# MARCHENKO-PASTUR LAW WITH RELAXED INDEPENDENCE CONDITIONS

#### JENNIFER BRYSON, ROMAN VERSHYNIN, AND HONGKAI ZHAO

ABSTRACT. We prove the Marchenko-Pastur law for the eigenvalues of  $p \times p$  sample covariance matrices in two new situations where the data does not have independent coordinates. In the first scenario – the block-independent model – the p coordinates of the data are partitioned into blocks in such a way that the entries in different blocks are independent but the entries from the same block may be dependent. In the second scenario – the random tensor model – the data is the homogeneous random tensor of order d, i.e. the coordinates of the data are all  $\binom{n}{d}$  different products of d variables chosen from a set of n independent random variables. We show that Marchenko-Pastur law holds for the block-independent model as long as the size of the largest block is o(p), and for the random tensor model as long as  $d = o(n^{1/3})$ . Our main technical tools are new concentration inequalities for quadratic forms in random variables with block-independent coordinates, and for random tensors.

### 1. Introduction

1.1. Marchenko-Pastur law. Consider a  $p \times m$  random matrix X with independent entries that have zero mean and unit variance. The limiting distribution of eigenvalues  $\lambda_i(W)$  of the sample covariance matrix  $W = \frac{1}{m}XX^{\mathsf{T}}$  is determined by the celebrated Marchenko-Pastur law [41]. This result is valid in the regime where the dimensions of X increase to infinity but the aspect ratio converges to a constant, i.e.  $p \to \infty$  and  $p/m \to \lambda \in (0, \infty)$ . Then, with probability 1, the empirical spectral distribution of the  $p \times p$  matrix W converges weakly to a deterministic distribution that is now called the Marchenko-Pastur law with parameter  $\lambda$ . More specifically, if  $\lambda \in (0,1)$ , then with probability 1 the following holds for each  $x \in \mathbb{R}$ :

$$F^{W}(x) := \frac{1}{p} \# \{ 1 \le i \le p : \lambda_{i}(W) \le x \} \to \int_{\infty}^{x} f_{\lambda}(t) dt$$

where  $f_{\lambda}$  is the Marchenko-Pastur density

(1.1) 
$$f_{\lambda}(x) = \frac{1}{2\pi\lambda x} \sqrt{\left[(\lambda_{+} - x)(x - \lambda_{-})\right]_{+}}, \quad \text{with} \quad \lambda_{\pm} = (1 \pm \sqrt{\lambda})^{2}.$$

A similar result also holds for  $\lambda > 1$ , but in that case the limiting distribution has an additional point mass of  $1-1/\lambda$  at the origin. A straightforward proof of the Marchenko-Pastur law using the Stieltjes transform is given in Chapter 3 of [18]. More extensive expositions of the

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, IRVINE, IRVINE, CA. 92697

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, IRVINE, IRVINE, CA. 92697

DEPARTMENT OF MATHEMATICS, DUKE UNIVERSITY, DURHAM, NC 27708

E-mail addresses: jabryson@uci.edu, rvershyn@uci.edu, zhao@math.duke.edu.

Jennifer Bryson was partially supported by NSF Graduate Research Fellowship Program DGE-1321846. Roman Vershynin is supported by USAF Grant FA9550-18-1-0031, NSF Grants DMS 1954233 and DMS 2027299, and U.S. Army Grant 76649-CS. Hongkai Zhao was partially supported by NSF grants DMS-1622490 and DMS-1821010.

Marchenko-Pastur law with proofs using both the moment method and the Stieltjes transform are given in [11, Chapter 3] and [9]. Furthermore, [9] includes a review of many existing works prior to 1999.

1.2. Relaxing independence. In many data sets it is natural to have independent columns, but not independent entries in the same column. For example, data collected from people, such as patient health information or personal movie ratings, will have independent columns since it is reasonable to assume each person's responses are independent of everyone else's responses. However, entries within a column are most likely not independent.

Several papers relaxing the independence within the columns already exist. Yin and Krishnaiah [59] required the independent columns  $X_k$ , to come from a spherically symmetric distribution; specifically, they require the distribution of  $X_k$  to be the same as that of  $PX_k$ where P is an orthogonal matrix. Aubrun [7] allowed  $X_k$  to be distributed uniformly on the  $l_n^m$ ball. That result was extended by Pajor and Pastur [45] for all isotropic log-concave measures. Hui and Pan [30] and Wei, Yang and Ying [55] considered independent columns  $X_k$  with m(k)dependent stationary entries as long as the length of  $X_k$  is  $O([m(k)]^4)$ . Hofmann-Credner and Stolz [29] and Friesen, Löwe and M. Stolz [20] assumed that the entries of X can be partitioned into independent subsets, while allowing the entries from the same subset to be dependent. Götze and Tikhomirov in [23] and [24] replace the independence assumptions entirely with certain martinale-type conditions. In a similar manner, Adamczak [2] showed that Marchenko-Pastur law holds if the Euclidean norms of the rows and columns of X concentrate around their means and the expectation of each entry of X conditioned on all other entries equals zero. Bai and Zhou [12] gave a sufficient condition in terms of concentration of quadratic forms, and Yao [58] used their condition to allow a time series dependence structure in X. Yaskov [56] gave a short proof with a slightly weaker condition on the concentration of quadratic forms than Bai and Zhou's result. O'Rourke [44] considered a class of random matrices with dependent entries where even the columns are not necessarily independent, but are uncorrelated; although columns that are far enough apart must be independent. Lastly, the papers [16], [14], and [39] consider structured matrices such as block Toeplitz, Hankel, and Markov matrices.

In this paper, we study two random matrix models with relaxed independence requirement. In our first model, we consider matrices with independent columns and each column is partitioned into blocks of the same size, and we only require the entries in different blocks to be independent.

**Definition 1.1** (Block-independent model). Consider a mean zero, isotropic<sup>1</sup> random vector  $x \in \mathbb{R}^p$ . Assume that the entries of x can be partitioned into blocks each of length  $d_k$ , in such a way that the entries in different blocks are independent. (The entries from the same block may be dependent.) Then we say that x follows the block-independent model.

The block-independent data structures arise naturally in many situations. For example Netflix's movie recommendation data set contains ratings of movies by many people. A single person's movie ratings are likely to have a block structure coming from different movie genres, i.e. someone who dislikes documentary movies will have a block of poor ratings, etc. Another example of such a block structure is the stock market. The Marchenko-Pastur law assuming

<sup>&</sup>lt;sup>1</sup>Isotropy means that the covariance matrix of x is identity, i.e.  $\mathsf{E} x x^\mathsf{T} = I$ . The isotropy assumption is convenient but not essential, and we show how to remove it in Section 1.6.

independence among all entries has been used as a comparison to the empirical spectral distribution of daily stock prices, see [34, 46]. However, a block structure is more realistic, since for each day the performance of stocks in the same sector of the market are likely to be correlated and stocks in different sectors can be considered to be independent.

In the second model we study, the independent columns of a random matrix are formed by vectorized independent symmetric random tensors.

**Definition 1.2** (Random tensor model). Consider an isotropic random vector  $x \in \mathbb{R}^n$  with independent entries. Let the random vector  $\mathbf{x} \in \mathbb{R}^{\binom{n}{d}}$  be obtained by vectorizing the symmetric tensor  $x^{\otimes d}$ . Thus, the entries of  $\mathbf{x}$  are indexed by d-element subsets  $\mathbf{i} \subset [n]$  and are defined as products of the entries of x over  $\mathbf{i}$ :

$$x_{\mathbf{i}} = \prod_{i \in \mathbf{i}} x_i = x_{i_1} x_{i_2} \cdots x_{i_d}, \quad \mathbf{i} = \{i_1, \dots, i_d\}.$$

Then we say that  $\boldsymbol{x}$  follows the random tensor model.

Although random tensors appear frequently in data science problems [6, 21, 31, 42, 54, 43, 47, 17, 15, 33, 26, 8, 5, 40, 52, 13, 61], a systematic theory of random tensors is still in its infancy.

1.3. **New results.** In this paper, we generalize Marchenko-Pastur law to the two models of random matrices described above. The following is our main result for the block-independence model (described in Definition 1.1 above).

**Theorem 1.3** (Marchenko-Pastur law for the block-independent model). Let  $X = X^{(p)}$ ,  $p = \sum_k d_k$ , be a sequence of  $p \times m$  random matrices, whose columns are independent and follow the block-independent model with blocks of sizes  $d_k = d_k(p)$ , the aspect ratio p/m converging to a number  $\lambda \in (0, \infty)$  and  $\max_k d_k = o(p)$  as  $p \to \infty$ . Assume that all entries of the random matrix X have uniformly bounded fourth moments. Then with probability 1 the empirical spectral distribution of the sample covariance matrix  $W = \frac{1}{m}XX^{\mathsf{T}}$  converges weakly in distribution to the Marchenko-Pastur distribution with parameter  $\lambda$ .

Remark 1.4. The requirement  $\max_k d_k = o(p)$  in Theorem 1.3 implies that the number of independent blocks grows to infinity. We will show that this condition is necessary in Section 1.7.

Our second main result is the Marchenko-Pastur law for the random tensor model (described in Definition 1.2).

**Theorem 1.5** (Marchenko-Pastur law for the random tensor model). Let  $X = X^{(p)}$ ,  $p = \binom{n}{d}$ ,  $n = 1, 2, \ldots$ , be a sequence of  $p \times m$  random matrices, whose columns are independent and follow the random tensor model with  $d = o(n^{1/3})$ , and the aspect ratio p/m converging to a number  $\lambda \in (0, \infty)$  as  $p \to \infty$ . Assume that the entries of the random vector x have uniformly bounded fourth moments. Then with probability 1 the empirical spectral distribution of the sample covariance matrix  $W = \frac{1}{m}XX^{\mathsf{T}}$  converges weakly in distribution to the Marchenko-Pastur distribution with parameter  $\lambda$ .

<sup>&</sup>lt;sup>2</sup>Note that for the random tensor model, the fourth moment assumption only concerns the entries of the random vector x. The fourth moments of the entries of the random tensor x, and thus of the entries of the random matrix X, can be very large. Indeed, if  $\mathsf{E}\,x_i^4 = K$  for all i, then  $\mathsf{E}\,x_i^4 = K^d$  by independence.

Remark 1.6. Better understood is the non-symmetric version of the random tensor model. Instead of considering the d-fold product  $x^{\otimes d}$  of a random vector  $x \in \mathbb{R}^n$ , consider d i.i.d. random vectors  $x_1, \ldots, x_d \in \mathbb{R}^n$  and consider their inner product  $x_1 \otimes \cdots \otimes x_d$ . Vectorizing this random tensor, we obtain a random vector in  $\mathbb{R}^{n^d}$ . Spectral properties of the non-symmetric random tensor model were studied in [5] in connection with physics and quantum information theory. Marchenko-Pastur law was proved for this model by Lytova [40] under the assumption d = o(n). The non-symmetric random tensor model is even more challenging as it is generated by fewer independent random variables.

1.4. Marchenko-Pastur law via concentration of quadratic forms. Our approach to both main results is based on concentration of quadratic forms. Starting with the original proof of Marchenko-Pastur law [41] via Stieltjes transform, many arguments in random matrix theory (e.g. [45, 25]), make crucial use of concentration of quadratic forms. Specifically, at the core of the proof of Marchenko-Pastur law lies the bound

$$(1.2) \qquad \qquad \mathsf{Var}(x^{\mathsf{T}} A x) = o(p^2)$$

where  $x \in \mathbb{R}^p$  is any column of the random matrix X and A is any deterministic  $p \times p$  matrix with  $||A|| \leq 1$ . If the entires of x have uniformly bounded fourth moments, one always has

$$Var(x^T A x) = E(x^T A x)^2 \le E ||x||_2^4 = O(p^2).$$

Thus, the requirement (1.2) is just a little stronger than the trivial bound.

Suppose the columns of the random matrix X are independent, but the entries of each columns may be dependent. Then for Marchenko-Pastur to hold for X, it is sufficient (but not necessary) to verify the concentration inequality (1.2). The sufficiency is given in the following result; the absence of necessity is noted in [2, Section 2.1, Example 3].

**Theorem 1.7** (Bai-Zhou [12]). Let  $X = X^{(p)}$ , p = 1, 2, ..., be a sequence of mean zero  $p \times m$  random matrices with independent columns. Assume the following as  $p \to \infty$ .

- 1. The aspect ratio p/m converges to a number  $\lambda \in (0, \infty)$  as  $p \to \infty$ .
- 2. For each p, all columns  $X_k$  of  $X^{(p)}$  have the same covariance matrix  $\Sigma = \Sigma^{(p)} = \mathsf{E}\, X_k X_k^\mathsf{T}$ . The spectral norm of the covariance matrix  $\Sigma^{(p)}$  is uniformly bounded, and the empirical spectral distribution of  $\Sigma^{(p)}$  converges to a deterministic distribution H.
- 3. For any deterministic  $p \times p$  matrices  $A = A^{(p)}$  with uniformly bounded spectral norm and for every column  $X_k$ , we have

$$\max_{k} \mathsf{Var}(X_k^{\mathsf{T}} A X_k) = o(p^2).$$

Then, with probability 1 the empirical spectral distribution of the sample covariance matrix  $W = \frac{1}{m} X X^{\mathsf{T}}$  converges weakly to a deterministic distribution whose Stieltjes transform satisfies

(1.3) 
$$s(z) = \int_0^\infty \frac{1}{t(1 - \lambda - \lambda zs) - z} dH(t), \quad z \in \mathbb{C}^+.$$

In the case where the entries of the columns are uncorrelated and have unit variance, we have  $\Sigma = I$  and Theorem 1.7 yields that the limiting distribution is the original Marchenko-Pastur law (1.1).

<sup>&</sup>lt;sup>3</sup>Although the main result in [40] is stated for fixed degree d, it can be allowed to grow as fast as d = o(n): see Lemma 3.3, Theorem 1.2, Definition 1.1, and Remark 4.1 in [40].

1.5. Concentration of quadratic forms: new results. Theorem 1.7 reduces proving Marchenko-Pastur law for our new models to the concentration of a quadratic form  $x^{T}Ax$ . If the random vector x has all independent entries, bounding the variance of this quadratic form is elementary. Moreover, in this case Hanson-Wright inequality (see e.g. [51, 48]) gives good probability tail bounds for the quadratic form.

But in our new models, the coordinates of the random vector are not independent. There seem to be no sufficiently powerful concentration inequalities available for such models. Known concentration inequalities for random chaoses [35, 37, 36, 3, 4, 22, 1] exhibit an unspecified (possibly exponential) dependence on the degree d, which is too bad for our purposes. An exception is the recent work [52] on concentration of random tensors with an optimal dependence on d. However, the results of [52] only apply for non-symmetric tensors and positive-semidefinite matrices A.

The following are new concentration inequalities for the block-independent model (Theorem 1.8) and the random tensor model (Theorem 1.9), which we will prove in Section 2 and Section 3 respectively.

**Theorem 1.8** (Variance of quadratic forms for block-independent model). Let  $x \in \mathbb{R}^p$  be a random vector that follows the block-independent model with blocks of sizes  $d_k$ . Then, for any fixed matrix  $A \in \mathbb{R}^{p \times p}$ , we have

$$\operatorname{Var}(x^{\mathsf{T}} A x) \le \|A\|^2 \Big( K \sum_{k} d_k^2 + 2p \Big).$$

Here K is the largest fourth moment of the entries of x.

This result combined with Theorem 1.7 immediately establishes Marchenko-Pastur law for the block-independence model:

Proof of Theorem 1.3. Apply Theorem 1.8 and simplify the conclusion using the bound  $\sum_k d_k^2 \le (\max_k d_k) \sum_k d_k = (\max_k d_k) p$ . We get

$$Var(x^{\mathsf{T}}Ax) \le p||A||^2 \left(K \max_k d_k + 2\right) = o(p^2),$$

if ||A|| = O(1), K = O(1), and  $\max_k d_k = o(p)$  as  $p \to \infty$ . This justifies condition 3 of Theorem 1.7. Applying this theorem with  $\Sigma = I$  we conclude Theorem 1.3.

**Theorem 1.9** (Variance of quadratic forms for random tensor model). There exist positive absolute constants C, c > 0 such that the following holds. Let  $\mathbf{x} \in \mathbb{R}^p$ ,  $p = \binom{n}{d}$ , be a random vector that follows the random tensor model. Then, for any fixed matrix  $A \in \mathbb{R}^{p \times p}$ , we have

$$Var(x^{T}Ax) \le C||A||^{2}p^{2}\left(\frac{K^{1/2}d}{n^{1/3}}\right)^{3/2},$$

if  $K^{1/2}d/n^{1/3} < c$ . Here K is the largest fourth moment of the entries of x.

This result combined with Theorem 1.7 immediately establishes Marchenko-Pastur law for the random tensor model:

Proof of Theorem 1.5. Theorem 1.9 yields

$$\mathsf{Var}(x^\mathsf{T} A x) = o(p^2)$$

whenever ||A|| = O(1), K = O(1), and  $d = o(n^{1/3})$ . This justifies condition 3 of Theorem 1.7. Applying this theorem with  $\Sigma = I$  we conclude Theorem 1.5.

1.6. Anisotropic block-independent model. In Definition 1.1 of the block-independent model we assumed for simplicity that the blocks are isotropic. Let us show how to remove this assumption and still obtain a version of Theorem 1.3; the limiting spectral distribution will then be the anisotropic Marchenko-Pastur law (1.3).

To see this, suppose all columns of our random matrix  $X = X^{(p)}$  have the same covariance matrix  $\Sigma = \Sigma^{(p)}$ . Assume that, as  $p \to \infty$ , we have  $\|\Sigma^{(p)}\| = O(1)$  and the empirical spectral distribution<sup>4</sup> of  $\Sigma^{(p)}$  converges to a deterministic distribution H.

Denoting as before by  $X_k$  the k-th column of X, we can represent it as  $X_k = \Sigma^{1/2} x_k$  where  $x_k$  is some isotropic random vector, i.e. one whose entries are uncorrelated and have unit variance. Then

$$\operatorname{Var}(X_k^{\mathsf{T}} A X_k) = \operatorname{Var}\left(x_k^{\mathsf{T}} \ \Sigma^{1/2} A \Sigma^{1/2} \ x_k\right).$$

Applying Theorem 1.8 for  $x = x_k$  and  $\Sigma^{1/2} A \Sigma^{1/2}$  instead of A, we conclude that  $\mathsf{Var}(X_k^\mathsf{T} A X_k) = o(p^2)$  if  $\|\Sigma\| = O(1)$ ,  $\|A\| = O(1)$ , K = O(1), and  $\max_k d_k = o(p)$ .

This justifies condition 3 of Theorem 1.7. Applying this theorem, we conclude that the limiting spectral distribution of  $W = \frac{1}{m}XX^{\mathsf{T}}$  converges to the anisotropic Marchenko-Pastur distribution (1.3).

1.7. **Optimality.** Here we show that the number of blocks in the block-independent model has to go  $\infty$ . Indeed, let  $X^{(p)}$  be a sequence of  $p \times m$  random matrices such that  $p/m \to \lambda > 0$  as  $p \to \infty$ , and whose columns are independent copies of an isotropic random vector  $x^{(p)} \in \mathbb{R}^p$ . According to a result of P. Yaskov [57, Theorem 2.1], a necessary condition for Marchenko-Pastur law is that

(1.4) 
$$\frac{1}{p} \|x^{(p)}\|_2^2 \to 1 \quad \text{in probability.}$$

This condition may fail if the number of independent blocks is O(1). To see this, take a random vector from the block-independent model with n equal length blocks (p = nd), and replace each block with a zero vector independently with probability 1/2. Multiply the result by  $\sqrt{2}$ . The resulting random vector  $x^{(p)}$  still follows the bock-independent model, but it equals zero with probability  $2^{-n}$ , a quantity that is bounded below by a positive constant if n = O(1). This violates the condition (1.4) and demonstrates that Marchenko-Pastur law fails in this case.

It is less clear whether our requirement on the degree  $d = o(n^{1/3})$  in Theorem 1.5 is optimal. In the light of (1.4), it seems that the optimal condition might be

$$d = o(n^{1/2}).$$

Indeed, consider a random vector  $x^{(p)} \in \mathbb{R}^p$ ,  $p = \binom{n}{d}$  obtained from a random vector  $x \in \mathbb{R}^n$  with i.i.d. coordinates, that follows the random tensor model. Then

$$U_p := \frac{1}{p} \| \boldsymbol{x}^{(p)} \|_2^2 = \frac{1}{\binom{n}{d}} \sum_{1 \le i_1 \le \dots \le i_d \le n} x_{i_1}^2 x_{i_2}^2 \cdots x_{i_d}^2$$

<sup>&</sup>lt;sup>4</sup>Since the blocks are independent, the covariance matrix  $\Sigma$  is block-diagonal. If  $\Sigma_j$  is the covariance matrix of the block j, then the spectral norm of  $\Sigma$  is the maximal spectral norm of  $\Sigma_j$ , and the empirical spectral distribution of  $\Sigma$  is the mixture of the empirical spectral distributions of all  $\Sigma_j$ .

is a U-statistic. According to a result of W. Hoeffding [28],

$$Var(U_p) \ge \frac{d^2}{n} Var(x_1^2).$$

Assume the variance of  $x_1^2$  is nonzero. If  $d \gtrsim n^{1/2}$  then  $\mathsf{Var}(U_p)$  does not converge to zero. This makes it plausible that the necessary condition (1.4) for Marchenko-Pastur law may be violated in this regime.

- 2. Quadratic forms in block-independent random vectors: Proof of Theorem 1.8
- 2.1. **Reductions.** Rearranging the entries of x, we can assume that the indices of the blocks are successive intervals, i.e. the kth block index set is  $I_k = \left\{\sum_{l=1}^{k-1} d_l + 1, \dots, \sum_{l=1}^k d_l\right\}$ . Since  $x^\mathsf{T} A x = \left(x^\mathsf{T} A x\right)^\mathsf{T} = x^\mathsf{T} A^\mathsf{T} x$ , the symmetric matrix  $\tilde{A} := (A + A^\mathsf{T})/2$  satisfies

$$x^{\mathsf{T}} A x = x^{\mathsf{T}} \tilde{A} x$$
 and  $\|\tilde{A}\| \le \frac{1}{2} (\|A\| + \|A^{\mathsf{T}}\|) = \|A\|.$ 

Therefore, it suffices to prove Theorem 1.8 for symmetric matrices A.

We will control the contribution of the diagonal and off-diagonal blocks of A separately. The diagonal blocks of A form the block-diagonal matrix  $D = (D_{ij})_{i,j=1}^{nd}$  defined as

$$D_{ij} = A_{ij}$$
 if  $i, j$  lie in the same block

and  $D_{ij} = 0$  otherwise. Now, decomposing  $x^{\mathsf{T}}Ax = x^{\mathsf{T}}Dx + x^{\mathsf{T}}(A-D)x$ , we have

$$(2.1) \qquad \operatorname{Var}(x^{\mathsf{T}} A x) \le 2 \operatorname{Var}(x^{\mathsf{T}} D x) + 2 \operatorname{Var}(x^{\mathsf{T}} (A - D) x).$$

Let us bound each of the two terms on the right hand side.

2.2. **Diagonal contribution.** The vector x can be decomposed into blocks  $\bar{x}_k := (x_i)_{i \in I_k}$ , and the matrix D consists of corresponding diagonal blocks  $\bar{D}_k := (D_{ij})_{i,j \in I_k}$ . Then  $x^\mathsf{T} D x = \sum_k \bar{x}_k^\mathsf{T} \bar{D}_k \bar{x}_k$ , and since  $\bar{x}_k$  are independent, this yields

$$\operatorname{Var}(x^{\mathsf{T}}Dx) = \sum_{k} \operatorname{Var}\left(\bar{x}_{k}^{\mathsf{T}}\bar{D}_{k}\bar{x}_{k}\right).$$

Now,

$$\mathsf{Var}\left(\bar{x}_{k}^{\mathsf{T}}\bar{D}_{k}\bar{x}_{k}\right) \leq \mathsf{E}\left(\bar{x}_{k}^{\mathsf{T}}\bar{D}_{k}\bar{x}_{k}\right)^{2} \leq \mathsf{E}\left(\|\bar{D}_{k}\|\,\|\bar{x}_{k}\|_{2}^{2}\right)^{2} \leq \|A\|^{2}\;\mathsf{E}\,\|\bar{x}_{k}\|_{2}^{4}.$$

Furthermore,

$$\mathsf{E} \, \|\bar{x}_k\|_2^4 = \sum_{i,j \in I_k} \mathsf{E} \, x_i^2 x_j^2 \le K d_k^2.$$

We conclude that

(2.2) 
$$\operatorname{Var}(x^{\mathsf{T}}Dx) \le K||A||^2 \sum_{k} d_k^2.$$

2.3. Off-diagonal contribution. By definition,

(2.3) 
$$\operatorname{Var}(x^{\mathsf{T}}(A-D)x) = \operatorname{E}\left(x^{\mathsf{T}}(A-D)x\right)^{2} - \left(\operatorname{E}x^{\mathsf{T}}(A-D)x\right)^{2}.$$

Denote by  $\mathcal{R}$  the set of all index pairs (i, j) such that i and j do not lie in the same block. Then

$$\mathsf{E}\left(x^\mathsf{T}(A-D)x\right)^2 = \mathsf{E}\left(\sum_{(i,j)\in\mathcal{R}}A_{ij}x_ix_j\right)^2 = \sum_{(i,j),(k,l)\in\mathcal{R}}A_{ij}A_{kl}\,\mathsf{E}\,x_ix_jx_kx_l$$

Consider any term  $\mathsf{E} x_i x_j x_k x_l$  that is nonzero. By the mean zero assumption and block-independence, none of the indices i, j, k or l may lie in their own block. This means that a pair of these indices lies in one block and another pair lies in a different block. By definition of  $\mathcal{R}$ , there there are only two ways to form such pairs: (i, k) in one block and (j, l) in another, or (i, l) in one block and (j, k) in another.

In the first scenario, block-independence yields

$$\mathsf{E}\,x_ix_jx_kx_l = \mathsf{E}\,x_ix_k\,\mathsf{E}\,x_jx_l.$$

By isotropy, this term equals 1 if i = k and j = l, and zero otherwise. In the second scenario, arguing similarly we get one if i = l and j = k, and zero otherwise. Therefore, breaking the sum according to the scenario and then using the symmetry of A, we obtain

$$\sum_{(i,j),(k,l)\in\mathcal{R}} A_{ij} A_{kl} \, \mathsf{E} \, x_i x_j x_k x_l = \sum_{(i,j)\in\mathcal{R}} A_{ij} A_{ij} + \sum_{(i,j)\in\mathcal{R}} A_{ij} A_{ji} = 2 \sum_{(i,j)\in\mathcal{R}} A_{ij}^2$$

$$\leq 2 \sum_{i,j} A_{ij}^2 \leq 2 \sum_{i=1}^p \sum_{j=1}^m A_{ij}^2 \leq 2p \|A\|^2.$$

We just bounded the first term in the right hand side of (2.3). The second term vanishes. Indeed,

$$\mathsf{E} \, x^{\mathsf{T}} (A - D) x = \sum_{(i,j) \in \mathcal{R}} A_{ij} \, \mathsf{E} \, x_i x_j = 0$$

since  $\mathsf{E}\,x_ix_j=0$  for all  $i\neq j$  by assumption. Summarizing, we bounded the off-diagonal contribution as follows:

$$\mathsf{Var}(x^{\mathsf{T}}(A-D)x) \le 2p\|A\|^2.$$

Combining this with the bound (2.2) on the diagonal contribution and substituting into (2.1), we conclude that

$$\operatorname{Var}(x^{\mathsf{T}}Ax) \leq \|A\|^2 \Big(K \sum_k d_k^2 + 2p\Big).$$

3. Quadratic forms in random tensors: Proof of Theorem 1.9

3.1. Reductions. Without loss of generality, we may assume that ||A|| = 1 by rescaling. Expanding  $\mathbf{x}^{\mathsf{T}}A\mathbf{x}$  as a double sum of terms  $A_{\mathbf{ij}}\mathbf{x_i}\mathbf{x_j}$ , and distinguishing the diagonal terms

 $(\mathbf{i} = \mathbf{j})$  and the off-diagonal terms  $(\mathbf{i} \neq \mathbf{j})$ , we have:

$$\operatorname{Var}(\boldsymbol{x}^{\mathsf{T}} A \boldsymbol{x}) = \mathsf{E}\left[|\boldsymbol{x}^{\mathsf{T}} A \boldsymbol{x} - \operatorname{tr} A|^{2}\right] \leq 2 \, \mathsf{E}\left[\left(\sum_{\mathbf{i}} A_{\mathbf{i}\mathbf{i}}(\boldsymbol{x}_{\mathbf{i}}^{2} - 1)\right)^{2}\right] + 2 \, \mathsf{E}\left[\left(\sum_{\mathbf{i} \neq \mathbf{j}} A_{\mathbf{i}\mathbf{j}} \boldsymbol{x}_{\mathbf{i}} \boldsymbol{x}_{\mathbf{j}}\right)^{2}\right]$$

$$=: 2S_{\operatorname{diag}} + 2S_{\operatorname{off}}.$$

Here we used the inequality  $(a+b)^2 \le 2a^2 + 2b^2$ .

3.2. **Diagonal contribution.** Expanding the square, we can express the diagonal contribution as

(3.2) 
$$S_{\text{diag}} = \sum_{\mathbf{i},\mathbf{k}} A_{\mathbf{i}\mathbf{i}} A_{\mathbf{k}\mathbf{k}} \, \mathsf{E}(\boldsymbol{x}_{\mathbf{i}}^2 - 1)(\boldsymbol{x}_{\mathbf{k}}^2 - 1).$$

Both meta-indices **i** and **k** range in all  $\binom{n}{d}$  subsets of [n] of cardinality d. Let v denote the overlap between these two subsets, i.e.

$$v := |\mathbf{i} \cap \mathbf{k}|.$$

If v = 0, the subsets are disjoint, the random variables  $x_i^2 - 1$  and  $x_k^2 - 1$  are independent and have mean zero, and thus

$$\mathsf{E}(\boldsymbol{x}_{\mathbf{i}}^2 - 1)(\boldsymbol{x}_{\mathbf{k}}^2 - 1) = 0.$$

Such terms do not contribute anything to the sum in (3.2).

If  $v \ge 1$ , the monomial  $\mathbf{x}_{\mathbf{i}}^2 \mathbf{x}_{\mathbf{k}}^2$  consists of v terms raised to the fourth power (coming from the indices that are both in  $\mathbf{i}$  and  $\mathbf{k}$ ) and 2(d-v) terms raised to the second power (coming from the symmetric difference of  $\mathbf{i}$  and  $\mathbf{k}$ ). Thus,

$$\left| \mathsf{E}(\boldsymbol{x}_{\mathbf{i}}^2 - 1)(\boldsymbol{x}_{\mathbf{k}}^2 - 1) \right| \le \mathsf{E}\,\boldsymbol{x}_{\mathbf{i}}^2 \boldsymbol{x}_{\mathbf{k}}^2 \le \max_{\alpha} \left( \mathsf{E}\,\boldsymbol{x}_{\alpha}^4 \right)^v \cdot \max_{\beta} \left( \mathsf{E}\,\boldsymbol{x}_{\beta}^2 \right)^{2(d-v)} \le K^v,$$

where we used the unit variance assumption.

There are  $\binom{n}{d}$  ways to choose **i**. Once we fix **i** and  $v \in \{1, \ldots, d\}$ , there are  $\binom{d}{v}\binom{n-d}{d-v}$  ways to choose **k**, since v indices must come from **i** and the remaining d-v indices must come from  $[n] \setminus \mathbf{i}$ . Therefore,

(3.3) 
$$S_{\text{diag}} \leq \binom{n}{d} \sum_{v=1}^{d} \binom{d}{v} \binom{n-d}{d-v} K^{v}.$$

To bound this sum, we can assume without loss of generality that K is a positive integer. Then the following elementary inequality holds:

$$\binom{d}{v}K^v \le \binom{Kd}{v},$$

and it can be quickly checked by writing the binomial coefficients in terms of factorials. Now, if we were summing v from zero as opposed from 1 in (3.3), we can use Vandermonde's identity and get

$$\sum_{v=0}^{d} {d \choose v} {n-d \choose d-v} K^v \le \sum_{v=0}^{d} {Kd \choose v} {n-d \choose d-v} = {n-d+Kd \choose d}.$$

Subtracting the zeroth term, we obtain

$$\sum_{v=1}^{d} {d \choose v} {n-d \choose d-v} K^{v} \le {n-d+Kd \choose d} - {n-d \choose d}.$$

Now use a stability property of binomial coefficients (Lemma 3.7), which tells us that

$$\binom{n-d+Kd}{d} - \binom{n-d}{d} \le \delta \binom{n-d}{d} \quad \text{where } \delta := \frac{2Kd^2}{n-2d+1},$$

as long as  $\delta \leq 1/2$ . According to our assumptions on the degree d, we do have  $\delta \leq 1/2$  when n is sufficiently large.

Summarizing, we have shown that

(3.4) 
$$S_{\text{diag}} \leq \binom{n}{d} \cdot \delta \binom{n-d}{d} \lesssim \binom{n}{d}^2 \cdot \frac{Kd^2}{n}.$$

3.3. Off-diagonal contribution: the cross moments. Expanding the square, we can express the off-diagonal contribution in (3.1) as

(3.5) 
$$S_{\text{off}} = \sum_{\mathbf{i} \neq \mathbf{j}} \sum_{\mathbf{k} \neq \mathbf{l}} A_{\mathbf{i}\mathbf{j}} A_{\mathbf{k}\mathbf{l}} \ \mathsf{E} \, \boldsymbol{x}_{\mathbf{i}} \boldsymbol{x}_{\mathbf{j}} \boldsymbol{x}_{\mathbf{k}} \boldsymbol{x}_{\mathbf{l}}.$$

Let us first bound the expectation of

$$x_i x_j x_k x_l = \prod_{i \in i} x_i \prod_{i \in i} x_j \prod_{k \in k} x_k \prod_{l \in l} x_l.$$

Without loss of generality, we can assume that this monomial of degree 4d has no linear factors, i.e. each of the factors  $x_{\alpha}$  of this monomial has degree at least 2, otherwise the expectation of the monomial is zero. Rearranging the factors, we can express the monomial as

(3.6) 
$$\mathbf{x_i} \mathbf{x_j} \mathbf{x_k} \mathbf{x_l} = \prod_{\alpha \in \Lambda_2} x_{\alpha}^2 \prod_{\beta \in \Lambda_3} x_{\beta}^3 \prod_{\gamma \in \Lambda_4} x_{\gamma}^4$$

for some disjoint sets  $\Lambda_2, \Lambda_3, \Lambda_4 \subset [n]$ . Thus,  $\Lambda_2$  consists of the indices that are covered by exactly two of the sets  $\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l}$ , and similarly for  $\Lambda_3$  and  $\Lambda_4$ . Since each of the four sets  $\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l}$  contains d indices, counting the indices with multiplicities gives

$$(3.7) 4d = 2|\Lambda_2| + 3|\Lambda_3| + 4|\Lambda_4|.$$

Since each index is covered at least by two of the four sets i, j, k, l, the cardinality of the set

$$\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l} = \Lambda_2 \sqcup \Lambda_3 \sqcup \Lambda_4$$

is at most 4d/2 = 2d. Let  $w \ge 0$  be the "defect" defined by

$$(3.9) |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}| = 2d - w.$$

Thus, w would be zero if every index is covered by exactly two sets, and w would be positive if there are triple or quadruple covered indices. From (3.8) and (3.9) we see that

$$2d - w = |\Lambda_2| + |\Lambda_3| + |\Lambda_4|.$$

Multiplying both sides of this equation by 2 and subtracting from (3.7), we get

$$(3.10) 2w = |\Lambda_3| + 2|\Lambda_4|,$$

a relation that will be useful in a moment.

Take expectation on both sides of (3.6). Using independence and the assumptions that  $\mathsf{E}\,x_\alpha^2=1$  and  $\mathsf{E}\,x_\alpha^4\leq K$  for each  $\alpha$ , we get

$$\mathsf{E}\,|\boldsymbol{x_i}\boldsymbol{x_j}\boldsymbol{x_k}\boldsymbol{x_l}| = \prod_{\beta\in\Lambda_3}\mathsf{E}\,|x_\beta|^3\cdot\prod_{\gamma\in\Lambda_4}\mathsf{E}\,x_\gamma^4 = \prod_{\beta\in\Lambda_3}\left(\,\mathsf{E}\,|x_\beta|^4\right)^{3/4}\cdot\prod_{\gamma\in\Lambda_4}\mathsf{E}\,x_\gamma^4 \leq K^{\frac{3}{4}|\Lambda_3|+|\Lambda_4|}.$$

Due to (3.10),

$$\frac{3}{4}|\Lambda_3| + |\Lambda_4| = \frac{3}{2}w - \frac{1}{2}|\Lambda_4| \le \frac{3}{2}w.$$

Thus we have shown that

$$\mathsf{E} | \boldsymbol{x}_{\mathbf{i}} \boldsymbol{x}_{\mathbf{i}} \boldsymbol{x}_{\mathbf{k}} \boldsymbol{x}_{\mathbf{l}} | < K^{3w/2}.$$

3.4. Sizes of intersections of meta-indices. Due to the last step, the off-diagonal contribution (3.5) can be bounded as follows:

$$(3.11) S_{\text{off}} \leq \sum_{\mathbf{i} \neq \mathbf{j}} \sum_{\mathbf{k} \neq \mathbf{l}} |A_{\mathbf{i}\mathbf{j}}| |A_{\mathbf{k}\mathbf{l}}| K^{3w/2},$$

where the sum only includes the sets  $\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l}$  that provide at least a *double cover*, i.e. such that every index from  $\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}$  must belong to at least two of these four sets. We quantified this property by the *defect*  $w \geq 0$ , which we defined by

$$|\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}| = |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| = 2d - w.$$

In preparation to bounding the double sum in (3.11), let us consider

$$|\mathbf{i} \cap \mathbf{j}| =: v, \quad |\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}| =: r,$$

and observe a few useful bounds involving w, v, and r.

**Lemma 3.1.** We have  $w \le v \le d - 1$ .

*Proof.* By definition,  $v = |\mathbf{i} \cap \mathbf{j}| \le |\mathbf{i}| = d$ . Moreover, v may not equal d, for this would mean that  $\mathbf{i} = \mathbf{j}$ , a possibility that is excluded in the double sum (3.11). This means that  $v \le d - 1$ . Next, we have

$$(3.12) |\mathbf{i} \cup \mathbf{j}| = |\mathbf{i}| + |\mathbf{j}| - |\mathbf{i} \cap \mathbf{j}| = 2d - v.$$

On the other hand,  $|\mathbf{i} \cup \mathbf{j}| \le |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| = 2d - w$ . Combining these two facts yields  $w \le v$ .

**Lemma 3.2.** We have  $r \leq v$  and  $r \leq 2w$ .

*Proof.* The first statement follows from definition. To prove the second statement, recall that the sets  $\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l}$  form at least a double cover of  $\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}$  and at least a triple cover of  $\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}$  (trivially). Since each of the four sets has d indices, counting the indices with multiplicities gives

$$4d \ge 2|\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}| + |\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}| = 2(2d - w) + r$$

by the definition of w and r. This yields  $r \leq 2w$ .

**Lemma 3.3.** We have  $r \leq d - v + w$ .

*Proof.* The sets  $\mathbf{i}$ ,  $\mathbf{j}$ ,  $\mathbf{k}$  obviously form at least a double cover of  $\mathbf{i} \cap \mathbf{j}$  and a triple cover of  $\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}$ . Since each of the three sets has d indices, counting the indices with multiplicities gives

$$3d \ge |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| + |\mathbf{i} \cap \mathbf{j}| + |\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}| = (2d - w) + v + r$$

by definition of w, v and r. Rearranging the terms completes the proof.

3.5. Number of choices of meta-indices. Let us fix w, v, and r, and estimate the number of possible choices for the sets i, j, k, l that conform to these w, v, and r. This would help us determining the number of terms in the double sum (3.11). Thus, we would like to know how many ways are there to choose four d-element sets  $\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l} \subset [n]$  that provide at least a double cover of  $\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}$ , and so that

(3.13) 
$$|\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| = |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}| = 2d - w, \quad |\mathbf{i} \cap \mathbf{j}| = v, \text{ and } |\mathbf{i} \cap \mathbf{j} \cap \mathbf{k}| = r.$$

Choosing i. This is easy: there are  $\binom{n}{d}$  ways to choose the d-element subset i from [n].

Choosing **j**. Recall that we need to obey  $|\mathbf{i} \cap \mathbf{j}| = v$ . Thus, for a fixed **i**, we have  $\binom{d}{v} \binom{n-d}{d-v}$  choices for  $\mathbf{j}$ , which is seen by first picking the v overlapping indices from  $\mathbf{i}$  and then the remaining d-v indices from  $\mathbf{i}^c$ .

Choosing k. Let us fix i and j. The set of all available indices [n], from which the indices of k can be chosen, can be partitioned into the three disjoint sets:

$$[n] = (\mathbf{i} \cap \mathbf{j}) \sqcup (\mathbf{i} \cup \mathbf{j})^c \sqcup (\mathbf{i} \triangle \mathbf{j}).$$

Let us see how many indices for k should come from each of these three sets.

As we see from (3.13), the v-element set  $\mathbf{i} \cap \mathbf{j}$  must contain exactly r indices of k, and these can be selected in  $\binom{v}{r}$  ways.

Next, we know from (3.12) that  $|(\mathbf{i} \cup \mathbf{j})^c| = n - (2d - v)$ , and

$$(3.15) \qquad |(\mathbf{i} \cup \mathbf{j})^c \cap \mathbf{k}| = |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| - |\mathbf{i} \cup \mathbf{j}| = (2d - w) - (2d - v) = v - w,$$

where we used (3.13) and (3.12). So, the set  $(\mathbf{i} \cup \mathbf{j})^c$  must contain exactly v - w indices of  $\mathbf{k}$ , and these can be selected in  $\binom{n-(2d-v)}{v-w}$  ways.<sup>5</sup>

Finally, by (3.12) and (3.13) we have

$$(3.16) |\mathbf{i} \triangle \mathbf{j}| = |\mathbf{i} \cup \mathbf{j}| - |\mathbf{i} \cap \mathbf{j}| = (2d - v) - v = 2(d - v).$$

We already allocated r + (v - w) indices of k to the first two sets on the right-hand side of (3.14). Thus, the number of indices for **k** that come from the third set,  $\mathbf{i} \triangle \mathbf{j}$ , must be

(3.17) 
$$|(\mathbf{i}\triangle\mathbf{j})\cap\mathbf{k}| = d - r - (v - w).$$

These indices can be selected in  $\binom{2(d-v)}{d-r-(v-w)}$  ways.<sup>6</sup> Summarizing, for fixed **i** and **j**, we have  $\binom{v}{r}\binom{n-(2d-v)}{v-w}\binom{2(d-v)}{d-r-(v-w)}$  choices for **k**.

<sup>&</sup>lt;sup>5</sup>Since the cardinality of any set is nonnegative, equation (3.15) provides an alternative proof of the bound  $w \leq v$  in Lemma 3.1.

 $<sup>\</sup>overline{^6}$ Since the cardinality of any set is nonnegative, equation (3.17) provides an alternative proof of Lemma 3.3.

Choosing l. Fix i, j and k. Recall that the sets i, j, k, l must form at least a double cover of  $\mathbf{i} \cup \mathbf{j} \cup \mathbf{k} \cup \mathbf{l}$ . This has two consequences. First, we must have

$$(3.18) l \subset \mathbf{i} \cup \mathbf{j} \cup \mathbf{k}$$

to avoid any single-covered indices in  $\mathbf{l}$ . Second,  $\mathbf{l}$  must contain all the *single indices*, i.e. those that belong to exactly one of the sets  $\mathbf{i}$ ,  $\mathbf{j}$ , or  $\mathbf{k}$ . The set of single indices, denoted  $\mathbf{s}$ , can be represented as

$$\mathbf{s} = (\mathbf{i}^c \cap \mathbf{j}^c \cap \mathbf{k}) \sqcup \left[ (\mathbf{i} \cap \mathbf{j}^c \cap \mathbf{k}^c) \sqcup (\mathbf{i}^c \cap \mathbf{j} \cap \mathbf{k}^c) \right] = \left[ (\mathbf{i} \cup \mathbf{j})^c \cap \mathbf{k} \right] \sqcup \left[ (\mathbf{i} \triangle \mathbf{j}) \cap \mathbf{k}^c \right].$$

At this stage, the sets  $\mathbf{i}$ ,  $\mathbf{j}$  and  $\mathbf{k}$  are all fixed, and so is  $\mathbf{s}$ .

To compute the cardinality of **s**, recall from (3.15) that  $|(\mathbf{i} \cup \mathbf{j})^c \cap \mathbf{k}| = v - w$ . Furthermore, using (3.16) and (3.17), we see that

$$|(\mathbf{i}\triangle\mathbf{j})\cap\mathbf{k}^c|=|(\mathbf{i}\triangle\mathbf{j})|-|(\mathbf{i}\triangle\mathbf{j})\cap\mathbf{k}|=2(d-v)-(d-r-(v-w))=d-w-v+r.$$

Thus, the number of single indices is

$$|\mathbf{s}| = (v - w) + (d - w - v + r) = d - 2w + r.$$

Since I must contain the set **s** of single indices, which is fixed, the only freedom in choosing I comes from selecting non-single indices. There are d - (d - 2w + r) = 2w - r of them,<sup>7</sup> and they must come from the set  $(\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}) \setminus \mathbf{s}$ , due to (3.18). Now, recalling (3.13), we have

$$|(\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}) \setminus \mathbf{s}| = |\mathbf{i} \cup \mathbf{j} \cup \mathbf{k}| - |\mathbf{s}| = (2d - w) - (d - 2w + r) = d + w - r.$$

Hence, for fixed **i**, **j** and **k**, we have  $\binom{d+w-r}{2w-r}$  choices for **l**.

3.6. Bounding the off-diagonal contribution by a binomial sum. We can now return to our bound (3.11) on the off-diagonal contribution. We can rewrite it as follows:

$$(3.19) S_{\text{off}} \leq \sum_{w,v,r} K^{3w/2} \sum_{\mathbf{i} \in \mathbf{I}} \sum_{\mathbf{j} \in \mathbf{J}(\mathbf{i})} \sum_{\mathbf{k} \in \mathbf{K}(\mathbf{i},\mathbf{j})} \sum_{\mathbf{l} \in \mathbf{L}(\mathbf{i},\mathbf{j},\mathbf{k})} |A_{\mathbf{i}\mathbf{j}}| |A_{\mathbf{k}\mathbf{l}}|.$$

The first sum is over all realizable v, w, and r, and the rest of the sums are over all possible choices for  $\mathbf{i}$ ,  $\mathbf{j}$ ,  $\mathbf{k}$  and  $\mathbf{l}$  that conform to the given v, w and r per (3.13). Thus, for instance,  $\mathbf{L}(\mathbf{i}, \mathbf{j}, \mathbf{k})$  consists of all possible choices for  $\mathbf{l}$  given  $\mathbf{i}, \mathbf{j}$  and  $\mathbf{k}$ . We observed various bounds on realizable v, w and r in Section 3.4, and we computed the cardinalities of the sets  $\mathbf{I}$ ,  $\mathbf{J}(\mathbf{i})$ ,  $\mathbf{K}(\mathbf{i}, \mathbf{j})$  and  $\mathbf{L}(\mathbf{i}, \mathbf{j}, \mathbf{k})$  in Section 3.5. This knowledge will help us to bound the five-fold sum in (3.19).

In order to do this, rewrite (3.19) as follows:

$$S_{\text{off}} \leq \sum_{w,v,r} K^{3w/2} \sum_{\mathbf{i} \in \mathbf{I}} \sum_{\mathbf{j} \in \mathbf{J}(\mathbf{i})} |A_{\mathbf{i}\mathbf{j}}| \sum_{\mathbf{k} \in \mathbf{K}(\mathbf{i},\mathbf{j})} \sum_{\mathbf{l} \in \mathbf{L}(\mathbf{i},\mathbf{j},\mathbf{k})} |A_{\mathbf{k}\mathbf{l}}|.$$

Note that  $|A_{\mathbf{k}\mathbf{l}}| \leq ||A|| = 1$  for all  $\mathbf{k}$  and  $\mathbf{l}$ , and

$$\sum_{\mathbf{i} \in \mathbf{J}(\mathbf{i})} |A_{\mathbf{i}\mathbf{j}}| \le |\mathbf{J}(\mathbf{i})|^{1/2} \left(\sum_{\mathbf{i} \in \mathbf{J}(\mathbf{i})} A_{\mathbf{i}\mathbf{j}}^2\right)^{1/2} \le |\mathbf{J}(\mathbf{i})|^{1/2} ||A|| = |\mathbf{J}(\mathbf{i})|^{1/2}.$$

<sup>&</sup>lt;sup>7</sup>Since the number of indices is non-negative, this provides an alternative proof of the bound  $r \leq 2w$  in Lemma 3.2.

Thus

$$S_{\text{off}} \leq \sum_{w,v,r} K^{3w/2} |\mathbf{I}| \cdot \max_{\mathbf{i}} |\mathbf{J}(\mathbf{i})|^{1/2} \cdot \max_{\mathbf{i},\mathbf{j}} |\mathbf{K}(\mathbf{i},\mathbf{j})| \cdot \max_{\mathbf{i},\mathbf{j},\mathbf{k}} |\mathbf{L}(\mathbf{i},\mathbf{j},\mathbf{k})|.$$

Now we can use the bounds we proved in Section 3.5 on the cardinalities of sets  $\mathbf{I}$ ,  $\mathbf{J}(\mathbf{i})$ ,  $\mathbf{K}(\mathbf{i}, \mathbf{j})$  and  $\mathbf{L}(\mathbf{i}, \mathbf{j}, \mathbf{k})$ , which are the number of choices for  $\mathbf{i}$ , for  $\mathbf{j}$  given  $\mathbf{i}$ , for  $\mathbf{k}$  given  $\mathbf{i}$ ,  $\mathbf{j}$ , and for  $\mathbf{l}$  given  $\mathbf{i}$ ,  $\mathbf{j}$ ,  $\mathbf{k}$ . We obtain

$$S_{\text{off}} \leq \sum_{w,v,r} K^{3w/2} \binom{n}{d} \binom{d}{v}^{1/2} \binom{n-d}{d-v}^{1/2} \binom{v}{r} \binom{n-(2d-v)}{v-w} \binom{2(d-v)}{d-r-(v-w)} \binom{d+w-r}{2w-r}$$

$$\leq \binom{n}{d} \sum_{w,v,r} K^{3w/2} B_1 B_2 B_3 B_4 B_5 B_6,$$
(3.20)

where  $B_m = B_m(n, d, w, v, r)$  denote the corresponding factors in this expression; for example  $B_2 = \binom{n-d}{d-v}^{1/2}$ .

3.7. The terms of the binomial sum. Let us observe a few bounds on the factors  $B_m$ . First,

$$(3.21) B_5 \le 2^{2(d-v)}$$

due to the inequality  $\binom{m}{k} \leq 2^m$ .

Next, since  $v \leq d + w - r$  by Lemma 3.3, we have  $B_3 = \binom{v}{r} \leq \binom{d+w-r}{r}$ . Combining this with  $B_6 = \binom{d+w-r}{2w-r}$ , we get

$$B_3B_6 \le \binom{d+w-r}{r} \binom{d+w-r}{2w-r} \le \binom{d+w-r}{w}^2.$$

Here we used the log-concavity property of binomial coefficients, see Lemma 3.5 in the appendix. Furthermore, we have  $w \le d$  by Lemma 3.1 and  $r \ge 0$ , so

(3.22) 
$$B_3 B_6 \le {2d \choose w}^2 \le (2ed)^{2w},$$

where we used an elementary bound from Lemma 3.4 in the last step.

Next, using the decay of the binomial coefficients (Lemma 3.6), we get

$$B_4 \le \binom{n - (2d - v)}{v - w} \le \left(\frac{v}{n - 2d + 1}\right)^w \binom{n - (2d - v)}{v}.$$

Now recall that  $v \leq d$  (Lemma 3.1) and note that our assumption on d with a sufficiently small constant c implies  $d \leq n/4$ . Thus

$$B_4 \le \left(\frac{2d}{n}\right)^w \binom{n - (2d - v)}{v}.$$

This expression can be conveniently combined with  $B_2^2 = \binom{n-d}{d-v}$ , since

$$B_2^2 B_4 \le \left(\frac{2d}{n}\right)^w \binom{n-d}{d-v} \binom{n-(2d-v)}{v} = \left(\frac{2d}{n}\right)^w \binom{d}{v} \binom{n-d}{d},$$

The last identity can be easily checked by expressing the binomial coefficients in terms of factorials. This expression in turn can be conveniently combined with  $B_1 = \binom{d}{v}^{1/2}$ , and we get

(3.23) 
$$B_1 B_2 B_4 = B_1 \cdot \frac{B_2^2 B_4}{B_2} = \left(\frac{2d}{n}\right)^w \binom{n-d}{d} \cdot \frac{\binom{d}{v}^{3/2}}{\binom{n-d}{d-v}^{1/2}}.$$

Now, using the elementary binomial bounds (Lemma 3.4), we obtain

$$\frac{\binom{d}{v}^{3/2}}{\binom{n-d}{d-v}^{1/2}} = \frac{\binom{d}{d-v}^{3/2}}{\binom{n-d}{d-v}^{1/2}} \le \left(\frac{e^{3/2}d^{3/2}}{(d-v)(n-d)^{1/2}}\right)^{d-v} \le \left(\frac{C_1d^{3/2}}{n^{1/2}}\right)^{d-v}.$$

In the last step we used that  $d - v \ge 1$  by Lemma 3.1 and that  $d \le n/2$ , which follows from our assumption on d if the constant c is chosen sufficiently small. Recall that by  $C_1$ ,  $C_2$ , etc. we denote suitable absolute constants. Returning to (3.23), we have shown that

(3.24) 
$$B_1 B_2 B_4 \le \left(\frac{2d}{n}\right)^w \binom{n}{d} \left(\frac{C_1 d^{3/2}}{n^{1/2}}\right)^{d-v}.$$

3.8. The final bound on the off-diagonal contribution. We can now combine our bounds (3.21), (3.22) and (3.24) on  $B_i$  and put them into (3.20). We obtain

$$S_{\text{off}} \le \binom{n}{d} \sum_{w,v,r} K^{3w/2} B_5 \cdot B_3 B_6 \cdot B_1 B_2 B_4 \le \binom{n}{d}^2 \sum_{w,v,r} \left( \frac{C_2 d^3 K^{3/2}}{n} \right)^w \left( \frac{C_3 d^{3/2}}{n^{1/2}} \right)^{d-v}.$$

Recall from Lemma 3.2 that  $0 \le r \le 2w$ , thus the sum over r includes at most 2w+1 terms. Similarly, Lemma 3.1 determines the ranges for the other two sums, namely  $0 \le w, v \le d-1$ . Hence

(3.25) 
$$S_{\text{off}} \le {n \choose d}^2 \sum_{w=0}^{d-1} (2w+1) \left(\frac{C_2 d^3 K^{3/2}}{n}\right)^w \cdot \sum_{v=0}^{d-1} \left(\frac{C_3 d^{3/2}}{n^{1/2}}\right)^{d-v}.$$

The sums over w and v in the right hand side of (3.25) can be easily estimated. To handle the sum over w, we can use the identity  $\sum_{k=0}^{\infty} kz^k = z/(1-z)^2$ , which is valid for all  $z \in (0,1)$ . Thus, the sum over w is bounded by an absolute constant, as long as  $C_2 d^3 K^{3/2}/n \le 1/2$ . The latter restriction holds by our assumption on d with a sufficiently small constant c.

Similarly, the sum over v in the right hand side of (3.25) is a partial sum of a geometric series. It is dominated by the leading term, i.e. the term where v = d - 1. Hence this sum is bounded by  $C_4 d^{3/2}/n^{1/2}$ , as long as  $C_3 d^{3/2}/n^{1/2} \le 1/2$ . The latter restriction holds by our assumption on d with a sufficiently small constant c.

Summarizing, we obtained the following bound on the off-diagonal contribution (3.5):

$$S_{\text{off}} \lesssim \binom{n}{d}^2 \frac{d^{3/2}}{n^{1/2}}.$$

Combining this with the bound (3.4) on the diagonal contribution and plugging into (3.2), we conclude that

$$\mathsf{E}\left[|\bm{x}^{\mathsf{T}} A \bm{x} - \operatorname{tr} A|^2\right] \lesssim \binom{n}{d}^2 \cdot \frac{K d^2}{n} + \binom{n}{d}^2 \frac{d^{3/2}}{n^{1/2}} \lesssim \binom{n}{d}^2 \cdot \frac{K^{3/4} d^{3/2}}{n^{1/2}}.$$

In the last step, we used the assumption that  $d \lesssim K^{-1/2} n^{1/3}$ . The proof of Theorem 1.9 is complete.

## Appendix. Elementary bounds on binomial coefficients

Here we record some bounds on binomial coefficients used throughout the paper.

**Lemma 3.4** (see e.g. Exercise 0.0.5 in [51]). For any integers  $1 \le d \le n$ , we have:

$$\left(\frac{n}{d}\right)^d \le \binom{n}{d} \le \sum_{k=0}^d \binom{n}{k} \le \left(\frac{en}{d}\right)^d.$$

**Lemma 3.5** (Log-concavity of binomial coefficients). We have

$$\binom{a}{b-c}\binom{a}{b+c} \le \binom{a}{b}^2.$$

for all positive integers a, b and c for which the binomial coefficients are defined.

*Proof.* Expressing the binomial coefficients in terms of factorials, we have

$$\frac{\binom{a}{b-c}\binom{a}{b+c}}{\binom{a}{b}^2} = \frac{b!/(b-c)!}{(b+c)!/b!} \cdot \frac{(a-b)!/(a-b-c)!}{(a-b+c)!/(a-b)!}$$

Examining the first fraction in the right hand side, we find that both the numerator and denominator consist of c terms. Each term in the numerator is bounded by the corresponding terms in the denominator. Thus the fraction is bounded by 1. We argue similarly for the second fraction, and thus the entire quantity is bounded by 1.

**Lemma 3.6** (Decay of binomial coefficients). For any positive integers  $s \leq t \leq m$ , we have

$$\binom{m}{t-s} \le \left(\frac{t}{m-t+1}\right)^s \binom{m}{t}.$$

*Proof.* The definition of binomial coefficients gives

$$\frac{\binom{m}{t-s}}{\binom{m}{t}} = \frac{t(t-1)\cdots(t-s+1)}{(m-t+s)(m-t+s-1)\cdots(m-t+1)} \le \frac{t^s}{(m-t+1)^s}.$$

**Lemma 3.7** (Stability of binomial coefficients). For any positive integers m, p and  $t \leq m$ , we have

$$\binom{m+p}{t} \le (1+\delta) \binom{m}{t} \quad where \ \delta := \frac{2tp}{m+1-t},$$

as long as  $\delta \leq 1/2$ .

*Proof.* The definition of binomial coefficients gives

$$\frac{\binom{m+p}{t}}{\binom{m}{t}} = \prod_{k=1}^{p} \left(1 + \frac{t}{m-t+k}\right) \le \left(1 + \frac{t}{m-t+1}\right)^{p}.$$

Now use the bound  $(1 + \epsilon)^p \le e^{\epsilon p} \le 1 + 2\epsilon p$ , which holds as long as  $\epsilon p \in [0, 1]$ .

### 4. Numerical Experiments

We present a few numerical experiments to verify that the empirical spectral densities for the block-independent model and the random tensor model tend to the Marchenko-Pastur laws. In all of our tests, the numerical results are computed from a single realization, i.e. we did not average over multiple trials.

Block-independent model experiments: In Figure 1 we show the empirical spectral densities for four experiments of block-independent matrices; in each case, they align very well with the corresponding Marchenko-Pastur density. In Figure 1a, the columns of  $X \in \mathbb{R}^{4000 \times 16000}$ consist of n = 2000 blocks, each of length d = 2 where the first entry of the block is  $z \sim N(0, 1)$ and the second entry is  $\frac{1}{\sqrt{2}}(z^2-1)$ . Thus the second entry is completely determined via a formula of the first entry. While this matrix has half the amount of randomness as an i.i.d. matrix of the same size, it still follows the same limiting distribution as the i.i.d. matrix. We see the densities match up very well even for these relatively small sized matrices. In Figure 1b, the columns of  $X \in \mathbb{R}^{1800 \times 12600}$  consist of n = 600 blocks each of length d = 3 where the first and second entry of the block are  $\pm \frac{1}{2}$  each with probability  $\frac{1}{2}$  and the third entry is a shifted XOR of the first and second (i.e. the third entry is  $\frac{1}{2}$  if the first and second entries have opposite signs and it is  $-\frac{1}{2}$  if the first and second entries have the same sign). In this case the variance of the entries is  $\frac{1}{4}$ , so it matches up with Marchenko-Pastur density with covariance matrix  $\Sigma = \frac{1}{4}I$  and  $\lambda = \frac{1}{7}$ . In Figure 1c, the columns of matrix  $X \in \mathbb{R}^{7000 \times 21000}$  have n = 10 blocks, where each block is length d = 700 and is of the form  $\pm \sqrt{d}e_i$  for i selected uniformly from [d], where  $\{e_i\}_{i=1}^d \in \mathbb{R}^d$  are the standard basis vectors in  $\mathbb{R}^d$ . This example shows that with the exchangeability criteria, it is possible for  $n \ll d$ . Additionally, we see the two densities agree very well, despite only having n = 10 blocks. Similar to Figure 1c, in Figure 1d the columns of matrix  $X \in \mathbb{R}^{6400 \times 12800}$  have n = 80 blocks, where each block is length d = 80 and is of the form  $\pm \sqrt{d}e_i$  for i selected uniformly from [d]. These figures and other experiments together suggest that having  $n \geq 10$  and dimensions in the low thousands is enough for the empirical spectral density of a block-independent model matrix to align quite well with the corresponding Marchenko-Pastur density.

Random tensor model experiments: In Figures 2 and 3, we look at vectorized 2-tensors and 3-tensors (d=2 and d=3 respectively). We see that the fourth moment of the entries appears to be important for the speed of convergence as  $n \to \infty$ . For both the 2-tensors and 3-tensors we consider three types of entries in the vector that we will tensor with itself: 1) the entries are Bernoulli  $\pm 1$  each with probability half - these entries have fourth moment of 1; 2) the entries are Uniform on  $[-\sqrt{3},\sqrt{3}]$  - these entries have fourth moment of  $\frac{9}{5}$ ; 3) the entries are standard normal - these entries have fourth moment of 3. In Figure 2 we compare the the empirical spectral density for 2-tensors with the corresponding Marchenko-Pastur density using n = 145. We see that the two densities match up quite well, and match up better when the entries had smaller fourth moments. We do the same experiments for 3-tensors in Figure 3 except now using n = 45, since n = 145 is too computationally costly as it would have  $\binom{145}{3} \approx 500,000$  rows. We see that the two densities match up quite well for the Bernoulli entry case, not very well for the uniform entry case, and very poorly for the standard normal case. These figures suggest there may even be a different limiting law for small values of n. Testing n = 100 does show (Figure 4) that the empirical densities are getting closer to the Marchenko-Pastur density as n increases. These experiments show that while the limiting density does tend to the Marchenko-Pastur density, they do not align very well for small values of n and the rate of convergence likely depends upon the largest fourth moment of the random vector.

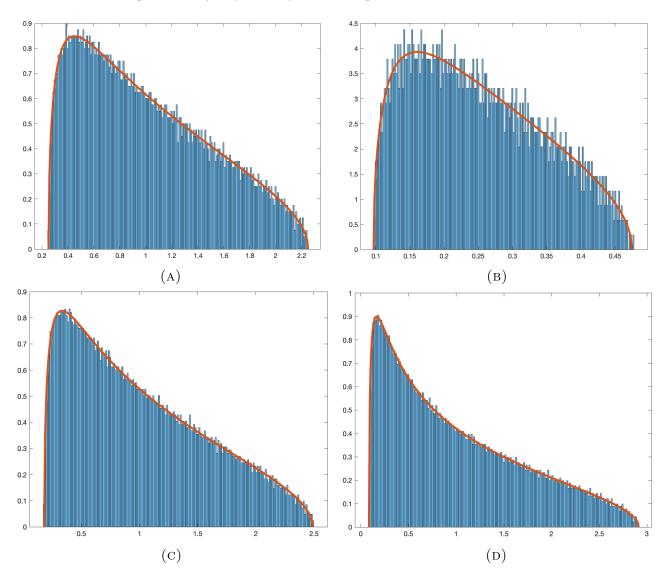


FIGURE 1. The Marchenko-Pastur density (red curve) vs. empirical spectral density for block-independent matrices described in Section 4.

## REFERENCES

- [1] R. Adamczak, Logarithmic Sobolev inequalities and concentration of measure for convex functions and polynomial chaoses, Bull. Pol. Acad. Sci. Math. 53 (2005), 221–238.
- [2] R. Adamczak, On the Marchenko-Pastur and circular law for some classes of random matrices with dependent entries, Electron. J. Prob. 16 (2011), no. 37, 1068-1095.
- [3] R. Adamczak, R. Latala, Tail and moment estimates for chaoses generated by symmetric random variables with logarithmically concave tails, Ann. Inst. Henri Poincaré Probab. Stat. 48 (2012), 1103–1136.
- [4] R. Adamczak, P. Wolff, Concentration inequalities for non-Lipschitz functions with bounded derivatives of higher order, Probability Theory and Related Fields 162 (2015), 531–586.

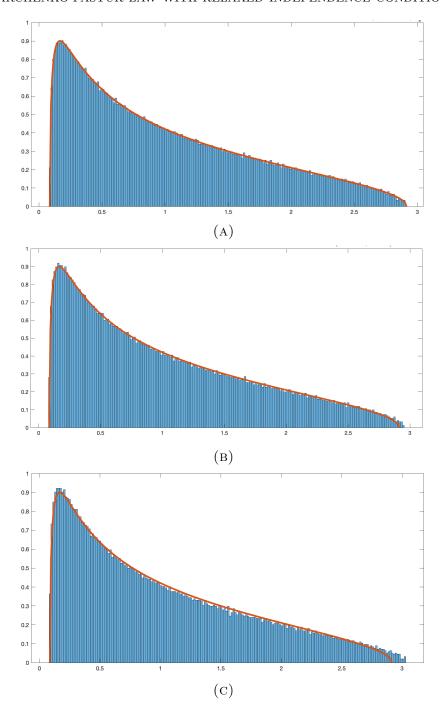


FIGURE 2. The Marchenko-Pastur density (red curve) vs. empirical spectral density for matrices in  $\mathbb{R}^{\binom{145}{2}\times 2\binom{145}{2}}$  whose columns are random 2-tensors as described in Section 4.

- [5] A. Ambainis, A. W. Harrow and M. B. Hastings, Random tensor theory: extending random matrix to mixtures of random product states, Communications in Mathematical Physics, 310 (2012), 25-74.
- [6] A. Anandkumar, R. Ge, D. Hsu, S. Kakade, M. Telgarsky, *Tensor decompositions for learning latent variable models*, The Journal of Machine Learning Research 15 (2014), 2773–2832.

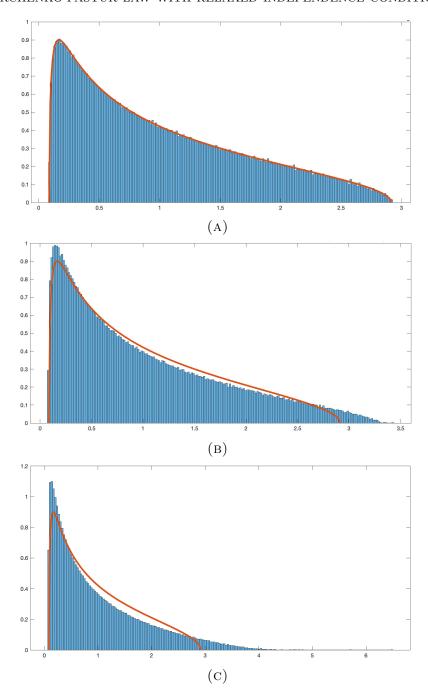


FIGURE 3. The Marchenko-Pastur density (red curve) vs. empirical spectral density for matrices in  $\mathbb{R}^{\binom{45}{3}\times 2\binom{45}{3}}$  whose columns are random 3-tensors as described in Section 4.

<sup>[7]</sup> G. Aubrun, Random points in the unit ball of  $l_p^n$ , Positivity 10 (2006), no. 4, 755–759.

<sup>[8]</sup> A. Auffinger, G. Ben Arous, J. Cerny, Random matrices and complexity of spin glasses, Communications on Pure and Applied Mathematics 66 (2013), 165–201.

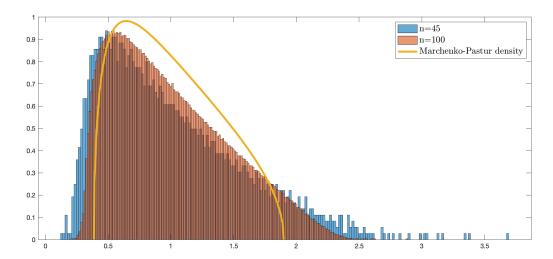


FIGURE 4. Empirical spectral density of  $\frac{1}{\binom{n}{3}}XX^T$ , where columns of  $X^T \in \mathbb{R}^{\binom{n}{3} \times \frac{1}{7}\binom{n}{3}}$  are 3-tensors of a random vector in  $\mathbb{R}^n$  with entries uniform on  $[-\sqrt{3}, \sqrt{3}]$  as described in Section 4.

- [9] Z. D. Bai, Methodologies in spectral analysis of large-dimensional random matrices, a review, Statistica Sinica 9 (1999), 611–677.
- [10] Z. D. Bai and J. W. Silverstein, On the empirical distribution of the eigenvalues of a class of large dimensional random matrices, J. of Multivariate Analysis 54 (1995), no. 2, 175–192.
- [11] Z. D. Bai and J. W. Silverstein, Spectral analysis of large dimensional random matrices. 2nd ed., Springer, 2010.
- [12] Z. Bai and W. Zhou, Large sample covariance matrices without independence structures in columns, Statistica Sinica 18 (2008), 425–442.
- [13] P. Baldi, R. Vershynin, *Polynomial threshold functions, hyperplane arrangements, and random tensors*, SIAM Journal on Mathematics of Data Science, to appear.
- [14] D. Banerjee and A. Bose, *Bulk behavior of some patterned block matrices*, Indian J. Pure Appl. Math. 47 (2016), no. 2, 273–289.
- [15] G. Ben Arous, S. Mei, A. Montanari, M. Nica, *The landscape of the spiked tensor model*, Communications on Pure and Applied Mathematics 72 (2019), 2282–2330.
- [16] W. Bryc, A. Dembo and T. Jiang, Spectral measure of large random Hankel, Markov and Toeplitz matrices, Ann. Prob. 34 (2006), no. 1, 1–38.
- [17] W.-K. Chen, Phase transition in the spiked random tensor with Rademacher prior, The Annals of Statistics 47 (2019), 2734–2756.
- [18] R. Couillet and M. Debbah, Random Matrix Methods for Wireless Communications. Cambridge University Press, 2011.
- [19] E. Dobriban, Efficient computation of limit spectra of sample covariance matrices, Random Matrices: Theory and Applications 04 (2015), no. 4, 1550019–1550055.
- [20] O. Friesen, M. Lowe, and M. Stolz, Gaussian fluctuations for sample covariance matrices with dependent data, Journal of Multivariate Analysis 114 (2013), 270–287.
- [21] R. Ge, F. Huang, C. Jin, Y. Yuan, Escaping from saddle points online stochastic gradient for tensor decomposition, COLT 2015 (Conference on Learning Theory), 797–842.
- [22] F. Götze, H. Sambale, A. Sinulis, Concentration inequalities for polynomials in  $\alpha$ -sub-exponential random variables, preprint (2019).
- [23] F. Götze, A. Tikhomirov. Limit theorems for spectra of positive random matrices under dependence, Zap. Nauchn. Sem. S.-Petersburg. Otdel. Mat. Inst. Steklov. (POMI), Vol. 311 (2004), Veroyatn. i Stat. 7, 92–123, 299.

- [24] F. Götze, A. Tikhomirov. Limit theorems for spectra of random matrices with martingale structure, Stein's Method and Applications, Singapore Univ. Press (2005), 181–195.
- [25] N. El Karoui, Concentration of measure and spectra of random matrices: applications to correlation matrices, elliptical distributions and beyond, Ann. Appl. Probab. 19 (2009), no.6, 2362–2405.
- [26] D. Ghoshdastidar and A. Dukkipati, Consistency of spectral hypergraph partitioning under planted partition model, The Annals of Statistics 45 (2017), 289–315.
- [27] U. Greander and J. W. Silverstein, Spectral analysis of networks with random topologies, SIAM J. Appl. Math. 32 (1977), 499–519.
- [28] W. Hoeffding, A class of statistics with asymptotically normal distribution, Annals of Mathematical Statistics 19 (1948), 293–325.
- [29] K. Hofmann-Credner, M. Stolz, Wigner theorems for random matrices with dependent entries: ensembles associated to symmetric spaces and sample covariance matrices, Electron. Commun. Probab. 13 (2008), 401–414.
- [30] J. Hui and G. M. Pan, Limiting spectral distribution for large sample covariance matrices with m-dependent elements, Commun. Stat. Theory Methods 39, (2010), 935–941.
- [31] P. Jain, S. Oh, Provable tensor factorization with missing data, NIPS 2014 (Advances in Neural Information Processing Systems), 1431–1439.
- [32] D. Jonnson, Some limit theorems for the eigenvalues of a sample covariance matrix, J. Multivariate Anal. 12 (1982), 1–38.
- [33] C. Kim, A. S. Bandeira, M. X. Goemans, Community detection in hypergraphs, spiked tensor models, and sum-of-squares, SampTA 2017 (International Conference on Sampling Theory and Applications), 124–128.
- [34] L. Laloux, P. Cizeau, M. Potters and J-P. Bouchaud, Random matrix theory and financial correlations, International Journal of Theoretical and Applied Finance 1 (2000), no. 03, pp.391–397.
- [35] R. Latala, Estimates of moments and tails of Gaussian chaoses, The Annals of Probability 34 (2006), 2315–2331.
- [36] R. Latala, R. Lochowski, Moment and tail estimates for multidimensional chaos generated by positive random variables with logarithmically concave tails, Stochastic inequalities and applications (2003), 77–92.
- [37] J. Lehec, *Moments of the Gaussian chaos*, Seminaire de Probabilites XLIII, 327–340, Lecture Notes in Math. (2006), Springer, Berlin, 2011.
- [38] J. Lei, K. Chen, B. Lynch, Consistent community detection in multi-layer network data, Biometrika, to appear (2019).
- [39] P. Loubaton, On the almost sure location of the singular values of certain Gaussian block-Hankel large random matrices, J. Theoretical Probab., submitted.
- [40] A. Lytova, Central Limit Theorem for Linear Eigenvalue Statistics for a Tensor Product Version of Sample Covariance Matrices, J. Theor. Prob. 31 (2018), 1024 –1057.
- [41] V. A. Marchenko and L. A. Pastur, Distribution of eigenvalues for some sets of random matrices, Math USSR Sbornik 1 (1967), 457–483.
- [42] A. Montanari, N. Sun, Spectral algorithms for tensor completion, Communications on Pure and Applied Mathematics 71 (2018), 2381–2425.
- [43] N. H. Nguyen, P. Drineas, T. D. Tran, Tensor sparsification via a bound on the spectral norm of random tensors, Information and Inference: A Journal of the IMA 4(2015), 195–229.
- [44] S. O'Rourke, A note on the Marchenko-Pastur law for a class of random matrices with dependent entries, Electron. Commun. Probab. 17 (2012), no. 28 ,1–13.
- [45] A. Pajor and L. Pastur, On the limiting empirical measure of eigenvalues of the sum of rank one matrices with log-concave distribution, Studia Math. 195 (2009), no. 1, 11–29.
- [46] V. Plerou, P. Gopikrishnan, B. Rosenow, L. Amaral, T. Guhr and H. Stanley, *Random matrix approach to cross correlations in financial data*, Phys. Rev. E 64 (2002), 66126-66144.
- [47] E. Richard, A. Montanari, A statistical model for tensor PCA, NIPS 2014 (Advances in Neural Information Processing Systems), 2897–2905.

- [48] M. Rudelson, R. Vershynin, *Hanson-Wright inequality and sub-gaussian concentration*, Electronic Communications in Probability 18 (2013), 1–9.
- [49] J. Silverstein, Strong convergence of the empirical distribution of eigenvalues of large dimensional random matrices, J. Multivariate Analysis 55 (1995), no. 2, 331–339.
- [50] T. Tao, *Topics in random matrix theory*. Graduate Studies in Mathematics, 132. American Mathematical Society, Providence, RI, 2012.
- [51] R. Vershynin, *High-dimensional probability: an introduction with applications in data science*. Cambridge University Press, 2018.
- [52] R. Vershynin, Concentration inequalities for random tensors, submitted (2019).
- [53] K. W. Wachter, The strong limits of random matrix spectra for sample matrices of independent elements, Ann. Probab. 6 (1978), 1-18.
- [54] Y. Wang, H.-Y. Tung, A. J. Smola and A. Anandkumar, Fast and guaranteed tensor decomposition via sketching, NIPS 2015 (Advances in Neural Information Processing Systems), 991–999.
- [55] M. Wei, G. Yang, and L. Ying, The limiting spectral distribution for large sample covariance matrices with unbounded m-dependent entries, Commun. Stat. Theory Methods 45 (2016), 6651–6662.
- [56] P. Yaskov, A short proof of the Marchenko-Pastur theorem, C. R. Math. Acad. Sci. Paris Ser. I, 354 (2016), 319–322.
- [57] P. Yaskov, Necessary and sufficient conditions for the Marchenko-Pastur theorem, Electronic Communications in Probability 21, no. 73 (2016), 1–8.
- [58] J. Yao, A note on a Marchenko-Pastur type theorem for time series, Statist. and Probab. Letters 82 (2012), 20-28.
- [59] Y. Q. Yin and P. R. Krishnaiah, Limit theorems for the eigenvalues of product of large-dimensional random matrices when the underlying distribution is isotropic, Teor. Veroyatnost. i Primenen. 31 (1986), 394-398.
- [60] Y. Q. Yin, Limiting spectral distribution for a class of random matrices, J. Multivariate Anal. 20 (1986), 50-68.
- [61] Z. Zhou, Y. Zhu, Sparse random tensors: concentration, regularization and applications, submitted (2019).