

Convergence of the largest eigenvalue of normalized sample covariance matrices when p and n both tend to infinity with their ratio converging to zero

B.B.Chen and G.M.Pan
School of Physical and Mathematical Sciences
Nanyang Technological University
Singapore

Abstract

Let $\mathbf{X}_p = (\mathbf{s}_1, \dots, \mathbf{s}_n) = (X_{ij})_{p \times n}$ where X_{ij} 's are independent and identically distributed (i.i.d) random variables with $EX_{11} = 0, EX_{11}^2 = 1$ and $EX_{11}^4 < \infty$. It is showed that the largest eigenvalue of the random matrix $\mathbf{A}_p = \frac{1}{2\sqrt{np}}(\mathbf{X}_p\mathbf{X}_p' - n\mathbf{I}_p)$ tends to 1 almost surely as $p \rightarrow \infty, n \rightarrow \infty$ with $p/n \rightarrow 0$.

Keyword empirical distribution, random matrices, maximum eigenvalue.

1. Introduction

Consider the sample covariance type matrix $\mathbf{S} = \frac{1}{n}\mathbf{X}_p\mathbf{X}_p'$, where $\mathbf{X}_p = (\mathbf{s}_1, \dots, \mathbf{s}_n) = (X_{ij})_{p \times n}$ and $X_{ij}, i = 1, \dots, p, j = 1, \dots, n$ are i.i.d. random variables with mean zero and variance 1. For such a matrix, much attention has been paid to asymptotic properties of its eigenvalues in the setting of $p/n \rightarrow c > 0$ as $p \rightarrow \infty$ and $n \rightarrow \infty$. For example, its empirical spectral distribution (ESD) function $F^{\mathbf{S}}(x)$ converges with probability one to the famous Marčenko and Pastur law (see [7] and [4]). Here the ESD for any matrix \mathbf{A} with real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p$ is defined by

$$F^{\mathbf{A}}(x) = \frac{1}{p} \#\{i : \lambda_i \leq x\},$$

where $\#\{\dots\}$ denotes the number of elements of the set. Also, with probability one its maximum eigenvalue and minimum eigenvalue converge, respectively, to the left end point and right end point of the support of Marčenko and Pastur's law (see [5] and [3]).

In contrast with asymptotic behaviors of \mathbf{S} in the case of $p/n \rightarrow c$, the asymptotic properties of \mathbf{S} have not been well understood when $p/n \rightarrow 0$. The first breakthrough was made in Bai and Yin(1988). They considered the normalized matrix

$$\mathbf{A}_p = \frac{1}{2\sqrt{np}}(\mathbf{X}_p\mathbf{X}'_p - nI_p)$$

and proved with probability one

$$F^{\mathbf{A}_p} \rightarrow F(x),$$

which is the so-called semicircle law with a density

$$F'(x) = \begin{cases} \frac{2}{\pi}\sqrt{1-x^2}, & \text{if } |x| \leq 1, \\ 0, & \text{if } |x| > 1. \end{cases}$$

One should note that the semicircle law is also the limit of the empirical spectral distribution of a symmetric random matrix whose diagonal are i.i.d. random variables and above diagonal elements are also i.i.d. (see [10]).

We investigate the maximum eigenvalue of \mathbf{A}_p under the setting of $p/n \rightarrow 0$ as $p \rightarrow \infty$ and $n \rightarrow \infty$ in this paper. The main results are presented in the following theorems.

Theorem 1 *Let $\mathbf{X}_p = (X_{ij})_{p \times n}$ where $\{X_{ij} : i = 1, 2, \dots, p; j = 1, 2, \dots, n\}$ are i.i.d. real random variables with $EX_{11} = 0, EX_{11}^2 = 1$ and $EX_{11}^4 < \infty$. Suppose that $n = n(p) \rightarrow \infty$ and $p/n \rightarrow 0$ as $p \rightarrow \infty$. Define*

$$\mathbf{A}_p = (A_{ij})_{p \times p} = \frac{1}{2\sqrt{np}}(\mathbf{X}_p\mathbf{X}'_p - nI_p).$$

Then as $p \rightarrow \infty$

$$\lambda_{max}(\mathbf{A}_p) \rightarrow 1 \quad a.s.,$$

where $\lambda_{max}(\mathbf{A}_p)$ represents the largest eigenvalue of \mathbf{A}_p .

Indeed, after truncation and normalization of the entries of the matrix \mathbf{A}_p , we may obtain a better result:

Theorem 2 Let $n = n(p) \rightarrow \infty$ and $p/n \rightarrow 0$ as $p \rightarrow \infty$. Define a $p \times p$ random matrix \mathbf{A}_p :

$$\mathbf{A}_p = (A_{ij})_{p \times p} = \frac{1}{2\sqrt{np}}(\mathbf{X}_p \mathbf{X}_p' - nI_p),$$

where $\mathbf{X}_p = (X_{ijp})_{p \times n}$. Suppose that X_{ijp} 's are i.i.d. real random variables and satisfy the following conditions

- 1) $EX_{11p} = 0, EX_{11p}^2 = 1, EX_{11p}^4 < \infty$ and
- 2) $|X_{ijp}| \leq \delta_p \sqrt[4]{np}$, where $\delta_p \downarrow 0$, but $\delta_p \sqrt[4]{np} \uparrow +\infty$, as $p \rightarrow \infty$.

Then, for any $\epsilon > 0, \ell > 0$

$$p(\lambda_{\max}(\mathbf{A}_p) \geq 1 + \epsilon) = o(p^{-\ell}).$$

So far we have considered the sample covariance type matrix \mathbf{S} . However a common used sample covariance matrix in statistics is

$$\mathbf{S}_1 = \frac{1}{n} \sum_{j=1}^n (\mathbf{s}_j - \bar{\mathbf{s}})(\mathbf{s}_j - \bar{\mathbf{s}})',$$

where

$$\bar{\mathbf{s}} = \frac{1}{n} \sum_{j=1}^n \mathbf{s}_j.$$

Similarly we re-normalize it as

$$\mathbf{A}_{p1} = \frac{1}{2} \sqrt{\frac{n}{p}} (\mathbf{S}_1 - I_p).$$

Theorem 3 Under assumptions of Theorem 1, as $p \rightarrow \infty$

$$\lambda_{\max}(\mathbf{A}_{p1}) \rightarrow 1 \quad a.s.,$$

where $\lambda_{\max}(\mathbf{A}_{p1})$ stands for the largest eigenvalues of \mathbf{A}_{p1} .

Estimating a population covariance matrix for high dimension data is a challenging task. Usually, one can not expect the sample covariance matrix to be a consistent estimate of a population covariance matrix when both p and n go to infinity, especially when the orders of p and n are very close to each other. In such circumstance, as argued in [8], operator norm consistent estimation of large population covariance matrix still has nice properties.

Suppose that Σ is a population covariance matrix, nonnegative definite symmetric matrix. Then $\Sigma^{1/2} \mathbf{s}_j, j = 1, \dots, n$ may be viewed as i.i.d. sample drawn from the population with covariance

matrix Σ , where $(\Sigma^{1/2})^2 = \Sigma$. The corresponding sample covariance matrix is

$$\mathbf{S}_2 = \frac{1}{n} \sum_{j=1}^n (\Sigma^{1/2} \mathbf{s}_j - \Sigma^{1/2} \bar{\mathbf{s}})(\Sigma^{1/2} \mathbf{s}_j - \Sigma^{1/2} \bar{\mathbf{s}})'$$

Theorem 3 indicates that the matrix \mathbf{S}_2 is an operator consistent estimation of Σ as long as $p/n \rightarrow 0$ when $p \rightarrow \infty$. Specifically, we have

Theorem 4 *In addition to the assumptions of Theorem 1, assume that $\|\Sigma\|$ is bounded. Then, as $p \rightarrow \infty$*

$$\|\mathbf{S}_2 - \Sigma\| = O\left(\sqrt{\frac{p}{n}}\right) \quad a.s.,$$

where $\|\cdot\|$ stands for the spectral norm of a matrix.

Remark 1 *Related work is [1], where the authors investigated quantitative estimates of the convergence of the empirical covariance matrix in the Log-concave ensemble. Here we obtain a convergence rate of the empirical covariance matrix when the sample vectors are in the form of $\Sigma^{1/2} \mathbf{s}_j$.*

1 Proof of Theorem 1

Throughout the paper, C denotes a constant whose value may vary from line to line. Also, all limits in the paper are taken as $p \rightarrow \infty$.

It follows from Theorem in [2] that

$$\liminf_{p \rightarrow \infty} \lambda_{\max}(\mathbf{A}_p) \geq 1 \quad a.s.. \quad (1)$$

Thus, it suffices to show that

$$\limsup_{p \rightarrow \infty} \lambda_{\max}(\mathbf{A}_p) \leq 1 \quad a.s.. \quad (2)$$

Let $\hat{\mathbf{A}}_p = \frac{1}{2\sqrt{np}}(\hat{\mathbf{X}}_p \hat{\mathbf{X}}_p' - nI_p)$, where $\hat{\mathbf{X}}_p = (\hat{X}_{ij})_{p \times n}$ and $\hat{X}_{ij} = X_{ij}I(|X_{ij}| \leq \delta_p \sqrt[4]{np})$ where δ_p is chosen as the larger of δ_p constructed as in (3) and δ_p as in (5). On the one hand, since $EX_{11}^4 < \infty$ for any $\delta > 0$ we have

$$\lim_{p \rightarrow \infty} \delta^{-4} E|X_{11}|^4 I(|X_{11}| > \delta \sqrt[4]{np}) = 0.$$

Since the above is true for arbitrary positive δ there exists a sequence of positive δ_p such that

$$\lim_{p \rightarrow \infty} \delta_p = 0, \quad \lim_{p \rightarrow \infty} \delta_p^{-4} E|X_{11}|^4 I(|X_{11}| > \delta_p \sqrt[4]{np}) = 0, \quad \delta_p \sqrt[4]{np} \uparrow +\infty. \quad (3)$$

On the other hand, since $EX_{11}^4 < \infty$ for any $\varepsilon > 0$

$$\sum_{k=1}^{\infty} 2^k P\left(|X_{11}| > \varepsilon 2^{\frac{k}{4}}\right) < \infty.$$

In view of the arbitrariness of ε there is a sequence of positive number ε_k such that

$$\varepsilon_k \rightarrow 0, \text{ as } k \rightarrow \infty, \quad \sum_{k=1}^{\infty} 2^k P\left(|X_{11}| > \varepsilon_k 2^{\frac{k}{4}}\right) < \infty. \quad (4)$$

For each k , let p_k be the maximum p such that $n(p)p \leq 2^k$. For $p(k-1) < p \leq p(k)$, set

$$\delta_p = 2\varepsilon_k. \quad (5)$$

Let $Z_t = X_{ij}, t = (i-1)n + j$ and obviously $\{Z_t\}$ are *i.i.d.* We then conclude from (4) and (5) that

$$\begin{aligned} P(\mathbf{A}_p \neq \hat{\mathbf{A}}_p, i.o) &\leq \lim_{K \rightarrow \infty} P\left(\bigcup_{k=K}^{\infty} \bigcup_{p_{k-1} < p \leq p_k} \bigcup_{i \leq p, j \leq n} \{|X_{ij}| > \delta_p \sqrt[4]{np}\}\right) \\ &\leq \lim_{K \rightarrow \infty} \sum_{k=K}^{\infty} P\left(\bigcup_{p_{k-1} < p \leq p_k} \bigcup_{t=1}^{2^k} \{|Z_t| > \varepsilon_k 2^{\frac{k}{4}}\}\right) \\ &= \lim_{K \rightarrow \infty} \sum_{k=K}^{\infty} P\left(\bigcup_{t=1}^{2^k} \{|Z_t| > \varepsilon_k 2^{\frac{k}{4}}\}\right) \\ &\leq \lim_{K \rightarrow \infty} \sum_{k=K}^{\infty} 2^k P\left(|Z_1| > \varepsilon_k 2^{\frac{k}{4}}\right) \\ &= 0 \quad a.s. \end{aligned}$$

It follows that $\lambda_{\max}(\mathbf{A}_p) - \lambda_{\max}(\hat{\mathbf{A}}_p) \rightarrow 0 \quad a.s.$ as $p \rightarrow \infty$.

From now on we write δ for δ_p to simplify notation. Moreover, set $\tilde{\mathbf{A}}_p = \frac{1}{2\sqrt{np}}(\tilde{\mathbf{X}}_p \tilde{\mathbf{X}}_p' - nI_p)$, where $\tilde{\mathbf{X}}_p = (\tilde{X}_{ij})_{p \times n}$ and $\tilde{X}_{ij} = \frac{\hat{X}_{ij} - E\hat{X}_{11}}{\sigma}$. Here $\sigma^2 = E(\hat{X}_{11} - E\hat{X}_{11})^2$ and $\sigma^2 \rightarrow 1$ as $p \rightarrow \infty$.

We obtain via (3)

$$|E\hat{X}_{11}| \leq \frac{E|X_{11}|^4 I(|X_{11}| > \delta_p \sqrt[4]{np})}{\delta^3 (np)^{3/4}} \leq \frac{C}{(np)^{3/4}}, \quad (6)$$

and

$$|\sigma^2 - 1| \leq CE|X_{11}|^2 I(|X_{11}| > \delta \sqrt[4]{np}) \leq \frac{E|X_{11}|^4 I(|X_{11}| > \delta \sqrt[4]{np})}{\delta^2 \sqrt{np}} = o\left(\frac{1}{\sqrt{np}}\right). \quad (7)$$

We conclude from the Rayleigh-Ritz theorem that

$$|\lambda_{\max}(\tilde{\mathbf{A}}_p) - \lambda_{\max}(\hat{\mathbf{A}}_p)|$$

$$\begin{aligned}
&\leq \frac{1}{2\sqrt{np}} \left| \sup_{\|\mathbf{z}\|=1} \left(\sum_{i \neq j} z_i z_j \sum_{k=1}^n \hat{X}_{ik} \hat{X}_{jk} + \sum_{i=1}^p z_i^2 \sum_{k=1}^n (\hat{X}_{ik}^2 - 1) \right) \right. \\
&\quad \left. - \sup_{\|\mathbf{z}\|=1} \left(\sum_{i \neq j} z_i z_j \sum_{k=1}^n \tilde{X}_{ik} \tilde{X}_{jk} + \sum_{i=1}^p z_i^2 \sum_{k=1}^n (\tilde{X}_{ik}^2 - 1) \right) \right| \\
&\leq \frac{1}{2\sqrt{np}} \left| 1 - \frac{1}{\sigma^2} \right| \sup_{\|\mathbf{z}\|=1} \left| \sum_{i \neq j} z_i z_j \frac{1}{\sqrt{np}} \sum_{k=1}^n \hat{X}_{ik} \hat{X}_{jk} + \sum_{i=1}^p z_i^2 \sum_{k=1}^n (\hat{X}_{ik}^2 - 1) \right| \\
&\quad + \frac{1}{2\sqrt{np}} \frac{2|EX_{11}|}{\sigma^2} \sup_{\|\mathbf{z}\|=1} \left| \sum_{i=1}^p \sum_{j=1}^p z_i z_j \sum_{k=1}^n \hat{X}_{ik} \right| \\
&\quad + \frac{1}{2\sqrt{np}} \frac{n|EX_{11}|^2}{\sigma^2} \sup_{\|\mathbf{z}\|=1} \left| \sum_{i=1}^p \sum_{j=1}^p z_i z_j \right| + \frac{n}{2\sqrt{np}} \left| 1 - \frac{1}{\sigma^2} \right| \\
&= A_1 + A_2 + A_3 + A_4.
\end{aligned}$$

By (7) and the strong law of large numbers, we have

$$\begin{aligned}
A_1 &= \frac{|\sigma^2 - 1|}{2\sqrt{np}\sigma^2} \sup_{\|\mathbf{z}\|=1} \left| \sum_{k=1}^n \left(\left(\sum_{i=1}^p z_i \hat{X}_{ik} \right)^2 - \sum_{i=1}^p z_i^2 \hat{X}_{ik}^2 \right) + \sum_{i=1}^p z_i^2 \sum_{k=1}^n (\hat{X}_{ik}^2 - 1) \right| \\
&\leq \frac{|\sigma^2 - 1|\sqrt{np}}{2\sigma^2} \cdot \frac{1}{np} \left(2 \left| \sum_{k=1}^n \sum_{i=1}^p \hat{X}_{ik}^2 \right| + \left| \sum_{i=1}^p \sum_{k=1}^n (\hat{X}_{ik}^2 - 1) \right| \right) \\
&\leq \frac{|\sigma^2 - 1|\sqrt{np}}{2\sigma^2} \cdot \frac{1}{np} \left(3 \left| \sum_{k=1}^n \sum_{i=1}^p X_{ik}^2 \right| + np \right) \\
&\rightarrow 0 \quad a.s..
\end{aligned}$$

Similarly, (6) and the strong law of large numbers yield

$$\begin{aligned}
A_2 &\leq \frac{1}{2\sqrt{np}} \cdot \frac{2|E\hat{X}_{11}|}{\sigma^2} \sup_{\|\mathbf{z}\|=1} \left| \sum_{j=1}^p z_j \left| \sum_{i=1}^p z_i \sum_{k=1}^n \hat{X}_{ik} \right| \right| \\
&\leq \frac{1}{2\sqrt{np}} \cdot \frac{2|E\hat{X}_{11}|}{\sigma^2} \cdot \sqrt{p} \cdot \left(\sum_{i=1}^p \left(\sum_{k=1}^n \hat{X}_{ik} \right)^2 \right)^{1/2} \\
&\leq \frac{C}{\sigma^2(np)^{1/4}} \left| \frac{1}{np} \sum_{i=1}^p \sum_{k=1}^n \hat{X}_{ik}^2 \right|^{1/2} \\
&\leq \frac{C}{\sigma^2(np)^{1/4}} \left| \frac{1}{np} \sum_{i=1}^p \sum_{k=1}^n X_{ik}^2 \right|^{1/2} \\
&\rightarrow 0 \quad a.s..
\end{aligned}$$

It is straightforward to conclude from (6) and (7) that

$$A_3 \rightarrow 0 \quad a.s. \quad A_4 \rightarrow 0 \quad a.s.$$

Thus, we have $\lambda_{\max}(\hat{\mathbf{A}}_p) - \lambda_{\max}(\tilde{\mathbf{A}}_p) \rightarrow 0 \quad a.s.$ By the above results, to prove (2), it is sufficient to show that $\limsup_{p \rightarrow \infty} \lambda_{\max}(\tilde{\mathbf{A}}_p) \leq 1 \quad a.s.$ To this end, we note that the matrix $\tilde{\mathbf{A}}_p$ satisfies all the assumptions in Theorem 2. Therefore we obtain (2) by Theorem 2 (whose argument is given in the next section). Together with (1), we finishes the proof of Theorem 1.

2 Proof of Theorem 2

Suppose that $\mathbf{z} = (z_1, \dots, z_p)$ is a unit vector. By the Rayleigh-Ritz theorem we then have

$$\begin{aligned} \lambda_{\max}(\mathbf{A}_p) &= \max_{\|z\|=1} \left(\sum_{i,j} z_i z_j A_{ij} \right) \\ &= \max_{\|z\|=1} \left(\sum_{i \neq j} z_i z_j A_{ij} + \sum_{i=1}^p z_i^2 A_{ii} \right) \\ &\leq \lambda_{\max}(\mathbf{B}_p) + \max_{i \leq p} |A_{ii}|, \end{aligned}$$

where $\mathbf{B}_p = (B_{ij})_{p \times p}$ with

$$B_{ij} = \begin{cases} 0, & \text{if } i = j, \\ \frac{1}{2\sqrt{np}} \sum_{k=1}^n X_{ikp} X_{jkp}, & \text{if } i \neq j. \end{cases}$$

To prove Theorem 2, it is sufficient to prove, for any $\epsilon > 0, \ell > 0$

$$P(\lambda_{\max}(\mathbf{B}_p) > 1 + \epsilon) = o(p^{-\ell}) \quad (8)$$

and

$$P\left(\max_{i \leq p} \frac{1}{\sqrt{np}} \left| \sum_{j=1}^n (X_{ij}^2 - 1) \right| > \epsilon\right) = o(p^{-\ell}). \quad (9)$$

We first prove (9). To simplify notation, let $Y_j = X_{1jp}^2 - 1$ and $C_1 = E|Y_1|^2$. Then $EY_j = 0$. Choose an appropriate sequence $h = h_p$ such that it satisfies, as $p \rightarrow \infty$

$$\begin{cases} h / \log p \rightarrow \infty \\ \delta^2 h / \log p \rightarrow 0 \\ \frac{\delta^4 p}{C_1} \geq \sqrt{p}. \end{cases} \quad (10)$$

We then have

$$\begin{aligned}
P(\max_{i \leq p} \frac{1}{\sqrt{np}} |\sum_{j=1}^n (X_{ij}^2 - 1)| > \epsilon) &\leq p \cdot P\left(|\sum_{j=1}^n (X_{1jp}^2 - 1)| > \epsilon \sqrt{np}\right) \\
&\leq \epsilon^{-h} p (\sqrt{np})^{-h} E|\sum_{j=1}^n Y_j|^h \\
&\leq \epsilon^{-h} p (\sqrt{np})^{-h} \sum_{m=1}^{h/2} \sum_{1 \leq j_1 < j_2 < \dots < j_m \leq n} \sum_{\substack{i_1+i_2+\dots+i_m=h \\ i_1 \geq 2, \dots, i_m \geq 2}} \frac{h!}{i_1! i_2! \dots i_m!} E|Y_{j_1}|^{i_1} E|Y_{j_2}|^{i_2} \dots E|Y_{j_m}|^{i_m} \\
&\leq \epsilon^{-h} p (\sqrt{np})^{-h} \sum_{m=1}^{h/2} \sum_{\substack{i_1+i_2+\dots+i_m=h \\ i_1 \geq 2, \dots, i_m \geq 2}} \frac{n!}{m!(n-m)!} \frac{h!}{i_1! i_2! \dots i_m!} E|Y_1|^{i_1} E|Y_1|^{i_2} \dots E|Y_1|^{i_m} \\
&\leq \epsilon^{-h} p \sum_{m=1}^{h/2} \sum_{\substack{i_1+i_2+\dots+i_m=h \\ i_1 \geq 2, \dots, i_m \geq 2}} n^m \frac{h!}{i_1! i_2! \dots i_m!} C_1^m (\delta^2 \sqrt{np})^{h-2m} \\
&\leq \epsilon^{-h} p \sum_{m=1}^{h/2} m^h \left(\frac{\delta^4 p}{C_1}\right)^{-m} \delta^{2h} \leq \epsilon^{-h} p \frac{h}{2} \cdot \left(\frac{\delta^2 h}{\log(\delta^4 p / C_1)}\right)^h \\
&\leq \left(\left(\frac{ph}{2}\right)^{1/h} \cdot \frac{2\delta^2 h}{\log p} \cdot \epsilon^{-1}\right)^h \leq \left(\frac{\xi}{\epsilon}\right)^h = o(p^{-\ell}),
\end{aligned}$$

where ξ is a constant satisfying $0 < \xi < \epsilon$. Below are some interpretations of the above inequalities:

a) The fifth inequality is because, $\frac{n!}{m!(n-m)!} < n^m$, $|Y_1| < \delta^2 \sqrt{np}$.

b) We use the fact $\sum_{\substack{i_1+i_2+\dots+i_m=h \\ i_1 \geq 2, \dots, i_m \geq 2}} \frac{h!}{i_1! i_2! \dots i_m!} < m^h$ in the sixth inequality.

c) The seventh inequality uses the fact that for any $a > 1, b > 0, t \geq 1$, $\frac{b}{\log a} > 1$, $a^{-t} t^b \leq \left(\frac{b}{\log a}\right)^b$.

d) The last two inequalities are due to (10).

e) With the facts that $\frac{\xi}{\epsilon} < 1, h/\log p \rightarrow \infty$, the last equality is true.

Thus (9) follows.

Next consider (8). For any $\epsilon > 0$, we have

$$\begin{aligned}
P(\lambda_{\max}(\mathbf{B}_p) \geq 1 + \epsilon) &\leq \frac{E\lambda_{\max}^k(\mathbf{B}_p)}{(1 + \epsilon)^k} \leq \frac{E\text{tr}(\mathbf{B}_p^k)}{(1 + \epsilon)^k} \\
&\leq \frac{1}{(1 + \epsilon)^k} \cdot \frac{1}{(2\sqrt{np})^k} \sum E(X_{i_1 j_1} X_{i_2 j_1} X_{i_2 j_2} X_{i_3 j_2} \dots X_{i_k j_k} X_{i_1 j_k}),
\end{aligned}$$

where $k = k_p$ satisfies, as $p \rightarrow \infty$

$$\left\{ \begin{array}{l} k/\log p \rightarrow \infty \\ \delta^{1/3}k/\log p \rightarrow 0 \quad , \\ \frac{\delta^2 \sqrt[4]{p}}{k^3} \geq 1 \end{array} \right.$$

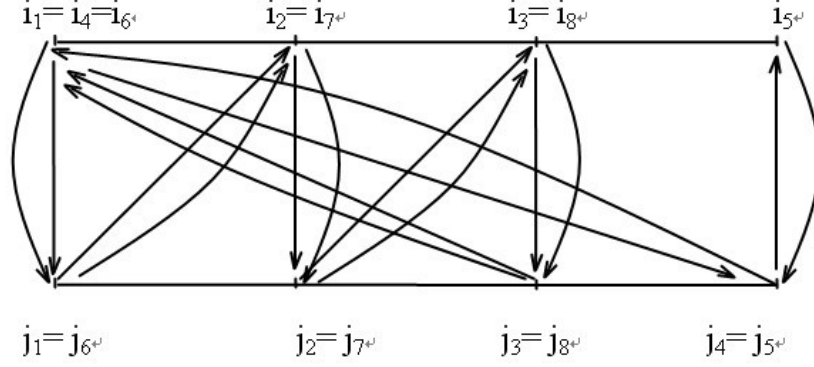
and the summation is taken with respect to j_1, j_2, \dots, j_k running over all integers in $\{1, 2, \dots, n\}$ and i_1, i_2, \dots, i_k running over all integers in $\{1, 2, \dots, p\}$ subject to the condition that $i_1 \neq i_2, i_2 \neq i_3, \dots, i_k \neq i_1$.

In order to get an up bound for $|\sum EX_{i_1 j_1} X_{i_2 j_1} \dots X_{i_k j_k} X_{i_1 j_k}|$, we need to construct a graph for given i_1, \dots, i_k and j_1, \dots, j_k , as in [5], [9] and [3]. We follow the presentation in [3] and [9] to introduce some fundamental concepts associated with the graph.

For the sequence (i_1, i_2, \dots, i_k) from $\{1, 2, \dots, p\}$ and the sequence (j_1, j_2, \dots, j_k) from $\{1, 2, \dots, n\}$, we define a directed graph as follows. Plot two parallel real lines, referred to as *I-line* and *J-line*, respectively. Draw $\{i_1, i_2, \dots, i_k\}$ on the *I-line*, called *I-vertices* and draw $\{j_1, j_2, \dots, j_k\}$ on the *J-line*, known as *J-vertices*. The vertices of the graph consist of the *I-vertices* and *J-vertices*. The edges of the graph are $\{e_1, e_2, \dots, e_{2k}\}$, where for $a = 1, \dots, k$, $e_{2a-1} = i_a j_a$ are called the column edges and $e_{2a} = j_a i_{a+1}$ are called row edges with the convention that $i_{2k+1} = i_1$. For each column edge e_{2a-1} , the vertices i_a and j_a are called the ends of the edge $i_a j_a$ and moreover i_a and j_a are, respectively, the initial and the terminal of the edge $i_a j_a$. Each row edge e_{2a} starts from the vertex j_b and ends with the vertex i_{b+1} .

Two vertices are said to coincide if they are both in the *I-line* or both in the *J-line* and they are identical. That is $i_a = i_b$ or $j_a = j_b$. Readers are also reminded that the vertices i_a and j_b are not coincident even if they have the same value because they are in different lines. We say that two edges are coincident if two edges have the same set of ends.

The graph constructed above is said to be a *W-graph* if each edge in the graph coincides with at least one other edge. Below is an example of a *W-graph*:



Two graphs are said to be isomorphic if one becomes another by an appropriate permutation on $\{1, 2, \dots, p\}$ of *I-vertices* and an appropriate permutation on $\{1, 2, \dots, n\}$ of *J-vertices*. A *W-graph* is called a canonical graph if $i_a \leq \max\{i_1, i_2, \dots, i_{a-1}\} + 1$ and $j_a \leq \max\{j_1, j_2, \dots, j_{a-1}\} + 1$ with $i_1 = j_1 = 1$, where $a = 1, 2, \dots, k$.

In the canonical graph, if $i_{a+1} = \max\{i_1, i_2, \dots, i_a\} + 1$, then the edge $j_a i_{a+1}$ is called a row innovation and if $j_a = \max\{j_1, j_2, \dots, j_{a-1}\} + 1$, then the edge $i_a j_a$ is called a column innovation. Apparently, a row innovation and a column innovation, respectively, lead to a new I-vertex and a new J-vertex except the first column innovation $i_1 j_1$ leading to a new I-vertex i_1 and a new J-vertex j_1 .

We now classify all edges into three types, T_1 , T_3 and T_4 . Let T_1 denote the set of all innovations including row innovations and column innovations. We further distinguish the column innovations as follows. An edge $i_a j_a$ is called a T_{11} edge if it is a column innovation and the edge $j_a i_{a+1}$ is a row innovation; An edge $i_b j_b$ is referred to as a T_{12} edge if it is a column innovation but $j_b i_{b+1}$ is not a row innovation. An edge e_j is said to be a T_3 edge if there is an innovation edge $e_i, i < j$ so that e_j is the first one to coincide with e_i . An edge is called a T_4 edge if it does not belong a T_1 edge or T_3 edge. The first appearance of a T_4 edge is referred to as a T_2 edge. There are two kinds of T_2 edges: (a) the first appearance of an edge that coincides with a T_3 edge, denoted by T_{21} edge; (b) the first appearance of a single non innovation, denoted by T_{22} edge.

We say that an edge e_i is single up to the edge $e_j, j \geq i$ if it does not coincide with any other edges among e_1, \dots, e_j except itself. A T_3 edge e_i is said to be regular if there are more than one innovation, adjacent to e_i and single up to e_{i-1} , among the edges $\{e_1, \dots, e_{i-1}\}$.

Corresponding to the above classification of the edges, we introduce the following notation and

list some useful facts.

1. Denote by l the total number of innovations.
2. Let r be the number of the row innovations. Moreover, let c denote the column innovations. We then have $r + c = l$.
3. Define r_1 to be the number of the T_{11} edges. Then $r_1 \leq r$ by the definition of a T_{11} edge. Also, the number of the T_{12} edges is $l - r - r_1$.
4. Let t be the number of the T_2 edges. Note that the number of the T_3 edges is the same as the number of the innovations and there are a total of $2k$ edges in the graph. It follows that the number of the T_4 edges is $2k - 2l$. On the other hand, each T_2 edge is also a T_4 edge. Therefore $t \leq 2k - 2l$.
5. Define μ to be the number of T_{21} edges. Obviously, $\mu \leq t$. The number of T_{22} edge is then $t - \mu$. Since each T_{21} edge coincides with one innovation, we let $n_i, i = 1, 2, \dots, \mu$ denote the number of T_4 edges which coincide with the i -th such innovation, $n_i \geq 0$.
6. $\mu_1 = \#\{i : n_i = 1, i = 1, 2, \dots, \mu\}$, where $\#\{\cdot\}$ denotes the cardinality of the set $\{\cdot\}$.
7. Let $m_j, j = 1, 2, \dots, t - \mu$ denote the number of T_4 edges which coincide with and include the j -th T_{22} edge. Note that $m_j \geq 2$.

We now claim that

$$\begin{aligned}
Etr(B_p^k) &\leq (2\sqrt{np})^{-k} \sum E(X_{i_1 j_1} X_{i_2 j_1} \dots X_{i_k j_k} X_{i_1 j_k}) \\
&= (2\sqrt{np})^{-k} \sum' \sum'' \sum''' \sum_* E(X_{i_1 j_1} X_{i_2 j_1} \dots X_{i_k j_k} X_{i_1 j_k}) \\
&\leq (2\sqrt{np})^{-k} \sum_{l=1}^k \sum_{r=1}^l \sum_{r_1=0}^r \sum_{t=0}^{2k-2l} \sum_{\mu=0}^t \sum_{\mu_1=0}^{\mu} \sum_* \binom{k}{r} \binom{r}{r_1} \binom{k-r_1}{l-r-r_1} \binom{2k-l}{l} \\
&\quad \times k^{3t} (t+1)^{6k-6l} (\delta \sqrt[4]{np})^{2k-2l-2t+\mu_1} p^{r+1} n^{l-r},
\end{aligned} \tag{11}$$

where the summation \sum' is with respect to different arrangements of three types of edges at the $2k$ different positions, the summation \sum'' over different canonical graphs with a given arrangement of the three types of edges for $2k$ positions, the third summation \sum''' with respect to all isomorphic graphs for a given canonical graph and the last notation \sum_* denotes the constraint that $i_1 \neq i_2, i_2 \neq i_3, \dots, i_k \neq i_1$.

Now, we explain why the above estimate is true:

- (i). The factor $(2\sqrt{np})^{-k}$ is obvious.
- (ii). If the graph is not a W -graph, which means there is a single edge in the graph, then the mean of the product of X_{ij} corresponding to this graph is zero (since $EX_{11} = 0$). Thus we have $l \leq k$. Moreover, the facts that $r \leq l$, $r_1 \leq r$, $t \leq 2k - 2l$, $\mu \leq t$ and $\mu_1 \leq \mu$ are easily obtained from the fact 1 to the fact 7 listed before.
- (iii). There are at most $\binom{k}{r}$ ways to choose r edges out of the k row edges to be the r row innovations. Subsequently, we consider how to select the column innovations. Observe that the definition of T_{11} edges, there are $\binom{r}{r_1}$ ways to select r_1 row innovations out of the total r row innovations so that the edge before each such r_1 row innovations is a T_{11} edge, column innovation. Moreover there are at most $\binom{k-r_1}{l-r-r_1}$ ways to choose $l-r-r_1$ edges out of the remaining $k-r_1$ column edges to be the $l-r-r_1$ T_{12} edges, the remaining column innovations.
- (iv). Given the position of the l innovations, there are at most $\binom{2k-l}{l}$ ways to select l edges out of the $2k-l$ edges to be T_3 edges. And the rest positions are for the T_4 edges. Therefore the first summation \sum' is bounded by $\sum_{l=1}^k \sum_{r=1}^l \sum_{r_1=0}^r \binom{k}{r} \binom{r}{r_1} \binom{k-r_1}{l-r-r_1} \binom{2k-l}{l}$.
- (v). By definition, each innovation (or each irregular T_3 edges) is uniquely determined by the subgraph prior to the innovation (or the irregular T_3). Moreover by Lemma 3.2 in [9] for each regular T_3 edge, there are at most $t+1$ innovations so that the regular T_3 edge coincides with one of them and by Lemma 3.3 in [9] there are at most $2t$ regular T_3 edges. Therefore there are at most $(t+1)^{2t} \leq (t+1)^{2(2k-2l)}$ ways to draw the regular T_3 edges.
- (vi). Once the positions of the innovations and the T_3 edges are fixed there are at most $\binom{(r+1)c}{t} \leq \binom{k^2}{t} \leq k^{2t}$ ways to arrange the t T_2 edges. After t positions of T_2 edges are determined there are at most t^{2k-2l} ways to distribute $2k-2l$ T_4 edges among the t positions. So there are at most $k^{2t} \cdot t^{2k-2l}$ ways to arrange T_4 edges. It follows that \sum'' is bounded by $\sum_{t=0}^{2k-2l} (t+1)^{2(2k-2l)} k^{2t} \cdot t^{2k-2l}$.
- (vii). The third summation \sum''' is bounded by $n^c p^{r+1}$ because the number of graphs in the isomorphic class for a given graph is $p(p-1) \cdots (p-r)n(n-1) \cdots (n-c+1)$.

(viii). Recalling the definitions of $l, r, t, \mu, \mu_1, n_i, m_i$ we have

$$EX_{i_1 j_1} X_{i_2 j_1} \dots X_{i_k j_k} X_{i_1 j_k} = (EX_{11}^2)^{l-\mu} \left(\prod_{i=1}^{\mu} EX_{11}^{n_i+2} \right) \left(\prod_{i=1}^{t-\mu} EX_{11}^{m_i} \right), \quad (12)$$

where $\sum_{i=1}^{\mu} n_i + \sum_{i=1}^{t-\mu} m_i = 2k - 2l$. Without loss of generality, we suppose $n_1 = n_2 = \dots = n_{\mu_1} = 1$ and $n_{\mu_1+1}, \dots, n_{\mu} \geq 2$ for convenience. It is easy to check that

$$E|X_{11}^s| \leq \begin{cases} M(\delta \sqrt[4]{np})^{s-4}, & \text{if } s \geq 4, M = \max\{EX_{11}^4, |EX_{11}^3|\} \\ (\delta \sqrt[4]{np})^{s-2}, & \text{if } s \geq 2. \end{cases}$$

Thus, (12) becomes

$$\begin{aligned} & |EX_{i_1 j_1} X_{i_2 j_1} \dots X_{i_k j_k} X_{i_1 j_k}| \\ & \leq \sum_{\mu=0}^t \sum_{\mu_1=0}^{\mu} |EX_{11}^3|^{\mu_1} |EX_{11}^4|^{t-\mu_1} (\delta \sqrt[4]{np})^{\sum_{i=\mu_1+1}^{\mu} n_i - 2(\mu - \mu_1)} (\delta \sqrt[4]{np})^{\sum_{i=1}^{t-\mu} m_i - 2(t-\mu)} \\ & \leq \sum_{\mu=0}^t \sum_{\mu_1=0}^{\mu} M^t (\delta \sqrt[4]{np})^{2k-2l-2t+\mu_1} \\ & \leq \sum_{\mu=0}^t \sum_{\mu_1=0}^{\mu} k^t (\delta \sqrt[4]{np})^{2k-2l-2t+\mu_1}, \quad \text{when } k \text{ is large enough.} \end{aligned} \quad (13)$$

The above points regarding the T_2 edges are discussed for $t > 0$, but they are still valid when $t = 0$ with the convention that $0^0 = 1$ in the term t^{2k-2l} , because in this case there are only T_1 edges and T_3 edges in the graph and thus $l = k$.

Consider the constraint \sum_{*} now. Note that for each T_{12} edge, say $i_a j_a$, it is a column innovation, but the next row edge $j_a i_{a+1}$ is not a row innovation. Since $i_{a+1} \neq i_a$, the edge $j_a i_{a+1}$ cannot coincide with the edge $i_a j_a$. Moreover, it also doesn't coincide with any edges before the edge $i_a j_a$ since j_a is a new vertex. So $j_a i_{a+1}$ must be a T_{22} edge. Thus, the number of the T_{12} edges can not exceed the number of the T_{22} edges. This implies $l - r - r_1 \leq t - \mu$. Moreover, note that $\mu_1 \leq \mu$. We then have

$$\begin{aligned} & n^{-k/2} p^{-k/2} n^{l-r} p^{r+1} (np)^{k/2-l/2-t/2+\mu_1/4} \\ & = (n/p)^{l/2} \cdot n^{-r-t/2+\mu_1/4} p^{r+1-t/2+\mu_1/4} \\ & \leq \left(\sqrt{\frac{p}{n}} \right)^{r-r_1} \cdot p^{-t/2} p. \end{aligned} \quad (14)$$

We thus conclude from (11) and (14) that

$$\begin{aligned}
Etr(B_p^k) &\leq 2^{-k} \sum_{l=1}^k \sum_{r=1}^l \sum_{r_1=0}^r \sum_{t=0}^{2k-2l} \sum_{\mu=0}^t \sum_{\mu_1=0}^{\mu} \binom{k}{r} \binom{r}{r_1} \binom{k-r_1}{l-r-r_1} \binom{2k-l}{l} \\
&\quad \times \left(\sqrt{\frac{p}{n}} \right)^{r-r_1} p^{-t/2} p k^{3t} (t+1)^{6k-6l} \delta^{2k-2l-2t+\mu_1}.
\end{aligned} \tag{15}$$

Moreover we claim that

$$\begin{aligned}
&p \left[2^{-k} \binom{k}{r} \right] \left[\binom{r}{r_1} \left(\sqrt{\frac{p}{n}} \right)^{r-r_1} \right] \left[\binom{k-r_1}{l-r-r_1} \delta^{l-r-r_1} \right] \\
&\quad \times \left[\binom{2k-l}{l} \left(\frac{\sqrt{p}\delta^3}{k^3} \right)^{-t} (t+1)^{6k-6l} \delta^{2k-2l} \right] \delta^{-(l-r-r_1)+3t-(2k-2l)} \cdot \delta^{2k-2l-2t+\mu_1} \\
&\leq p^2 \left(1 + \sqrt{\frac{p}{n}} \right)^k (1+\delta)^k \left(1 + \frac{24^3 k^3 \delta}{\log^3 p} \right)^{2k}.
\end{aligned} \tag{16}$$

Indeed, the above claim is based on the following five facts.

$$1) \quad 2^{-k} \binom{k}{r} \leq 2^{-k} \sum_{r=0}^k \binom{k}{r} = 1.$$

$$2) \quad \binom{r}{r_1} \left(\sqrt{\frac{p}{n}} \right)^{r-r_1} = \binom{r}{r-r_1} \left(\sqrt{\frac{p}{n}} \right)^{r-r_1} \leq \sum_{s=0}^r \binom{r}{s} \left(\sqrt{\frac{p}{n}} \right)^s = \left(1 + \sqrt{\frac{p}{n}} \right)^r \leq \left(1 + \sqrt{\frac{p}{n}} \right)^k.$$

$$3) \quad \binom{k-r_1}{l-r-r_1} \delta^{l-r-r_1} \leq \sum_{s=0}^{k-r_1} \binom{k-r_1}{s} \delta^s = (1+\delta)^{k-r_1} \leq (1+\delta)^k.$$

4) By the fact that $\binom{2k-l}{l} \leq \binom{2k}{2l}$, and the inequality $a^{-t}(t+1)^b \leq a \left(\frac{b}{\log a} \right)^b$, for $a > 1$, $b > 0$, $t > 0$ and $\frac{\delta^2 \sqrt{p}}{k^3} \geq \sqrt[4]{p}$, we have

$$\begin{aligned}
&\binom{2k-l}{l} \left(\frac{\sqrt{p}\delta^3}{k^3} \right)^{-t} (t+1)^{6k-6l} \delta^{2k-2l} \\
&\leq \binom{2k}{2l} \frac{\sqrt{p}\delta^3}{k^3} \left(\frac{6k-6l}{\log(\sqrt{p}\delta^3/k^3)} \right)^{6k-6l} \delta^{2k-2l} \\
&\leq p \binom{2k}{2l} \left(\frac{24k}{\log p} \right)^{6k-6l} \delta^{2k-2l} \\
&\leq p \binom{2k}{2l} \left(\frac{24^3 k^3 \delta}{\log^3 p} \right)^{2k-2l} \\
&\leq p \sum_{s=0}^{2k} \binom{2k}{s} \left(\frac{24^3 k^3 \delta}{\log^3 p} \right)^{2k-s} \\
&= p \left(1 + \frac{24^3 k^3 \delta}{\log^3 p} \right)^{2k}.
\end{aligned}$$

5) When p is large enough, $\delta^{-(l-r-r_1)+3t-(2k-2l)} \cdot \delta^{2k-2l-2t+\mu_1} = \delta^{t-(l-r-r_1)} \cdot \delta^{\mu_1} \leq 1$, since $\delta \rightarrow 0$ and $l - r - r_1 \leq t - \mu$.

Summarizing (15) and (16) we obtain that

$$\begin{aligned}
Etr(B_p^k) &\leq k \cdot l \cdot r \cdot (2k - 2l) \cdot t \cdot \mu \cdot p^2 \left(1 + \sqrt{\frac{p}{n}}\right)^k (1 + \delta)^k \left(1 + \frac{24^3 l^3 \delta}{\log^3 p}\right)^{2k} \\
&\leq 8k^6 p^2 \left(1 + \sqrt{\frac{p}{n}}\right)^k (1 + \delta)^k \left(1 + \frac{24^3 l^3 \delta}{\log^3 p}\right)^{2k} \\
&= \left((8k^6)^{1/k} p^{2/k} \left(1 + \sqrt{\frac{p}{n}}\right) (1 + \delta) \left(1 + \frac{24^3 k^3 \delta}{\log^3 p}\right)^2 \right)^k \\
&\leq \eta^k,
\end{aligned}$$

where η is a constant satisfying $1 < \eta < 1 + \epsilon$. Here the last inequality uses the facts below:

- i. $(p^2)^{1/k} \rightarrow 1$, because $k/\log p \rightarrow \infty$,
- ii. $(8k^6)^{1/k} \rightarrow 1$, because $k \rightarrow \infty$,
- iii. $\left(1 + \sqrt{\frac{p}{n}}\right) \rightarrow 1$, because $p/n \rightarrow 0$,
- iv. $(1 + \delta) \rightarrow 1$, because $\delta \rightarrow 0$,
- v. $\frac{24^3 \cdot k^3 \delta}{\log p^3} \rightarrow 0$, because $\frac{\delta^{1/3} k}{\log p} \rightarrow 0$.

It follows that

$$P(\lambda_{max}(\mathbf{B}_p) > 1 + \epsilon) \leq \left(\frac{\eta}{1 + \epsilon}\right)^k = o(p^{-\ell})$$

since $k/\log p \rightarrow \infty$ and $\frac{\eta}{1+\epsilon} < 1$. The proof is complete.

3 Proof of Theorem 3

Note that

$$\mathbf{S}_1 = \mathbf{S} - \bar{s}\bar{s}' \tag{17}$$

By the Fan inequality

$$\sup_x |F^{\mathbf{A}_{p1}}(x) - F^{\mathbf{A}_p}(x)| \leq \frac{1}{p}.$$

Thus from theorem in [2] we see that

$$F^{\mathbf{A}_{p1}}(x) \xrightarrow{a.s.} F(x),$$

specified in the introduction. It follows that

$$\liminf_{p \rightarrow \infty} \lambda_{\max}(\mathbf{A}_{p1}) \geq 1.$$

Let \mathbf{z} be a unit vector. In view of (17) we obtain

$$\mathbf{z}' \mathbf{A}_{p1} \mathbf{z} = \mathbf{z}' \mathbf{A}_p \mathbf{z} - \frac{1}{2} \sqrt{\frac{n}{p}} \mathbf{z}' \bar{\mathbf{S}} \mathbf{z} \leq \mathbf{z}' \mathbf{A}_p \mathbf{z},$$

which implies that

$$\lambda_{\max}(\mathbf{A}_{p1}) \leq \lambda_{\max}(\mathbf{A}_p).$$

This, together with Theorem 1, finishes the proof of Theorem 3.

4 Proof of Theorem 4

Theorem 4 follows from Theorem 3 and the fact that

$$\|\mathbf{S}_2 - \Sigma\| = \|\Sigma^{1/2}(\mathbf{S}_1 - \mathbf{I}_p)\Sigma^{1/2}\| \leq \|\mathbf{S}_1 - \mathbf{I}_p\| \|\Sigma\|.$$

References

- [1] Adamczak, R., Litvak, A, Pajor, A. and Tomczak-Jaegermann. N. (2010) Quantitative estimates of the convergence of the empirical covariance matrix in Log-concave ensembles. *J. Amer. Math. Soc.* **23**, 535-561.
- [2] Bai,Z.D. and Yin,Y.Q. 1988. Convergence to the semicircle law. *Ann. Probab.* **16**, 863-875.
- [3] Bai,Z.D. and Yin,Y.Q. (1993) limit of the smallest eigenvalue of a large dimensional sample covariance matrix. *Ann. Probab.* **21** , 1275-1294.
- [4] Jonsson, D. (1982) Some limit theorems of the eigenvalues of sample covariance matrix. *J. Multivariate Anal.*, **12**, 1-38.
- [5] Geman, S. (1980) A limit theorem for the norm of random matrices. *Ann. Probab.*, **8**, 252-261.
- [6] Fan, K. (1951) Maximum properties and inequalities for the eigenvalues of completely continuous operators. *Proc. Nat. Acad. Sci. U.S.A.* , **37**, 760-766.

- [7] Marčenko, V. A. & Pastur, L. A., Distribution for some sets of random matrices. *Math. USSR-Sb.*, **1** (1967), 457-483.
- [8] Karoui, E. N. Operator norm consistent estimation of large dimensional sparse covariance matrices., *Ann.Stat.*, **6**. 2008, 2717-2756.
- [9] Yin, Y.Q., Bai, Z.D. and Krishnaiah, P.R. 1988. On the Limit of the Largest Eigenvalue of the Large Dimensional Sample Covariance Matrix. *Prob theory and its related fields.* **20**, 50-68.
- [10] Wigner, E. P. (1958). On the distributions of the roots fo certain symmetric matrices, *Ann. Math.* **67** 325-327.