

# Uniform Hanson-Wright type concentration inequalities for unbounded entries via the entropy method

Yegor Klochkov\*

Humboldt Universität zu Berlin  
klochkoy@hu-berlin.de

Nikita Zhivotovskiy†

Higher School of Economics, Now at Google Research, Brain Team.  
nikita.zhivotovskiy@phystech.edu

August 9, 2019

## Abstract

This paper is devoted to uniform versions of the Hanson-Wright inequality for a random vector  $X \in \mathbb{R}^n$  with independent subgaussian components. The core technique of the paper is based on the entropy method combined with truncations of both gradients of functions of interest and of the components of  $X$  itself. Our results recover, in particular, the classic uniform bound of [Talagrand \(1996\)](#) for Rademacher chaoses and the more recent uniform result of [Adameczak \(2015\)](#) which holds under certain rather strong assumptions on the distribution of  $X$ . We provide several applications of our techniques: we establish a version of the standard Hanson-Wright inequality, which is tighter in some regimes. Extending our results we show a version of the dimension-free matrix Bernstein inequality that holds for random matrices with a subexponential spectral norm. We apply the derived inequality to the problem of covariance estimation with missing observations and prove an almost optimal high probability version of the recent result of [Lounici \(2014\)](#). Finally, we show a uniform Hanson-Wright-type inequality in the Ising model under Dobrushin’s condition. A closely related question was posed by [Marton \(2003\)](#).

*Keywords:* concentration inequalities, modified logarithmic Sobolev inequalities, uniform Hanson-Wright inequalities, Rademacher chaos, matrix Bernstein inequality

## 1 Introduction

The concentration properties of quadratic forms of random variables is a classic topic in probability. A well-known result is due to Hanson and Wright (we refer to the form of this inequality presented in [Rudelson and Vershynin \(2013\)](#)), which claims that if  $A$  is an  $n \times n$  real matrix and

---

\*Financial support from the German Research Foundation (DFG) via the International Research Training Group 1792 “High Dimensional Nonstationary Time Series” in Humboldt-Universität zu Berlin is gratefully acknowledged.

†Parts of this work were done while the author was a postdoctoral fellow at Technion. Nikita Zhivotovskiy was supported by RSF grant No. 18-11-00132.

$X = (X_1, \dots, X_n)$  is a random vector in  $\mathbb{R}^n$  with independent centered components satisfying  $\max_i \|X_i\|_{\psi_2} \leq K$  (we will recall the definition of  $\|\cdot\|_{\psi_2}$  below), then for all  $t \geq 0$

$$\mathbb{P}(|X^\top AX - \mathbb{E}X^\top AX| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|A\|_{\text{HS}}^2}, \frac{t}{K^2 \|A\|}\right\}\right), \quad (1.1)$$

for some absolute  $c > 0$ , where  $\|A\|_{\text{HS}} = \sqrt{\sum_{i,j} A_{i,j}^2}$  defines the Hilbert-Schmidt norm and  $\|A\|$  is the operator norm of  $A$ . An important extension of these results is when instead of just one matrix  $A$  we have a family of matrices  $\mathcal{A}$  and want to understand the behaviour of random quadratic forms simultaneously for all matrices in the family. As a concrete example we consider an order-2 Rademacher chaos: given a family  $\mathcal{A} \subset \mathbb{R}^{n \times n}$  of  $n \times n$  real symmetric matrices with zero diagonal, that is for all  $A \in \mathcal{A}$  we have  $A_{ii} = 0$  for all  $i = 1, \dots, n$ , one wants to study the following random variable

$$Z_{\mathcal{A}}(\varepsilon) = \sup_{A \in \mathcal{A}} \sum_{i,j=1}^n A_{ij} \varepsilon_i \varepsilon_j = \sup_{A \in \mathcal{A}} \varepsilon^\top A \varepsilon,$$

where  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$  is a sequence of independent Rademacher signs taking values  $\pm 1$  with equal probabilities. In the celebrated paper [Talagrand \(1996\)](#) it was shown, in particular, that there is an absolute constant  $c > 0$ , such that for any  $t \geq 0$

$$\mathbb{P}(|Z_{\mathcal{A}}(\varepsilon) - \mathbb{E}Z_{\mathcal{A}}(\varepsilon)| \geq t) \leq 2 \exp\left(-c \min\left(\frac{t^2}{(\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|)^2}, \frac{t}{\sup_{A \in \mathcal{A}} \|A\|}\right)\right). \quad (1.2)$$

Similar inequalities in the Gaussian case follow from the results in [Borell \(1984\)](#) and [Arcones and Gine \(1993\)](#). Apart from the new techniques that were used to prove (1.2), the significance of this result is that previously (see, for example, [Ledoux and Talagrand \(2013\)](#)) similar bounds were one-sided and had a multiplicative constant greater than 1 before  $\mathbb{E}Z_{\mathcal{A}}(\varepsilon)$ . Results with a multiplicative factor not equal to 1 are usually called *deviation inequalities* in contrast to *concentration bounds* of the form (1.2) that are studied below. A simplified proof of the upper tail of (1.2), that is the upper bound on  $\mathbb{P}(Z_{\mathcal{A}}(\varepsilon) - \mathbb{E}Z_{\mathcal{A}}(\varepsilon) \geq t)$ , appeared later in [Boucheron et al. \(2003\)](#). We will refer to inequalities of this form as (one-sided) *concentration inequalities*.

It is worth mentioning in advance that our main results are one-sided concentration inequalities. This is because the entropy method, used extensively in our proofs, is known to have some limitations when applied to prove lower tail inequalities (see the discussions in [Ledoux \(2001\)](#); [Boucheron et al. \(2013\)](#)). It would be interesting for future work to consider similar bounds for the lower tails.

Observe that when for every  $A \in \mathcal{A}$  the diagonal elements are zero, the corresponding quadratic forms are centered, that is  $\mathbb{E}\varepsilon^\top A \varepsilon = 0$ . In the general situation we will be interested in the analysis of

$$Z_{\mathcal{A}}(X) = \sup_{A \in \mathcal{A}} (X^\top AX - \mathbb{E}X^\top AX), \quad (1.3)$$

for a random vector  $X$  taking its values in  $\mathbb{R}^n$ . The analysis of both the expectation and the concentration/deviation properties of this random variable has appeared recently in many papers. To name several deviation inequalities: [Krahmer et al. \(2014\)](#) study  $\mathbb{E}Z_{\mathcal{A}}(X)$  and deviations of  $Z_{\mathcal{A}}(X)$  for classes of positive semidefinite matrices with applications to compressive sensing, [Dicker and Erdogdu \(2017\)](#) prove deviation inequalities for  $\sup_{A \in \mathcal{A}} (X^\top AX - \mathbb{E}X^\top AX)$  and subgaussian vectors  $X$  under some extra assumptions. Additionally, a recent paper [Adamczak et al. \(2018b\)](#) studies deviation bounds for  $Z = \|X^\top AX - \mathbb{E}X^\top AX\|$  with Banach space-valued matrices  $A$  and Gaussian variables, providing upper and lower bounds for the moments. The deviation inequality for general subgaussian vectors and a single positive semi-definite matrix was obtained in [Hsu et al. \(2012\)](#). Returning to concentration inequalities, it was shown in [Adamczak \(2015\)](#) that if  $X$  satisfies the so-called *concentration property* with constant  $K$ , that is for every 1-Lipschitz function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  and any  $t \geq 0$  we have  $\mathbb{E}|\varphi(X)| < \infty$  and

$$\mathbb{P}(|\varphi(X) - \mathbb{E}\varphi(X)| \geq t) \leq 2 \exp(-t^2/2K^2), \quad (1.4)$$

then the following bound, similar to (1.2), holds for every  $t \geq 0$ ,

$$\mathbb{P}(|Z_{\mathcal{A}}(X) - \mathbb{E}Z_{\mathcal{A}}(X)| \geq t) \leq 2 \exp \left( -c \min \left( \frac{t^2}{K^2 (\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|)^2}, \frac{t}{K^2 \sup_{A \in \mathcal{A}} \|A\|} \right) \right). \quad (1.5)$$

This result has application to covariance estimation and recovers another recent concentration result of [Koltchinskii and Lounici \(2017\)](#); this is discussed further in Section 2. The drawback of (1.5) is that the required concentration property places strong restrictions on the distribution of  $X$ : while it is satisfied by the standard Gaussian distribution as well as by some log-concave distributions (see [Ledoux \(2001\)](#)), it is not known whether the concentration property holds for general subgaussian entries and even in the simplest case of Rademacher random vectors.

In this paper we extend the aforementioned results in two directions. We extend the result of [Boucheron et al. \(2003\)](#) for bounded variables by allowing non-zero diagonal values of the matrices and unbounded subgaussian variables  $X_i$ . First, let us recall the following definition. For  $\alpha > 0$  denote the  $\psi_\alpha$ -norm of a random variable  $Y$  by

$$\|Y\|_{\psi_\alpha} = \inf \left\{ t \geq 0 : \mathbb{E} \exp \left( \frac{|Y|^\alpha}{t^\alpha} \right) \leq 2 \right\},$$

which is a proper norm whenever  $\alpha \geq 1$ . A random variable  $Y$  with  $\|Y\|_{\psi_1} < \infty$  is referred to as subexponential and  $\|Y\|_{\psi_2} < \infty$  is referred to as subgaussian and the corresponding norm is usually named a subgaussian norm. We also use the  $L_p(P)$  norm. For  $p \geq 1$  we set  $\|Y\|_{L_p} = (\mathbb{E}|Y|^p)^{\frac{1}{p}}$ . One of our main contributions is the following upper-tail bound.

**Theorem 1.1.** *Suppose that the components of  $X = (X_1, \dots, X_n)$  are independent centered random variables and  $\mathcal{A}$  is a finite family of  $n \times n$  real symmetric matrices. Denote  $M = \|\max_i |X_i|\|_{\psi_2}$ . Then, for any  $t \geq \max\{M \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|, M^2 \sup_{A \in \mathcal{A}} \|A\|\}$  we have*

$$\mathbb{P}(Z_{\mathcal{A}}(X) - \mathbb{E}Z_{\mathcal{A}}(X) \geq t) \leq \exp \left( -c \min \left( \frac{t^2}{M^2 (\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|)^2}, \frac{t}{M^2 \sup_{A \in \mathcal{A}} \|A\|} \right) \right), \quad (1.6)$$

where  $c > 0$  is an absolute constant and  $Z_{\mathcal{A}}(X)$  is defined by (1.3).

**Remark 1.1.** In Theorem 1.1 and below we assume that all  $A \in \mathcal{A}$  are symmetric. This was done only for convenience of presentation and in fact, the analysis may be performed for general square matrices. The only difference will be that in many places  $A$  should be replaced by  $\frac{1}{2}(A + A^T)$ .

**Remark 1.2.** Notice that even though the above result is stated for finite sets  $\mathcal{A}$ , it also holds for arbitrary bounded sets. Indeed, for a bounded set of matrices  $\mathcal{A}$ , since these matrices are finite dimensional we can consider an increasing sequence  $\mathcal{A}_1 \subset \mathcal{A}_2 \subset \dots \subset \mathcal{A}$  of finite epsilon-nets of  $\mathcal{A}$  such that the pointwise convergence  $Z_{\mathcal{A}_k}(X) \rightarrow Z_{\mathcal{A}}(X)$  holds. This pointwise convergence implies convergence in probability, in particular,

$$\lim_{k \rightarrow \infty} \mathbb{P}(Z_{\mathcal{A}_k}(X) - \mathbb{E}Z_{\mathcal{A}_k}(X) \geq t) = \mathbb{P}(Z_{\mathcal{A}}(X) - \mathbb{E}Z_{\mathcal{A}}(X) \geq t).$$

Since for a subset  $\mathcal{A}_k \subset \mathcal{A}$  the values  $\mathbb{E} \sup_{A \in \mathcal{A}_k} \|AX\|^2$  and  $\sup_{A \in \mathcal{A}_k} \|A\|$  are not greater than those for the original set  $\mathcal{A}$ , we obtain the bound (1.6) for arbitrary bounded sets. For the sake of simplicity, we only consider finite sets below.

In particular, Theorem 1.1 recovers the right-tail of the result of [Talagrand \(1.2\)](#) up to absolute constants, since in this case we obviously have  $\|\max_i |\varepsilon_i|\|_{\psi_2} \lesssim 1$ . Furthermore, the result of Theorem 1.1 works without the assumption used in [Talagrand \(1996\)](#) and [Boucheron et al. \(2003\)](#) that diagonals of all matrices in  $\mathcal{A}$  are zero. Moreover, it is also applicable in some situations when the concentration property (1.4) holds: indeed, if  $X$  is a standard normal vector in  $\mathbb{R}^n$

then it is well known (see [Ledoux and Talagrand \(2013\)](#)) that  $M = \|\max_i |X_i|\|_{\psi_2} \sim \sqrt{\log n}$ . If moreover the identity matrix  $I_n \in \mathcal{A}$  then  $\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \geq \mathbb{E} \|X\| \gtrsim \sqrt{n}$ . Therefore, in this case the factor  $M$  is only of at most logarithmic order when compared to  $\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|$ .

In the special case that  $\mathcal{A}$  consists of just one matrix, our bound recovers the bound that is similar to the original Hanson-Wright inequality. On the one hand, our bound may have an extra logarithmic factor that depends on the dimension  $n$ . On the other hand, the original term  $\max_i \|X_i\|_{\psi_2} \|A\|_{\text{HS}}$  is replaced by the better term  $\mathbb{E} \|AX\|$ . We discuss this phenomenon below. The core of the proof of the Hanson-Wright inequality in [Rudelson and Vershynin \(2013\)](#) is based on the decoupling technique which may be used (at least in a straightforward way) to prove the deviation inequality—but not the concentration inequality—for  $\sup_{A \in \mathcal{A}} (X^\top AX - \mathbb{E} X^\top AX)$  in the case that  $\mathcal{A}$  consists of more than one matrix.

A natural question to ask is whether one may improve [Theorem 1.1](#) and replace  $M = \|\max_i |X_i|\|_{\psi_2}$  by  $K = \max_i \|X_i\|_{\psi_2}$ . In [Section 2](#) we discuss that in the deviation version of [Theorem 1.1](#) this replacement is not possible in some cases. This is quite unexpected in light of the fact that  $\|\max_i |X_i|\|_{\psi_2}$  does not appear in the original Hanson-Wright inequality. Therefore, we believe that the form of our result is close to optimal. We also provide the following extension of [Theorem 1.1](#) which may be better in some cases.

**Proposition 1.2.** *Suppose that the components of  $X = (X_1, \dots, X_n)$  are independent centered random variables. Suppose also that the variables  $X_i$  are distributed symmetrically ( $X_i$  has the same distribution as  $-X_i$ ). Let  $\mathcal{A}$  be a finite family of  $n \times n$  real symmetric matrices. Denote  $M = \|\max_i |X_i|\|_{\psi_2}$  and  $K = \max_i \|X_i\|_{\psi_2}$  and let  $G$  be a standard Gaussian vector in  $\mathbb{R}^n$ . Then, for any  $t \geq \max\{MK \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\|, MK \sup_{A \in \mathcal{A}} \|A\|\}$  we have*

$$\mathbb{P}(Z_{\mathcal{A}}(X) - \mathbb{E} Z_{\mathcal{A}}(X) \geq t) \leq \exp \left( -c \min \left( \frac{t^2}{M^2 K^2 (\mathbb{E} \sup_{A \in \mathcal{A}} \|AG\|)^2}, \frac{t}{MK \sup_{A \in \mathcal{A}} \|A\|} \right) \right),$$

where  $c > 0$  is an absolute constant and  $Z_{\mathcal{A}}(X)$  is defined by [\(1.3\)](#).

**Remark 1.3.** [Proposition 1.2](#) is closer to the standard Hanson-Wright inequality [\(1.1\)](#). Indeed, in the case that  $\mathcal{A} = \{A\}$  we have  $\mathbb{E} \|AG\| \sim \|A\|_{\text{HS}}$ . The difference is that  $K^4$  and  $K^2$  are replaced by  $M^2 K^2$  and  $MK$  respectively.

We proceed with some notation that will be used below. For a non-negative random variable  $Y$ , define its *entropy* as

$$\text{Ent}(Y) = \mathbb{E} Y \log Y - \mathbb{E} Y \log \mathbb{E} Y.$$

Instead of the concentration property [\(1.4\)](#), we also discuss the following closely related property:

**Assumption 1.** *We say that a random vector  $X$  taking value in  $\mathbb{R}^n$  satisfies the logarithmic Sobolev inequality with constant  $K > 0$  if for any continuously differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  we have*

$$\text{Ent}(f^2) \leq 2K^2 \mathbb{E} \|\nabla f(X)\|^2, \tag{1.7}$$

whenever both sides of the inequality are not infinite.

One of the technical contributions of this paper is that we use a similar scheme to prove [Theorem 1.1](#) and to recover [\(1.5\)](#) under the logarithmic Sobolev [Assumption 1](#). The application of logarithmic Sobolev inequalities requires computation of the gradient of the function of interest, that is, in our case, the gradient of  $Z_{\mathcal{A}}(X) = \sup_{A \in \mathcal{A}} (X^\top AX - \mathbb{E} X^\top AX)$ . In the analysis that we present, there is a need to control the behaviour of  $\nabla Z_{\mathcal{A}}(X)$  (or its analogs) and, as in [Boucheron et al. \(2003\)](#) and [Adamczak \(2015\)](#), we use a truncation argument to do so. However, in both cases our proofs make use of the *entropy variational formula* of [Boucheron et al. \(2013\)](#), that states that for random variables  $Y, W$  with  $\mathbb{E} \exp(W) < \infty$  we have

$$\mathbb{E}(W \exp(\lambda Y)) \leq \mathbb{E} \exp(\lambda Y) \log(\mathbb{E} \exp(W)) + \text{Ent}(\exp(\lambda Y)). \tag{1.8}$$

Doing so allows us to shorten the proofs and avoid some technicalities appearing in previous papers. Finally, to prove Theorem 1.1 we use a second truncation argument: this argument is based on the Hoffman-Jørgensen inequality (see Ledoux and Talagrand (2013)). We also present two lemmas which are used several times in the text. Both results have short proofs and may be of independent interest.

**Lemma 1.3.** *Suppose that for random variables  $Z, W$  and any  $\lambda > 0$  we have*

$$\text{Ent}(e^{\lambda Z}) \leq \lambda^2 \mathbb{E} W e^{\lambda Z} \quad \text{and} \quad \mathbb{P}(W > L + \theta t) \leq e^{-t}, \quad (1.9)$$

where  $\theta, L$  are positive constants. Then, the following concentration result holds

$$\mathbb{P}(Z - \mathbb{E}Z > t) \leq \exp\left(-c \min\left\{\frac{t^2}{L + \theta}, \frac{t}{\sqrt{\theta}}\right\}\right), \quad (1.10)$$

where  $c > 0$  is an absolute constant. If, moreover, (1.9) holds for any  $\lambda \leq 0$ , we have

$$\mathbb{P}(|Z - \mathbb{E}Z| > t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{L + \theta}, \frac{t}{\sqrt{\theta}}\right\}\right).$$

The second technical result is a version of the convex concentration inequality of Talagrand (1996) which does not require the boundedness of the components of  $X$ .

**Lemma 1.4.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex,  $L$ -Lipschitz function with respect to the Euclidean norm on  $\mathbb{R}^n$  and  $X = (X_1, \dots, X_n)$  be a random vector with independent components. Then, for any  $t \geq CL \|\max_i |X_i|\|_{\psi_2}$  we have*

$$\mathbb{P}(|f(X) - \mathbb{E}f(X)| > t) \leq \exp\left(-c \frac{t^2}{L^2 \|\max_i |X_i|\|_{\psi_2}^2}\right),$$

where  $c, C > 0$  are absolute constants.

Despite generalizing existing results on convex concentration, the result of Lemma 1.4 follows easily from the truncation approach combined with the Hoffman-Jørgensen inequality. As another application of this technique we provide a version of the matrix Bernstein inequality that holds for random matrices with subexponential spectral norm. For clarity of presentation, this inequality is first presented in Section 4. Finally, the same argument showing that it is not possible to replace  $M = \|\max_i |X_i|\|_{\psi_2}$  by  $K = \max_i \|X_i\|_{\psi_2}$  in Theorem 1.1 is used to show that the same is not possible in Lemma 1.4.

We sum up the structure of the paper:

- Section 2 is devoted to applications and discussions and consists of several parts. At first, we give a simple proof of the uniform bound of Adamczak (2015) under the logarithmic Sobolev assumption. The second paragraph is devoted to improvements of the non-uniform Hanson-Wright inequality (1.1) in the subgaussian regime. Furthermore, we apply our techniques to obtain a uniform concentration result similar to Theorem 1.1 in a particular case of non-independent components. We consider the Ising model under Dobrushin's condition, a setting that has been studied recently by Adamczak et al. (2018a) and Götze et al. (2018). The question we study was raised by Marton (2003) in a closely related scenario. Finally, we show that it is not possible in general to replace  $\|\max_i |X_i|\|_{\psi_2}$  with  $\max_i \|X_i\|_{\psi_2}$  in Theorem 1.1 by providing an appropriate counterexample.
- In Section 3 we present our proof of Theorem 1.1. While doing so, we prove Lemma 1.3 and Lemma 1.4. Finally, we give a proof of Proposition 1.2.
- In Section 4 we formulate and prove the dimension-free matrix Bernstein inequality that holds for random matrices with subexponential spectral norm. We demonstrate how our Bernstein inequality can be used in the context of covariance estimation for subgaussian observations improving the state-of-the-art result of Lounici (2014) for covariance estimation with missing observations.

## 2 Some applications and discussions

We begin with some notation that will be used throughout the paper. For a random vector  $X$  taking its values in  $\mathbb{R}^n$  let  $X_1, \dots, X_n$  denote its components. When all components of  $X$  are independent let  $X'_i$  denote an independent copy of the component  $X_i$ . Throughout the paper  $C, c > 0$  are absolute constants that may change from line to line. We write  $a \lesssim b$  if  $a \leq Cb$  for some absolute constant  $C > 0$ . Moreover, if  $a \lesssim b$  and  $b \lesssim a$  we write  $a \sim b$ .

Furthermore, for a square matrix  $A$ , denote by  $\text{Diag}(A)$  the diagonal matrix that has the same elements on the diagonal as  $A$ . The off-diagonal part of  $A$  is defined by  $\text{Off}(A) = A - \text{Diag}(A)$ ; we define  $\text{diag}(\mathbf{a})$  as a  $n \times n$  diagonal matrix with diagonal elements  $\mathbf{a} \in \mathbb{R}^n$ . Finally, for two symmetric (Hermitian) matrices  $A, B$  we write  $A \prec B$  if  $B - A$  is positive-definite and  $A \preceq B$  if  $B - A$  is positive-semidefinite. In what follows we also use the following equivalent formulations of tail inequalities. Assume that for a random variable  $Y$  and some  $a, b > 0$  we have that for any  $t \geq 1$ ,

$$\mathbb{P}(Y > \max(a\sqrt{t}, bt)) \leq e^{-t}.$$

The last inequality implies for any  $u \geq \max(a, b)$ ,

$$\mathbb{P}(Y > u) \leq \exp\left(-\min\left\{\frac{u^2}{a^2}, \frac{u}{b}\right\}\right),$$

and vice versa.

### Uniform Hanson-Wright inequalities under the logarithmic Sobolev condition

In this paragraph we recover a result of [Adamczak \(2015\)](#) under Assumption 1. Consider a random variable  $Z_{\mathcal{A}}(X)$  defined by (1.3), a function of  $X$  that satisfies logarithmic Sobolev assumption (1.7).

Following [Adamczak \(2015\)](#) we assume without loss of generality, that  $\mathcal{A}$  is a finite set of matrices. Then  $Z_{\mathcal{A}}$  is Lebesgue-a.e. differentiable and

$$\|\nabla Z_{\mathcal{A}}(X)\| \leq 2 \sup_{A \in \mathcal{A}} \|AX\|,$$

bounded by a Lipschitz function of  $X$  with good concentration properties.

**Remark 2.1.** Note that Assumption 1 applies only for smooth functions, so that a standard smoothing argument should be used (see e.g. [Ledoux \(2001\)](#)). For the sake of completeness we recall this argument in Section A. In what follows in this section we assume that none of these potential technical problems appear.

In particular, since  $X$  satisfies the logarithmic Sobolev condition with constant  $K$ , we have by Theorem 5.3 in [Ledoux \(2001\)](#) that

$$\mathbb{P}\left(\sup_{A \in \mathcal{A}} \|AX\| \geq \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + K\sqrt{t} \sup_{A \in \mathcal{A}} \|A\|\right) \leq e^{-t}.$$

Taking squares and using  $(a + b)^2 \leq 2a^2 + 2b^2$  we get

$$\mathbb{P}\left(\sup_{A \in \mathcal{A}} \|AX\|^2 \geq 2\left(\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|\right)^2 + 2K^2 \sup_{A \in \mathcal{A}} \|A\|^2 t\right) \leq e^{-t}.$$

Furthermore, the logarithmic Sobolev condition implies for any  $\lambda \in \mathbb{R}$

$$\text{Ent}(e^{\lambda Z_{\mathcal{A}}(X)}) \leq 4K^2 \lambda^2 \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|^2 e^{\lambda Z_{\mathcal{A}}(X)}.$$

Therefore, by Lemma 1.3 it holds for any  $t \geq 0$  that

$$\mathbb{P}\left(|Z_{\mathcal{A}}(X) - \mathbb{E} Z_{\mathcal{A}}(X)| > C\left(K \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \sqrt{t} + K^2 \sup_{A \in \mathcal{A}} \|A\| t\right)\right) \leq 2e^{-t},$$

which coincides with (1.5) for  $K$ -concentrated vectors up to absolute constant factors.

**Remark 2.2.** This result may be used directly to prove the concentration for  $\|\hat{\Sigma} - \Sigma\|$ , where  $\hat{\Sigma}$  is the sample covariance defined as  $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N X_i X_i^\top$  and  $X_1, \dots, X_N$  are centered Gaussian vectors with the covariance matrix  $\Sigma$  (see Theorem 4.1 in Adamczak (2015)). We return to the covariance estimation problem in Section 4.

**Remark 2.3.** We note some additional connections between the convex concentration property (1.4) and Assumption 1. It is known that (1.4) follows from the logarithmic Sobolev inequality by taking  $f(X) = \exp(\lambda(\varphi(X) - \mathbb{E}\varphi(X))/2)$  for  $\lambda > 0$  which implies

$$\text{Ent}(\exp(\lambda(\varphi(X) - \mathbb{E}\varphi(X)))) \leq \frac{K^2 \lambda^2}{2} \mathbb{E} \exp(\lambda(\varphi(X) - \mathbb{E}\varphi(X))).$$

This immediately implies (1.4) via the standard Herbst argument, see Boucheron et al. (2013). Moreover, the last inequality is equivalent to the concentration property. Indeed, from the concentration property we know that  $\|\varphi(X) - \mathbb{E}\varphi(X)\|_{\psi_2} \lesssim K$  and this implies (see van Handel (2016)) that for all  $\lambda \in \mathbb{R}$

$$\text{Ent}(\exp(\lambda(\varphi(X) - \mathbb{E}\varphi(X)))) \lesssim K^2 \lambda^2 \mathbb{E} \exp(\lambda(\varphi(X) - \mathbb{E}\varphi(X))).$$

### Improving Hanson-Wright inequality in the subgaussian regime

Our analysis implies, in particular, an improved version of Hanson-Wright inequality (1.1) in some cases. We consider a centered random vector  $X = (X_1, \dots, X_n)$  with independent subgaussian components and set  $K = \max_i \|X_i\|_{\psi_2}$ ,  $M = \|\max_i |X_i|\|_{\psi_2}$ . In this case (1.1) implies that with probability at least  $1 - 2e^{-t}$  we have

$$X^\top A X - \mathbb{E} X^\top A X \lesssim K^2 \left( \|A\|_{\text{HS}} \sqrt{t} + \|A\| t \right). \quad (2.1)$$

At the same time, Theorem 1.1 for a single matrix  $\mathcal{A} = \{A\}$  implies with the same probability

$$X^\top A X - \mathbb{E} X^\top A X \lesssim M \mathbb{E} \|AX\| \sqrt{t} + M^2 \|A\| t. \quad (2.2)$$

Observe that when  $|X_i| \leq L$  almost surely for every  $i \leq n$ , we have  $M \lesssim \min\{K\sqrt{\log n}, L\}$ . The following example illustrates the difference between these two bounds.

**Example 2.1.** Assume,  $\delta_1, \dots, \delta_n$  are independent Bernoulli random variables with the same mean  $\delta$  and let  $\delta \leq \frac{1}{4}$ . For  $X = (\delta_1 - \delta, \dots, \delta_n - \delta)$  we easily get

$$\mathbb{E} \|AX\| \leq \sqrt{\mathbb{E} X^\top A^2 X} \leq \sqrt{\delta} \|A\|_{\text{HS}}.$$

On the other hand, for  $\delta \leq \frac{1}{4}$  we have

$$\begin{aligned} \|X_1\|_{\psi_2}^2 &= \|\delta_1 - \delta\|_{\psi_2}^2 \sim \sup_{\lambda \in \mathbb{R}} \frac{\log(\mathbb{E} \exp(\lambda(\delta_1 - \delta)))}{\lambda^2} \\ &= \sup_{\lambda \in \mathbb{R}} \frac{\log(\delta \exp(\lambda(1 - \delta)) + (1 - \delta) \exp(-\lambda\delta))}{\lambda^2} = \frac{1 - 2\delta}{4 \log((1 - \delta)/\delta)} \sim \frac{1}{|\log \delta|}, \end{aligned}$$

where the last line follows directly from Theorem 1.1 in Schlemm (2016) (a result equivalent to Theorem 1.1 was also obtained in Berend and Kontorovich (2013)). Therefore, the standard Hanson-Wright inequality implies that with probability at least  $1 - e^{-t}$  we have

$$X^\top A X - \mathbb{E} X^\top A X \lesssim \frac{1}{|\log \delta|} \left( \|A\|_{\text{HS}} \sqrt{t} + \|A\| t \right),$$

while (2.2) and  $M \lesssim \min\{K\sqrt{\log n}, 1\}$  imply that for  $t \geq 1$  and  $\delta \leq \frac{1}{4}$  it holds with probability at least  $1 - 2e^{-t}$  that

$$X^\top A X - \mathbb{E} X^\top A X \lesssim \min \left\{ \sqrt{\frac{\delta \log n}{|\log \delta|}}, \sqrt{\delta} \right\} \|A\|_{\text{HS}} \sqrt{t} + \min \left\{ \frac{\log n}{|\log \delta|}, 1 \right\} \|A\| t. \quad (2.3)$$

It is easy to verify that  $\lim_{\delta \rightarrow 0+} \sqrt{\delta} |\log \delta| = 0$ , thus the inequality (2.3) is better than Hanson-Wright inequality for this  $X$  in the subgaussian regime (when the  $t$ -term is dominated by the  $\sqrt{t}$ -term).

### Uniform concentration results in the Ising model

Consider a random vector  $\sigma \in \{-1, 1\}^n$  with the distribution defined by

$$\pi(\sigma) = \frac{1}{Z'} \exp \left( \sum_{i,j=1}^n J_{ij} \sigma_i \sigma_j - \sum_{i=1}^n h_i \sigma_i \right),$$

where  $Z'$  is a normalizing factor. This distribution defines the *Ising model* with parameters  $J = (J_{ij})_{i,j=1}^n$  and  $\mathbf{h} = (h_i)_{i=1}^n$ . For an arbitrary function  $f$  on  $\{-1, 1\}^n$  denote a difference operator,

$$|\mathfrak{d}f|^2(\sigma) = \frac{1}{2} \sum_{i=1}^n (f(\sigma) - f(T_i \sigma))^2 \pi(-\sigma_i \mid \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots),$$

where the operator  $T_i \sigma = (\sigma_1, \dots, \sigma_{i-1}, -\sigma_i, \sigma_{i+1}, \dots)$  flips the sign of the  $i$ -th component, and  $\pi(\cdot \mid \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots)$  is conditional distribution of the  $i$ -th component given the rest of the elements. The following recent result provides the logarithmic Sobolev inequality for  $\sigma$  under Dobrushin-type conditions.

**Theorem 2.1** (Proposition 1.1, Götze et al. (2018)). *Suppose,  $\|\mathbf{h}\|_\infty \leq \alpha$  and  $J$  satisfies  $J_{ii} = 0$  and*

$$\|J\|_{1 \rightarrow 1} = \max_{i=1, \dots, n} \sum_{j=1}^n |J_{ij}| \leq 1 - \rho \quad (2.4)$$

There is a constant  $C = C(\alpha, \rho)$ , such that for an arbitrary function  $f$  on  $\{-1, 1\}^n$  we have

$$\text{Ent}(f^2) \leq 2C \mathbb{E} |\mathfrak{d}f|^2.$$

**Remark 2.4.** Following Götze et al. (2018) the condition (2.4) will be called *Dobrushin's condition*.

We may obtain the following uniform concentration result which is a simple outcome of our Lemma 1.3 and Theorem 2.1.

**Proposition 2.2.** *Let  $\mathcal{A}$  be a finite set of symmetric matrices with zero diagonal. It holds in the Ising model under Dobrushin's condition and  $\|\mathbf{h}\|_\infty \leq \alpha$  that for any  $t \geq 0$*

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \sigma^\top A \sigma - \mathbb{E} \sup_{A \in \mathcal{A}} \sigma^\top A \sigma \geq t \right) \leq \exp \left( -C \min \left( \frac{t^2}{\left( \mathbb{E} \sup_{A \in \mathcal{A}} \|A \sigma\| + \sup_{A \in \mathcal{A}} \|A\| \right)^2}, \frac{t}{\sup_{A \in \mathcal{A}} \|A\|} \right) \right), \quad (2.5)$$

where  $C$  depends only on  $\alpha, \rho$ .

*Proof.* Let  $\sigma'_{(i)} = (\sigma_1, \dots, \sigma_{i-1}, \sigma'_i, \sigma_{i+1}, \dots)$ , where given all but the  $i$ -th element of  $\sigma$ , the variables  $\sigma_i$  and  $\sigma'_i$  are independent and are distributed according to  $\pi(\cdot \mid \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots)$ . Obviously, we may have all  $\sigma_1, \dots, \sigma_i$  and  $\sigma'_1, \dots, \sigma'_n$  defined on the same discrete probability space, and thus we will use the notation  $\pi(\cdot)$  and  $\pi(\cdot \mid \cdot)$  for the distribution and the conditional distribution. Therefore, we have

$$\begin{aligned} \mathbb{E} |\mathfrak{d}f|^2(\sigma) &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} (f(\sigma) - f(T_i \sigma))^2 \pi(-\sigma_i \mid \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots) \\ &= \sum_{i=1}^n \sum_{\sigma \in \{-1, 1\}^n} \pi(\sigma) \sum_{\sigma'_i \in \{-1, 1\}} (f(\sigma) - f(\sigma'_{(i)}))_+^2 \pi(\sigma'_i \mid \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots) \end{aligned}$$

where we switched from  $\frac{1}{2}(f(\sigma) - f(\sigma'_{(i)}))^2$  to  $(f(\sigma) - f(\sigma'_{(i)}))_+^2$  due to the symmetry between  $\sigma_i$  and  $\sigma'_i$ .



Denoting  $\sigma^{-i} = (\sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots, \sigma_n)$  and using the independence of  $\sigma_i$  and  $\sigma'_i$  given  $\sigma^{-i}$  we observe that  $\pi(\sigma_i, \sigma'_i | \sigma^{-i}) = \pi(\sigma_i | \sigma^{-i})\pi(\sigma'_i | \sigma^{-i})$ . Moreover, it follows from the definition of conditional probability that

$$\begin{aligned} \pi(\sigma)\pi(\sigma'_i | \sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots) &= \pi(\sigma^{-i})\pi(\sigma_i | \sigma^{-i})\pi(\sigma'_i | \sigma^{-i}) \\ &= \pi(\sigma^{-i})\pi(\sigma_i, \sigma'_i | \sigma^{-i}) = \pi(\sigma'_i, \sigma_i, \sigma^{-i}). \end{aligned}$$

Finally, we get

$$\mathbb{E}|\mathfrak{d}f|^2(\sigma) = \sum_{i=1}^n \sum_{(\sigma, \sigma'_i) \in \{-1, 1\}^{n+1}} (f(\sigma) - f(\sigma'_i))_+^2 \pi(\sigma, \sigma'_i) = \sum_{i=1}^n \mathbb{E}(f(\sigma) - f(\sigma'_i))_+^2.$$

Now we want to consider the function

$$Z_{\mathcal{A}}(\sigma) = \sup_{A \in \mathcal{A}} \sigma^\top A \sigma, \quad (2.6)$$

where  $\mathcal{A}$  is a given finite set of symmetric matrices with zero diagonal (the diagonal is not important here, since  $\sigma_i^2 = 1$ ). Let us apply Theorem 2.1 to  $f(\sigma) = e^{\lambda Z_{\mathcal{A}}(\sigma)/2}$ . Since for  $x \geq y$  and  $\lambda \geq 0$  we have  $(e^{\lambda x/2} - e^{\lambda y/2})^2 = e^{\lambda x}(1 - e^{-\lambda(x-y)/2})^2 \leq \frac{\lambda^2}{4} e^{\lambda x}(x-y)^2$ , it holds that

$$\begin{aligned} \mathbb{E}|\mathfrak{d}f|^2(\sigma) &= \mathbb{E} \sum_{i=1}^n (f(\sigma) - f(\sigma'_i))_+^2 = \mathbb{E} e^{\lambda Z_{\mathcal{A}}(\sigma)} \sum_{i=1}^n (1 - e^{-\lambda(Z_{\mathcal{A}}(\sigma) - Z_{\mathcal{A}}(\sigma'_i))/2})_+^2 \\ &\leq \frac{\lambda^2}{4} \mathbb{E} e^{\lambda Z_{\mathcal{A}}(\sigma)} \sum_{i=1}^n (Z_{\mathcal{A}}(\sigma) - Z_{\mathcal{A}}(\sigma'_i))_+^2, \end{aligned}$$

where for  $\tilde{A}$  (maximizer of (2.6)) we have,

$$\begin{aligned} \sum_{i=1}^n (Z_{\mathcal{A}}(\sigma) - Z_{\mathcal{A}}(\sigma'_i))_+^2 &\leq \sum_{i=1}^n (\sigma^\top \tilde{A} \sigma - [\sigma'_i]^\top \tilde{A} \sigma'_i)_+^2 = \sum_{i=1}^n \left( 2(\sigma_i - \sigma'_i) \sum_{j=1}^n \tilde{A}_{ij} \sigma_j \right)_+^2 \\ &\leq 16 \sup_{A \in \mathcal{A}} \|A\sigma\|^2. \end{aligned}$$

Note that concentration for  $\sup_{A \in \mathcal{A}} \|A\sigma\|$  is implied by the same result. Indeed, we have

$$\begin{aligned} \sum_{i=1}^n \left( \sup_{A \in \mathcal{A}, \gamma \in S^{n-1}} \gamma^\top A \sigma - \sup_{A \in \mathcal{A}, \gamma \in S^{n-1}} \gamma^\top A \sigma'_i \right)_+^2 &\leq \sum_{i=1}^n (\tilde{w}^\top \sigma - \tilde{w}^\top \sigma'_i)_+^2 \\ &= \sum_{i=1}^n (\tilde{w}_i(\sigma_i - \sigma'_i))_+^2 \leq 4 \sup_{A \in \mathcal{A}} \|A\|, \end{aligned}$$

where  $\tilde{w}^\top = \gamma^\top A$  is such that  $\sup_{A \in \mathcal{A}} \|A\sigma\| = \tilde{w}^\top \sigma$ . Thus, the expectation of the corresponding difference operator is bounded by  $4 \sup_{A \in \mathcal{A}} \|A\|$ . Therefore, due to the standard Herbst argument (Proposition 6.1 in Boucheron et al. (2013)) Theorem 2.1 implies

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \|A\sigma\| > \mathbb{E} \sup_{A \in \mathcal{A}} \|A\sigma\| + C \sup_{A \in \mathcal{A}} \|A\| \sqrt{t} \right) \leq e^{-t}.$$

To sum up, by Theorem 2.1 we have

$$\text{Ent}(e^{\lambda Z_{\mathcal{A}}(\sigma)}) \leq \lambda^2 \mathbb{E}(4 \sup_{A \in \mathcal{A}} \|A\sigma\|) e^{\lambda Z_{\mathcal{A}}(\sigma)}.$$

It is left to apply Lemma 1.3 which finishes the proof of the following inequality

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \sigma^\top A \sigma - \mathbb{E} \sup_{A \in \mathcal{A}} \sigma^\top A \sigma > C(\sqrt{t} \mathbb{E} \sup_{A \in \mathcal{A}} \|A\sigma\| + (\sqrt{t} + t) \sup_{A \in \mathcal{A}} \|A\|) \right) \geq 1 - e^{-t}, \quad (2.7)$$

where  $C$  only depends on  $\alpha, \rho$  from Theorem 2.1. The claim follows.  $\square$

**Remark 2.5.** In the case that  $\mathcal{A} = \{A\}$  our result implies the upper tail of the recent concentration inequality proved in Adamczak et al. (2018a) (see Theorem 2.2 and Example 2.5). To show this fact (denoting  $\bar{\sigma} = \sigma - \mathbb{E}\sigma$ ) we observe that

$$\mathbb{E}\|A\sigma\| \leq \mathbb{E}\|A\bar{\sigma}\| + \|A\mathbb{E}\sigma\| = \mathbb{E}\|A\bar{\sigma}\| + \left(\sum_{i=1}^n \left(\sum_{j=1}^n A_{i,j}\mathbb{E}\sigma_j\right)^2\right)^{\frac{1}{2}}.$$

Now, it is well known that  $\text{Ent}(f^2) \leq 2c\mathbb{E}|\mathfrak{d}f|^2$  implies the Poincare inequality  $\text{Var}(f) \leq c\mathbb{E}|\mathfrak{d}f|^2$ . Therefore, we have

$$\|\mathbb{E}\bar{\sigma}\bar{\sigma}^\top\| = \sup_{u \in S^{n-1}} \text{Var}(u^\top \bar{\sigma}) \leq (c(\alpha, \rho)/2) \sup_{u \in S^{n-1}} 4\|u\|^2 = 2c(\alpha, \rho).$$

This implies,

$$\mathbb{E}\|A\bar{\sigma}\|^2 = \mathbb{E}\text{Tr}(A^2\bar{\sigma}\bar{\sigma}^\top) \leq \|A\|_{HS}^2 \|\mathbb{E}\bar{\sigma}\bar{\sigma}^\top\| \leq 2c(\rho, \alpha)\|A\|_{HS}^2,$$

where we used that  $\text{Tr}(BD) \leq \text{Tr}(B)\|D\|$  which holds for any pair of symmetric and nonnegative matrices  $B, D$ . Finally, we have

$$\mathbb{E}\|A\sigma\| \leq C(\rho, \alpha)\|A\|_{HS} + \left(\sum_{i=1}^n \left(\sum_{j=1}^n A_{i,j}\mathbb{E}\sigma_j\right)^2\right)^{\frac{1}{2}}.$$

The right-hand side term appears instead of  $\mathbb{E}\|A\sigma\|$  in aforementioned Example 2.5.

### Replacing $\|\max_i |X_i|\|_{\psi_2}$ with $\max_i \|X_i\|_{\psi_2}$ in Theorem 1.1 and Lemma 1.4

Here we show that it is essentially not possible to substitute  $\|\max_i |X_i|\|_{\psi_2}$  with  $\max_i \|X_i\|_{\psi_2}$  in Theorem 1.1 by presenting a concrete counterexample which was kindly suggested by Radosław Adamczak. Suppose the opposite: there is an absolute constant  $C > 0$  such that for any set of matrices  $\mathcal{A}$  and any subgaussian random variables  $X_1, \dots, X_n$  it holds with probability at least  $1 - e^{-t}$  that

$$Z_{\mathcal{A}}(X) \leq C \left( \mathbb{E}Z_{\mathcal{A}}(X) + \max_i \|X_i\|_{\psi_2} \sqrt{t} \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \max_i \|X_i\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\|t \right), \quad (2.8)$$

which implies that for some other constant  $C' > 0$  we have

$$\mathbb{E}^{1/2} Z_{\mathcal{A}}(X)^2 \leq C' \left( \mathbb{E}Z_{\mathcal{A}}(X) + \max_i \|X_i\|_{\psi_2} \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \max_i \|X_i\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\| \right).$$

Notice that here we allow a multiplicative constant not equal to 1 in front of the expectation. Let us take  $\mathcal{A} = \{A^{(1)}, \dots, A^{(n)}\}$  with  $A^{(i)}$  having only one nonzero element  $A_{ii}^{(i)} = 1$ . For the sake of simplicity we take i.i.d.  $X_1, \dots, X_n$  with  $\mathbb{E}X_i^2 = 1$ . This implies

$$Z_{\mathcal{A}}(X) = \max_{i \leq n} (X_i^2 - 1), \quad \sup_{A \in \mathcal{A}} \|AX\| = \max_{i \leq n} |X_i|, \quad \sup_{A \in \mathcal{A}} \|A\| = 1.$$

Assuming that  $\|X_1\|_{\psi_2} \leq 4$  we have

$$\left\| \max_{i \leq n} X_i^2 - 1 \right\|_{L_2} \leq C' \left( \mathbb{E} \max_{i \leq n} (X_i^2 - 1) + 4\mathbb{E} \max_{i \leq n} |X_i| + 16 \right),$$

which, since  $\|\max_{i \leq n} X_i^2\|_{L_1} \geq \max_{i \leq n} \|X_i\|_{L_2} = 1$ , implies

$$\left\| \max_{i \leq n} X_i^2 \right\|_{L_2} \leq 1 + C' (\|\max_{i \leq n} X_i^2\|_{L_1} + 4\mathbb{E} \max_{i \leq n} |X_i| + 15) \leq (1 + 20C') \left\| \max_{i \leq n} X_i^2 \right\|_{L_1}.$$

Note that this inequality also holds if we rescale  $X'_i = \alpha X_i$  for an arbitrary  $\alpha > 0$ . Therefore, if  $\|X_1\|_{\psi_2} \leq 4\|X_1\|_{L_2}$ , we can always rescale our random variables to have  $\|X_1\|_{L_2} = 1$  and  $\|X_1\|_{\psi_2} \leq 4$ , so that the above inequality still holds.

Taking the latter into account we conclude that there is a constant  $D > 0$ , such that if a centered random  $X_1$  satisfies  $\|X_1\|_{\psi_2} \leq 4\|X_1\|_{L_2}$ , then for any  $n \geq 1$  the following inequality holds

$$\|\max_{i \leq n} X_i^2\|_{L_2} \leq D \|\max_{i \leq n} X_i^2\|_{L_1}. \quad (2.9)$$

It is known that such hypercontractivity of maxima implies certain regularity of tails of  $X_1^2$ . In this case by Theorem 4.6 in [Hitczenko et al. \(1998\)](#) for any  $\rho, \varepsilon > 0$  there is another constant  $A = A(D, \rho, \varepsilon) > 1$  such that for every  $t \geq t_0 = \rho\|X_1\|_{L_1}$  we have

$$A\mathbb{P}(X_1^2 > At) \leq \varepsilon\mathbb{P}(X_1^2 > t),$$

so that taking  $\rho = \varepsilon = 1$ , there is  $A = A(D) > 1$  such that for all  $t \geq \|X_1\|_{L_1}$  we have

$$\mathbb{P}(X_1^2 > At) \leq \frac{1}{A}\mathbb{P}(X_1^2 > t). \quad (2.10)$$

The latter does not have to hold for every subgaussian random variable  $X_1$ . For instance, taking a symmetric random variable  $X_1$  with  $\mathbb{P}(|X_1| = 1) = 1 - e^{-r}$  and  $\mathbb{P}(|X_1| = \sqrt{r}) = e^{-r}$  for  $r \geq 4 > 4 \log 2$  we have  $\mathbb{E} \exp\left(\frac{|X_1|^2}{2}\right) = e^{\frac{1}{2}}(1 - e^{-r}) + e^{-r + \frac{r}{2}} \leq e^{\frac{1}{2}} + e^{-\frac{r}{2}} \leq 2$ , which implies  $\|X_1\|_{\psi_2} \leq \sqrt{2}$ . Moreover, for  $r \geq 4$  we also have  $\mathbb{E}X_1^2 \geq 1 - e^{-\frac{r}{2}} \geq \frac{1}{2}$ , thus  $\|X_1\|_{L_2} \geq 1/\sqrt{2}$  and the conditions of (2.9) are satisfied. But for large enough  $r > At$  and for  $t = t_0$ , we have

$$\mathbb{P}(X_1^2 > At) = \mathbb{P}(X_1^2 > t) = e^{-r},$$

therefore breaking the tail regularity (2.10). Therefore, it is impossible to establish an inequality of the form (2.8). We note that it is also possible to prove that (2.9) may not hold for  $X_1$  defined above via some direct calculations.

For the same reason it is not possible to replace  $\|\max_{i \leq n} |X_i|\|_{\psi_2}$  with  $\max_{i \leq n} \|X_i\|_{\psi_2}$  in Lemma 1.4. Indeed, suppose that for any convex  $L$ -Lipschitz function  $f$  we have

$$\mathbb{P}\left(f(X) \leq C(\mathbb{E}f(X) + L \max_{i \leq n} \|X_i\|_{\psi_2} \sqrt{t})\right) \leq e^{-t}.$$

Taking  $f(X) = \max_{i \leq n} |X_i|$ , which is convex and 1-Lipschitz, we get

$$\|\max_{i \leq n} X_i^2\|_{L_2} = \|\max_{i \leq n} |X_i|\|_{L_4} \leq C' \left( \mathbb{E} \max_{i \leq n} |X_i| + \max_{i \leq n} \|X_i\|_{\psi_2} \right).$$

The same choice of  $X_1$  implies (2.9) and leads to a contradiction.

### 3 Proof of Theorem 1.1

In this section we assume that the components of  $X$  are independent. We recall that  $X'_i$  denotes an independent copy of the component  $X_i$ . The main tool of the proof is the modified logarithmic Sobolev inequality (see Theorem 2 in [Boucheron et al. \(2003\)](#) or Theorem 6.15 in [Boucheron et al. \(2013\)](#)). For the sake of brevity we denote  $Z = Z_{\mathcal{A}}(X)$  in this section. Let us set

$$Z'_i = Z_{\mathcal{A}}(X^{(i)}), \quad X^{(i)} = (X_1, \dots, X_{i-1}, X'_i, X_i, \dots, X_n).$$

Then by the symmetrized version of the inequality we have that for any  $\lambda$ ,

$$\text{Ent}(e^{\lambda Z}) \leq \sum_{i=1}^n \mathbb{E} e^{\lambda Z} \tau(-\lambda(Z - Z'_i)_+),$$

where  $\tau(x) = x(e^x - 1)$ . Since  $\tau(x) \leq x^2$  for  $x \leq 0$ , we have for all  $\lambda \geq 0$ ,

$$\text{Ent}(e^{\lambda Z}) \leq \lambda^2 \mathbb{E} V_+ e^{\lambda Z}, \quad V_+ := \mathbb{E}' \sum_{i=1}^n (Z - Z'_i)_+^2. \quad (3.1)$$

The right-hand side of the inequality can be “decoupled” by the variational entropy formula (1.8), as it is done in the proof of Lemma 1.3 which was presented in the introduction.

*Proof of Lemma 1.3.* We have

$$\text{Ent}(e^{\lambda Z}) \leq \lambda^2 L \mathbb{E} e^{\lambda Z} + \lambda^2 \mathbb{E} (W - L)_+ e^{\lambda Z}.$$

Due to the deviation bound for  $W$  it holds for some absolute constant  $C > 0$  that

$$\mathbb{E} \exp\left(\frac{(W - L)_+}{C\theta}\right) \leq e.$$

Therefore, by (1.8) we have

$$\mathbb{E} (W - L)_+ / (C\theta) e^{\lambda Z} \leq \mathbb{E} e^{\lambda Z} + \text{Ent}(e^{\lambda Z}),$$

which implies

$$(1 - C\theta\lambda^2) \text{Ent}(e^{\lambda Z}) \leq \lambda^2 (L + C\theta) \mathbb{E} e^{\lambda Z}.$$

By the standard Herbst argument (see e.g., Proposition 6.1 in Boucheron et al. (2013)) we have for any  $0 < \lambda \leq (2C\theta)^{-1/2}$ ,

$$\log \mathbb{E} \exp(\lambda(Z - \mathbb{E}Z)) \leq 2(L + C\theta)\lambda^2.$$

This moment generating function bound is known to immediately imply the right-tail concentration bound (see the properties of subgamma random variables in Boucheron et al. (2013)). Finally, if (1.9) holds for all  $\lambda \in \mathbb{R}$ , the two sided inequality can be derived in the same way.  $\square$

**Remark 3.1.** Note, there is as well a moment version of the modified logarithmic Sobolev inequality, see e.g., Theorem 2 in Boucheron et al. (2005). By this theorem it holds for all  $q \geq 2$  that

$$\|(Z - \mathbb{E}Z)_+\|_{L_q} \leq \sqrt{2\kappa q} \|\sqrt{V_+}\|_{L_q},$$

where  $\kappa < 2$  is an absolute constant. Then if we have

$$\|\sqrt{V_+}\|_{L_q} \leq \sqrt{L} + \sqrt{\theta q}, \quad \forall q \geq 2, \quad (3.2)$$

which is equivalent to the second inequality in (1.9) up to absolute constant factors, then it holds for any  $q \geq 2$

$$\|(Z - \mathbb{E}Z)_+\|_{L_q} \leq \sqrt{4Lq} + \sqrt{4\theta q}.$$

The last inequality implies (1.10) up to absolute constant factors. We note that similar moment computations were used in Boucheron et al. (2005) to analyze the Rademacher chaos. Similarly, one can introduce the moment analog of the logarithmic Sobolev inequality (see equation 3 in Adamczak and Wolff (2015)):

$$\|Z(X) - \mathbb{E}Z(X)\|_{L_q} \leq K\sqrt{q} \|\nabla Z(X)\|_{L_q},$$

where  $K > 0$  is a constant,  $\|\cdot\|$  stands for the standard Euclidean norm and  $q \geq 2$ . Now, if it holds (which in some cases may be derived by the second application of the moment analog of the logarithmic Sobolev inequality)

$$\|\nabla Z(X)\|_{L_q} \leq \mathbb{E} |\nabla Z(X)| + \| |\nabla Z(X)| - \mathbb{E} |\nabla Z(X)| \|_{L_q} \leq \sqrt{L} + K\sqrt{\theta q}, \quad \forall q \geq 2$$

then

$$\|Z - \mathbb{E}Z\|_{L_q} \leq K(\sqrt{Lq} + K\sqrt{\theta q}),$$

which implies the result similar to (1.10).

Finally, we establish a version of our result that requires neither that  $X_i$  is centered nor that  $X_i$  has variance one. It can happen that  $\mathbb{E}X^\top AX \neq \text{Tr}(A)$ , but in fact, the value we subtract does not really affect the concentration properties. In general we can consider the random variable

$$Z = \sup_{A \in \mathcal{A}} (X^\top AX - g(A)), \quad (3.3)$$

where  $g : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  is an arbitrary function.

**Lemma 3.1.** *Suppose that the components  $X_i$  are independent but not necessarily centered, and  $|X_i| \leq K$  almost surely. Then for  $Z$  defined by (3.3) and for any  $t \geq 1$  it holds with probability at least  $1 - e^{-t}$  that*

$$Z - \mathbb{E}Z \leq C \left( K \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\| \right) \sqrt{t} + K^2 \sup_{A \in \mathcal{A}} \|A\| t,$$

where  $C$  is an absolute constant.

*Proof.* Let  $\tilde{A}$  be the matrix that maximizes  $Z(X)$  given  $X$ . We have

$$\begin{aligned} \sum_{i \leq n} (Z - Z'_i)_+^2 &\leq \sum_{i \leq n} \left( X^\top \tilde{A} X - [X^{(i)}]^\top \tilde{A} X^{(i)} \right)^2 \\ &= \sum_{i \leq n} \left( 2(X_i - X'_i) \sum_{j \neq i} \tilde{a}_{ij} X_j + \tilde{a}_{ii}(X_i^2 - X_i'^2) \right)^2 \\ &= \sum_{i \leq n} (X_i - X'_i)^2 \left( 2 \sum_{j \neq i} \tilde{a}_{ij} X_j + \tilde{a}_{ii}(X_i + X'_i) \right)^2 \\ &\leq (2K)^2 \sum_{i \leq n} \left( 2 \sum_j \tilde{a}_{ij} X_j + \tilde{a}_{ii}(X'_i - X_i) \right)^2, \end{aligned}$$

where the last line follows from  $|X_i - X'_i| \leq 2K$ . The factor 2 appears in the second line because  $\tilde{A}$  is symmetric and thus  $X'_i$  is counted twice. Applying the triangle inequality we get

$$V_+ = \mathbb{E}' \sum_{i \leq n} (Z - Z'_i)_+^2 \leq (2K)^2 \mathbb{E}' \sup_{A \in \mathcal{A}} (2\|AX\| + \|\text{Diag}(A)X\| + \|\text{Diag}(A)X'\|)^2,$$

where  $\mathbb{E}'[\cdot] = \mathbb{E}[\cdot | X]$  denotes the expectation with respect to the variables  $X'_1, \dots, X'_n$  only. Thus,

$$V_+ \leq 12K^2 \left( 4 \sup_{A \in \mathcal{A}} \|AX\|^2 + \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|^2 + \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|^2 \right),$$

where we used  $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$ . Since  $|X_i| \leq K$ , we have by convex concentration for Lipschitz functions (see e.g. Theorem 6.10 in [Boucheron et al. \(2013\)](#))

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \|AX\| > \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + 2\sqrt{2}K \sup_{A \in \mathcal{A}} \|A\| \sqrt{t} \right) \leq e^{-t}. \quad (3.4)$$

Using  $(a + b)^2 \leq 2a^2 + 2b^2$  we have

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \|AX\|^2 > 2 \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \right)^2 + 16K^2 \sup_{A \in \mathcal{A}} \|A\|^2 t \right) \leq e^{-t}. \quad (3.5)$$

Similar inequality holds for the term  $\sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|^2$ . Moreover, by the Poincaré inequality (Theorem 3.17 in [Boucheron et al. \(2013\)](#)) we have

$$\begin{aligned} \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|^2 &= \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\| \right)^2 + \text{Var} \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\| \right) \\ &\leq \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\| \right)^2 + (2K)^2 \sup_{A \in \mathcal{A}} \|\text{Diag}(A)\|^2. \end{aligned}$$

Since  $\|\text{Diag}(A)\| \leq \|A\|$ , we get that for  $L \sim K^2 (\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|)^2$  and  $\theta \sim K^4 (\sup_{A \in \mathcal{A}} \|A\|)^2$  we have

$$\mathbb{P}(V_+ > L + \theta t) \leq e^{-t}.$$

Therefore, due to the modified logarithmic Sobolev inequality (3.1) we can use Lemma 1.3. This provides us with the inequality

$$\mathbb{P}(Z - \mathbb{E}Z > C(\sqrt{L} + \theta\sqrt{t} + \sqrt{\theta t})) \leq e^{-t},$$

where we can neglect the  $\theta$  in front of  $\sqrt{t}$  when  $t \geq 1$ .  $\square$

Note that our bound has the term  $\mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\|$  which can be avoided in the case of centered variables  $X_i$ . Therefore, we obtain the bound matching the previous results (1.5) and (1.2).

**Corollary 3.2.** *Suppose that  $|X_i| \leq K$  almost surely and  $\mathbb{E}X_i = 0$ . Then for any  $t \geq 1$  it holds with probability at least  $1 - e^{-t}$  that*

$$Z - \mathbb{E}Z \leq C \left( K \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \sqrt{t} + K^2 \sup_{A \in \mathcal{A}} \|A\| t \right),$$

where  $C > 0$  is an absolute constant.

In the next two lemmas we show how to get rid of the diagonal term. This finishes the proof of the corollary above.

**Lemma 3.3.** *Suppose that  $Y \in \mathbb{R}^n$  has the i.i.d. components with symmetric distribution and let  $\mathcal{B}$  be a finite set of  $n \times n$  positive-semidefinite symmetric matrices. Then we have*

$$\mathbb{E} \sup_{B \in \mathcal{B}} Y^\top \text{Diag}(B)Y \leq \mathbb{E} \sup_{B \in \mathcal{B}} Y^\top BY.$$

*Proof.* Since any  $B \in \mathcal{B}$  is positive-semidefinite,  $\sup_{B \in \mathcal{B}} x^\top Bx$  is a convex function of  $x \in \mathbb{R}^n$ . Moreover,  $Y \stackrel{d}{=} \text{diag}(\varepsilon)Y$  for an independent Rademacher vector  $\varepsilon \in \{1, -1\}^n$ . Therefore, by Jensen's inequality

$$\begin{aligned} \mathbb{E} \sup_{B \in \mathcal{B}} Y^\top BY &= \mathbb{E} \mathbb{E}_\varepsilon \sup_{B \in \mathcal{B}} Y^\top \text{diag}(\varepsilon)B \text{diag}(\varepsilon)Y \\ &\geq \mathbb{E} \sup_{B \in \mathcal{B}} \mathbb{E}_\varepsilon Y^\top \text{diag}(\varepsilon)B \text{diag}(\varepsilon)Y \\ &= \mathbb{E} \sup_{B \in \mathcal{B}} Y^\top \text{Diag}(B)Y, \end{aligned}$$

where  $\mathbb{E}_\varepsilon$  denotes the expectation with respect to  $\varepsilon$  given  $Y$ .  $\square$

**Lemma 3.4.** *For  $X$  with the components that are independent and mean zero we have*

$$\mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)X\| \leq C \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|,$$

where  $C > 0$  is an absolute constant.

*Proof.* Let  $X'$  be an independent copy of  $X$ . By the standard symmetrization argument together with Jensen's inequality and the triangle inequality we have

$$\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \leq \mathbb{E} \sup_{A \in \mathcal{A}} \|A(X - X')\| \leq 2 \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|. \quad (3.6)$$

Observe that  $X - X' \stackrel{d}{=} (X - X') \text{diag}(\varepsilon) = \text{diag}(X - X')\varepsilon$  where  $\varepsilon \in \{1, -1\}^n$  is an independent Rademacher vector. Therefore, we have

$$\mathbb{E} \sup_{A \in \mathcal{A}} \|A(X - X')\| = \mathbb{E} \mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}} \|A \text{diag}(X - X')\varepsilon\|,$$

where  $\mathbb{E}_\varepsilon$  denotes the expectation with respect to  $\varepsilon$ . Conditionally on  $(X - X')$  set  $\mathcal{A}_{X, X'} = \{A \text{diag}(X - X') : A \in \mathcal{A}\}$ . Let  $\mathbf{a}_1, \dots, \mathbf{a}_n$  be the columns of  $A$ . Notice that for any matrix  $A$  we have  $\text{Diag}(A^\top A) = \text{diag}(\|\mathbf{a}_1\|^2, \dots, \|\mathbf{a}_n\|^2) \succeq \text{diag}(A_{11}^2, \dots, A_{nn}^2) = \text{Diag}(A)^2$ . Therefore, by Lemma 3.3 we have

$$\mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}_{X, X'}} \|\text{Diag}(A)\varepsilon\|^2 \leq \mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}_{X, X'}} \|A\varepsilon\|^2. \quad (3.7)$$

We now want to get rid of the squares in (3.7). Let  $\mathcal{B}$  be an arbitrary set of symmetric  $n \times n$  matrices and let us fix some  $B \in \mathcal{B}$ . We have  $\mathbb{E}\|B\varepsilon\|^2 = \|B\|_{HS}^2$  and by Khinchin's inequality we have

$$\mathbb{E}\|B\varepsilon\| \geq \frac{1}{\sqrt{2}}\|B\|_{HS},$$

with the optimal constant due to Szarek (1976). Thus, we have

$$\mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\| \geq \sup_{B \in \mathcal{B}} \mathbb{E}\|B\varepsilon\| \geq \frac{1}{\sqrt{2}} \sup_{B \in \mathcal{B}} \|B\|.$$

Furthermore, by the convex Poincare inequality (Theorem 3.17, Boucheron et al. (2013)) we have,

$$\text{Var}(\sup_{B \in \mathcal{B}} \|B\varepsilon\|) = \mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\|^2 - \left( \mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\| \right)^2 \leq 4 \sup_{B \in \mathcal{B}} \|B\|^2.$$

Therefore,  $\mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\|^2 \leq (\mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\|)^2 + 4 \sup_{B \in \mathcal{B}} \|B\|^2 \leq 9 (\mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\|)^2$  and we get

$$\left( \mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\| \right)^2 \leq \mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\|^2 \leq 9 \left( \mathbb{E} \sup_{B \in \mathcal{B}} \|B\varepsilon\| \right)^2.$$

The last inequality combined with (3.7) implies

$$\mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}_{X, X'}} \|\text{Diag}(A)\varepsilon\| \leq \left( \mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}_{X, X'}} \|\text{Diag}(A)\varepsilon\|^2 \right)^{\frac{1}{2}} \leq 3 \mathbb{E}_\varepsilon \sup_{A \in \mathcal{A}_{X, X'}} \|A\varepsilon\|.$$

Now, taking the expectation with respect to  $X, X'$  and applying (3.6) again we finish the proof.  $\square$

### 3.1 Truncation for unbounded variables

In this section we finish the proof of Theorem 1.1. In order to apply the bounded version of our inequality, we want to truncate each variable  $X_i$ , which can be done by the approach from Adamczak (2008) (see references therein for more details on various applications of this method), where it was used in the context of Talagrand's concentration inequality. Suppose that  $\|\max_i |X_i|\|_{\psi_2} < \infty$  and set

$$Y_i = X_i \mathbf{1}(|X_i| \leq M), \quad W_i = X_i - Y_i, \quad (3.8)$$

with  $M = 8\mathbb{E} \max |X_i|$ . We have,

$$\begin{aligned} Z_{\mathcal{A}}(X) &= \sup_{A \in \mathcal{A}} (Y^\top AY - \mathbb{E}X^\top AX + W^\top AX + W^\top AY) \\ &\leq \sup_{A \in \mathcal{A}} (Y^\top AY - \mathbb{E}X^\top AX) + \sup_{A \in \mathcal{A}} |W^\top AX| + \sup_{A \in \mathcal{A}} |W^\top AY| \\ &\leq \sup_{A \in \mathcal{A}} (Y^\top AY - \mathbb{E}X^\top AX) + \|W\| \sup_{A \in \mathcal{A}} \|AX\| + \|W\| \sup_{A \in \mathcal{A}} \|AY\|. \end{aligned} \quad (3.9)$$

The variables  $Y_i$  are now bounded by the value  $M$ . Therefore, the first term of the last line can be analyzed by Lemma 3.1.

To bound the rest we need to control the deviations of  $\|W\|$ . We have  $\|W\|^2 = W_1^2 + \dots + W_n^2$  is a sum of independent random variables with bounded  $\psi_1$ -norm. Thus, we can control it's

expectation via the Hoffman-Jørgensen inequality. Due to the choice of the truncation level we have by Markov's inequality

$$\mathbb{P}\left(\max_{i \leq n} W_i^2 > 0\right) = \mathbb{P}\left(\max_{i \leq n} |X_i| > M\right) \leq \frac{\mathbb{E} \max_{i \leq n} |X_i|}{M} \leq \frac{1}{8}.$$

Denoting  $S_k = W_1^2 + \dots + W_k^2$  we have  $\|W\|^2 = S_n$ . Then,

$$P\left(\max_{k \leq n} |S_k| > 0\right) \leq P\left(\max_{i \leq n} W_i^2 > 0\right) \leq \frac{1}{8}.$$

Therefore, by Proposition 6.8 in [Ledoux and Talagrand \(2013\)](#) we have

$$\mathbb{E}\|W\|^2 = \mathbb{E}S_n \leq 8\mathbb{E} \max_{i \leq n} W_i^2 \lesssim \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2}^2,$$

where the latter holds since  $\left\| \max_{i \leq n} W_i^2 \right\|_{\psi_1} \leq \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2}^2$ . Furthermore, by Theorem 6.21 in [Ledoux and Talagrand \(2013\)](#) we have

$$\begin{aligned} \left\| \sum_{i=1}^n W_i^2 - \mathbb{E}W_i^2 \right\|_{\psi_1} &\leq K_1 \left( \mathbb{E}\|W\|^2 - \mathbb{E}\|W\|^2 + \left\| \max_{i \leq n} W_i^2 - \mathbb{E}W_i^2 \right\|_{\psi_1} \right) \\ &\leq 2K_1 \left( \mathbb{E}\|W\|^2 + \left\| \max_{i \leq n} W_i^2 \right\|_{\psi_1} \right) \\ &\lesssim \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2}^2, \end{aligned}$$

where  $K_1$  is an absolute constant. Given the bound on the expectation of  $\|W\|^2$  we have

$$\| \|W\| \|_{\psi_2} \lesssim \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2}.$$

Finally, we obtain the deviation bound: for every  $t > 0$  we have

$$\mathbb{P}\left(\|W\| \geq C\sqrt{t} \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2}\right) \leq 2e^{-t}. \quad (3.10)$$

Now we apply Lemma 3.1 to the bounded variables  $Y$ . Notice that our theorem does not require the variables to be centered. This assumption is only used in Corollary 3.2. Taking this into account, Lemma 3.1 can be applied to the variables  $Y$  as follows. Set  $g(A) = \mathbb{E}X^\top AX$  and  $Z_{\mathcal{A}}(Y) = \sup_{A \in \mathcal{A}} (Y^\top AY - g(A))$ . By Lemma 3.1 we have

$$Z_{\mathcal{A}}(Y) - \mathbb{E}Z_{\mathcal{A}}(Y) \lesssim M\sqrt{t} \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AY\| + \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)Y\| \right) + M^2t \sup_{A \in \mathcal{A}} \|A\| \quad (3.11)$$

with probability at least  $1 - e^{-t}$ . Finally, we have to replace the expectations  $\mathbb{E}Z_{\mathcal{A}}(Y)$ ,  $\mathbb{E} \sup_{A \in \mathcal{A}} \|AY\|$  and  $\mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)Y\|$  in (3.11) by their counterparts, taken with respect to  $X$ , as in the original formulation of the result.

First, we want to provide a concentration bound for the convex function  $\sup_{A \in \mathcal{A}} \|AX\|$  that accounts for unbounded variables. As a matter of fact, we prove the following Lemma which is even slightly stronger than Lemma 1.4.

**Lemma 3.5.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be separately convex<sup>1</sup>  $L$ -Lipschitz with respect to the Euclidean norm in  $\mathbb{R}^n$  and  $X = (X_1, \dots, X_n)$  be a random vector with the independent components. Then it holds for  $t \geq 1$  that*

$$\mathbb{P}\left(f(X) > \mathbb{E}f(X) + C \left\| \max_{i \leq n} |X_i| \right\|_{\psi_2} L\sqrt{t}\right) \leq e^{-t},$$

<sup>1</sup>This means that for every  $i = 1, \dots, n$  it is a convex function of  $i$ -th variable if the rest of the variables are fixed. Any convex function is separately convex.



where  $C > 0$  is an absolute constant. Additionally, if  $f$  is convex and  $L$ -Lipschitz, then for any  $t > 0$ ,

$$\mathbb{P}\left(|f(X) - \mathbb{E}f(X)| > C\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} L\sqrt{t}\right) \leq 4e^{-t}.$$

*Proof.* By convex concentration (Theorem 6.10 in [Boucheron et al. \(2013\)](#)) for bounded  $Y_i$  defined by (3.8) we have that for any  $t > 0$ ,

$$\mathbb{P}\left(f(Y) > \mathbb{E}f(Y) + C\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} L\sqrt{t}\right) \leq e^{-t}.$$

Moreover, due to the Lipschitz assumption and (3.10) we have

$$|f(X) - f(Y)| \leq L\|W\| \lesssim L\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \sqrt{1+t},$$

where the latter holds with probability at least  $1 - e^{-t}$ . Integrating these two bounds we also get

$$|\mathbb{E}f(X) - \mathbb{E}f(Y)| \lesssim L\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2}, \quad (3.12)$$

which together implies that with probability at least  $1 - e^{-t}$  we have

$$\begin{aligned} f(X) - \mathbb{E}f(X) &\leq f(Y) - \mathbb{E}f(Y) + |f(X) - f(Y)| + |\mathbb{E}f(X) - \mathbb{E}f(Y)| \\ &\lesssim L\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \sqrt{t}. \end{aligned}$$

The proof of the lower tail bound follows from Theorem 7.12 in [Boucheron et al. \(2013\)](#) and the standard relation between the median and the expectation which holds in our case.  $\square$

From Lemma 3.5 due to the fact that  $\sup_{A \in \mathcal{A}} \|AX\|$  is  $\sup_{A \in \mathcal{A}} \|A\|$ -Lipschitz we have

$$\mathbb{P}\left(\sup_{A \in \mathcal{A}} \|AX\| > \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + C\left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\| \sqrt{t}\right) \leq 2e^{-t}. \quad (3.13)$$

Moreover, similar to (3.12) we have

$$\left|\mathbb{E} \sup_{A \in \mathcal{A}} \|AY\| - \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|\right| \lesssim \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\|. \quad (3.14)$$

Next, we bound the difference between  $\mathbb{E}Z_{\mathcal{A}}(X)$  and  $\mathbb{E}Z_{\mathcal{A}}(Y)$ .

**Lemma 3.6.** *We have*

$$|\mathbb{E}Z_{\mathcal{A}}(Y) - \mathbb{E}Z_{\mathcal{A}}(X)| \lesssim \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\|.$$

*Proof.* Similarly to (3.9) we have

$$\begin{aligned} |\mathbb{E}Z_{\mathcal{A}}(Y) - \mathbb{E}Z_{\mathcal{A}}(X)| &\leq \mathbb{E}\|W\| \sup_{A \in \mathcal{A}} \|AX\| + \mathbb{E}\|W\| \sup_{A \in \mathcal{A}} \|AY\| \\ &\leq \mathbb{E}^{1/2}\|W\|^2 (\mathbb{E}^{1/2} \sup_{A \in \mathcal{A}} \|AX\|^2 + \mathbb{E}^{1/2} \sup_{A \in \mathcal{A}} \|AY\|^2), \end{aligned} \quad (3.15)$$

where by (3.10)  $\mathbb{E}^{1/2}\|W\|^2 \lesssim \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2}$  and by (3.13),

$$\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|^2 \lesssim \left(\mathbb{E} \sup_{A \in \mathcal{A}} \|AX\|\right)^2 + \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\|^2.$$

Taking the square root we get

$$\mathbb{E}^{1/2} \sup_{A \in \mathcal{A}} \|AX\|^2 \lesssim \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \left\|\max_{i \leq n} |X_i|\right\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\|.$$

Similarly and using (3.14) we have,

$$\begin{aligned} \mathbb{E}^{1/2} \sup_{A \in \mathcal{A}} \|AY\|^2 &\lesssim \mathbb{E} \sup_{A \in \mathcal{A}} \|AY\| + \|\max_{i \leq n} |X_i|\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\| \\ &\lesssim \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \|\max_{i \leq n} |X_i|\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\|. \end{aligned}$$

Plugging it in (3.15) we get the required inequality.  $\square$

Therefore, in (3.11) we can use Lemma 3.6 to get

$$\mathbb{E}Z_{\mathcal{A}}(Y) \leq \mathbb{E}Z_{\mathcal{A}}(X) + C \left( \|\max_{i \leq n} |X_i|\|_{\psi_2} \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \|\max_{i \leq n} |X_i|\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\| \right), \quad (3.16)$$

and by Lemma 3.14 (neglecting the diagonal term for centered  $X$  due to Lemma 3.4)

$$\mathbb{E} \sup_{A \in \mathcal{A}} \|AY\| + \mathbb{E} \sup_{A \in \mathcal{A}} \|\text{Diag}(A)Y\| \leq C \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \|\max_{i \leq n} |X_i|\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\| \right). \quad (3.17)$$

Finally, with probability at least  $1 - e^{-t}$  for  $t \geq 1$  we have from (3.9), (3.14) and (3.13)

$$\begin{aligned} |Z_{\mathcal{A}}(X) - Z_{\mathcal{A}}(Y)| &\leq \|W\| \sup_{A \in \mathcal{A}} \|AY\| + \|W\| \sup_{A \in \mathcal{A}} \|AX\| \\ &\lesssim \|W\| \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| + \|W\| \|\max_{i \leq n} |X_i|\|_{\psi_2} \sup_{A \in \mathcal{A}} \|A\| \sqrt{t}, \end{aligned}$$

which using (3.10) turns into

$$|Z_{\mathcal{A}}(X) - Z_{\mathcal{A}}(Y)| \lesssim \|\max_{i \leq n} |X_i|\|_{\psi_2} \mathbb{E} \sup_{A \in \mathcal{A}} \|AX\| \sqrt{t} + \|\max_{i \leq n} |X_i|\|_{\psi_2}^2 \sup_{A \in \mathcal{A}} \|A\| t.$$

Combining the last inequality together with (3.16) and (3.17) we finish the proof of Theorem 1.1.

### 3.2 Proof of Proposition 1.2

The proof is essentially based on the application of the next standard deviation bound instead of the concentration bound of (3.13) in the proof of Theorem 1.1. Since we did not find an exact reference, we derive this inequality here.

**Lemma 3.7.** *Suppose that  $X_1, \dots, X_n$  are independent centered random variables and  $\mathcal{A}$  is a finite set of symmetric matrices. Let  $G$  be a standard normal vector in  $\mathbb{R}^n$ . Then it holds with probability at least  $1 - Ce^{-t}$  that*

$$\sup_{A \in \mathcal{A}} \|AX\| \lesssim \|\max_{i \leq n} |X_i|\|_{\psi_2} \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\| + \sup_{A \in \mathcal{A}} \|A\| \sqrt{t} \right),$$

where  $C > 0$  is an absolute constant.

*Proof.* At first, we observe that  $\sup_{A \in \mathcal{A}} \|AX\| = \sup_{A \in \mathcal{A}, \gamma \in S^{n-1}} \gamma^T AX$ . Consider the metric  $\rho$  defined by  $\rho(a, b) = \|a - b\| \|\max_{i \leq n} |X_i|\|_{\psi_2}$  for any  $a, b \in \mathbb{R}^n$ . By Theorem 2.2.26 in Talagrand (2014) it holds for  $t \geq 0$  and an absolute constant  $C > 0$  that with probability at least  $1 - C \exp(-t)$ ,

$$\sup_{A \in \mathcal{A}, \gamma \in S^{n-1}} \gamma^T AX \lesssim \text{diam}(\mathcal{A}S^{n-1}, \rho) \sqrt{t} + \gamma_2(\mathcal{A}S^{n-1}, \rho),$$

where  $\text{diam}(\mathcal{A}S^{n-1}) = \sup_{x, y \in \mathcal{A}S^{n-1}} \|x - y\| \|\max_{i \leq n} |X_i|\|_{\psi_2} \leq 2 \sup_{A \in \mathcal{A}} \|A\| \|\max_{i \leq n} |X_i|\|_{\psi_2}$  and the functional  $\gamma_2$  is also defined in Talagrand (2014). For the sake of brevity, we will not introduce its definition here. Finally, applying Talagrand's majorizing measure theorem (Theorem 2.4.1 in Talagrand (2014)) we have

$$\gamma_2(\mathcal{A}S^{n-1}, \rho) \lesssim \|\max_{i \leq n} |X_i|\|_{\psi_2} \mathbb{E} \sup_{x \in \mathcal{A}S^{n-1}} x^T G = \|\max_{i \leq n} |X_i|\|_{\psi_2} \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\|.$$

The claim follows.  $\square$

Setting  $M = 8\mathbb{E} \max_i |X_i|$  and  $K = \max_i \|X_i\|_{\psi_2}$  consider the truncation scheme that is used in (3.8). Due to the assumption that all  $X_i$  are symmetrically distributed, we have  $\mathbb{E}Y_i = 0$ . Therefore, Lemma 3.7 implies

$$\mathbb{P} \left( \sup_{A \in \mathcal{A}} \|AY\| > CK \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\| + \sup_{A \in \mathcal{A}} \|A\| \sqrt{t} \right) \right) \leq e^{-t},$$

which can be used instead of the convex concentration inequality (3.4) when dealing with the modified logarithmic Sobolev inequality. Following this proof and using the fact that  $\max_i |Y_i| \leq M$  almost surely, we end up with the concentration bound

$$Z_{\mathcal{A}}(Y) - \mathbb{E}Z_{\mathcal{A}}(Y) \lesssim MK \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\| \sqrt{t} + \sup_{A \in \mathcal{A}} \|A\| t \right),$$

which holds with probability at least  $1 - e^{-t}$  for any  $t > 1$ . Furthermore, we slightly modify the derivations of the previous section by using Lemma 3.7 instead of (3.13). In particular, we get with probability at least  $1 - e^{-t}$  for any  $t > 1$ ,

$$|Z_{\mathcal{A}}(X) - Z_{\mathcal{A}}(Y)| \lesssim MK \left( \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\| \sqrt{t} + \sup_{A \in \mathcal{A}} \|A\| t \right),$$

and taking expectation we also get  $|\mathbb{E}Z_{\mathcal{A}}(X) - \mathbb{E}Z_{\mathcal{A}}(Y)| \lesssim MK \mathbb{E} \sup_{A \in \mathcal{A}} \|AG\|$ . The claim follows from (3.9).

## 4 The matrix Bernstein inequality in the subexponential case

As we mentioned above, one of the prominent applications of the uniform Hanson-Wright inequalities is a recent concentration result in the Gaussian covariance estimation problem. It is known that covariance estimation may be alternatively approached by the matrix Bernstein inequality, see e.g. Wei and Minsker (2017); Lounici (2014). Following the truncation approach, which was taken above, we provide a version of matrix Bernstein inequality that does not require uniformly bounded matrices. The standard version of the inequality (see Tropp (2012) and reference therein) may be formulated as follows: consider random independent matrices  $X_1, \dots, X_N \in \mathbb{R}^{n \times n}$ , such that almost surely  $\max_i \|X_i\| \leq L$ . It holds

$$\mathbb{P} \left( \left\| \sum_{i=1}^N X_i - \mathbb{E}X_i \right\| > u \right) \leq n \exp \left( -c \left( \frac{u^2}{\sigma^2} \wedge \frac{u}{L} \right) \right),$$

where  $c$  is an absolute constant and  $\sigma^2 = \left\| \mathbb{E} \sum_{i=1}^N (X_i - \mathbb{E}X_i)^2 \right\|$ . The first problem with this result is that it does not hold in general cases when  $\max_i \|X_i\|_{\psi_1}$  or  $\max_i \|X_i\|_{\psi_2}$  are bounded. The second problem is the bound depends on the dimension  $n$ . This does not allow to apply this result to operators in infinite-dimensional Hilbert spaces.

For a positive-definite real square matrix  $A$  we define the *effective rank* as  $\tilde{\mathbf{r}}(A) = \frac{\text{Tr}(A)}{\|A\|}$ . We show the following bound.

**Proposition 4.1.** *Consider the set of random independent Hermitian matrices  $X_1, \dots, X_N \in \mathbb{C}^{n \times n}$  such that  $\|X_i\|_{\psi_1} < \infty$ . Set  $M = \|\max_{i \leq N} X_i\|_{\psi_1}$  and let the positive definite Hermitian matrix  $R$  be such that  $\mathbb{E} \sum_{i=1}^N X_i^2 \preceq R$ . Finally, set  $\sigma^2 = \|R\|$ . There are absolute constants  $c, C, c_1 > 0$  such that for any  $u \geq c_1 \max\{M, \sigma\}$  we have*

$$\mathbb{P} \left( \left\| \sum_{i=1}^N X_i - \mathbb{E}X_i \right\| > u \right) \leq C \tilde{\mathbf{r}}(R) \exp \left( -c \left( \frac{u^2}{\sigma^2} \wedge \frac{u}{M} \right) \right).$$

**Remark 4.1.** Using the well known bound for the maximum of subexponential random variables (see [Ledoux and Talagrand \(2013\)](#)) we have

$$\left\| \max_{i \leq N} \|X_i\| \right\|_{\psi_1} \lesssim \log N \max_{i \leq N} \left\| \|X_i\| \right\|_{\psi_1}.$$

Therefore, up to absolute constant factors we may state the bound for  $M = \log N \max_{i \leq N} \left\| \|X_i\| \right\|_{\psi_1}$ . When  $n = 1$  the effective rank plays no role and our bound recovers the version of classical Bernstein inequality which is due to [Adamczak \(2008\)](#). In this paper, it is also shown that the  $\log N$  factor cannot be removed in general. This means that  $M = \left\| \max_{i \leq N} \|X_i\| \right\|_{\psi_1}$  can not be replaced by  $\max_{i \leq N} \left\| \|X_i\| \right\|_{\psi_1}$  in general.

*Proof.* Fix  $U > 0$  and consider the decomposition

$$X_i = Y_i + Z_i, \quad Y_i = X_i \mathbf{1}(\|X_i\| \leq U), \quad Z_i = X_i \mathbf{1}(\|X_i\| > U),$$

so that the matrices  $Y_i$  are uniformly bounded by  $U$  in the operator norm. By the triangle inequality and the union bound we have

$$\mathbb{P} \left( \left\| \sum_{i=1}^N X_i - \mathbb{E}X_i \right\| > 2u \right) \leq \mathbb{P} \left( \left\| \sum_{i=1}^N Y_i - \mathbb{E}Y_i \right\| > u \right) + \mathbb{P} \left( \left\| \sum_{i=1}^N Z_i - \mathbb{E}Z_i \right\| > u \right).$$

Therefore, two parts can be treated separately. Throughout this proof  $c > 0$  is an absolute constant which may change from line to line. It is known that uniformly bounded random matrices satisfy the Bernstein-type inequality (see Theorem 3.1 in [Minsker \(2017\)](#)) for  $u \geq \frac{1}{6}(U + \sqrt{U^2 + 36\sigma^2})$

$$\mathbb{P} \left( \left\| \sum_{i=1}^N Y_i - \mathbb{E}Y_i \right\| > u \right) \leq 14\tilde{\mathfrak{r}} \left( \mathbb{E} \sum_{i=1}^N (Y_i - \mathbb{E}Y_i)^2 \right) \exp \left( - \frac{cu^2}{\left\| \sum_{i=1}^N \mathbb{E}(Y_i - \mathbb{E}Y_i)^2 \right\| + Uu} \right),$$

where we used  $\|Y_i\| \leq U$ . However, since we want to present this bound in terms of  $X_i$  and not  $Y_i$ , we need the following modification of the proof of Minsker's theorem. Using the notation of his proof, it follows from Lemma 3.1 in [Minsker \(2017\)](#):

$$\log \mathbb{E} \exp(\theta(Y_i - \mathbb{E}Y_i)) \leq \frac{\phi(\theta U)}{U^2} \mathbb{E}(Y_i - \mathbb{E}Y_i)^2 \leq \frac{\phi(\theta U)}{U^2} 2\mathbb{E}Y_i^2 \leq \frac{\phi(\theta U)}{U^2} 2\mathbb{E}X_i^2,$$

where  $\phi(t) = e^t - t - 1$ . Now, using the same lines of the proof, instead of formula (3.4) we have

$$\mathbb{E} \text{Tr} \phi \left( \theta \sum_{i=1}^N (Y_i - \mathbb{E}Y_i) \right) \leq \text{Tr} \left( \exp \left( \frac{\phi(\theta U)}{U^2} 2 \sum_{i=1}^N \mathbb{E}X_i^2 \right) - I_d \right),$$

and lines (3.5) with the condition  $\sum_{i=1}^n \mathbb{E}X_i^2 \leq R$  imply

$$\exp \left( \frac{\phi(\theta U)}{U^2} 2 \sum_{i=1}^N \mathbb{E}X_i^2 \right) - I_d \leq \exp \left( \frac{2\phi(\theta U)}{U^2} R \right) - I_d \leq \frac{R}{\sigma^2} \exp \left( \frac{2\phi(\theta U)}{U^2} \sigma^2 \right),$$

where  $\sigma^2 = \|R\|$ . Following the last lines of the proof of Theorem 3.1, we finally have

$$\mathbb{P} \left( \left\| \sum_{i=1}^N Y_i - \mathbb{E}Y_i \right\| > u \right) \leq 14\tilde{\mathfrak{r}}(R) \exp \left( - \frac{cu^2}{\sigma^2 + Uu} \right), \quad (4.1)$$

for  $u \geq C \max\{U, \sigma\}$ .

We proceed with the analysis of  $Z_i$ . Set  $U = 8\mathbb{E}\max_{i \leq N} \|X_i\|$ , then we have by Markov's inequality

$$\mathbb{P}\left(\max_{k \leq N} \left\| \sum_{i=1}^k Z_i \right\| > 0\right) \leq \mathbb{P}\left(\max_{i \leq N} \|Z_i\| > 0\right) = \mathbb{P}\left(\max_{i \leq N} \|X_i\| > U\right) \leq 1/8.$$

Thus, we can apply Proposition 6.8 from [Ledoux and Talagrand \(2013\)](#) to  $Z_i$  taking its values in the Banach space  $(\mathbb{C}^{n \times n}, \|\cdot\|)$  equipped with the spectral norm. We have

$$\mathbb{E} \left\| \sum_{i=1}^N Z_i \right\| \leq 8\mathbb{E} \max_{i \leq N} \|Z_i\|,$$

which implies that for some absolute constant  $K > 0$ ,

$$\mathbb{E} \left\| \sum_{i=1}^N Z_i - \mathbb{E}Z_i \right\| \leq 2\mathbb{E} \left\| \sum_{i=1}^N Z_i \right\| \leq 16\mathbb{E} \max_{i \leq N} \|Z_i\| \leq K \left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1}.$$

Using Theorem 6.21 from [Ledoux and Talagrand \(2013\)](#) in  $(\mathbb{C}^{n \times n}, \|\cdot\|)$  we have,

$$\begin{aligned} \left\| \left\| \sum_{i=1}^N Z_i - \mathbb{E}Z_i \right\| \right\|_{\psi_1} &\leq K_1 \left( \mathbb{E} \left\| \sum_{i=1}^N Z_i - \mathbb{E}Z_i \right\| + \left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1} \right) \\ &\leq K_2 \left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1}, \end{aligned}$$

where  $K_1, K_2 > 0$  are absolute constants. This implies that for any  $u \geq \left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1}$  we have

$$\mathbb{P}\left(\left\| \sum_{i=1}^N Z_i - \mathbb{E}Z_i \right\| > u\right) \leq \exp\left(-\frac{cu}{\left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1}}\right),$$

where  $c > 0$  is an absolute constant. Combining it with (4.1) and that for some absolute  $C > 0$  we have  $U \leq C \left\| \max_{i \leq N} \|X_i\| \right\|_{\psi_1}$  and  $\left\| \max_{i \leq N} \|Z_i\| \right\|_{\psi_1} \leq \left\| \max_{i \leq N} \|X_i\| \right\|_{\psi_1}$ , we prove the claim.  $\square$

To the best of our knowledge, the Proposition 4.1 is the first to combine two important properties: it simultaneously captures the effective rank instead of the dimension  $n$  and is valid for matrices with subexponential operator norm (the matrix Bernstein inequality in the unbounded case was previously granted under the so-called Bernstein moment condition; we refer to [Tropp \(2012\)](#) and the references therein). We should also compare our results with Proposition 2 of [Koltchinskii \(2011\)](#). His inequality has the same form as our bound, but instead of the effective rank, the original dimension  $n$  is used and  $M = \left\| \max_{i \leq N} \|X_i\| \right\|_{\psi_1}$  is replaced by  $\max_{i \leq N} \left\| \|X_i\| \right\|_{\psi_1} \log\left(N \left(\max_{i \leq N} \left\| \|X_i\| \right\|_{\psi_1}\right)^2 / \sigma^2\right)$ .

### An application to covariance estimation with missing observations

Now we turn to the problem studied in [Koltchinskii and Lounici \(2017\)](#) and [Lounici \(2014\)](#). Suppose, we want to estimate the covariance structure of a random subgaussian vector  $X \in \mathbb{R}^n$  (which will be assumed centered) based on  $N$  i.i.d. observations  $X_1, \dots, X_N$ . For the sake of brevity, we work with the finite-dimensional case, while as in [Koltchinskii and Lounici \(2017\)](#) our results do not depend explicitly on the dimension  $n$ . Recall that a centered random vector  $X \in \mathbb{R}^n$  is *subgaussian* if for all  $u \in \mathbb{R}^n$  we have

$$\|\langle X, u \rangle\|_{\psi_2} \lesssim (\mathbb{E}\langle X, u \rangle^2)^{\frac{1}{2}}. \quad (4.2)$$

Observe that this definition does not require any independence of the components of  $X$ .

In what follows we discuss a more general framework suggested by [Lounici \(2014\)](#). Let  $\delta_{i,j}$ ,  $i \leq N, j \leq n$  be independent Bernoulli random variables with the same mean  $\delta$ . We assume that instead of observing  $X_1, \dots, X_N$  we observe vectors  $Y_1, \dots, Y_N$  which are defined as  $Y_i^j = \delta_{i,j} X_i^j$ . This means that some components of vectors  $X_1, \dots, X_N$  are missing (replaced by zero), each with probability  $1 - \delta$ . Since  $\delta$  can be easily estimated, we assume it to be known. Following [Lounici \(2014\)](#), denote

$$\hat{\Sigma}^{(\delta)} = \frac{1}{N} \sum_{i=1}^N Y_i Y_i^\top.$$

It can be easily shown that the estimator

$$\hat{\Sigma} = (\delta^{-1} - \delta^{-2}) \text{Diag}(\hat{\Sigma}^{(\delta)}) + \delta^{-2} \hat{\Sigma}^{(\delta)}$$

is an unbiased estimator of  $\Sigma = \mathbb{E} X_i X_i^\top$ . In particular,

$$\Sigma = (\delta^{-1} - \delta^{-2}) \text{Diag}(\mathbb{E} Y_i Y_i^\top) + \delta^{-2} \mathbb{E} Y_i Y_i^\top. \quad (4.3)$$

**Theorem 4.2.** *Under the assumptions defined above, it holds with probability at least  $1 - e^{-t}$  for  $t \geq 1$  that*

$$\|\hat{\Sigma} - \Sigma\| \lesssim \|\Sigma\| \max \left( \sqrt{\frac{\tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma)}{N\delta^2}}, \sqrt{\frac{t}{N\delta^2}}, \frac{\tilde{\mathbf{r}}(\Sigma) (\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N\delta^2} \right).$$

**Remark 4.2.** The upper bound above provides some important improvements upon Proposition 3 in [Lounici \(2014\)](#), which is

$$\|\hat{\Sigma} - \Sigma\| \lesssim \|\Sigma\| \max \left( \sqrt{\frac{\tilde{\mathbf{r}}(\Sigma) \log n}{N\delta^2}}, \sqrt{\frac{\tilde{\mathbf{r}}(\Sigma) t}{N\delta^2}}, \frac{\tilde{\mathbf{r}}(\Sigma) (\log n + t) (\log N + t)}{N\delta^2} \right) \quad (4.4)$$

The bound (4.4) depends on  $n$  and therefore is not applicable in the infinite dimensional scenarios. It also contains a term proportional to  $t^2$ , which appeared due to a straightforward truncation of each observation. Moreover, this result has an unnecessary factor  $\tilde{\mathbf{r}}(\Sigma)$  in the term  $\sqrt{\frac{\tilde{\mathbf{r}}(\Sigma) t}{N\delta^2}}$ . Finally, when  $\delta = 1$  tighter results may be obtained using high probability generic chaining bounds for quadratic processes. In particular, Theorem 9 in [Koltchinskii and Lounici \(2017\)](#) implies with probability at least  $1 - e^{-t}$ ,

$$\|\hat{\Sigma} - \Sigma\| \lesssim \|\Sigma\| \max \left( \sqrt{\frac{\tilde{\mathbf{r}}(\Sigma)}{N}}, \sqrt{\frac{t}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma)}{N}, \frac{t}{N} \right). \quad (4.5)$$

Unfortunately, this analysis may not be implied for  $\delta < 1$  in general, since the assumption (4.2) does not hold for the vector  $Y$ , defined by  $Y_i^j = \delta_{i,j} X_i^j$ . Therefore, our technique is a reasonable alternative that works for general  $\delta$  and is almost as tight as (4.5) when  $\delta = 1$ . We also remark that for  $\delta = 1$  even sharper versions of (4.5) were obtained in [Mendelson and Zhivotovskiy \(2018\)](#). However, their statistical procedure differs from the sample covariance matrix  $\hat{\Sigma}$ .

To prove Theorem 4.2 we need the following technical lemma, parts of which may as well be found in [Lounici \(2014\)](#).

**Lemma 4.3.** *Let  $X \in \mathbb{R}^n$  satisfy (4.2) with covariance matrix  $\Sigma$ , and  $Y = (\delta_1 X^1, \dots, \delta_n X^n)$  where  $\delta_i$ ,  $i \leq n$  are independent Bernoulli random variables with the same mean  $\delta$ . We have*

$$\|\text{Diag}(YY^\top)\|_{\psi_1} \lesssim \tilde{\mathbf{r}}(\Sigma) \|\Sigma\|, \quad \|\text{Off}(YY^\top)\|_{\psi_1} \lesssim \tilde{\mathbf{r}}(\Sigma) \|\Sigma\|.$$

Additionally, it holds for some absolute constant  $C > 0$  that

$$\mathbb{E} \text{Off}(YY^\top)^2 \preceq C\delta^2 \text{Tr}(\Sigma) (\Sigma + \text{Diag}(\Sigma)), \quad \text{and} \quad \mathbb{E} \text{Diag}(YY^\top)^2 \preceq C\delta \text{Tr}(\Sigma) \text{Diag}(\Sigma). \quad (4.6)$$

*Proof.* Observe that  $\|\text{Diag}(YY^\top)\| \leq \|Y\|^2$  and  $\|\text{Off}(YY^\top)\| \leq \|YY^\top\| + \|\text{Diag}(YY^\top)\| \leq 2\|Y\|^2$ . Therefore,

$$\|\|\text{Off}(YY^\top)\|\|_{\psi_1} \leq 2\|\|Y\|\|_{\psi_2}^2 \leq 2\|\|X\|\|_{\psi_2}^2 \lesssim \text{Tr}(\Sigma),$$

and the same bound holds for  $\|\|\text{Diag}(YY^\top)\|\|_{\psi_1}$ .

Let  $A$  be an arbitrary symmetric matrix and let us calculate  $\mathbb{E}(A \odot \boldsymbol{\delta}\boldsymbol{\delta}^\top)^2$  where  $\odot$  denotes the Hadamard product and  $\boldsymbol{\delta} = (\delta_1, \dots, \delta_n)$  is a vector with independent components having Bernoulli distribution with the same mean  $\delta$ . We have,

$$\left[\mathbb{E}(A \odot \boldsymbol{\delta}\boldsymbol{\delta}^\top)^2\right]_{ii} = \mathbb{E} \sum_k A_{ik} \delta_i \delta_k A_{ki} \delta_i \delta_k = \sum_k A_{ik} A_{ik} \mathbb{E} \delta_i^2 \delta_k^2 = \delta^2 [A^2]_{ii} + (\delta - \delta^2) A_{ii}^2.$$

If  $i \neq j$  we have for the  $i, j$ -th element

$$\begin{aligned} \left[\mathbb{E}(A \odot \boldsymbol{\delta}\boldsymbol{\delta}^\top)^2\right]_{ij} &= \mathbb{E} \sum_k A_{ik} \delta_i \delta_k A_{kj} \delta_j \delta_k = \sum_k A_{ik} A_{kj} \mathbb{E} \delta_i \delta_j \delta_k^2 \\ &= \delta^3 [A^2]_{ij} + (\delta^2 - \delta^3) (A_{ii} A_{ij} + A_{ij} A_{jj}). \end{aligned}$$

This can be put together in the following expression,

$$\begin{aligned} \mathbb{E}(\boldsymbol{\delta}\boldsymbol{\delta}^\top \odot A)^2 &= \delta^3 A^2 + (\delta^2 - \delta^3) [\text{Diag}(A^2) + \text{Off}(A)\text{Diag}(A) + \text{Diag}(A)\text{Off}(A)] \\ &\quad + (\delta - \delta^2) \text{Diag}(A)^2. \end{aligned}$$

Note that all of these matrices are positive semi-definite, apart from the term  $\text{Off}(A)\text{Diag}(A) + \text{Diag}(A)\text{Off}(A)$ , which we can obviously bound by  $\frac{1}{2}(\text{Off}(A) + \text{Diag}(A))^2 = A^2/2$ . Taking into account  $\delta \leq 1$ , we have the following

$$\begin{aligned} \mathbb{E}(\boldsymbol{\delta}\boldsymbol{\delta}^\top \odot A)^2 &\preceq \frac{1}{2}(\delta^3 + \delta^2)A^2 + (\delta^2 - \delta^3)\text{Diag}(A^2) + (\delta - \delta^2)\text{Diag}(A)^2 \\ &\preceq \delta^2(A^2 + \text{Diag}(A^2)) + \delta\text{Diag}(A)^2. \end{aligned}$$

Recall that  $Y = \text{diag}(\boldsymbol{\delta})X$ . Therefore, we have  $\text{Off}(YY^\top) = \boldsymbol{\delta}\boldsymbol{\delta}^\top \odot \text{Off}(XX^\top)$ . Since the latter has zero diagonal, the term with  $\delta$  in the formula above disappears. Therefore,

$$\mathbb{E}\text{Off}(YY^\top)^2 \preceq \delta^2 [\mathbb{E}\text{Off}(XX^\top)^2 + \text{Diag}(\mathbb{E}\text{Off}(XX^\top)^2)]. \quad (4.7)$$

It holds  $\mathbb{E}\text{Off}(XX^\top)^2 \preceq 2\mathbb{E}(XX^\top)^2 + 2\mathbb{E}\text{Diag}(XX^\top)^2$ , and we also have from [Lounici \(2014\)](#) that  $\mathbb{E}(XX^\top)^2 \preceq C\text{Tr}(\Sigma)\Sigma$ . Additionally, due to [\(4.2\)](#) we immediately have  $\mathbb{E}X_i^4 \lesssim \Sigma_{ii}^2$ . Finally, the following bound holds

$$\mathbb{E}\text{Diag}(XX^\top)^2 \preceq C\text{Diag}(\Sigma)^2 \preceq C\text{Tr}(\Sigma)\text{Diag}(\Sigma).$$

Plugging these bounds into [\(4.7\)](#) we get the second inequality. As for the diagonal case we have for  $A = \text{Diag}(XX^\top)$ ,

$$\mathbb{E}\text{Diag}(YY^\top)^2 \preceq C\delta\mathbb{E}\text{Diag}(XX^\top)^2 \preceq C\delta\text{Tr}(\Sigma)\text{Diag}(\Sigma).$$

□

**Lemma 4.4.** *For  $Y$  as in [Lemma 4.3](#) and any unit vector  $u \in \mathbb{R}^n$  we have*

$$\mathbb{E}(u^\top \text{Off}(YY^\top)u)^2 \lesssim \delta^2 \|\Sigma\|^2, \quad \mathbb{E}(u^\top \text{Diag}(YY^\top)u)^2 \lesssim \delta \|\Sigma\|^2.$$

*Proof.* Let  $v \in \mathbb{R}^n$  be an arbitrary unit vector. First, we want to check that

$$\|u^\top \text{Diag}(XX^\top)v\|_{L_4} \lesssim \|\Sigma\|, \quad \|u^\top \text{Off}(XX^\top)v\|_{L_4} \lesssim \|\Sigma\|. \quad (4.8)$$

Obviously,  $\|u^\top XX^\top v\|_{L_4} \leq \|u^\top X\|_{L_8} \|v^\top X\|_{L_8} \lesssim \|\Sigma\|$ , so it is enough to check that only for the diagonal. Let us apply the symmetrization argument. Let  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$  denote independent Rademacher variables. Then,

$$u^\top \text{Diag}(XX^\top)v = \mathbb{E}_\varepsilon \varepsilon^\top \text{diag}(u)XX^\top \text{diag}(v)\varepsilon = \mathbb{E}_\varepsilon u_\varepsilon XX^\top v_\varepsilon,$$

where  $u_\varepsilon = (u_1\varepsilon_1, \dots, u_n\varepsilon_n)^\top$  (and similarly for  $v_\varepsilon$ ) and  $\mathbb{E}_\varepsilon$  denotes the conditional expectation with respect to  $\varepsilon$  given  $X$ . Then, by Jensen's and Hölder's inequalities,

$$\mathbb{E} (u^\top \text{Diag}(XX^\top)v)^4 \leq \mathbb{E} (u_\varepsilon^\top XX^\top v_\varepsilon)^4 = \mathbb{E}_\varepsilon \mathbb{E}^{1/2}[(u_\varepsilon^\top X)^8 \mid \varepsilon] \mathbb{E}^{1/2}[(v_\varepsilon^\top X)^8 \mid \varepsilon] \lesssim \|\Sigma\|^4,$$

thus implying (4.8).

Consider two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ . We show the following bound,

$$\mathbb{E} \left( \sum_{i \neq j} \delta_i \delta_j a_i b_j \right)^2 \leq 18\delta^2 \|\mathbf{a}\|^2 \|\mathbf{b}\|^2 + 2\delta^4 \left( \sum_i a_i \right)^2 \left( \sum_i b_i \right)^2. \quad (4.9)$$

First, using  $\mathbb{E}Z^2 = \mathbb{E}(Z - \mathbb{E}Z)^2 + (\mathbb{E}Z)^2$  and the fact that  $\mathbb{E}\delta_i = \delta$  we have,

$$\mathbb{E} \left( \sum_{i \neq j} \delta_i \delta_j a_i b_j \right)^2 = \mathbb{E} \left( \sum_{i \neq j} (\delta_i - \delta)(\delta_j - \delta) a_i b_j \right)^2 + \left( \sum_{i \neq j} \delta^2 a_i b_j \right)^2.$$

To the first term we apply the decoupling inequality (Theorem 6.1.1 in Vershynin (2016)). Namely, defining  $\delta'_1, \dots, \delta'_n$  as independent copies of  $\delta_1, \dots, \delta_n$  we have,

$$\begin{aligned} \mathbb{E} \left( \sum_{i \neq j} (\delta_i - \delta)(\delta_j - \delta) a_i b_j \right)^2 &\leq 16\mathbb{E} \left( \sum_{i \neq j} (\delta_i - \delta)(\delta'_j - \delta) a_i b_j \right)^2 \\ &= 16 \sum_{i \neq j} \sum_{k \neq l} \mathbb{E}(\delta_i - \delta)(\delta'_j - \delta)(\delta_k - \delta)(\delta'_l - \delta) a_i b_j a_k b_l, \end{aligned}$$

where in the last sum only the terms with  $k = i$  and  $l = j$  do not vanish. Since  $\mathbb{E}(\delta_i - \delta)^2 = \delta(1 - \delta)$ , we have

$$\mathbb{E} \left( \sum_{i \neq j} (\delta_i - \delta)(\delta_j - \delta) a_i b_j \right)^2 \leq 16 \sum_{i \neq j} a_i^2 b_j^2 \delta^2 (1 - \delta)^2 \leq 16\delta^2 \|\mathbf{a}\|^2 \|\mathbf{b}\|^2.$$

It remains to bound the second term. Using  $(x + y)^2 \leq 2x^2 + 2y^2$  together with the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \left( \sum_{i \neq j} \delta^2 a_i b_j \right)^2 &\leq 2\delta^4 \left( \sum_i a_i b_i \right)^2 + 2\delta^4 \left( \sum_{i,j} a_i b_j \right)^2 \\ &\leq 2\delta^4 \|\mathbf{a}\|^2 \|\mathbf{b}\|^2 + 2\delta^4 \left( \sum_i a_i \right)^2 \left( \sum_i b_i \right)^2. \end{aligned}$$

Putting these two bounds together and using  $\delta \leq 1$  we get the required inequality. Since  $u^\top \text{Off}(YY^\top)v = \boldsymbol{\delta}^\top \text{diag}(u) \text{Off}(XX^\top) \text{diag}(v) \boldsymbol{\delta}$ , we can apply (4.9) with  $\mathbf{a} = \text{diag}(u)X$  and  $\mathbf{b} = \text{diag}(u)X$ . This implies

$$\begin{aligned} \mathbb{E} (u^\top \text{Off}(YY^\top)v)^2 &\lesssim \delta^2 \mathbb{E} \|\text{diag}(u)X\|^2 \|\text{diag}(v)X\|^2 + \delta^4 \mathbb{E} (u^\top X)^2 (v^\top X)^2 \\ &\lesssim \delta^2 \mathbb{E}^{1/2} \|\text{diag}(u)X\|^4 \mathbb{E}^{1/2} \|\text{diag}(v)X\|^4 + \delta^4 \mathbb{E}^{1/2} (u^\top X)^4 \mathbb{E}^{1/2} (v^\top X)^4. \end{aligned}$$



Due to (4.2) we have  $\mathbb{E}^{1/4}(u^\top X)^4 \lesssim \|\Sigma\|^{1/2}$ . Moreover, the vector  $\text{diag}(u)X$  also satisfies the subgaussian assumption (4.2) and has the covariance matrix  $\text{diag}(u)\Sigma\text{diag}(u)$ . Therefore, we have

$$\mathbb{E}^{1/2}\|\text{diag}(u)X\|^4 \lesssim \text{Tr}(\text{diag}(u)\Sigma\text{diag}(u)) \lesssim \sum_i u_i^2 \Sigma_{ii} \lesssim \max_i \Sigma_{ii} \lesssim \|\Sigma\|,$$

where we used that  $\|u\| = 1$ . Similar inequalities hold for the vector  $v$ . Therefore, we conclude that

$$\mathbb{E}(u^\top \text{Off}(YY^\top)u)^2 \lesssim \delta^2 \|\Sigma\|^2.$$

Finally, we have for the diagonal term

$$\begin{aligned} \mathbb{E}(u^\top \text{Diag}(YY^\top)v)^2 &= \mathbb{E}\left(\sum_i \delta_i u_i v_i X_i^2\right)^2 = \delta^2 \mathbb{E}(u^\top \text{Diag}(XX^\top)v)^2 + (\delta - \delta^2) \sum_i u_i^2 v_i^2 \mathbb{E}X_i^4 \\ &\lesssim \delta^2 \|\Sigma\|^2 + (\delta - \delta^2) \sum_i u_i^2 v_i^2 \|\Sigma\|^2 \lesssim \delta \|\Sigma\|^2. \end{aligned}$$

□

Before we present the proof of the deviation bound, let us recall the following version of Talagrand's concentration inequality for empirical processes. Remarkably, the following result can be proven using very similar techniques: at first, one may use the modified logarithmic Sobolev inequality to prove a version of Talagrand's concentration inequality in the bounded case and then use the same truncation argument as in the proof of Theorem 1.1 to get the result in the unbounded case.

**Theorem 4.5** (Theorem 4 in Adamczak (2008)). *Let  $X_1, \dots, X_N \in \mathcal{X}$  be a sample of independent random variables and  $\mathcal{F}$  be a countable class of measurable functions  $\mathcal{X} \mapsto \mathbb{R}$  such that  $\sup_{f \in \mathcal{F}} \|f(X_i)\|_{\psi_1} < \infty$ . Define*

$$Z_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^N f(X_i) - \mathbb{E}f(X_i) \right| \quad (4.10)$$

and  $\sigma^2 = \sup_{f \in \mathcal{F}} \sum_{i=1}^N \mathbb{E}f^2(X_i)$ . There is an absolute constant  $C > 0$  such that

$$\mathbb{P}(Z_{\mathcal{F}} > 2\mathbb{E}Z_{\mathcal{F}} + t) \leq \exp\left(-\frac{t^2}{4\sigma^2}\right) + 3 \exp\left(-\frac{t}{C \|\max_i \sup_f \|f(X_i)\|_{\psi_1}\|}\right).$$

*Proof of Theorem 4.2.* At first, using (4.3) we have

$$\|\hat{\Sigma} - \Sigma\| \lesssim \delta^{-1} \left\| \text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)}) \right\| + \delta^{-2} \left\| \text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)}) \right\|.$$

Let us apply Proposition 4.1 to the term  $N\text{Off}(\hat{\Sigma}^{(\delta)}) = \sum_{i=1}^N \text{Off}(Y_i Y_i^\top)$ , where

$$R = CN\delta^2 \text{Tr}(\Sigma)(\Sigma + \text{Diag}(\Sigma)).$$

We have  $\tilde{\mathbf{r}}(R) \leq 2\tilde{\mathbf{r}}(\Sigma)$  and  $\|R\| \lesssim N\delta^2 \text{Tr}(\Sigma)\|\Sigma\|$ . Therefore, with probability at least  $1 - e^{-t}$ ,

$$\begin{aligned} \|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\| &\lesssim \max\left(\sqrt{\frac{\delta^2 \text{Tr}(\Sigma)\|\Sigma\|(\log \tilde{\mathbf{r}}(\Sigma) + t)}{N}}, \frac{\text{Tr}(\Sigma)(\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N}\right) \\ &= \|\Sigma\| \max\left(\sqrt{\frac{\delta^2 \tilde{\mathbf{r}}(\Sigma)(\log \tilde{\mathbf{r}}(\Sigma) + t)}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma)(\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N}\right). \end{aligned} \quad (4.11)$$

Integrating this bound (see e.g. Theorem 2.3 in [Boucheron et al. \(2013\)](#)) we easily get

$$\mathbb{E}\|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\| \lesssim \|\Sigma\| \max\left(\sqrt{\frac{\delta^2 \tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma)}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma) \log N}{N}\right).$$

Finally, we apply Theorem 4.5 to the set of functions  $\mathcal{F}$  indexed by  $\gamma \in S^{n-1}$  and defined by

$$f_\gamma(X_i) = \gamma^\top \text{Off}(Y_i Y_i^\top) \gamma,$$

so that  $Z_{\mathcal{F}} = N\|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\|$  in (4.10). Then, by Lemma 4.4 we have  $\sigma^2 \lesssim \delta^2 N \|\Sigma\|^2$  and by Lemma 4.3  $\|\max_i \sup_f |f(X_i)|\|_{\psi_1} = \|\max_i \|\text{Off}(Y_i Y_i^\top)\|_{\psi_1}\| \lesssim \tilde{\mathbf{r}}(\Sigma) \|\Sigma\| \log N$ , so that with probability at least  $1 - e^{-t}$  for  $t \geq 1$ ,

$$\begin{aligned} \|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\| &\lesssim \mathbb{E}\|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\| + \delta \|\Sigma\| \sqrt{\frac{t}{N}} + \|\Sigma\| \frac{\tilde{\mathbf{r}}(\Sigma) t \log N}{N} \\ &\lesssim \|\Sigma\| \max\left(\sqrt{\frac{\delta^2 \tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma)}{N}}, \sqrt{\frac{\delta^2 t}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma) (\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N}\right). \end{aligned}$$

We proceed with the diagonal term. Applying Proposition 4.1 to the sum  $N \text{Diag}(\hat{\Sigma}^{(\delta)}) = \sum_{i=1}^N \text{Diag}(Y_i Y_i^\top)$  with  $R = CN\delta \text{Tr}(\Sigma) \text{Diag}(\Sigma)$  we have  $\tilde{\mathbf{r}}(R) \lesssim \tilde{\mathbf{r}}(\Sigma)$  and  $\|R\| \lesssim N\delta \text{Tr}(\Sigma) \|\Sigma\|$ . Therefore, with probability at least  $1 - e^{-t}$  we have

$$\|\text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)})\| \lesssim \|\Sigma\| \max\left(\sqrt{\frac{\delta \tilde{\mathbf{r}}(\Sigma) (\log \tilde{\mathbf{r}}(\Sigma) + t)}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma) (\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N}\right). \quad (4.12)$$

Again, integrating this inequality we get

$$\mathbb{E}\|\text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)})\| \lesssim \|\Sigma\| \max\left(\sqrt{\frac{\delta \tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma)}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma) \log N}{N}\right).$$

We have  $\mathbb{E}(u^\top \text{Diag}(Y_i Y_i^\top) u)^2 \lesssim \delta \|\Sigma\|^2$  and  $\|\max_i \|\text{Off}(Y_i Y_i^\top)\|_{\psi_1}\| \lesssim \tilde{\mathbf{r}}(\Sigma) \|\Sigma\| \log N$  by Lemma 4.4 and Lemma 4.3 respectively. By Theorem 4.5 we have with probability at least  $1 - e^{-t}$ ,

$$\begin{aligned} \|\text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)})\| &\lesssim \mathbb{E}\|\text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)})\| + \|\Sigma\| \sqrt{\frac{\delta t}{N}} + \|\Sigma\| \frac{\tilde{\mathbf{r}}(\Sigma) t \log N}{N} \\ &\lesssim \|\Sigma\| \max\left(\sqrt{\frac{\delta \tilde{\mathbf{r}}(\Sigma) \log \tilde{\mathbf{r}}(\Sigma)}{N}}, \sqrt{\frac{\delta t}{N}}, \frac{\tilde{\mathbf{r}}(\Sigma) (\log \tilde{\mathbf{r}}(\Sigma) + t) \log N}{N}\right). \end{aligned}$$

Finally, we combine the off-diagonal and diagonal bounds via the triangle inequality and get

$$\|\hat{\Sigma} - \Sigma\| \leq \delta^{-2} \|\text{Off}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Off}(\hat{\Sigma}^{(\delta)})\| + \delta^{-1} \|\text{Diag}(\hat{\Sigma}^{(\delta)}) - \mathbb{E}\text{Diag}(\hat{\Sigma}^{(\delta)})\|.$$

□

## Acknowledgement

We are indebted to Radosław Adamczak for his very useful feedback at several stages of this paper. We are especially grateful for his suggestion to study the question of Marton in the Ising model and for providing us with an important example in Section 2.

## References

- Adamczak, R. (2008). A tail inequality for suprema of unbounded empirical processes with applications to markov chains. *Electronic Journal of Probability*, 13:1000–1034.
- Adamczak, R. (2015). A note on the Hanson-Wright inequality for random vectors with dependencies. *Electronic Communications in Probability*, 20.
- Adamczak, R., Kotowski, M., Polaczyk, B., and Strzelecki, M. (2018a). A note on concentration for polynomials in the Ising model. *arXiv:1809.03187*.
- Adamczak, R., Latała, R., and Meller, R. (2018b). Hanson-Wright inequality in Banach spaces. *arXiv:1811.00353*.
- Adamczak, R. and Wolff, P. (2015). Concentration inequalities for non-Lipschitz functions with bounded derivatives of higher order. *Probability Theory and Related Fields*, 162:531–586.
- Arcones, M. and Gine, E. (1993). On decoupling, series expansions, and tail behavior of chaos processes. *Journal of Theoretical Probability*, 6:101–122.
- Berend, D. and Kontorovich, A. (2013). On the concentration of the missing mass. *Electron. Commun. Probab.*, 18,(3):7 pp.
- Borell, C. (1984). On the Taylor series of a Wiener polynomial. Seminar Notes on multiple stochastic integration, polynomial chaos and their integration. *Case Western Reserve Univ., Cleveland*.
- Boucheron, S., Bousquet, O., Lugosi, G., and Massart, P. (2005). Moment inequalities for functions of independent random variables. *The Annals of Probability*, 33(2):514–560.
- Boucheron, S., Lugosi, G., and Massart, P. (2003). Concentration inequalities using the entropy method. *The Annals of Probability*, 31(3):1583–1614.
- Boucheron, S., Lugosi, G., and Massart, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- Dicker, L. H. and Erdogdu, M. (2017). Flexible results for quadratic forms with applications to variance components estimation. *The Annals of Statistics*, 45(1):386–414.
- Götze, F., Sambale, H., and Sinulis, A. (2018). Higher order concentration for functions of weakly dependent random variables. *arXiv:1801.06348*.
- Hitczenko, P., Kwapien, S., Li, W., Schechtman, G., Schlumprecht, T., and Zinn, J. (1998). Hypercontractivity and Comparison of Moments of Iterated Maxima and Minima of Independent Random Variables. *Electronic Journal of Probability*, 3.
- Hsu, D., Kakade, S., and Zhang, T. (2012). A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.*, 17(52):6 pp.
- Koltchinskii, V. (2011). Von Neumann entropy penalization and low-rank matrix estimation. *Annals of Statistics*, 39(6):2936–2973.
- Koltchinskii, V. and Lounici, K. (2017). Concentration inequalities and moment bounds for sample covariance operators. *Bernoulli*, 23(1):110–133.
- Krahmer, F., Mendelson, S., and Rauhut., H. (2014). Suprema of Chaos Processes and the Restricted Isometry Property. *Communications in Pure and Applied Mathematics*, 67:1877–1904.
- Ledoux, M. (2001). *The concentration of measure phenomenon*, volume 89 of *Mathematical surveys and Monographs*. American Mathematical Society.

- Ledoux, M. and Talagrand, M. (2013). *Probability in Banach Spaces: Isoperimetry and Processes*. Springer-Verlag Berlin Heidelberg.
- Lounici, K. (2014). High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, 20(3):1029–1058.
- Marton, K. (2003). Measure concentration and strong mixing. *Studia scientiarum mathematicarum hungarica*, 40:95–113.
- Mendelson, S. and Zhivotovskiy, N. (2018). Robust covariance estimation under  $L_4 - L_2$  norm equivalence. *To appear in Annals of Statistics*.
- Minsker, S. (2017). On Some Extensions of Bernstein’s Inequality for Self-adjoint Operators. *Statistics and Probability Letters*, 127:111–119.
- Rudelson, M. and Vershynin, R. (2013). Hanson-Wright inequality and sub-gaussian concentration. *Electron. Commun. Probab.*, 18.
- Schlemm, E. (2016). The Kearns–Saul Inequality for Bernoulli and Poisson-Binomial Distributions. *Journal of Theoretical Probability*, 29:48–62.
- Szarek, S. (1976). On the best constants in the Khinchin inequality. *Studia Mathematica*, 58(2):197–208.
- Talagrand, M. (1996). New concentration inequalities in product spaces. *Inventiones mathematicae*, 126:505–563.
- Talagrand, M. (2014). *Upper and lower bounds for stochastic processes: modern methods and classical problems*, volume 60. Springer Science & Business Media.
- Tropp, J. (2012). User-Friendly Tail Bounds for Sums of Random Matrices. *Foundations of Computational Mathematics*, 12:389–434.
- van Handel, R. (2016). Probability in High Dimension. *Lecture Notes, Princeton University*.
- Vershynin, R. (2016). *High-Dimensional Probability: An Introduction with Applications*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press.
- Wei, X. and Minsker, S. (2017). Estimation of the covariance structure of heavy-tailed distributions. In *Advances in Neural Information Processing Systems*, pages 2859–2868.

## A An approximation argument for non-smooth functions

In order to apply the logarithmic Sobolev assumption (1.7) rigorously we need to take smooth approximations of the function

$$Z_{\mathcal{A}}(X) = \sup_{A \in \mathcal{A}} (X^{\top} A X - \mathbb{E} X^{\top} A X).$$

Notice that we have,

$$|Z_{\mathcal{A}}(X) - Z_{\mathcal{A}}(Y)| \leq \|X - Y\| \left( \sup_{A \in \mathcal{A}} \|A X\| + \sup_{A \in \mathcal{A}} \|A Y\| \right).$$

The following simple lemma shows how to apply the logarithmic Sobolev inequality to non-differentiable functions.

**Lemma A.1.** *Suppose that  $X$  satisfies Assumption 1. Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be such that*

$$|f(x) - f(y)| \leq |x - y| \max(L(x), L(y)),$$

*for some continuous non-negative function  $L(x)$ . Then for some absolute constant  $C > 0$  and any  $\lambda \in \mathbb{R}$  it holds*

$$\text{Ent}(e^{\lambda f}) \leq CK^2 \lambda^2 \mathbb{E}L(x)^2 e^{\lambda f}$$

*Proof.* Set  $h(x) = x^2(1-x)_+^2$  and consider the smoothing kernel supported on the unit ball defined by

$$\phi(u) = \frac{1}{Z_h} h(\|u\|^2), \quad Z_h = \int h(\|u\|^2) du = S_{n-1} \int_0^\infty h(r^2) dr,$$

where  $S_{n-1}$  is a surface area of the unit sphere in  $\mathbb{R}^n$ . Note that since  $\phi$  is radial,  $\nabla\phi(u) = -\nabla\phi(-u)$  and also,

$$\int \|u\| \|\nabla\phi(u)\| du = \frac{2S_{n-1}}{Z_h} \int_0^\infty r^2 |g'(r)| dr = \frac{2 \int_0^\infty r^2 |h'(r)| dr}{\int_0^\infty h(r^2) dr} = C_h.$$

Setting  $\phi_m(u) = m^{-1}\phi(u/m)$  for  $m \in \mathbb{N}$  we have  $\nabla\phi_m(u) = m^{-2}(\nabla\phi)(u/m)$ . Therefore, we have

$$\int \|u\| \|\nabla\phi_m(u)\| du = \int \left\| \frac{u}{m} \right\| \left\| (\nabla\phi) \left( \frac{u}{m} \right) \right\| d\frac{u}{m} = C_h.$$

Take  $F(x) = e^{\lambda f(x)/2}$  and let us consider a sequence of smooth approximations  $F_m(x) = \int \phi_m(x-u) F(u) du$ , so that  $F_m(x)$  tends to  $F$  pointwise. This is possible due to the fact that  $F$  is continuous. Moreover, we have due to the symmetry

$$\begin{aligned} \nabla F_m(x) &= \int (\nabla\phi_m)(x-u) F(u) du = \int (\nabla\phi_m)(u) F(x-u) du \\ &= \frac{1}{2} \int (\nabla\phi_m)(u) [F(x-u) - F(x+u)] du. \end{aligned}$$

Since  $\phi_m(u)$  vanishes for  $\|u\| \geq 1/m$ , we have

$$\begin{aligned} \|\nabla F_m(x)\| &\leq \frac{1}{2} \sup_{\|u\| \leq m^{-1}} \frac{|F(x-u) - F(x+u)|}{\|u\|} \int \|u\| \|\nabla\phi_m(u)\| du \\ &= C_h \sup_{\|u\| \leq m^{-1}} \frac{|F(x-u) - F(x+u)|}{2\|u\|}. \end{aligned}$$

It is easy to see that

$$|F(x) - F(y)| = |e^{\lambda f(x)/2} - e^{\lambda f(y)/2}| \leq \frac{\lambda}{2} \|x - y\| \max(e^{\lambda f(x)/2}, e^{\lambda f(y)/2}) \max(L(x), L(y)),$$

and therefore,

$$\|\nabla F_m(x)\| \leq \frac{\lambda C_h}{2} \tilde{F}_m(x) \times L_m(x),$$

where we set  $L_m(x) = \sup_{y: \|x-y\| \leq m^{-1}} L(y)$  and  $\tilde{F}_m(x) = \sup_{\|x-y\| \leq m^{-1}} e^{\lambda f(y)/2}$  that tend pointwise to  $L(x)$  and  $F(x)$ , respectively, as  $m \rightarrow \infty$ . Since each function  $f_m$  is smooth, we have by Assumption 1 that

$$\text{Ent}(F_m^2) \leq K^2 \mathbb{E} \|\nabla F_m(x)\|^2 \leq \frac{\lambda^2 C_h^2}{4} K^2 \mathbb{E} L_m^2(x) \tilde{F}_m(x)^2.$$

Taking the limit  $m \rightarrow \infty$  we prove the the required inequality.  $\square$