

Spectral distributions of high-dimensional sample correlation matrices under infinite variance

Johannes Heiny

Ruhr-University Bochum

Joint work with

Jianfeng Yao (HKU),

Thomas Mikosch and Jorge Yslas (Copenhagen).

Random Matrices and Complex Data Analysis Workshop,
December 10-12, 2019, Shanghai

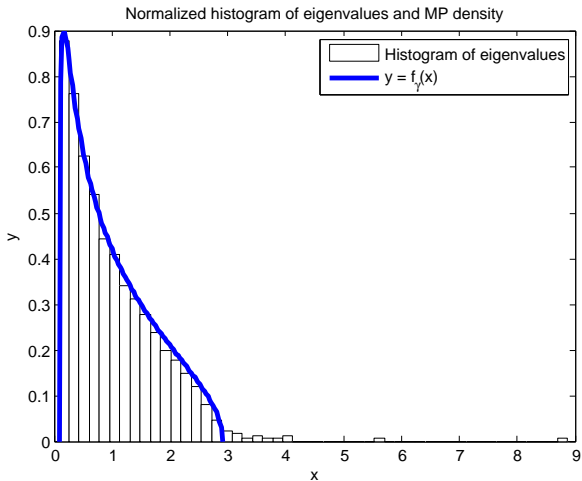


Figure: These are **NOT** spikes!

Setup for the picture

- **Data matrix** $\mathbf{X} = \mathbf{X}_n$: $p \times n$ matrix with iid centered entries and generic variable $X \stackrel{d}{=} X_{11}$.

$$\mathbf{X} = (X_{it})_{i=1,\dots,p;t=1,\dots,n}$$

- **Sample covariance matrix** $\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}'$
- **Ordered eigenvalues** of \mathbf{S}

$$\lambda_1(\mathbf{S}) \geq \lambda_2(\mathbf{S}) \geq \dots \geq \lambda_p(\mathbf{S})$$

- **Sample correlation matrix**

$$\mathbf{R} = (\text{diag}(\mathbf{S}))^{-1/2} \mathbf{S} (\text{diag}(\mathbf{S}))^{-1/2} .$$

- **Regular variation** with index $\alpha > 0$:

$$\mathbb{P}(|X| > x) = x^{-\alpha} L(x),$$

where L is a slowly varying function.

- This implies $\mathbb{E}[|X|^{\alpha+\varepsilon}] = \infty$ for any $\varepsilon > 0$.
- **Normalizing sequence** (a_{np}^2) such that

$$np \mathbb{P}(X^2 > a_{np}^2 x) \rightarrow x^{-\alpha/2}, \quad \text{as } n \rightarrow \infty \text{ for } x > 0.$$

Then $a_{np} = (np)^{1/\alpha} \ell(np)$ for a slowly varying function ℓ .

Diagonal

\mathbf{X} with iid regularly varying entries $\alpha \in (0, 4)$ and $p = n^\beta$ with $\beta \in [0, 1]$. We have

$$a_{np}^{-2} \|\mathbf{X}\mathbf{X}' - \text{diag}(\mathbf{X}\mathbf{X}')\| \xrightarrow{\mathbb{P}} 0,$$

where $\|\cdot\|$ denotes the spectral norm.

$$(\mathbf{X}\mathbf{X}')_{ij} = \sum_{t=1}^n X_{it}X_{jt}.$$

- **Weyl's inequality**

$$\max_{i=1,\dots,p} |\lambda_i(\mathbf{A} + \mathbf{B}) - \lambda_i(\mathbf{A})| \leq \|\mathbf{B}\|.$$

- Choose $\mathbf{A} + \mathbf{B} = \mathbf{X}\mathbf{X}'$ and $\mathbf{A} = \text{diag}(\mathbf{X}\mathbf{X}')$ to obtain

$$a_{np}^{-2} \max_{i=1,\dots,p} |\lambda_i(\mathbf{X}\mathbf{X}') - \lambda_i(\text{diag}(\mathbf{X}\mathbf{X}'))| \xrightarrow{\mathbb{P}} 0, \quad n \rightarrow \infty.$$

- **Note:** Limit theory for $(\lambda_i(\mathbf{S}))$ reduced to (S_{ii}) .

Theorem (Heiny and Mikosch, 2016)

\mathbf{X} with iid regularly varying entries $\alpha \in (0, 4)$ and $p_n = n^\beta \ell(n)$ with $\beta \in [0, 1]$.

- ① If $\beta \in [0, 1]$, then

$$a_{np}^{-2} \max_{i=1, \dots, p} |\lambda_i(\mathbf{X}\mathbf{X}') - \lambda_i(\text{diag}(\mathbf{X}\mathbf{X}'))| \xrightarrow{\mathbb{P}} 0.$$

- ② If $\beta \in ((\alpha/2 - 1)_+, 1]$, then

$$a_{np}^{-2} \max_{i=1, \dots, p} |\lambda_i(\mathbf{X}\mathbf{X}') - X_{(i), np}^2| \xrightarrow{\mathbb{P}} 0.$$

Example: Eigenvalues

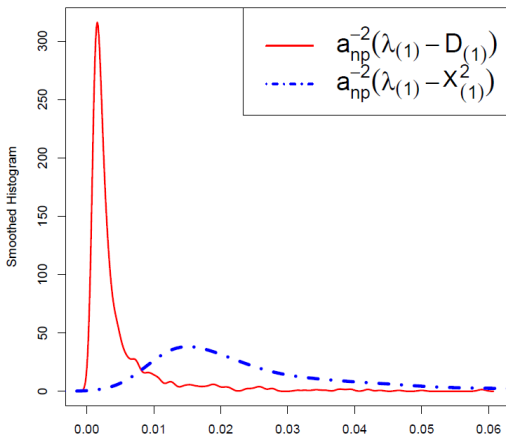


Figure: Smoothed histogram based on 20000 simulations of the approximation error for the normalized eigenvalue $a_{np}^{-2}\lambda_1(\mathcal{S})$ for entries X_{it} with $\alpha = 1.6$, $\beta = 1$, $n = 1000$ and $p = 200$.

- \mathbf{v}_k unit eigenvector of \mathbf{S} associated to $\lambda_k(\mathbf{S})$
- Unit eigenvectors of $\text{diag}(\mathbf{S})$ are canonical basisvectors \mathbf{e}_j .

Eigenvectors

\mathbf{X} with iid regularly varying entries with index $\alpha \in (0, 4)$ and $p_n = n^\beta \ell(n)$ with $\beta \in [0, 1]$. Then for any fixed $k \geq 1$,

$$\|\mathbf{v}_k - \mathbf{e}_{L_k}\|_{\ell_2} \xrightarrow{\mathbb{P}} 0, \quad n \rightarrow \infty.$$

Localization vs. Delocalization

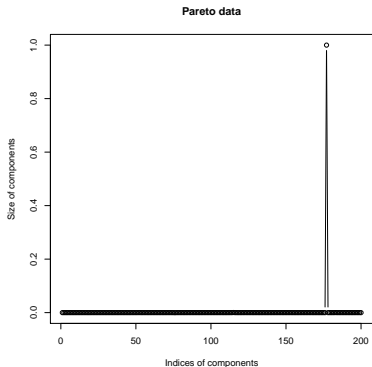


Figure: $X \sim \text{Pareto}(0.8)$

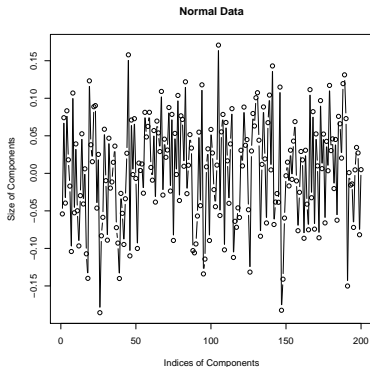


Figure: $X \sim N(0, 1)$

Components of eigenvector \mathbf{v}_1 . $p = 200$, $n = 1000$.

Point process convergence

$$N_n = \sum_{i=1}^p \delta_{a_n^{-2} \lambda_i(\mathbf{X}\mathbf{X}')} \xrightarrow{d} \sum_{i=1}^{\infty} \delta_{\Gamma_i^{-2/\alpha}} = N$$

The limit is a PRM on $(0, \infty)$ with mean measure $\mu(x, \infty) = x^{-\alpha/2}, x > 0$, and

$$\Gamma_i = E_1 + \cdots + E_i, \quad (E_i) \text{ iid standard exponential.}$$

- **Limiting distribution:** For $k \geq 1$,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(a_{np}^{-2} \lambda_k \leq x) &= \lim_{n \rightarrow \infty} \mathbb{P}(N_n(x, \infty) < k) = \mathbb{P}(N(x, \infty) < k) \\ &= \sum_{s=0}^{k-1} \frac{(x^{-\alpha/2})^s}{s!} e^{-x^{-\alpha/2}}, \quad x > 0. \end{aligned}$$

- **Limiting distribution:** For $k \geq 1$,

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{P}(a_{np}^{-2} \lambda_k \leq x) &= \lim_{n \rightarrow \infty} \mathbb{P}(N_n(x, \infty) < k) = \mathbb{P}(N(x, \infty) < k) \\ &= \sum_{s=0}^{k-1} \frac{(x^{-\alpha/2})^s}{s!} e^{-x^{-\alpha/2}}, \quad x > 0.\end{aligned}$$

- **Largest eigenvalue**

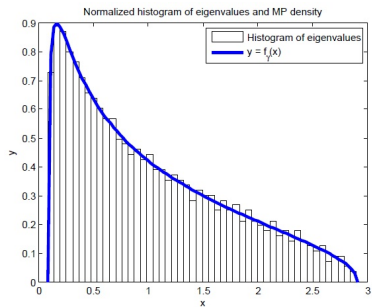
$$\frac{n}{a_{np}^2} \lambda_1(\mathbf{S}) \xrightarrow{d} \Gamma_1^{-\alpha/2},$$

where the limit has a *Fréchet distribution* with parameter $\alpha/2$.

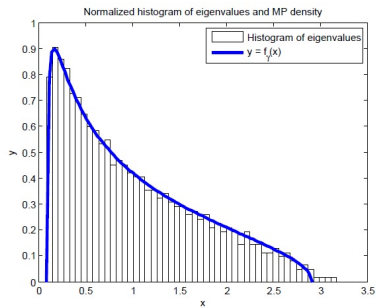
Soshnikov (2006), Auffinger et al. (2009), Auffinger and Tang (2016),

Davis et al. (2014, 2016²), JH and Mikosch (2016)

$$\alpha = 3.99$$

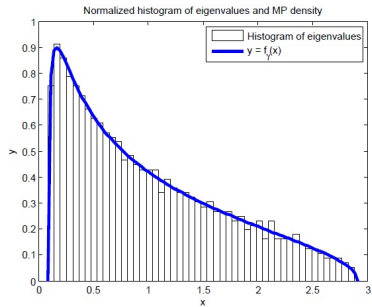


(a) Sample correlation

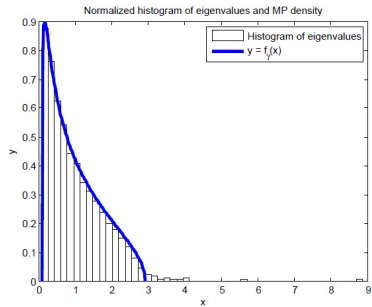


(b) Sample covariance

$$\alpha = 3.99, n = 2000, p = 1000$$



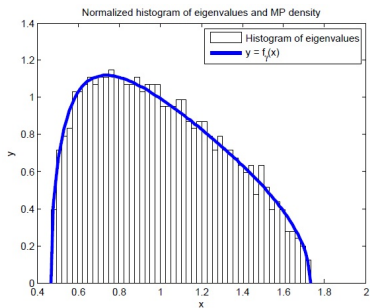
(a) Sample correlation



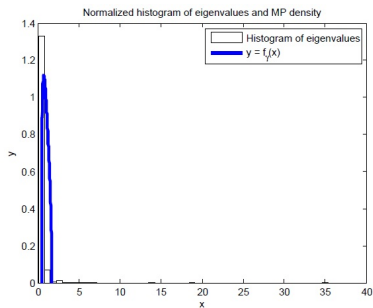
(b) Sample covariance

$$\alpha = 3, n = 2000, p = 1000$$

$$\alpha = 2.1$$



(a) Sample correlation



(b) Sample covariance

$$\alpha = 2.1, n = 10000, p = 1000$$

Limiting spectral distribution of $(\mathbf{X}\mathbf{X}')$ under $\mathbb{E}[X^2] = \infty$:

- Regular variation with $\alpha < 2$:

$$F_{a_{n+p}^{-2}} \mathbf{X}\mathbf{X}' \rightarrow G_{\alpha}^{\gamma} \text{ weakly,}$$

whose density g_{α}^{γ} satisfies

$$g_{\alpha}^{\gamma}(x) \sim c x^{-1-\alpha/2}, \quad x \rightarrow \infty.$$

Ben Arous and Guionnet (2008), Belinschi et al. (2009)

- Assumption: X symmetric and regularly varying with index $\alpha \in (0, 2)$.
- **Goal:** For $k \geq 1$, find the limit of

$$\mathbb{E}\left[\int x^k F_{\mathbf{R}}(dx)\right] = \frac{1}{p}\mathbb{E}[\text{tr}(\mathbf{R}^k)]$$

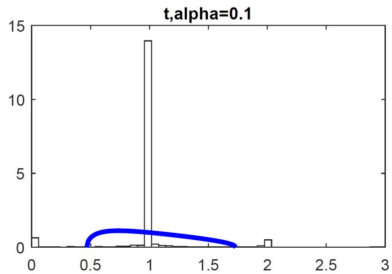
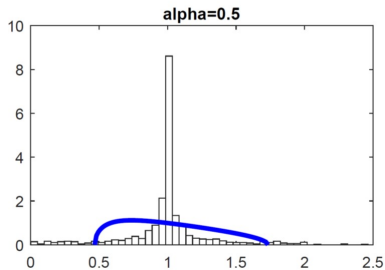
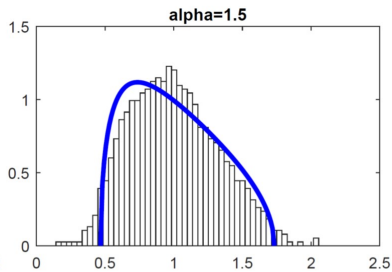
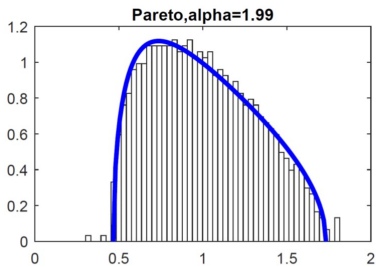
One has

$$\mathbb{E}[\text{tr}(\mathbf{R}^k)] = \sum_{i_1, \dots, i_k=1}^p \underbrace{\sum_{t_1, \dots, t_k=1}^n \mathbb{E}[Y_{i_1 t_1} Y_{i_2 t_1} \cdots Y_{i_k t_k} Y_{i_1 t_k}]}_{:=F(i_1, \dots, i_k)}.$$

Assumption: X symmetric $\Rightarrow Y_{ij}$ symmetric

$$Y_{ij} = \frac{X_{ij}}{\sqrt{\sum_{t=1}^n X_{it}^2}}$$

$$\begin{aligned}
 \frac{1}{p} \mathbb{E}[\text{tr}(\mathbf{R}^k)] &\rightarrow \beta_k(\gamma) + \frac{2}{\alpha} \sum_{r=2}^{k-2} \gamma^{r-1} \sum_{q=0}^{r-2} (\Gamma(1 - \alpha/2))^{-r+q+1} \\
 &\sum_{I \in \mathcal{C}_{r,k}^{(q)}} \sum_{s=1}^{t^*(\tilde{I})} \left(\frac{\alpha/2}{\Gamma(1 - \alpha/2)} \right)^s \sum_{T \in \mathcal{C}_{s,|\tilde{I}|}(\tilde{I})} \left(\prod_{i=1}^{r-q} \frac{\Gamma(d_i(\tilde{I}, T))}{\Gamma(N_i(\tilde{I}))} \right) \\
 &\prod_{(i,t) \in \Delta(\tilde{I}, T)} \Gamma\left(\frac{m_{it}(\tilde{I}, T) - \alpha}{2} \right).
 \end{aligned}$$





- **Random walk**

$$S_n = X_1 + \cdots + X_n, \quad n \geq 1.$$

- ① (X_i) are iid random variables with generic element X .
- ② $\mathbb{E}[X] = 0$ and $\mathbb{E}[X^2] = 1$.

- **Dimension** $p = p_n \rightarrow \infty$
- Consider iid copies $(S_n^{(i)})_{i \leq p}$ of S_n and define the **point process**

$$N_n = \sum_{i=1}^p \delta_{d_p(S_n^{(i)}/\sqrt{n}-d_p)}.$$

We want to prove:

$$N_n = \sum_{i=1}^p \delta_{d_p(S_n^{(i)}/\sqrt{n}-d_p)} \xrightarrow{d} N, \quad n \rightarrow \infty,$$

where N is a Poisson random measure with mean measure $\mu(x, \infty) = e^{-x}$, $x \in \mathbb{R}$, and

$$d_p = \sqrt{2 \log p} - \frac{\log \log p + \log 4\pi}{2(2 \log p)^{1/2}}.$$

We want to prove:

$$N_n = \sum_{i=1}^p \delta_{d_p(S_n^{(i)}/\sqrt{n}-d_p)} \xrightarrow{d} N, \quad n \rightarrow \infty,$$

where N is a Poisson random measure with mean measure $\mu(x, \infty) = e^{-x}$, $x \in \mathbb{R}$, and

$$d_p = \sqrt{2 \log p} - \frac{\log \log p + \log 4\pi}{2(2 \log p)^{1/2}}.$$

- Note: d_p is the centering and normalizing sequence for the maximum of p iid standard normals.
- By [Resnick \(2007\)](#), this is equivalent to

$$p \mathbb{P}(d_p(S_n/\sqrt{n} - d_p) > x) \rightarrow e^{-x}, \quad x \in \mathbb{R}.$$

H., Mikosch, Yslas (2019+)

Assume that the sequence (p_n) satisfies the following conditions:

(C1) $p = O(n^{(s-2)/2})$ for $s > 2$ if $\mathbb{E}[|X|^s] < \infty$.

(C2) $p = \exp(o(n^{1/3}))$ if $\mathbb{E}[\exp(h|X|)] < \infty$ for some $h > 0$.

Then

$$p \mathbb{P}(d_p(S_n/\sqrt{n} - d_p) > x) \rightarrow e^{-x}, \quad x \in \mathbb{R}.$$

H., Mikosch, Yslas (2019+)

Assume that the sequence (p_n) satisfies the following conditions:

(C1) $p = O(n^{(s-2)/2})$ for $s > 2$ if $\mathbb{E}[|X|^s] < \infty$.

(C2) $p = \exp(o(n^{1/3}))$ if $\mathbb{E}[\exp(h|X|)] < \infty$ for some $h > 0$.

Then

$$p \mathbb{P}(d_p(S_n/\sqrt{n} - d_p) > x) \rightarrow e^{-x}, \quad x \in \mathbb{R}.$$

Precise large deviation bounds of the type

$$\sup_{0 \leq y \leq \gamma_n} \left| \frac{\mathbb{P}(S_n/\sqrt{n} > y)}{\overline{\Phi}(y)} - 1 \right| \rightarrow 0, \quad n \rightarrow \infty,$$

Under (C1): $\gamma_n = \sqrt{(s-2) \log n}$, Michel (1974)

Under (C2): $\gamma_n = o(n^{1/6})$, Petrov (1972)

- **Data matrix** $\mathbf{X} = \mathbf{X}_n$: $p \times n$ matrix with iid entries with generic element X .

$$\mathbf{X} = (X_{it})_{i=1,\dots,p;t=1,\dots,n}$$

- **Sample covariance matrix**

$$\mathbf{S} = \mathbf{X}\mathbf{X}'$$

Dependent random walks

$$S_{ij} = \sum_{t=1}^n X_{it}X_{jt}, \quad i < j.$$

- **Off-diagonal point process:**

$$N_n^S = \sum_{1 \leq i < j \leq p} \delta_{\tilde{d}_p(S_{ij}/\sqrt{n} - \tilde{d}_p)},$$

where $\tilde{d}_p = d_{p(p-1)/2}$.

$$N_n^S = \sum_{1 \leq i < j \leq p} \delta_{\tilde{d}_p(S_{ij}/\sqrt{n} - \tilde{d}_p)}$$

Theorem: H., Mikosch, Yslas (2019+)

Assume that the sequence (p_n) satisfies:

- $p = O(n^{(s-2)/4})$ for $s > 2$ if $\mathbb{E}[|X|^s] < \infty$.
- $p = \exp(o(n^{1/3}))$ if $\mathbb{E}[\exp(h |X_{11}X_{12}|)] < \infty$ for some $h > 0$.

Then

$$N_n^S \xrightarrow{d} N.$$

Remark: Entries of \mathbf{X} do not have to be identically distributed.

- Note that

$$N = \sum_{i=1}^{\infty} \delta_{-\log \Gamma_i},$$

where $\Gamma_i = E_1 + \cdots + E_i$, $i \geq 1$, and (E_i) is iid standard exponential.

- Note that

$$N = \sum_{i=1}^{\infty} \delta_{-\log \Gamma_i},$$

where $\Gamma_i = E_1 + \dots + E_i$, $i \geq 1$, and (E_i) is iid standard exponential.

- For fixed k ,

$$\tilde{d}_p(S_{(i)}/\sqrt{n} - \tilde{d}_p)_{i=1,\dots,k} \xrightarrow{d} (-\log \Gamma_i)_{i=1,\dots,k}.$$

- In particular, **Jiang (2004)**

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\tilde{d}_p(S_{(1)}/\sqrt{n} - \tilde{d}_p) \leq x\right) = \exp(-e^{-x}).$$

- **Fang Han's** talk on Tuesday, **Songxi Chen's** talk on Wednesday

Sample correlation matrix

$$\mathbf{R} = (\text{diag}(\mathbf{S}))^{-1/2} \mathbf{S} (\text{diag}(\mathbf{S}))^{-1/2}$$

Sample correlation matrix

$$\mathbf{R} = (\text{diag}(\mathbf{S}))^{-1/2} \mathbf{S} (\text{diag}(\mathbf{S}))^{-1/2}$$

Theorem: H., Mikosch, Yslas (2019+)

Assume that the sequence (p_n) satisfies:

- $p = O(n^{(s-2)/4})$ for $s > 2$ if $\mathbb{E}[|X|^s] < \infty$.
- $p = \exp(o(n^{1/3}))$ if $\mathbb{E}[\exp(h |X_{11} X_{12}|)] < \infty$ for some $h > 0$.

Then

$$N_n^R = \sum_{1 \leq i < j \leq p} \delta_{\tilde{d}_p(\sqrt{n}R_{ij} - \tilde{d}_p)} \xrightarrow{d} N.$$

Thank you!

(Z_{it}) : iid field of regularly varying random variables.

- **Stochastic volatility model:**

$$\mathbf{X} = (Z_{it} \sigma_{it}^{(n)})_{p \times n}$$

(Z_{it}) : iid field of regularly varying random variables.

- **Stochastic volatility model:**

$$\mathbf{X} = (Z_{it} \sigma_{it}^{(n)})_{p \times n}$$

- **Generate deterministic covariance structure \mathbf{A} :**

$$\mathbf{X} = \mathbf{A}^{1/2} \mathbf{Z}$$

Davis et al. (2014)

Heavy Tails and Dependence

(Z_{it}) : iid field of regularly varying random variables.

- **Dependence among rows and columns:**

$$X_{it} = \sum_{l=0}^{\infty} \sum_{k=0}^{\infty} h_{kl} Z_{i-k,t-l}$$

with some constants h_{kl} . Davis et al. (2016)

Heavy Tails and Dependence

(Z_{it}) : iid field of regularly varying random variables.

- **Dependence among rows and columns:**

$$X_{it} = \sum_{l=0}^{\infty} \sum_{k=0}^{\infty} h_{kl} Z_{i-k,t-l}$$

with some constants h_{kl} . Davis et al. (2016)

- **Relation to iid case:**

$$\mathbf{X}\mathbf{X}' = \sum_{l_1, l_2=0}^{\infty} \sum_{k_1, k_2=0}^{\infty} h_{k_1 l_1} h_{k_2 l_2} \mathbf{Z}(k_1, l_1) \mathbf{Z}'(k_2, l_2),$$

where

$$\mathbf{Z}(k, l) = (Z_{i-k,t-l})_{i=1,\dots,p;t=1,\dots,n}, \quad l, k \in \mathbb{Z}.$$

Heavy Tails and Dependence

(Z_{it}) : iid field of regularly varying random variables.

- **Dependence among rows and columns:**

$$X_{it} = \sum_{l=0}^{\infty} \sum_{k=0}^{\infty} h_{kl} Z_{i-k,t-l}$$

with some constants h_{kl} . Davis et al. (2016)

- **Relation to iid case:**

$$\mathbf{X}\mathbf{X}' = \sum_{l_1, l_2=0}^{\infty} \sum_{k_1, k_2=0}^{\infty} h_{k_1 l_1} h_{k_2 l_2} \mathbf{Z}(k_1, l_1) \mathbf{Z}'(k_2, l_2),$$

where

$$\mathbf{Z}(k, l) = (Z_{i-k,t-l})_{i=1, \dots, p; t=1, \dots, n}, \quad l, k \in \mathbb{Z}.$$

- **Location of squares:**

$$M_{ij} = \sum_{l \in \mathbb{Z}} h_{il} h_{jl}, \quad i, j \in \mathbb{Z}.$$

- For $s \geq 0$,

$$\mathbf{X}_n(s) = (X_{i,t+s})_{i=1,\dots,p; t=1,\dots,n}, \quad n \geq 1.$$

Then $\mathbf{X}_n = \mathbf{X}_n(0)$.

- **Autocovariance matrix** for lag s

$$\mathbf{X}_n(0)\mathbf{X}_n(s)'$$

- Limit theory for **singular values** of such matrices.

- **Autocovariance matrix** for lag s

$$\mathbf{C}_n(s) = \begin{cases} \mathbf{X}_n(0)\mathbf{X}_n(s)', & \text{if } \alpha < 2(1 + \beta), \\ \mathbf{X}_n(0)\mathbf{X}_n(s)' - \mathbb{E}[\mathbf{X}_n(0)\mathbf{X}_n(s)'], & \text{if } \alpha > 2(1 + \beta), \end{cases}$$

- Consider

$$\mathbf{P}_n(s_1, s_2) = \sum_{s=s_1}^{s_2} \mathbf{C}_n(s)\mathbf{C}_n(s)' \quad \text{for fixed } 0 \leq s_1 \leq s_2.$$



$$(\mathbf{M}(s))_{ij} = \sum_{l \in \mathbb{Z}} h_{i,l} h_{j,l+s}, \quad i, j \in \mathbb{Z}.$$

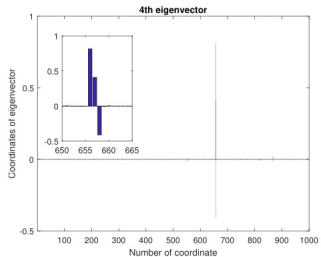
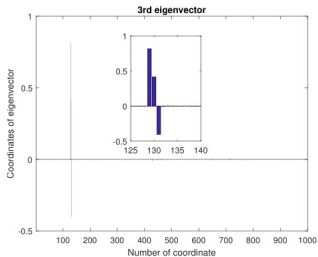
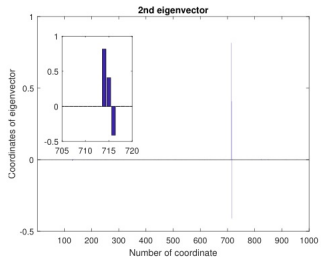
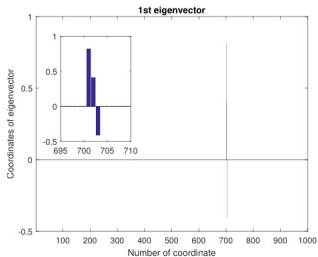
For $0 \leq s_1 \leq s_2 < \infty$, we define the positive semi-definite matrix

$$\mathbf{K}(s_1, s_2) = \sum_{s=s_1}^{s_2} \mathbf{M}(s) \mathbf{M}(s)'$$

- Eigenvector approximation

$$\|\mathbf{y}_i(s_1, s_2) - \mathbf{u}_{b(i)}^{a(i)}(s_1, s_2)\|_{\ell_2} \xrightarrow{\mathbb{P}} 0, \quad n \rightarrow \infty.$$

Autocovariance eigenvectors



Autocovariance eigenvectors

