0, \sigma, \bar{\Delta}, \Delta, M, R)$. Thus, we arrive at $\frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 \geq \lambda (\|\tilde{\boldsymbol{\theta}}\|_1 - \|\boldsymbol{\theta}^*\|_1)$.

For $\mathbf{a} \in \mathbb{R}^d$, $\mathcal{J} \subset [d]$ we obtain $\mathbf{a}_{\mathcal{J}} \in \mathbb{R}^d$ by keeping entries of \mathbf{a} in \mathcal{J} while setting others to zero. Let \mathcal{A} be the support of $\boldsymbol{\theta}^*$, $\mathcal{A}^c = [d] \setminus \mathcal{A}$, then we have

$$\begin{aligned} \frac{1}{2} \|\tilde{\mathbf{Y}}\|_1 &\geq \|\boldsymbol{\theta}^*\|_1 + \|\tilde{\mathbf{Y}}\|_1 - \|\boldsymbol{\theta}^*\|_1 \\ &= \|\boldsymbol{\theta}^*\|_1 + \|\tilde{\mathbf{Y}}_{\mathcal{A}} + \tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1 - \|\boldsymbol{\theta}^*\|_1 \\ &\geq \|\boldsymbol{\theta}^*\|_1 + \|\tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1 - \|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1 - \|\boldsymbol{\theta}^*\|_1 \\ &= \|\tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1 - \|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1. \end{aligned} \quad (40)$$

Further use $\frac{1}{2} \|\tilde{\mathbf{Y}}\|_1 \leq \frac{1}{2} \|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1 + \frac{1}{2} \|\tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1$, we obtain $\|\tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1 \leq 3\|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1$. Hence, we have $\|\tilde{\mathbf{Y}}\|_1 \leq \|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1 + \|\tilde{\mathbf{Y}}_{\mathcal{A}^c}\|_1 \leq 4\|\tilde{\mathbf{Y}}_{\mathcal{A}}\|_1 \leq 4\sqrt{s}\|\tilde{\mathbf{Y}}\|_2$. Now, we further substitute this into (39) and obtain

$$\begin{aligned} \frac{1}{2} \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 &\leq \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\boldsymbol{\theta}^*\|_1 - \|\tilde{\boldsymbol{\theta}}\|_1) \\ &\leq \frac{3\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 \leq 6\lambda \sqrt{s} \|\tilde{\mathbf{Y}}\|_2. \end{aligned}$$

Thus, we arrive at the desired error bound for ℓ_2 -norm

$$\|\tilde{\mathbf{Y}}\|_2 \lesssim \mathcal{L} \sqrt{\frac{\delta s \log d}{n}}, \text{ with } \mathcal{L} := \frac{(\sigma + \bar{\Delta})^2 (\Delta + M^{1/(2l)})}{\kappa_0^{3/2}}.$$

We simply use $\|\tilde{\mathbf{Y}}\|_1 \leq 4\sqrt{s}\|\tilde{\mathbf{Y}}\|_2$ again to establish the bound for $\|\tilde{\mathbf{Y}}\|_1$. The proof is complete. \square

Proof of Theorem 10: *Proof.* The proof is divided into two steps. Compared to the last proof, due to the heavy-tailedness of \mathbf{x}_k , the step of ‘‘verifying RSC’’ reduces to the simpler argument in (42).

Step 1. Proving $\lambda \geq \beta \|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty$ for some pre-specified $\beta > 2$

Recall that (\mathbf{Q}, \mathbf{b}) are constructed from the quantized data as $\mathbf{Q} = \frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d$ and $\mathbf{b} = \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k$. Thus, our main aim in this step is to prove that

$$\lambda = C_1 (R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

suffices to ensure

$$\lambda \geq \beta \left\| \frac{1}{n} \sum_{k=1}^n \left(\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d \right) \boldsymbol{\theta}^* - \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k \right\|_\infty$$

with the promised probability and any pre-specified $\beta > 2$. We let $\hat{\mathbf{x}}_k = \tilde{\mathbf{x}}_k + \boldsymbol{\tau}_k + \mathbf{w}_k$, $\hat{y}_k = \tilde{y}_k + \phi_k + \vartheta_k$ with quantization errors \mathbf{w}_k and ϑ_k . Analogously to (34), we have $\mathbb{E}[\hat{y}_k \hat{\mathbf{x}}_k] = \mathbb{E}[\tilde{y}_k \tilde{\mathbf{x}}_k]$ and $\mathbb{E}[\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top \boldsymbol{\theta}^*] = \frac{\bar{\Delta}^2}{4} \boldsymbol{\theta}^* + \mathbb{E}[\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top \boldsymbol{\theta}^*]$. Thus, the term we want to bound can be first decomposed into two concentration terms (I_1, I_3) and one bias term (I_2):

$$\begin{aligned} & \left\| \left(\frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \frac{\bar{\Delta}^2}{4} \mathbf{I}_d \right) \boldsymbol{\theta}^* - \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k \right\|_\infty \\ & \leq \left\| \frac{1}{n} \sum_{k=1}^n \hat{y}_k \hat{\mathbf{x}}_k - \mathbb{E}[\hat{y}_k \hat{\mathbf{x}}_k] \right\|_\infty + \left\| \mathbb{E}[\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top \boldsymbol{\theta}^*] - \mathbb{E}[\tilde{y}_k \tilde{\mathbf{x}}_k] \right\|_\infty \\ & + \left\| \left(\frac{1}{n} \sum_{k=1}^n \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \mathbb{E}[\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top] \right) \boldsymbol{\theta}^* \right\|_\infty := I_1 + I_2 + I_3, \quad (41) \end{aligned}$$

Step 1.1. Bounding I_1

Denote the i -th entry of $\mathbf{x}_k, \tilde{\mathbf{x}}_k, \hat{\mathbf{x}}_k, \boldsymbol{\tau}_k, \mathbf{w}_k$ by $x_{ki}, \tilde{x}_{ki}, \hat{x}_{ki}, \tau_{ki}, w_{ki}$, respectively. Since $|\hat{y}_k| \leq |\tilde{y}_k| + |\phi_k| + |\vartheta_k| \leq |\tilde{y}_k| + \Delta$, $\hat{x}_{ki} \leq |\tilde{x}_{ki}| + |\tau_{ki}| + |w_{ki}| \leq |\tilde{x}_{ki}| + \frac{3}{2}\bar{\Delta}$, for any integer $q \geq 2$ we can bound the moments as

$$\begin{aligned} & \sum_{k=1}^n \mathbb{E} \left| \frac{\hat{y}_k \hat{x}_{ki}}{n} \right|^q \leq \frac{(\zeta_x + \frac{3}{2}\bar{\Delta})^{q-2} (\zeta_y + \Delta)^{q-2}}{n^q} \sum_{k=1}^n \mathbb{E} |\hat{y}_k \hat{x}_{ki}|^2 \\ & \leq \frac{[(\zeta_x + \frac{3}{2}\bar{\Delta})(\zeta_y + \Delta)]^{q-2}}{n^q} \sum_{k=1}^n \sqrt{\mathbb{E} |\hat{y}_k|^4 \mathbb{E} |\hat{x}_{ki}|^4} \\ & \lesssim \left(\frac{(\zeta_x + \frac{3}{2}\bar{\Delta})(\zeta_y + \Delta)}{n} \right)^{q-2} \left(\frac{M + \Delta^4 + \bar{\Delta}^4}{n} \right), \end{aligned}$$

and in the last inequality we use $\mathbb{E}|\tilde{y}_k| \leq \mathbb{E}|y_k|^4 \leq M$ and $\mathbb{E}|\tilde{x}_{ki}|^4 \leq \mathbb{E}|x_{ki}|^4 \leq M$. With these preparations, we apply Bernstein's inequality (Lemma 1) and a union bound, yielding that

$$\begin{aligned} \mathbb{P} \left(I_1 \geq C_5 \left\{ \sqrt{\frac{(M + \Delta^4 + \bar{\Delta}^4)t}{n}} + \frac{(\zeta_x + \frac{3}{2}\bar{\Delta})(\zeta_y + \Delta)t}{n} \right\} \right) \\ \leq 2d \cdot \exp(-t). \end{aligned}$$

Set $t = \delta \log d$, we obtain that

$$I_1 \lesssim (\sqrt{M} + \Delta^2 + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 2d^{1-\delta}$.

Step 1.2. Bounding I_2

Noting that $\mathbb{E}(y_k \mathbf{x}_k) = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top \boldsymbol{\theta}^*) + \mathbb{E}(\epsilon_k \mathbf{x}_k) = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top \boldsymbol{\theta}^*)$, we could further decompose I_2 as

$$\begin{aligned} I_2 & \leq \|\mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top \boldsymbol{\theta}^*) - \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top \boldsymbol{\theta}^*)\|_\infty \\ & + \|\mathbb{E}(y_k \mathbf{x}_k) - \mathbb{E}(\tilde{y}_k \tilde{\mathbf{x}}_k)\|_\infty := I_{21} + I_{22} \end{aligned}$$

To bound I_{21} , we note that the assumption and truncation procedure for \mathbf{x}_k are the same as in Theorem 2; thus, Step 2 in the proof of Theorem 2 can yield that

$$\|\mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top)\|_\infty \leq \sqrt{\frac{\delta M \log d}{n}}.$$

Thus, we have

$$I_{21} \leq \|\mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top - \mathbf{x}_k \mathbf{x}_k^\top)\|_\infty \|\boldsymbol{\theta}^*\|_1 \leq R\sqrt{M} \sqrt{\frac{\delta \log d}{n}}.$$

To bound I_{22} , we estimate the i -th entry

$$\begin{aligned} & |\mathbb{E}(y_k x_{ki} - \tilde{y}_k \tilde{x}_{ki})| \\ & = |\mathbb{E}(y_k x_{ki} - \tilde{y}_k \tilde{x}_{ki}) (\mathbb{1}(|y_k| > \zeta_y) + \mathbb{1}(|x_{ki}| \geq \zeta_x))| \\ & \leq \mathbb{E}(|y_k x_{ki}| \mathbb{1}(|y_k| > \zeta_y)) + \mathbb{E}(|y_k x_{ki}| \mathbb{1}(|x_{ki}| \geq \zeta_x)) \\ & \stackrel{(i)}{\leq} M \left(\frac{1}{\zeta_x^2} + \frac{1}{\zeta_y^2} \right) \lesssim \sqrt{\frac{\delta M \log d}{n}}, \end{aligned}$$

where (i) is because

$$\begin{aligned} & \mathbb{E}(|y_k x_{ki}| \mathbb{1}(|y_k| > \zeta_y)) \\ & \leq [\mathbb{E}|y_k^2 x_{ki}^2|]^{1/2} \sqrt{\mathbb{P}(|y_k| > \zeta_y)} \\ & \leq (\mathbb{E}y_k^4)^{1/4} (\mathbb{E}x_{ki}^4)^{1/4} \sqrt{\frac{\mathbb{E}y_k^4}{\zeta_y^4}} \leq \frac{M}{\zeta_y^2} \end{aligned}$$

and applying similar treatment to $\mathbb{E}(|y_k x_{ki}| \mathbb{1}(|x_{ki}| > \zeta_x))$. Overall, we have

$$I_2 \lesssim R\sqrt{M} \sqrt{\frac{\delta \log d}{n}}.$$

Step 1.3. Bounding I_3

We first note that

$$I_3 = \|(\mathbf{Q} - \boldsymbol{\Sigma}^*) \boldsymbol{\theta}^*\|_\infty \leq \|\mathbf{Q} - \boldsymbol{\Sigma}^*\|_\infty \|\boldsymbol{\theta}^*\|_1 \leq R \|\mathbf{Q} - \boldsymbol{\Sigma}^*\|_\infty.$$

By Theorem 2, we know that

$$\|\mathbf{Q} - \boldsymbol{\Sigma}^*\|_\infty \lesssim (\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 2d^{2-\delta}$, which leads to

$$I_3 \leq \|\mathbf{Q} - \boldsymbol{\Sigma}^*\|_\infty \|\boldsymbol{\theta}^*\|_1 \lesssim R(\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}.$$

Thus, by combining everything, we obtain that

$$\|\mathbf{Q}\boldsymbol{\theta}^* - \mathbf{b}\|_\infty \lesssim (R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

holds with probability at least $1 - 4d^{2-\delta}$. Compared to our choice of λ , the claim of this step follows.

Step 2. Bounding the Estimation Error

We are now ready to bound the estimation error of any $\tilde{\theta}$ satisfying (14). Set $\theta = \theta^*$ in (41), it gives

$$\langle Q\tilde{\theta} - \mathbf{b} + \lambda \cdot \partial\|\tilde{\theta}\|_1, \tilde{\mathbf{Y}} \rangle \leq 0.$$

Recall that we can assume

$$\|Q - \Sigma^*\|_\infty \leq C_6(\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

with probability at least $1 - 2d^{2-\delta}$, which leads to

$$\begin{aligned} \tilde{\mathbf{Y}}^\top Q \tilde{\mathbf{Y}} &= \tilde{\mathbf{Y}}^\top \Sigma^* \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}^\top (\Sigma^* - Q) \tilde{\mathbf{Y}} \\ &\geq \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 - C_6(\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}} \|\tilde{\mathbf{Y}}\|_1^2. \end{aligned} \quad (42)$$

Thus, it follows that

$$\begin{aligned} 0 &\geq \langle Q\tilde{\theta} - \mathbf{b} + \lambda \cdot \partial\|\tilde{\theta}\|_1, \tilde{\mathbf{Y}} \rangle \\ &= \langle Q\theta^* - \mathbf{b}, \tilde{\mathbf{Y}} \rangle + \tilde{\mathbf{Y}}^\top Q \tilde{\mathbf{Y}} + \lambda \langle \partial\|\tilde{\theta}\|_1, \tilde{\mathbf{Y}} \rangle \\ &\stackrel{(i)}{\geq} C_0 \sqrt{M} \|\tilde{\mathbf{Y}}\|_2^2 - C_6(\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}} \|\tilde{\mathbf{Y}}\|_1^2 \\ &\quad - \|Q\theta^* - \mathbf{b}\|_\infty \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1) \\ &\stackrel{(ii)}{\geq} C_0 \sqrt{M} \|\tilde{\mathbf{Y}}\|_2^2 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1) \\ &\quad - \left(2C_6 R(\sqrt{M} + \bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}} + \|Q\theta^* - \mathbf{b}\|_\infty \right) \|\tilde{\mathbf{Y}}\|_1 \\ &\stackrel{(iii)}{\geq} C_0 \sqrt{M} \|\tilde{\mathbf{Y}}\|_2^2 - \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1). \end{aligned} \quad (43)$$

Note that (i) is due to (42) and $\|\theta^*\|_1 - \|\tilde{\theta}\|_1 \geq \langle \partial\|\tilde{\theta}\|_1, -\tilde{\mathbf{Y}} \rangle$, in (ii) we use $\|\tilde{\mathbf{Y}}\|_1 \leq \|\tilde{\theta}\|_1 + \|\theta^*\|_1 \leq 2R$, and from Step 1 (iii) holds when

$$\lambda = C_2(R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2) \sqrt{\frac{\delta \log d}{n}}$$

with sufficiently large C_2 . Therefore, we arrive at $\frac{1}{2} \|\tilde{\mathbf{Y}}\|_1 \geq \|\tilde{\theta}\|_1 - \|\theta^*\|_1$. Similar to Step 3 in the proof of Theorem 9, we can show $\|\tilde{\mathbf{Y}}\|_1 \leq 4\sqrt{s} \|\tilde{\mathbf{Y}}\|_2$. Applying (43) again, it yields

$$\kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 \leq \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda \|\tilde{\mathbf{Y}}\|_1 \leq \frac{3\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 \leq 6\sqrt{s}\lambda \|\tilde{\mathbf{Y}}\|_2.$$

Thus, we obtain $\|\tilde{\mathbf{Y}}\|_2 \lesssim \mathcal{L} \sqrt{\frac{\delta s \log d}{n}}$ with $\mathcal{L} := \frac{R\sqrt{M} + \Delta^2 + R\bar{\Delta}^2}{\kappa_0}$. The proof can be concluded by further applying $\|\tilde{\mathbf{Y}}\|_1 \leq 4\sqrt{s} \|\tilde{\mathbf{Y}}\|_2$ to derive the bound for ℓ_1 -norm. \square

Proof of Proposition 1:

Proof. We let $\theta = \theta^*$ in (14), then proceeds as the proof of Theorem 10:

$$\begin{aligned} 0 &\geq \langle Q\tilde{\theta} - \mathbf{b} + \lambda \cdot \partial\|\tilde{\theta}\|_1, \tilde{\mathbf{Y}} \rangle \\ &= \tilde{\mathbf{Y}}^\top \Sigma^* \tilde{\mathbf{Y}} + \langle Q\theta^* - \mathbf{b}, \tilde{\mathbf{Y}} \rangle \\ &\quad + \tilde{\mathbf{Y}}^\top (Q - \Sigma^*) \tilde{\mathbf{Y}} + \lambda \langle \partial\|\tilde{\theta}\|_1, \tilde{\mathbf{Y}} \rangle \\ &\stackrel{(i)}{\geq} \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 - \|Q\theta^* - \mathbf{b}\|_\infty \|\tilde{\mathbf{Y}}\|_1 \\ &\quad - \|Q - \Sigma^*\|_\infty \|\tilde{\mathbf{Y}}\|_1^2 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1) \\ &\stackrel{(ii)}{\geq} \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1) \\ &\quad - (\|Q\theta^* - \mathbf{b}\|_\infty + 2R \cdot \|Q - \Sigma^*\|_\infty) \|\tilde{\mathbf{Y}}\|_1 \\ &\stackrel{(iii)}{\geq} \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 - \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\tilde{\theta}\|_1 - \|\theta^*\|_1), \end{aligned} \quad (44)$$

where in (i) we use $\lambda_{\min}(\Sigma^*) \geq \kappa_0$ and $\|\theta^*\|_1 - \|\tilde{\theta}\|_1 \geq \langle \partial\|\tilde{\theta}\|_1, -\tilde{\mathbf{Y}} \rangle$, (ii) is by $\|\tilde{\mathbf{Y}}\|_1 \leq \|\tilde{\theta}\|_1 + \|\theta^*\|_1 \leq 2R$, in (iii) we use the assumption (15). Thus, by $\kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 \geq 0$ we obtain $2(\|\tilde{\theta}\|_1 - \|\theta^*\|_1) \leq \|\tilde{\mathbf{Y}}\|_1$. Similarly to Step 3 in the proof of Theorem 9, we can show $\|\tilde{\mathbf{Y}}\|_1 \leq 4\sqrt{s} \|\tilde{\mathbf{Y}}\|_2$. Again we use (44), it gives

$$\begin{aligned} \kappa_0 \|\tilde{\mathbf{Y}}\|_2^2 &\leq \frac{\lambda}{2} \|\tilde{\mathbf{Y}}\|_1 + \lambda (\|\theta^*\|_1 - \|\tilde{\theta}\|_1) \\ &\leq \frac{3}{2} \lambda \|\tilde{\mathbf{Y}}\|_1 \leq 6\lambda\sqrt{s} \|\tilde{\mathbf{Y}}\|_2, \end{aligned}$$

thus displaying $\|\tilde{\mathbf{Y}}\|_2 \leq \frac{6\lambda\sqrt{s}}{\kappa_0}$. The proof can be concluded by using $\|\tilde{\mathbf{Y}}\|_1 \leq 4\sqrt{s} \|\tilde{\mathbf{Y}}\|_2$. \square

Proof of Theorem 11:

Proof. To invoke Proposition 1, it is enough to verify (15). Recalling $\Sigma^* = \mathbb{E}(\mathbf{x}_k \mathbf{x}_k^\top)$, we first invoke [22, Thm. 1] and obtain $\|Q - \Sigma^*\|_\infty = O(\sigma^2 \sqrt{\frac{\delta \log d (\log n)^2}{n}})$ holds with probability at least $1 - 2d^{2-\delta}$. This confirms $\lambda \gtrsim R \cdot \|Q - \Sigma^*\|_\infty$. Then, it remains to upper bound $\|Q\theta^* - \mathbf{b}\|_\infty$:

$$\begin{aligned} \|Q\theta^* - \mathbf{b}\|_\infty &\leq \|(Q - \Sigma^*)\theta^*\|_\infty + \|\Sigma^*\theta^* - \mathbf{b}\|_\infty \\ &\leq \|Q - \Sigma^*\|_\infty \|\theta^*\|_1 + \left\| \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_{k1} - \mathbb{E}(y_k \mathbf{x}_k) \right\|_\infty \\ &\stackrel{(i)}{\lesssim} \sigma^2 (R+1) \sqrt{\frac{\delta \log d (\log n)^2}{n}}, \end{aligned}$$

where in (i) we use a known estimate from [22, Eq. (A.31)]:

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{\mathbf{x}}_{k1} - \mathbb{E}(y_k \mathbf{x}_k) \right\|_\infty \lesssim \sigma^2 \sqrt{\frac{\delta \log d (\log n)^2}{n}} \right) \\ \geq 1 - 2d^{1-\delta}. \end{aligned}$$

Thus, by setting

$$\lambda = C_1 \sigma^2 R \sqrt{\frac{\delta \log d (\log n)^2}{n}},$$

(15) can be satisfied with probability at least $1 - 4d^{2-\delta}$, hence using Proposition 1 concludes the proof. \square

Proof of Theorem 12:

Proof. The proof is again based on Proposition 1 and some ingredients from [22]. From [22, Thm. 4],

$$\|\mathbf{Q} - \Sigma^*\|_\infty \lesssim \left(\frac{M^2 \delta \log d}{n}\right)^{1/4}$$

holds with probability at least $1 - 2d^{2-\delta}$, thus confirming $\lambda \gtrsim R \cdot \|\mathbf{Q} - \Sigma^*\|_\infty$ with the same probability. Moreover,

$$\begin{aligned} \|\mathbf{Q}\theta^* - \mathbf{b}\|_\infty &\leq \|(\mathbf{Q} - \Sigma^*)\theta^*\|_\infty + \|\Sigma^*\theta^* - \mathbf{b}\|_\infty \\ &\leq \|\mathbf{Q} - \Sigma^*\|_\infty \|\theta^*\|_1 + \left\| \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{x}_{k1} - \mathbb{E}(y_k \mathbf{x}_k) \right\|_\infty \\ &\stackrel{(i)}{\lesssim} \sqrt{M}(R+1) \left(\frac{\delta \log d}{n}\right)^{1/4}, \end{aligned}$$

where (i) is due to a known estimate from [22, Eq. (A.34)]:

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{\gamma_x \gamma_y}{n} \sum_{k=1}^n \dot{y}_k \dot{x}_{k1} - \mathbb{E}(y_k \mathbf{x}_k) \right\|_\infty \lesssim \sqrt{M} \left(\frac{\delta \log d}{n}\right)^{1/4} \right) \\ \geq 1 - 2d^{1-\delta} \end{aligned}$$

Thus, with probability at least $1 - 4d^{2-\delta}$, (15) holds if

$$\lambda = C_1 \sqrt{M} R \left(\frac{\delta \log d}{n}\right)^{1/4}$$

with sufficiently large C_1 . The proof can be concluded by invoking Proposition 1. \square

2) **Uniform Recovery Guarantee:** We need some auxiliary results to support the proof. The first one is a concentration inequality for product process due to Mendelson [67]; the following statement can be directly adapted from [40, Thm. 8] by specifying the pseudo-metrics as ℓ_2 -distance.

Lemma 9. (Concentration of Product Process). *Let $\{g_a\}_{a \in \mathcal{A}}$ and $\{h_b\}_{b \in \mathcal{B}}$ be stochastic processes indexed by two sets $\mathcal{A} \subset \mathbb{R}^p$, $\mathcal{B} \subset \mathbb{R}^q$, both defined on a common probability space $(\Omega, \mathcal{A}, \mathbb{P})$. We assume that there exist $K_A, K_B, r_A, r_B \geq 0$ such that*

$$\begin{aligned} \|g_a - g_{a'}\|_{\psi_2} &\leq K_A \|a - a'\|_2, \|g_a\|_{\psi_2} \leq r_A, \forall a, a' \in \mathcal{A}; \\ \|h_b - h_{b'}\|_{\psi_2} &\leq K_B \|b - b'\|_2, \|h_b\|_{\psi_2} \leq r_B, \forall b, b' \in \mathcal{B}. \end{aligned}$$

Finally, let X_1, \dots, X_m be independent copies of a random variable $X \sim \mathbb{P}$, then for every $u \geq 1$ the following holds with probability at least $1 - 2 \exp(-cu^2)$

$$\begin{aligned} \sup_{\substack{a \in \mathcal{A} \\ b \in \mathcal{B}}} \frac{1}{n} \left| \sum_{i=1}^n g_a(X_i) h_b(X_i) - \mathbb{E}[g_a(X_i) h_b(X_i)] \right| \\ \leq C \left(\frac{(K_A \cdot \omega(\mathcal{A}) + u \cdot r_A) \cdot (K_B \cdot \omega(\mathcal{B}) + u \cdot r_B)}{n} \right. \\ \left. + \frac{r_A \cdot K_B \cdot \omega(\mathcal{B}) + r_B \cdot K_A \cdot \omega(\mathcal{A}) + u \cdot r_A r_B}{\sqrt{n}} \right), \end{aligned}$$

where $\omega(\mathcal{A}) = \mathbb{E} \sup_{a \in \mathcal{A}} (\mathbf{g}^\top a)$ with $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_p)$ is the Gaussian width of $\mathcal{A} \subset \mathbb{R}^p$, and similarly, $\omega(\mathcal{B})$ is the Gaussian width of \mathcal{B} .

We will use the following result that can be found in [61, Thm. 8].

Lemma 10. *Let $(X_u)_{u \in \mathcal{T}}$ be a random process indexed by points in a bounded set $\mathcal{T} \subset \mathbb{R}^n$. Assume that the process has sub-Gaussian increments, i.e., there exists $M > 0$ such that $\|X_u - X_v\|_{\psi_2} \leq M \|u - v\|_2$ holds for any $u, v \in \mathcal{T}$. Then for every $t > 0$, the event*

$$\sup_{u, v \in \mathcal{T}} |X_u - X_v| \leq CM \cdot (\omega(\mathcal{T}) + t \cdot \text{diam}(\mathcal{T}))$$

holds with probability at least $1 - \exp(-t^2)$, where $\text{diam}(\mathcal{T}) := \sup_{x, y \in \mathcal{T}} \|x - y\|_2$ denotes the diameter of \mathcal{T} .

Proof of Theorem 13:

Proof. We start from the optimality

$$\sum_{k=1}^n (\dot{y}_k - \mathbf{x}_k^\top \hat{\theta})^2 \leq \sum_{k=1}^n (\dot{y}_k - \mathbf{x}_k^\top \theta^*)^2.$$

By substituting $\hat{\theta} = \theta^* + \hat{\Upsilon}$ and performing some algebra, we obtain

$$\sum_{k=1}^n (\mathbf{x}_k^\top \hat{\Upsilon})^2 \leq 2 \sum_{k=1}^n (\dot{y}_k - \mathbf{x}_k^\top \theta^*) \mathbf{x}_k^\top \hat{\Upsilon}.$$

Due to the constraint we have $\|\theta^* + \hat{\Upsilon}\|_1 \leq \|\theta^*\|_1$, then similar to (40) we can show $\|\hat{\Upsilon}\|_1 \leq 2\sqrt{s} \|\hat{\Upsilon}\|_2$ holds. Thus, we let $\mathcal{V} = \{v : \|v\|_2 = 1, \|v\|_1 \leq 2\sqrt{s}\}$, then the following holds uniformly for all $\theta^* \in \Sigma_{s, R_0}$

$$\begin{aligned} \|\hat{\Upsilon}\|_2^2 \cdot \inf_{v \in \mathcal{V}} \sum_{k=1}^n (\mathbf{x}_k^\top v)^2 \\ \leq 2 \|\hat{\Upsilon}\|_2 \cdot \sup_{v \in \mathcal{V}} \sum_{k=1}^n (\dot{y}_k - \mathbf{x}_k^\top \theta^*) \mathbf{x}_k^\top v. \end{aligned} \quad (45)$$

Similarly to previous developments, our strategy is to lower bound the left-hand side while upper bound the right hand side, but with the difference that the bounds must be valid uniformly for all $\theta^* \in \Sigma_{s, R_0}$.

Step 1. Bounding the Left-Hand Side From Below

Letting $\bar{\Delta} = 0$ and restricting v to \mathcal{V} , we use (38) in the proof of Theorem 9, then for some constant $c(\kappa_0, \sigma)$ depending on κ_0, σ , with probability at least $1 - d^{-\delta}$

$$\inf_{v \in \mathcal{V}} \frac{1}{\sqrt{n}} \left[\sum_{k=1}^n (\mathbf{x}_k^\top v)^2 \right]^{1/2} \geq \sqrt{\kappa_0} - c(\kappa_0, \sigma) \cdot \sqrt{\frac{s \log d}{n}}.$$

Thus, if $n \geq \frac{4c^2(\kappa_0, \sigma)}{\kappa_0} s \log d$, then it holds that

$$\inf_{v \in \mathcal{V}} \sum_{k=1}^n (\mathbf{x}_k^\top v)^2 \geq \frac{1}{4} \kappa_0 n.$$

Step 2. Bounding the Right-Hand Side Uniformly

To pursue the uniformity over $\theta^* \in \Sigma_{s, R_0}$, we take a supremum by replacing specific θ^* with $\sup_{\theta \in \Sigma_{s, R_0}}$, then we consider the upper bound on

$$I := \sup_{\theta \in \Sigma_{s, R_0}} \sup_{v \in \mathcal{V}} \sum_{k=1}^n (\dot{y}_k - \mathbf{x}_k^\top \theta) \mathbf{x}_k^\top v, \quad (46)$$

where $\dot{y}_k = \mathcal{Q}_\Delta(\tilde{y}_k + \tau_k)$, $\tilde{y}_k = \mathcal{T}_{\zeta_y}(y_k) := \text{sign}(y_k) \min\{|y_k|, \zeta_y\}$, $y_k = \mathbf{x}_k^\top \boldsymbol{\theta} + \epsilon_k$; note that $\dot{y}_k, \tilde{y}_k, y_k$ depend on $\boldsymbol{\theta}$, and we will use notation $\dot{y}_{\boldsymbol{\theta},k}, \tilde{y}_{\boldsymbol{\theta},k}, y_{\boldsymbol{\theta},k}$ to indicate such dependence when necessary. In this proof, the ranges of $\boldsymbol{\theta}$ and \mathbf{v} (e.g., in supremum), if omitted, are respectively $\boldsymbol{\theta} \in \Sigma_{s,R_0}$ and $\mathbf{v} \in \mathcal{V}$. Now let the quantization noise be $\xi_k = \dot{y}_k - \tilde{y}_k$, observing that $\mathbb{E}(y_k \mathbf{x}_k^\top \mathbf{v}) = \mathbb{E}(\boldsymbol{\theta}^\top \mathbf{x}_k \mathbf{x}_k^\top \mathbf{v}) + \mathbb{E}(\epsilon_k \mathbf{x}_k^\top \mathbf{v}) = \mathbb{E}(\boldsymbol{\theta}^\top \mathbf{x}_k \mathbf{x}_k^\top \mathbf{v})$, then we can first decompose I as

$$\begin{aligned} I &\leq \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n \xi_k \mathbf{x}_k^\top \mathbf{v} + \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n (\tilde{y}_k \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\tilde{y}_k \mathbf{x}_k^\top \mathbf{v}]) \\ &\quad + \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n \mathbb{E}((\tilde{y}_k - y_k) \mathbf{x}_k^\top \mathbf{v}) \\ &\quad + \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n (\boldsymbol{\theta}^\top \mathbf{x}_k \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\boldsymbol{\theta}^\top \mathbf{x}_k \mathbf{x}_k^\top \mathbf{v}]) \\ &:= I_0 + I_1 + I_2 + I_3, \end{aligned} \quad (47)$$

where I_0 is the term arising from quantization, I_1 is the concentration term involving truncation of heavy-tailed data for which we develop some new machinery to bound it, I_2 is the bias term, I_3 is a more regular concentration term that can be bounded via Lemma 9. In the remainder of the proof, we will bound I_1, I_2, I_3 separately and finally deal with I_0 .

Step 2.1. Bounding I_1

Using $\tilde{y}_k = \mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta} + \epsilon_k)$, I_1 is concerned with the concentration of the product process

$$\left\{ \sum_{k=1}^n \mathcal{T}_{\zeta_y}(\boldsymbol{\theta}^\top \mathbf{x}_k + \epsilon_k) \mathbf{x}_k^\top \mathbf{v} \right\}_{\boldsymbol{\theta}, \mathbf{v}}$$

about its mean. It is natural to apply Lemma 9 towards this end, but we lack good bound on $\|\mathcal{T}_{\zeta_y}(\boldsymbol{\theta}^\top \mathbf{x}_k + \epsilon_k)\|_{\psi_2}$ because of the heavy-tailedness of ϵ_k (on the other hand, the bound $O(\zeta_y)$ is just too crude to yield a sharp rate). Our strategy is already introduced in the mainbody — we introduce $\tilde{z}_k := \tilde{y}_k - \mathcal{T}_{\zeta_y}(\epsilon_k)$ and decompose I_1 as

$$\begin{aligned} I_1 &\leq \underbrace{\sup_{\mathbf{v}, \boldsymbol{\theta}} \sum_{k=1}^n (\tilde{z}_k \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\tilde{z}_k \mathbf{x}_k^\top \mathbf{v}])}_{:= I_{11}} \\ &\quad + \underbrace{\sup_{\mathbf{v}} \sum_{k=1}^n (\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k^\top \mathbf{v} - \mathbb{E}[\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k^\top \mathbf{v}])}_{:= I_{12}}. \end{aligned}$$

Thus, it suffices to bound I_{11} and I_{12} .

Step 2.1.1. Bounding I_{11}

We use Lemma 9 to bound I_{11} . For any $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{V}$, it is evident that we have $\|\mathbf{x}_k^\top \mathbf{v}\|_{\psi_2} \leq \|\mathbf{x}_k\|_{\psi_2} \leq \sigma$ and $\|\mathbf{x}_k^\top \mathbf{v}_1 - \mathbf{x}_k^\top \mathbf{v}_2\|_{\psi_2} \leq \sigma \|\mathbf{v}_1 - \mathbf{v}_2\|_2$. Regarding

$$\tilde{z}_k = \tilde{z}_{\boldsymbol{\theta},k} := \mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta} + \epsilon_k) - \mathcal{T}_{\zeta_y}(\epsilon_k)$$

indexed by $\boldsymbol{\theta} \in \Sigma_{s,R_0}$, the 1-Lipschitzness of $\mathcal{T}_{\zeta_y}(\cdot)$ gives $|\tilde{z}_k| \leq |\mathbf{x}_k^\top \boldsymbol{\theta}|$, and then the definition of sub-Gaussian norm yields

$$\|\tilde{z}_k\|_{\psi_2} \leq \|\mathbf{x}_k^\top \boldsymbol{\theta}\|_{\psi_2} \leq \|\mathbf{x}_k\|_{\psi_2} \|\boldsymbol{\theta}\|_2 \leq R_0 \sigma$$

(this addresses the aforementioned issue). Further, for any $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Sigma_{s,R_0}$ we verify the sub-Gaussian increments

$$\begin{aligned} &|\tilde{z}_{\boldsymbol{\theta}_1,k} - \tilde{z}_{\boldsymbol{\theta}_2,k}| \\ &= |\mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta}_1 + \epsilon_k) - \mathcal{T}_{\zeta_y}(\epsilon_k) \\ &\quad - \mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta}_2 + \epsilon_k) + \mathcal{T}_{\zeta_y}(\epsilon_k)| \\ &= |\mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta}_1 + \epsilon_k) - \mathcal{T}_{\zeta_y}(\mathbf{x}_k^\top \boldsymbol{\theta}_2 + \epsilon_k)| \\ &\leq |\mathbf{x}_k^\top (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)|, \end{aligned} \quad (48)$$

which leads to

$$\begin{aligned} \|\tilde{z}_{\boldsymbol{\theta}_1,k} - \tilde{z}_{\boldsymbol{\theta}_2,k}\|_{\psi_2} &\leq \|\mathbf{x}_k^\top (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)\|_{\psi_2} \\ &\leq \|\mathbf{x}_k\|_{\psi_2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2 \leq \sigma \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2 \end{aligned}$$

With these preparations, we can invoke Lemma 9 use the well-known estimates $\omega(\Sigma_{s,R_0}), \omega(\mathcal{V}) = O(\sqrt{s \log d})^{17}$ to obtain that, with probability at least $1 - 2 \exp(-cu^2)$ we have

$$\begin{aligned} I_{11} &\lesssim \sigma^2 \left[\sqrt{n} (\omega(\Sigma_{s,R_0}) + \omega(\mathcal{V}) + u) \right. \\ &\quad \left. + (\omega(\Sigma_{s,R_0}) + u) \cdot (\omega(\mathcal{V}) + u) \right] \\ &\lesssim \sigma^2 \left[\sqrt{n} (\sqrt{s \log d} + u) \right. \\ &\quad \left. + (\sqrt{s \log d} + u) \cdot (\sqrt{s \log d} + u) \right]. \end{aligned} \quad (49)$$

Therefore, we can set $u = \sqrt{\delta s \log d}$ in (49), under the scaling of $n \gtrsim \delta s \log d$ it provides

$$\mathbb{P}(I_{11} \lesssim \sigma^2 \sqrt{n \delta s \log d}) \geq 1 - 2d^{-\delta \Omega(s)}. \quad (50)$$

Step 2.1.2. Bounding I_{12}

By $\|\mathbf{v}\|_1 \leq 2\sqrt{s}$ we have

$$I_{12} \leq 2\sqrt{s} \left\| \sum_{k=1}^n (\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k - \mathbb{E}[\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k]) \right\|_\infty.$$

Then to apply Bernstein's inequality, for integer $q \geq 2$ and $i \in [d]$, analogously to (35) in the proof of Theorem 9, we can bound that

$$\begin{aligned} &\sum_{k=1}^n \mathbb{E} \left| \frac{\mathcal{T}_{\zeta_y}(\epsilon_k) x_{ki}}{n} \right|^q \\ &\leq \left(\frac{\zeta_y}{n} \right)^{q-2} \frac{1}{n^2} \sum_{k=1}^n \mathbb{E} |\mathcal{T}_{\zeta_y}^2(\epsilon_k) x_{ki}^q| \\ &\leq \left(\frac{\sigma \zeta_y}{n} \right)^{q-2} \left(\frac{\sigma^2 M^{\frac{1}{l}}}{n} \right) (Cq)^{\frac{q}{2}} \leq \frac{q!}{2} v_0 c_0^{q-2}, \end{aligned}$$

for some $v_0 = O(\frac{\sigma^2 M^{1/l}}{n})$, $c_0 = O(\frac{\sigma \zeta_y}{n})$. Then we use Lemma 1 to obtain that, with probability at least $1 - 2 \exp(-t)$ we have

$$\begin{aligned} &\left| \frac{1}{n} \sum_{k=1}^n (\mathcal{T}_{\zeta_y}(\epsilon_k) x_{ki} - \mathbb{E}[\mathcal{T}_{\zeta_y}(\epsilon_k) x_{ki}]) \right| \\ &\leq C \sigma \left(M^{\frac{1}{2l}} \sqrt{\frac{t}{n}} + \frac{\zeta_y t}{n} \right) \end{aligned}$$

¹⁷In fact, we have the tighter estimate $\omega(\Sigma_{s,R_0}), \omega(\mathcal{V}) = O(\sqrt{s \log(\frac{ed}{s})})$ (e.g., [76]) but we simply put $\sqrt{s \log d}$ to be consistent with earlier results concerning unconstrained Lasso.

Then we use $\zeta_y \asymp (\sigma + M^{\frac{1}{2l}}) \sqrt{\frac{n}{\delta \log d}}$, set $t \asymp \delta \log d$, and take a union bound over $i \in [d]$ to obtain that,

$$\begin{aligned} & \left\| \sum_{k=1}^n \frac{1}{n} (\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k - \mathbb{E}[\mathcal{T}_{\zeta_y}(\epsilon_k) \mathbf{x}_k]) \right\|_{\infty} \\ & \lesssim \sigma(M^{1/(2l)} + \sigma) \sqrt{\frac{\delta \log d}{n}} \end{aligned}$$

holds with probability at least $1 - 2d^{1-\delta}$, which implies the following under the same probability

$$I_{12} \lesssim \sigma(M^{\frac{1}{2l}} + \sigma) \sqrt{ns\delta \log d}. \quad (51)$$

Therefore, combining (50) and (51), we obtain that

$$I_1 \lesssim \sigma(M^{\frac{1}{2l}} + \sigma) \sqrt{ns\delta \log d}$$

with the promised probability.

Step 2.2. Bounding I_2

For this bias term the supremum does not make things harder. We begin with

$$\begin{aligned} I_2 &= n \cdot \sup_{\boldsymbol{\theta}, \mathbf{v}} \mathbb{E}((\tilde{y}_k - y_k) \mathbf{x}_k^{\top} \mathbf{v}) \\ &\leq 2n\sqrt{s} \cdot \sup_{\boldsymbol{\theta}} \|\mathbb{E}(\tilde{y}_k - y_k) \mathbf{x}_k\|_{\infty}. \end{aligned}$$

Fix any $\boldsymbol{\theta} \in \Sigma_{s, R_0}$, we have

$$\mathbb{E}|y_k|^{2l} \lesssim \mathbb{E}|\mathbf{x}_k^{\top} \boldsymbol{\theta}|^{2l} + \mathbb{E}|\epsilon_k|^{2l} \lesssim M + \sigma^{2l}.$$

Then following arguments similarly to (36) we obtain

$$\|\mathbb{E}(\tilde{y}_k - y_k) \mathbf{x}_k\|_{\infty} \lesssim \frac{\sigma M^{1/l}}{\zeta_y} \lesssim \sigma(M^{1/(2l)} + \sigma) \sqrt{\frac{\delta \log d}{n}},$$

which implies $I_2 \lesssim \sigma M^{\frac{1}{2l}} \sqrt{ns\delta \log d}$.

Step 2.3. Bounding I_3

It is evident that we can apply Lemma 9 with $(g_{\boldsymbol{\theta}}(\mathbf{x}_k), h_{\mathbf{v}}(\mathbf{x}_k)) = (\boldsymbol{\theta}^{\top} \mathbf{x}_k, \mathbf{v}^{\top} \mathbf{x}_k)$, $(\mathcal{A}, \mathcal{B}) = (\Sigma_{s, R_0}, \mathcal{V})$, and hence with $K_{\mathcal{A}}, r_{\mathcal{A}}, K_{\mathcal{B}}, r_{\mathcal{B}} = O(\sigma)$. Along with $\omega(\Sigma_{s, R_0}), \omega(\mathcal{V}) \lesssim \sqrt{s \log d}$, we obtain that the following holds with probability at least $1 - 2 \exp(-cu^2)$:

$$I_3 \leq \sigma^2 \left[(\sqrt{s \log d} + u)^2 + \sqrt{n}(\sqrt{s \log d} + u) \right].$$

By taking $u \asymp \sqrt{s\delta \log d}$, under the scaling $n \gtrsim \delta s \log d$, it follows that $I_3 \lesssim \sigma^2 \sqrt{ns\delta \log d}$ with probability at least $1 - 2d^{-\delta\Omega(s)}$.

Step 2.4. Bounding I_0

It remains to bound $I_0 = \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n \xi_k \mathbf{x}_k^{\top} \mathbf{v}$. Bounding I_0 is similar to establishing the ‘‘limited projection distortion (LPD)’’ property in [94], but the key distinction is that $\boldsymbol{\theta}$ and \mathbf{v} in I_0 take value in different spaces.

The main difficulty associated with ‘‘sup $_{\boldsymbol{\theta}}$ ’’ lies in the discontinuity of the quantization noise $\xi_k := \mathcal{Q}_{\Delta}(\tilde{y}_k + \tau_k) - \tilde{y}_k$, which we overcome by a covering argument and some machinery developed in [94, Prop. 6.1]. However, the essential difference from [94] is that we use Lemma 10 to handle ‘‘sup $_{\mathbf{v}}$ ’’, while [94] again used covering argument for \mathbf{v} to strengthen their Proposition 6.1 to their Proposition 6.2, which is unfortunately insufficient in our setting because the covering

number of \mathcal{V} significantly increases under smaller covering radius (on the other hand, using covering argument for \mathbf{v} suffices for the analyses in [94] regarding a different estimator named *projected back projection*).

Let us first construct a ρ -net of Σ_{s, R_0} denoted by $\mathcal{G} = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N\}$, so that for any $\boldsymbol{\theta} \in \Sigma_{s, R_0}$ we can pick $\boldsymbol{\theta}' \in \mathcal{G}$ satisfying $\|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2 \leq \rho$; here, the covering radius ρ is to be chosen later, and we assume that $N \leq \left(\frac{9d}{\rho s}\right)^s$ [76, Lemma 3.3]. As is standard in a covering argument, we first control I_0 over the net \mathcal{G} (by replacing ‘‘sup $_{\boldsymbol{\theta} \in \Sigma_{s, R_0}}$ ’’ with ‘‘sup $_{\boldsymbol{\theta} \in \mathcal{G}}$ ’’), and then bound the approximation error induced by such replacement.

Step 2.4.1. Bounding I_0 over \mathcal{G}

In this step, we want to bound $I_{0, \mathcal{G}} := \sup_{\boldsymbol{\theta} \in \mathcal{G}} \sup_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n \xi_k \mathbf{x}_k^{\top} \mathbf{v}$. First let us consider a fixed $\boldsymbol{\theta} \in \Sigma_{s, R_0}$. Then since $|\xi_k| \leq \Delta$, we have $\|\xi_k \mathbf{x}_k\|_{\psi_2} \lesssim \Delta \sigma$. Because $\{\xi_k \mathbf{x}_k : k \in [n]\}$ are independent zero mean, by [92, Prop. 2.6.1] we have $\left\| \sum_{k=1}^n \xi_k \mathbf{x}_k \right\|_{\psi_2} \lesssim \sqrt{n} \Delta \sigma$. Define $\mathcal{V}' = \mathcal{V} \cup \{0\}$, then for any $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{V}'$ we have

$$\left\| \left(\sum_{k=1}^n \xi_k \mathbf{x}_k \right)^{\top} \mathbf{v}_1 - \left(\sum_{k=1}^n \xi_k \mathbf{x}_k \right)^{\top} \mathbf{v}_2 \right\|_{\psi_2} \leq C\sqrt{n} \Delta \sigma \|\mathbf{v}_1 - \mathbf{v}_2\|_2.$$

Thus, by Lemma 10, it holds with probability at least $1 - \exp(-u^2)$ that

$$\begin{aligned} & \sup_{\mathbf{v} \in \mathcal{V}} \left(\sum_{k=1}^n \xi_k \mathbf{x}_k \right)^{\top} \mathbf{v} \\ & \leq \sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{V}'} \left| \left(\sum_{k=1}^n \xi_k \mathbf{x}_k \right)^{\top} \mathbf{v} - \left(\sum_{k=1}^n \xi_k \mathbf{x}_k \right)^{\top} \mathbf{v}' \right| \\ & \leq C\sqrt{n} \Delta \sigma (\omega(\mathcal{V}') + u) \leq C_1 \sqrt{n} \Delta \sigma \left(\sqrt{s \log \frac{9d}{s}} + u \right). \end{aligned}$$

Moreover, by a union bound over \mathcal{G} , we obtain that

$$I_{0, \mathcal{G}} \lesssim \sigma \Delta \sqrt{n} \left(\sqrt{s \log \frac{9d}{s}} + u \right)$$

holds with probability at least $1 - \exp(s \log \frac{9d}{\rho s} - u^2)$. We set $u \asymp \sqrt{s\delta \log \left(\frac{9d}{\rho s}\right)}$ and arrive at

$$\mathbb{P} \left(I_{0, \mathcal{G}} \lesssim \sigma \Delta \sqrt{ns\delta \log \frac{9d}{\rho s}} \right) \geq 1 - \left(\frac{9d}{\rho s}\right)^{-\Omega(\delta s)}. \quad (52)$$

Step 2.4.2. Bounding the Approximation Error

From now on we indicate the dependence of ξ_k on $\boldsymbol{\theta}$ by using the notation $\xi_{\boldsymbol{\theta}, k} := \mathcal{Q}_{\Delta}(\tilde{y}_{\boldsymbol{\theta}, k} + \tau_k) - \tilde{y}_{\boldsymbol{\theta}, k}$ where $\tilde{y}_{\boldsymbol{\theta}, k} = \mathcal{T}_{\zeta_y}(\mathbf{x}_k^{\top} \boldsymbol{\theta} + \epsilon_k)$. For any $\boldsymbol{\theta} \in \Sigma_{s, R_0}$ we pick one $\boldsymbol{\theta}' \in \mathcal{G}$ such that $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 \leq \rho$; we fix such correspondence and remember that from now on every $\boldsymbol{\theta} \in \Sigma_{s, R_0}$ is associated with some $\boldsymbol{\theta}' \in \mathcal{G}$, (which of course depends on $\boldsymbol{\theta}$ but our notation omits such dependence). Thus, we can bound $I_0 = \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n \xi_{\boldsymbol{\theta}, k} \mathbf{x}_k^{\top} \mathbf{v}$ as

$$\begin{aligned} I_0 &\leq \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n \xi_{\boldsymbol{\theta}', k} \mathbf{x}_k^{\top} \mathbf{v} + \sup_{\boldsymbol{\theta}, \mathbf{v}} \sum_{k=1}^n (\xi_{\boldsymbol{\theta}, k} - \xi_{\boldsymbol{\theta}', k}) \mathbf{x}_k^{\top} \mathbf{v} \\ &\leq I_{0, \mathcal{G}} + I_{01}. \end{aligned} \quad (53)$$

Note that the bound on $I_{0,\mathcal{G}}$ is available in (52), so it remains to bound $I_{01} := \sup_{\theta, \mathbf{v}} \sum_{k=1}^n (\xi_{\theta,k} - \xi_{\theta',k}) \mathbf{x}_k^\top \mathbf{v}$, which can be understood as the approximation error of the net \mathcal{G} regarding the empirical process of interest. To facilitate the presentation we switch to the more compact notations — let $\mathbf{X} \in \mathbb{R}^{n \times d}$ with rows \mathbf{x}_k^\top be the sensing matrix, $\xi_{\theta} = [\xi_{\theta,k}] \in \mathbb{R}^n$ be the quantization error indexed by θ , $\tau = [\tau_k] \in \mathbb{R}^n$ be the random dither vector, $\epsilon = [\epsilon_k] \in \mathbb{R}^n$ be the heavy-tailed noise vector, $\mathbf{y}_{\theta} = [y_{\theta,k}] = \mathbf{X}\theta + \epsilon \in \mathbb{R}^n$ and $\tilde{\mathbf{y}}_{\theta} = [\tilde{y}_{\theta,k}] = \mathcal{T}_{\zeta_y}(\mathbf{y}_{\theta})$ be the measurement vector and truncated measurement vector, respectively. With these conventions we can write $I_{01} = \sup_{\theta, \mathbf{v}} (\xi_{\theta} - \xi_{\theta'})^\top \mathbf{X}\mathbf{v}$. Recall that a specific θ' has been specified for each $\theta \in \Sigma_{s,R_0}$, so defining $\Psi_{\theta} := \xi_{\theta} - \xi_{\theta'}$ allows us to write $I_{01} = \sup_{\theta, \mathbf{v}} \Psi_{\theta}^\top \mathbf{X}\mathbf{v}$. Further, we define $\tilde{\Psi}_{\theta} := \tilde{\mathbf{y}}_{\theta} - \tilde{\mathbf{y}}_{\theta'}$, $\hat{\Psi}_{\theta} := \mathcal{Q}_{\Delta}(\tilde{\mathbf{y}}_{\theta} + \tau) - \mathcal{Q}_{\Delta}(\tilde{\mathbf{y}}_{\theta'} + \tau)$ and make the following observation

$$\Psi_{\theta} = \xi_{\theta} - \xi_{\theta'} = \hat{\Psi}_{\theta} - \tilde{\Psi}_{\theta}. \quad (54)$$

We pause to establish a property of \mathbf{X} that holds w.h.p.. Specifically, we restrict (37) to $\mathbf{v} \in \mathcal{V}$ (recall that $\tilde{\Delta} = 0$ and so $\tilde{\mathbf{X}} = \mathbf{X}$ and $\tilde{\Sigma} = \Sigma^*$), then under the promised probability it holds for some $c(\kappa_0, \sigma)$ that

$$\sup_{\mathbf{v} \in \mathcal{V}} \left| \frac{\|\mathbf{X}\mathbf{v}\|_2}{\sqrt{n}} - \|\sqrt{\Sigma^*}\mathbf{v}\|_2 \right| \leq c(\kappa_0, \sigma) \sqrt{\frac{\delta s \log d}{n}}.$$

Thus, when $n \gtrsim \delta s \log d$ with for large enough hidden constant depending on (κ_0, σ) , it holds that

$$\sup_{\mathbf{v} \in \mathcal{V}} \frac{\|\mathbf{X}\mathbf{v}\|_2}{\sqrt{n}} \leq \sup_{\mathbf{v} \in \mathcal{V}} \|\sqrt{\Sigma^*}\mathbf{v}\|_2 + \sqrt{\kappa_1} \leq 2\sqrt{\kappa_1}. \quad (55)$$

We proceed by assuming we are on this event, which allows us to bound I_{01} as

$$\begin{aligned} I_{01} &= \sup_{\theta, \mathbf{v}} \Psi_{\theta}^\top \mathbf{X}\mathbf{v} \leq \sup_{\theta} \|\Psi_{\theta}\|_2 \sup_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_2 \\ &\leq 2\sqrt{\kappa_1 n} \cdot \sup_{\theta} \|\Psi_{\theta}\|_2. \end{aligned} \quad (56)$$

To bound $\sup_{\theta} \|\Psi_{\theta}\|_2$, motivated by (54), we will investigate $\tilde{\Psi}_{\theta}$ and $\hat{\Psi}_{\theta}$ more carefully. We pick a threshold $\eta \in (0, \frac{\Delta}{2})$ (that is to be chosen later), and by the 1-Lipschitzness of $\mathcal{T}_{\zeta_y}(\cdot)$ we have

$$\begin{aligned} \sup_{\theta} \|\tilde{\Psi}_{\theta}\|_2 &= \sup_{\theta} \|\tilde{\mathbf{y}}_{\theta} - \tilde{\mathbf{y}}_{\theta'}\|_2 \leq \sup_{\theta} \|\mathbf{y}_{\theta} - \mathbf{y}_{\theta'}\|_2 \\ &= \sup_{\theta} \|\mathbf{X}(\theta - \theta')\|_2 \leq 2\sqrt{\kappa_1 n} \rho, \end{aligned} \quad (57)$$

where the last inequality is because $\theta - \theta'$ is $2s$ -sparse, hence (55) implies $\|\mathbf{X}(\theta - \theta')\|_2 \leq 2\sqrt{\kappa_1 n} \|\theta - \theta'\|_2 \leq 2\sqrt{\kappa_1 n} \rho$.

To proceed, we will define for specific θ the index vectors $\mathbf{J}_{\theta,1}, \mathbf{J}_{\theta,2} \in \{0,1\}^n$ and use $|\mathbf{J}_{\theta,1}|$ to denote the number of 1s in $\mathbf{J}_{\theta,1}$ (similar meaning for $|\mathbf{J}_{\theta,2}|$). Specifically, using the entry-wise notation $\tilde{\Psi}_{\theta} = [\tilde{\Psi}_{\theta,k}]$ we define $\mathbf{J}_{\theta,1} = [\mathbb{1}(|\tilde{\Psi}_{\theta,k}| \geq \eta)]$. Recall that (57) gives $\sup_{\theta} \|\tilde{\Psi}_{\theta}\|_2^2 \leq 4\kappa_1 n \rho^2$; combined with the simple observation $\sup_{\theta} \|\tilde{\Psi}_{\theta}\|_2^2 \geq \sup_{\theta} \eta^2 |\mathbf{J}_{\theta,1}|$, we obtain a uniform bound on $|\mathbf{J}_{\theta,1}|$ as

$$\sup_{\theta \in \Sigma_{s,R_0}} |\mathbf{J}_{\theta,1}| \leq \frac{4\kappa_1 n \rho^2}{\eta^2}. \quad (58)$$

Next, we define the index vector $\mathbf{J}_{\theta,2}$ for $\theta \in \mathcal{G}$: first let $\mathcal{E}_{\theta,i} = \{\mathcal{Q}_{\Delta}(\tilde{y}_{\theta,i} + \tau_i + t) \text{ is discontinuous in } t \in [-\eta, \eta]\}$, and then we define $\mathbf{J}_{\theta,2} := [\mathbb{1}(\mathcal{E}_{\theta,i})]$. Then by Lemma 11 that we prove later, we have

$$\begin{aligned} \mathbb{P}\left(\sup_{\theta \in \mathcal{G}} |\mathbf{J}_{\theta,2}| \leq \frac{Cn\eta}{\Delta}\right) \\ \geq 1 - \exp\left(-\frac{cn\eta}{\Delta} + s \log \frac{9d}{\rho s}\right) := 1 - \mathcal{P}_1. \end{aligned} \quad (59)$$

Note that $\mathcal{E}_{\theta,i}$ does not happen (i.e., $\mathbb{1}(\mathcal{E}_{\theta,i}) = 0$) means that $\mathcal{Q}_{\Delta}(\tilde{y}_{\theta,i} + \tau_i + t)$ is continuous in $t \in [-\eta, \eta]$; combined with the definition of $\mathcal{Q}_{\Delta}(\cdot)$, this is also equivalent to the statement that “ $\mathcal{Q}_{\Delta}(\tilde{y}_{\theta,i} + \tau_i + t)$ remains constant in $t \in [-\eta, \eta]$.” Thus, given a fixed $\theta \in \Sigma_{s,R_0}$ and its associated θ' , suppose that the i -th entry of $\mathbf{J}_{\theta',2}$ is zero (meaning that “ $\mathcal{Q}_{\Delta}(\tilde{y}_{\theta',i} + \tau_i + t)$ remains constant in $t \in [-\eta, \eta]$ ”), if additionally i -th entry of $\mathbf{J}_{\theta,1}$ is zero (i.e., $|\tilde{\Psi}_{\theta,i}| < \eta$), then the i -th entry of $\hat{\Psi}_{\theta} = \mathcal{Q}_{\Delta}(\tilde{\mathbf{y}}_{\theta} + \tau) - \mathcal{Q}_{\Delta}(\tilde{\mathbf{y}}_{\theta'} + \tau)$ vanishes:

$$\begin{aligned} \hat{\Psi}_{\theta,i} &= \mathcal{Q}_{\Delta}(\tilde{y}_{\theta,i} + \tau_i) - \mathcal{Q}_{\Delta}(\tilde{y}_{\theta',i} + \tau_i) \\ &= \mathcal{Q}_{\Delta}(\tilde{y}_{\theta',i} + \tilde{\Psi}_{\theta,i} + \tau_i) - \mathcal{Q}_{\Delta}(\tilde{y}_{\theta',i} + \tau_i) = 0; \end{aligned}$$

combining with (54), this implies $\Psi_{\theta,i} = -\tilde{\Psi}_{\theta,i}$. Recall from (56) that we want to bound $\sup_{\theta} \|\Psi_{\theta}\|_2$. Write $\mathbf{J}_{\theta,1}^c = \mathbf{1} - \mathbf{J}_{\theta,1}$ and $\mathbf{J}_{\theta',2}^c = \mathbf{1} - \mathbf{J}_{\theta',2}$, then denoting hadamard product by \odot and using the decomposition $\mathbf{1} = \max\{\mathbf{J}_{\theta,1}, \mathbf{J}_{\theta',2}\} + \min\{\mathbf{J}_{\theta,1}^c, \mathbf{J}_{\theta',2}^c\}$ and (54) we have the display (60), where (i) is because entries of Ψ_{θ} equal to those of $-\tilde{\Psi}_{\theta}$ if the index corresponds to $\min\{\mathbf{J}_{\theta,1}^c, \mathbf{J}_{\theta',2}^c\} = 1$, and in (ii) we use the simple bound $\|\Psi_{\theta}\|_{\infty} \leq \|\xi_{\theta}\|_{\infty} + \|\xi_{\theta'}\|_{\infty} \leq 2\Delta$ and the derived bounds on $\sup_{\theta} |\mathbf{J}_{\theta,1}|$, $\sup_{\theta \in \mathcal{G}} |\mathbf{J}_{\theta,2}|$, $\sup_{\theta} \|\tilde{\Psi}_{\theta}\|_2$ in (58), (59) and (57), respectively.

Step 2.4.3. Concluding the Bound on I_0

We are ready to put pieces together, specify ρ, η , and conclude the bound on I_0 . Overall, with probability at least $1 - \mathcal{P}_1 - \mathcal{P}_2$ for \mathcal{P}_1 defined in (59) and some \mathcal{P}_2 within the promised probability, combining (52), (53), (56) and (60) we obtain

$$I_0 \lesssim \sigma \Delta \sqrt{ns \delta \log \frac{9d}{\rho s} + \frac{\kappa_1 n \Delta \rho}{\eta}} + n \sqrt{\kappa_1 \eta \Delta} + \kappa_1 \rho n.$$

Thus, we take the (near-optimal) choice of (ρ, η) as $\rho \asymp \frac{\Delta}{\sqrt{\kappa_1}} \left(\frac{s\delta}{n}\right)^{3/2}$ and $\eta \asymp \frac{\delta \Delta s}{n} \log \frac{9d}{\rho s}$, under which we obtain that, with the promised probability (as \mathcal{P}_1 is also sufficiently small), we obtain the bound on I_0 as

$$I_0 \lesssim \sigma \Delta \sqrt{ns \delta \log \left(\frac{\kappa_1 d^2 n^3}{\Delta^2 s^5 \delta^3}\right)} \quad (61)$$

We can conclude the proof with all the works above. Substituting $\inf_{\mathbf{v} \in \mathcal{V}} \sum_{k=1}^n (\mathbf{x}_k^\top \mathbf{v})^2 \geq \frac{1}{4} \kappa_0 n$ and the definition of I in (46) into (45), then we obtain $\frac{1}{4} \kappa_0 n \|\hat{\mathbf{Y}}\|_2^2 \leq 2\|\hat{\mathbf{Y}}\|_2 I$ that holds uniformly for all $\theta \in \Sigma_{s,R_0}$, which implies $\sup_{\theta} \|\hat{\mathbf{Y}}\|_2 \leq \frac{8I}{\kappa_0 n}$. Substituting the derived bounds on I_1, I_2, I_3, I_0 into (47), with the promised probability we have

$$I \lesssim \sigma \left(\sigma + M^{\frac{1}{2i}}\right) \sqrt{ns \delta \log d} + \sigma \Delta \sqrt{ns \delta \log \left(\frac{\kappa_1 d^2 n^3}{\Delta^2 s^5 \delta^3}\right)},$$

$$\begin{aligned}
\sup_{\theta} \|\Psi_{\theta}\|_2 &= \sup_{\theta} \|\Psi_{\theta} \odot \max\{\mathbf{J}_{\theta,1}, \mathbf{J}_{\theta',2}\} + \Psi_{\theta} \odot \min\{\mathbf{J}_{\theta,1}^c, \mathbf{J}_{\theta',2}^c\}\|_2 \\
&\leq \sup_{\theta} \|\Psi_{\theta} \odot \max\{\mathbf{J}_{\theta,1}, \mathbf{J}_{\theta',2}\}\|_2 + \sup_{\theta} \|\Psi_{\theta} \odot \min\{\mathbf{J}_{\theta,1}^c, \mathbf{J}_{\theta',2}^c\}\|_2 \\
&\stackrel{(i)}{\leq} \sup_{\theta} \|\Psi_{\theta}\|_{\infty} \cdot \sqrt{|\mathbf{J}_{\theta,1}| + |\mathbf{J}_{\theta',2}|} + \sup_{\theta} \|\tilde{\Psi}_{\theta} \odot \min\{\mathbf{J}_{\theta,1}^c, \mathbf{J}_{\theta',2}^c\}\|_2 \\
&\stackrel{(ii)}{\leq} C \left\{ \Delta \sqrt{n} \left(\frac{\sqrt{\kappa_1} \rho}{\eta} + \sqrt{\frac{\eta}{\Delta}} \right) + \rho \sqrt{\kappa_1 n} \right\},
\end{aligned} \tag{60}$$

so the uniform bound on $\|\hat{\mathbf{Y}}\|_2$ follows immediately. Further using $\|\hat{\mathbf{Y}}\|_1 \leq 2\sqrt{s}\|\hat{\mathbf{Y}}\|_2$ completes the proof. \square

Lemma 11. (Bounding $\sup_{\theta \in \mathcal{G}} |\mathbf{J}_{\theta,2}|$). *Along the proof of Theorem 13, it holds that*

$$\mathbb{P}\left(\sup_{\theta \in \mathcal{G}} |\mathbf{J}_{\theta,2}| \geq \frac{Cn\eta}{\Delta}\right) \leq \exp\left(-\frac{cn\eta}{\Delta} + s \log \frac{9d}{\rho s}\right). \tag{62}$$

Proof. Notation and details in the proof of Theorem 13 will be used. We first consider a fixed $\theta \in \mathcal{G}$, and by a simple shifting $\mathcal{E}_{\theta,i}$ happens if and only if $\mathcal{Q}_{\Delta}(\cdot)$ is discontinuous in $[\tilde{y}_{\theta,i} + \tau_i - \eta, \tilde{y}_{\theta,i} + \tau_i + \eta]$, which is also equivalently to

$$[\tilde{y}_{\theta,i} + \tau_i - \eta, \tilde{y}_{\theta,i} + \tau_i + \eta] \cap (\Delta \cdot \mathbb{Z}) = \emptyset.$$

Because $\tau_i \sim \mathcal{U}([-\frac{\Delta}{2}, \frac{\Delta}{2}])$ and $\eta < \frac{\Delta}{2}$, $\mathbb{P}(\mathcal{E}_{\theta,i}) = \frac{2\eta}{\Delta}$ is valid independent of the location of $[\tilde{y}_{\theta,i} - \eta, \tilde{y}_{\theta,i} + \eta]$. Thus, for fixed θ , by conditioning on (\mathbf{X}, ϵ) , $|\mathbf{J}_{\theta,2}|$ follows the binomial distribution with n trials and probability of success $p := \frac{2\eta}{\Delta}$. This allows us to write $|\mathbf{J}_{\theta,2}| = \sum_{k=1}^n J_k$ with J_k i.i.d. following Bernoulli distribution with success probability $\mathbb{E}J_k = p$. Then for any integer $q \geq 2$ we have

$$\begin{aligned}
\sum_{k=1}^n \mathbb{E}|J_k - \mathbb{E}J_k|^q &\leq \sum_{k=1}^n \mathbb{E}|J_k - \mathbb{E}J_k|^2 \\
&\leq np(1-p) \leq \frac{q!}{2} np.
\end{aligned}$$

Now we invoke Bernstein's inequality (Lemma 1) to obtain that for any $t > 0$,

$$\mathbb{P}(|\mathbf{J}_{\theta,2}| - np \geq \sqrt{2npt} + t) \leq \exp(-t).$$

We let $t = cnp$ and take a union bound over $\theta \in \mathcal{G}$; this yields the desired claim since $|\mathcal{G}| \leq \left(\frac{9d}{\rho s}\right)^s$. \square

Michael K. Ng (Senior Member, IEEE) received the B.Sc. and M.Phil. degrees from The University of Hong Kong, Hong Kong, in 1990 and 1992, respectively, and the Ph.D. degree from The Chinese University of Hong Kong, Hong Kong, in 1995. From 1995 to 1997, he was a Research Fellow with the Computer Sciences Laboratory, The Australian National University, Canberra, ACT, Australia. He was an Assistant Professor/Associate Professor with The University of Hong Kong from 1997 to 2005. He was a Professor/Chair Professor (2005-2019) with the Department of Mathematics, Hong Kong Baptist University, Hong Kong, Chair Professor (2019-2023) with the Department of Mathematics, The University of Hong. He is currently a Chair Professor in Mathematics and Chair Professor in Data Science at Hong Kong Baptist University. His research interests include applied and computational mathematics, machine learning and artificial intelligence, and data science. Dr. Ng serves as an editorial board member of several international journals. He was selected for the 2017 Class of Fellows of the Society for Industrial and Applied Mathematics. He received the Feng Kang Prize for his significant contributions to scientific computing.

Di Wang (Member, IEEE) received the Ph.D. degree in computer science from The State University of New York (SUNY) at Buffalo. He is currently an Assistant Professor of computer science and a Faculty Member of statistics with the Division of Computer, Electrical and Mathematical Sciences and Engineering (CEMSE), King Abdullah University of Science and Technology (KAUST). He is also the Principal Investigator (PI) of the Provable Responsible AI and Data Analytics (PRADA) Laboratory, a member of the Computational Bioscience Research Center (CBRC), and an affiliated Faculty Member with the SDAIA-KAUST Center of Excellence in Data Science and Artificial Intelligence (SDAIA-KAUST AI). His research interests include trustworthy machine learning, machine learning theory, and AI for science.

Junren Chen received the B.Sc. degree in Mathematics and Applied Mathematics from Sun Yat-sen University. He is currently pursuing the Ph.D. degree with the Department of Mathematics, The University of Hong Kong. His research interests include compressed sensing, high-dimensional statistics, signal and image processing, quantization, and optimization. He received the Hong Kong Ph.D. Fellowship from the Hong Kong Research Grants Council for supporting his Ph.D. study.