

Eigenvectors of some large sample covariance matrix ensembles

*Random Matrix Workshop, Télécom ParisTech
Monday, October 11th 2010*

Olivier Ledoit – Sandrine Péché

`oledoit@iew.uzh.ch` - `sandrine.peche@ujf-grenoble.fr`

University of Zürich – Université Grenoble 1

Sample Covariance Matrix

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

- N = number of variables

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

- N = number of variables
- p = sample size

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

- N = number of variables
- p = sample size
- N and p go to infinity together with $N/p \rightarrow c \in (0, +\infty)$

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

- N = number of variables
- p = sample size
- N and p go to infinity together with $N/p \rightarrow c \in (0, +\infty)$
- X_N = real or complex iid random variables mean 0, variance 1, bounded 12^{th} moment

Sample Covariance Matrix

$$S_N = \frac{1}{p} \Sigma_N^{1/2} X_N X_N^* \Sigma_N^{1/2}$$

- N = number of variables
- p = sample size
- N and p go to infinity together with $N/p \rightarrow c \in (0, +\infty)$
- X_N = real or complex iid random variables mean 0, variance 1, bounded 12^{th} moment
- Σ_N = population covariance matrix

Population Covariance Matrix Σ_N

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N
- Eigenvalues: $\tau_1 \leq \dots \leq \tau_N$

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N
- Eigenvalues: $\tau_1 \leq \dots \leq \tau_N$
- Eigenvectors: v_1, \dots, v_N

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N
- Eigenvalues: $\tau_1 \leq \dots \leq \tau_N$
- Eigenvectors: v_1, \dots, v_N
- Empirical Spectral Distribution (e.s.d.):

$$H_N(\tau) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\tau_i, +\infty)}(\tau)$$

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N
- Eigenvalues: $\tau_1 \leq \dots \leq \tau_N$
- Eigenvectors: v_1, \dots, v_N
- Empirical Spectral Distribution (e.s.d.):
$$H_N(\tau) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\tau_i, +\infty)}(\tau)$$
- $H_N(\tau) \rightarrow H(\tau)$ at all points of continuity of H

Population Covariance Matrix Σ_N

- Hermitian positive definite matrix
- Independent of X_N
- Eigenvalues: $\tau_1 \leq \dots \leq \tau_N$
- Eigenvectors: v_1, \dots, v_N
- Empirical Spectral Distribution (e.s.d.):
$$H_N(\tau) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\tau_i, +\infty)}(\tau)$$
- $H_N(\tau) \rightarrow H(\tau)$ at all points of continuity of H
- $\text{Supp}(H)$ bounded away from 0 and $+\infty$

Spectral Decomposition of S_N

Spectral Decomposition of S_N

■ eigenvalues: $\lambda_1 \leq \dots \leq \lambda_N$

Spectral Decomposition of S_N

- eigenvalues: $\lambda_1 \leq \dots \leq \lambda_N$
- eigenvectors: u_1, \dots, u_N

Spectral Decomposition of S_N

- eigenvalues: $\lambda_1 \leq \dots \leq \lambda_N$
- eigenvectors: u_1, \dots, u_N
- e.s.d: $F_N(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$

Spectral Decomposition of S_N

- eigenvalues: $\lambda_1 \leq \dots \leq \lambda_N$
- eigenvectors: u_1, \dots, u_N
- e.s.d: $F_N(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$

Marčenko and Pastur (1967), Silverstein (1995):

$$\exists F \quad \text{s.t.} \quad F_N(\lambda) \xrightarrow{\text{a.s.}} F(\lambda)$$

at all points of continuity of F

Stieltjes Transform

Stieltjes Transform

$$\forall z \in \mathbb{C}^+ \quad m_F(z) = \int_{-\infty}^{+\infty} \frac{1}{\lambda - z} dF(\lambda)$$

Stieltjes Transform

$$\forall z \in \mathbb{C}^+ \quad m_F(z) = \int_{-\infty}^{+\infty} \frac{1}{\lambda - z} dF(\lambda)$$

$$m_{F_N}(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} = \frac{1}{N} \text{Tr} \left[(S_N - zI)^{-1} \right]$$

Stieltjes Transform

$$\forall z \in \mathbb{C}^+ \quad m_F(z) = \int_{-\infty}^{+\infty} \frac{1}{\lambda - z} dF(\lambda)$$

$$m_{F_N}(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} = \frac{1}{N} \text{Tr} \left[(S_N - zI)^{-1} \right]$$

Inversion formula: if F is continuous at a and b :

$$F(b) - F(a) = \lim_{\eta \rightarrow 0^+} \frac{1}{\pi} \int_a^b \text{Im} [m_F(\xi + i\eta)] d\xi$$

MP67/Silverstein (1995) Equation

MP67/Silverstein (1995) Equation

$\forall z \in \mathbb{C}^+$, $m = m_F(z)$ is the unique solution in $\{m \in \mathbb{C} : \frac{c-1}{z} + cm \in \mathbb{C}^+\}$ to

$$m = \int_{-\infty}^{+\infty} \frac{1}{\tau(1-c-czm) - z} dH(\tau)$$

Extension to the Real Line

Extension to the Real Line

Silverstein and Choi (1995):

Extension to the Real Line

Silverstein and Choi (1995):

- $\forall \lambda \in \mathbb{R} - \{0\}$, $\lim_{z \in \mathbb{C}^+ \rightarrow \lambda} m_F(z) \equiv \check{m}_F(\lambda)$
exists

Extension to the Real Line

Silverstein and Choi (1995):

- $\forall \lambda \in \mathbb{R} - \{0\}$, $\lim_{z \in \mathbb{C}^+ \rightarrow \lambda} m_F(z) \equiv \check{m}_F(\lambda)$ exists
- F has continuous derivative $F' = \frac{1}{\pi} \text{Im} [\check{m}_F]$ on $\mathbb{R} - \{0\}$

Intuition for MP67/S95 Equation

Intuition for MP67/S95 Equation

- Sample eigenvalues are a reflection of population eigenvalues

Intuition for MP67/S95 Equation

- Sample eigenvalues are a reflection of population eigenvalues
- They are noisier, more diffuse, like spreading butter

Intuition for MP67/S95 Equation

- Sample eigenvalues are a reflection of population eigenvalues
- They are noisier, more diffuse, like spreading butter
- The higher the c , the more spreading there is

Intuition for MP67/S95 Equation

- Sample eigenvalues are a reflection of population eigenvalues
- They are noisier, more diffuse, like spreading butter
- The higher the c , the more spreading there is
- (Not obvious) Large eigenvalues get more spread out than small ones

Generalization of MP67/S95 Equation

Generalization of MP67/S95 Equation

g : piecewise continuous function

Generalization of MP67/S95 Equation

g : piecewise continuous function

$$m_{F_N}(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times 1$$

Generalization of MP67/S95 Equation

g : piecewise continuous function

$$m_{F_N}(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times 1$$

$$\Theta_N^g(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

Generalization of MP67/S95 Equation

g : piecewise continuous function

$$m_{F_N}(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times 1$$

$$\Theta_N^g(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

$$g(\tau) \equiv 1 \quad \iff \quad \Theta_N^g = m_{F_N}$$

Generalization of MP67/S95 Equation

Generalization of MP67/S95 Equation

$$m_{F_N}(z) = \frac{1}{N} \text{Tr} \left[(S_N - zI)^{-1} \right]$$

Generalization of MP67/S95 Equation

$$m_{F_N}(z) = \frac{1}{N} \text{Tr} \left[(S_N - zI)^{-1} \right]$$
$$\Theta_N^g(z) = \frac{1}{N} \text{Tr} \left[(S_N - zI)^{-1} g(\Sigma_N) \right]$$

Generalization of MP67/S95 Equation

Generalization of MP67/S95 Equation

$$\exists \Theta^g : \quad \forall z \in \mathbb{C}^+ \quad \Theta_N^g(z) \xrightarrow{\text{a.s.}} \Theta^g(z)$$

Generalization of MP67/S95 Equation

$$\exists \Theta^g : \quad \forall z \in \mathbb{C}^+ \quad \Theta_N^g(z) \xrightarrow{\text{a.s.}} \Theta^g(z)$$

$$\Theta^g(z) = \int_{-\infty}^{+\infty} g(\tau) \times \frac{1}{\tau [1 - c - czm_F(z)] - z} dH(\tau)$$

Generalization of MP67/S95 Equation

$$\exists \Theta^g : \quad \forall z \in \mathbb{C}^+ \quad \Theta_{N}^g(z) \xrightarrow{\text{a.s.}} \Theta^g(z)$$

$$\Theta^g(z) = \int_{-\infty}^{+\infty} g(\tau) \times \frac{1}{\tau [1 - c - czm_F(z)] - z} dH(\tau)$$

Same integration kernel!

Extension to the Real Line

Extension to the Real Line

$$\Theta_N^g(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

Extension to the Real Line

$$\Theta_N^g(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\lambda_i, +\infty)}(\lambda) \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

Extension to the Real Line

$$\Theta_N^g(z) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\lambda_i, +\infty)}(\lambda) \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j)$$

$$\Omega_N^g(\lambda) \xrightarrow{\text{a.s.}} \Omega^g(\lambda) = \lim_{\eta \rightarrow 0^+} \frac{1}{\pi} \int_{-\infty}^{\lambda} \text{Im} [\Theta^g(l + i\eta)] dl$$

wherever Ω^g is continuous

Sample Eigenvectors

Sample Eigenvectors

Fix τ and take $g = \mathbf{1}_{(-\infty, \tau)}$

Sample Eigenvectors

Fix τ and take $g = \mathbf{1}_{(-\infty, \tau)}$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N |u_i^* v_j|^2 \mathbf{1}_{[\lambda_i, +\infty)}(\lambda) \times \mathbf{1}_{[\tau_j, +\infty)}(\tau)$$

$$\rightarrow \int_{-\infty}^{\lambda} \int_{-\infty}^{\tau} \frac{c t}{|t [1 - c - c m_F(l)] - l|^2} dH(t) dF(l)$$

Sample Eigenvectors

Fix τ and take $g = \mathbf{1}_{(-\infty, \tau)}$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N |u_i^* v_j|^2 \mathbf{1}_{[\lambda_i, +\infty)}(\lambda) \times \mathbf{1}_{[\tau_j, +\infty)}(\tau)$$

$$\rightarrow \int_{-\infty}^{\lambda} \int_{-\infty}^{\tau} \frac{c t}{|t [1 - c - c l m_F(l)] - l|^2} dH(t) dF(l)$$

$$N |u_i^* v_j|^2 \approx \frac{c \lambda_i \tau_j}{|\tau_j [1 - c - c \lambda_i \check{m}_F(\lambda_i)] - \lambda_i|^2}$$

Eigenvalues with Multiplicity

Eigenvalues with Multiplicity

Σ_N has K distinct eigenvalues t_1, \dots, t_K
with multiplicities n_1, \dots, n_K

Eigenvalues with Multiplicity

Σ_N has K distinct eigenvalues t_1, \dots, t_K
with multiplicities n_1, \dots, n_K

$P_k =$ projection onto k^{th} eigenspace

Eigenvalues with Multiplicity

Σ_N has K distinct eigenvalues t_1, \dots, t_K
with multiplicities n_1, \dots, n_K

P_k = projection onto k^{th} eigenspace

$$|P_k u_i|^2 \approx \frac{n_k c \lambda_i t_k}{N |t_k [1 - c - c \lambda_i \check{m}_F(\lambda_i)] - \lambda_i|^2}$$

Estimating the Covariance Matrix (1)

Estimating the Covariance Matrix (1)

Frobenius norm: $\|A\| = \sqrt{\text{Tr}(AA^*)}$

Estimating the Covariance Matrix (1)

Frobenius norm: $\|A\| = \sqrt{\text{Tr}(AA^*)}$

U_N : matrix of eigenvectors of S_N

Estimating the Covariance Matrix (1)

Frobenius norm: $\|A\| = \sqrt{\text{Tr}(AA^*)}$

U_N : matrix of eigenvectors of S_N

Find matrix closest to Σ_N among those that have eigenvectors U_N

Estimating the Covariance Matrix (1)

Frobenius norm: $\|A\| = \sqrt{\text{Tr}(AA^*)}$

U_N : matrix of eigenvectors of S_N

Find matrix closest to Σ_N among those that have eigenvectors U_N

$$\min_{D_N \text{ diagonal}} \|U_N D_N U_N^* - \Sigma_N\|$$

Estimating the Covariance Matrix (1)

Frobenius norm: $\|A\| = \sqrt{\text{Tr}(AA^*)}$

U_N : matrix of eigenvectors of S_N

Find matrix closest to Σ_N among those that have eigenvectors U_N

$$\min_{D_N \text{ diagonal}} \|U_N D_N U_N^* - \Sigma_N\|$$

Solution:

$$\tilde{D}_N = \text{Diag}(\tilde{d}_1, \dots, \tilde{d}_N) \quad \text{where} \quad \tilde{d}_i = u_i^* \Sigma_N u_i$$

Estimating the Covariance Matrix (2)

Estimating the Covariance Matrix (2)

Take $g(\tau) = \tau$

Estimating the Covariance Matrix (2)

Take $g(\tau) = \tau$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N u_i^* \Sigma_N u_i \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$$

$$\rightarrow \int_{-\infty}^{\lambda} \frac{l}{|1 - c - clm_F(l)|^2} dF(l)$$

Estimating the Covariance Matrix (2)

Take $g(\tau) = \tau$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N u_i^* \Sigma_N u_i \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$$

$$\rightarrow \int_{-\infty}^{\lambda} \frac{l}{|1 - c - clm_F(l)|^2} dF(l)$$

$$u_i^* \Sigma_N u_i \approx \frac{\lambda_i}{|1 - c - c\lambda_i \check{m}_F(\lambda_i)|^2}$$

Oracle Estimator

Oracle Estimator

Keep same eigenvectors as those of S_n ,

Oracle Estimator

Keep same eigenvectors as those of S_n , divide i^{th} sample eigenvalue by $|1 - c - c\lambda_i\check{m}_F(\lambda_i)|^2$

Oracle Estimator

Keep same eigenvectors as those of S_n , divide i^{th} sample eigenvalue by $|1 - c - c\lambda_i\check{m}_F(\lambda_i)|^2$
→ *oracle estimator* \tilde{S}_N

Oracle Estimator

Keep same eigenvectors as those of S_n , divide i^{th} sample eigenvalue by $|1 - c - c\lambda_i\check{m}_F(\lambda_i)|^2$
→ *oracle estimator* \tilde{S}_N

Percentage Relative Improvement in Average Loss:

$$PRIAL = 100 \times \left[1 - \frac{\mathbb{E} \left\| \tilde{S}_N - U_N \tilde{D}_N U_N^* \right\|^2}{\mathbb{E} \left\| S_N - U_N \tilde{D}_N U_N^* \right\|^2} \right]$$

Monte-Carlo Simulations

Monte-Carlo Simulations

10,000 simulations

Monte-Carlo Simulations

10,000 simulations
 $c=1/2$

Monte-Carlo Simulations

10,000 simulations

$c=1/2$

Population eigenvalues:

- 20% equal to 1
- 40% equal to 3
- 40% equal to 10

Monte-Carlo Simulations

10,000 simulations

$c=1/2$

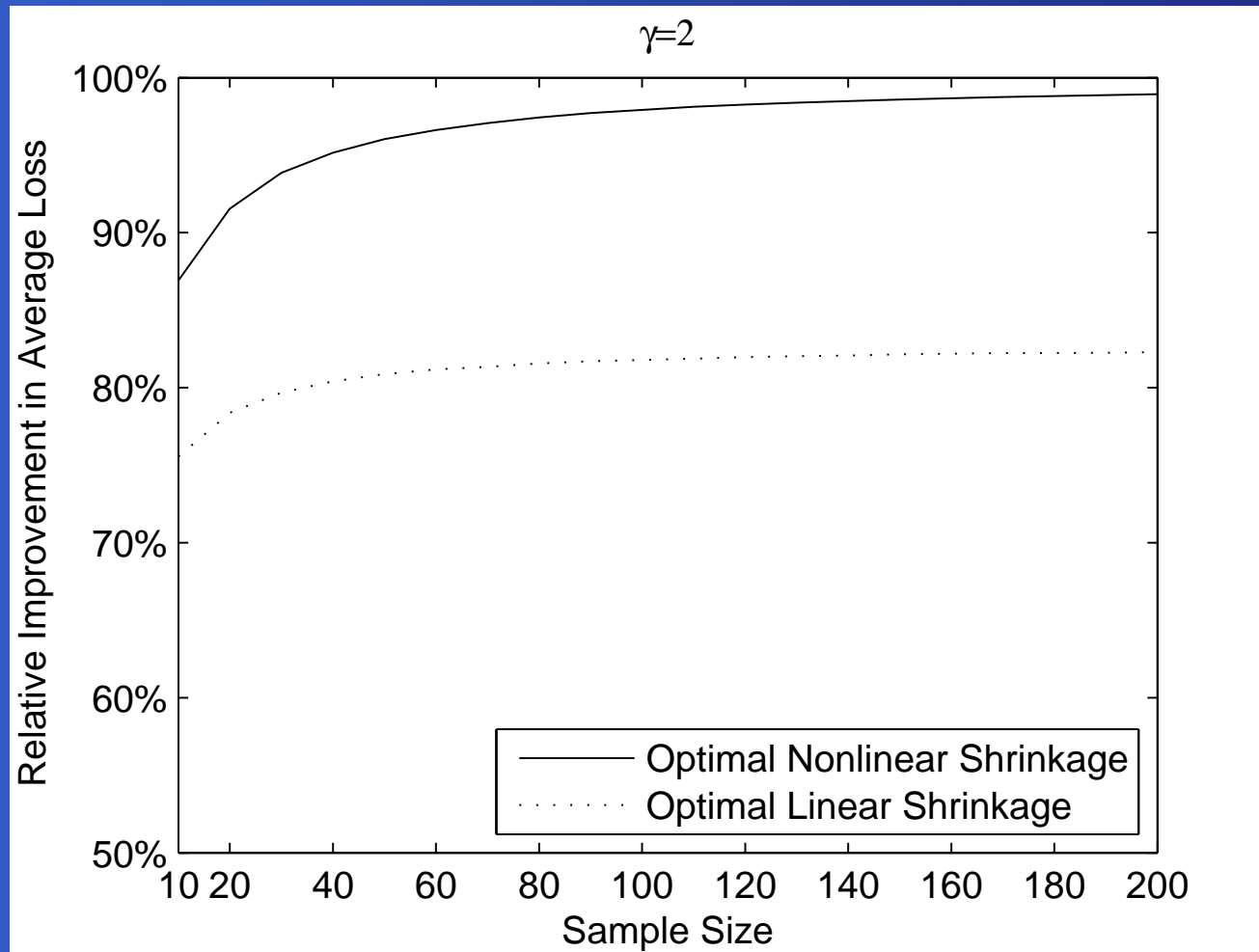
Population eigenvalues:

- 20% equal to 1
- 40% equal to 3
- 40% equal to 10

Compare with Ledoit-Wolf (2004) *linear* shrinkage estimator

Simulation Results

Simulation Results



Inverse of the Covariance Matrix (1)

Inverse of the Covariance Matrix (1)

Find matrix closest to Σ_N^{-1} among those that have eigenvectors U_N

Inverse of the Covariance Matrix (1)

Find matrix closest to Σ_N^{-1} among those that have eigenvectors U_N

$$\min_{\Delta_N \text{ diagonal}} \|U_N \Delta_N U_N^* - \Sigma_N^{-1}\|$$

Inverse of the Covariance Matrix (1)

Find matrix closest to Σ_N^{-1} among those that have eigenvectors U_N

$$\min_{\Delta_N \text{ diagonal}} \|U_N \Delta_N U_N^* - \Sigma_N^{-1}\|$$

Solution:

$$\tilde{\Delta}_N = \text{Diag}(\tilde{\delta}_1, \dots, \tilde{\delta}_N) \quad \text{where} \quad \tilde{\delta}_i = u_i^* \Sigma_N^{-1} u_i$$

Inverse of the Covariance Matrix (1)

Find matrix closest to Σ_N^{-1} among those that have eigenvectors U_N

$$\min_{\Delta_N \text{ diagonal}} \|U_N \Delta_N U_N^* - \Sigma_N^{-1}\|$$

Solution:

$$\tilde{\Delta}_N = \text{Diag}(\tilde{\delta}_1, \dots, \tilde{\delta}_N) \quad \text{where} \quad \tilde{\delta}_i = u_i^* \Sigma_N^{-1} u_i$$

$$u_i^* \Sigma_N^{-1} u_i \geq (u_i^* \Sigma_N u_i)^{-1}$$

Inverse of the Covariance Matrix (2)

Inverse of the Covariance Matrix (2)

Take $g(\tau) = \frac{1}{\tau}$

Inverse of the Covariance Matrix (2)

Take $g(\tau) = \frac{1}{\tau}$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N u_i^* \Sigma_N^{-1} u_i \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$$

$$\rightarrow \int_{-\infty}^{\lambda} \frac{1 - c - 2cl \operatorname{Re}[m_F(l)]}{l} dF(l)$$

Inverse of the Covariance Matrix (2)

Take $g(\tau) = \frac{1}{\tau}$

$$\Omega_N^g(\lambda) = \frac{1}{N} \sum_{i=1}^N u_i^* \Sigma_N^{-1} u_i \mathbf{1}_{[\lambda_i, +\infty)}(\lambda)$$

$$\rightarrow \int_{-\infty}^{\lambda} \frac{1 - c - 2cl \operatorname{Re}[m_F(l)]}{l} dF(l)$$

$$u_i^* \Sigma_N^{-1} u_i \approx \frac{1 - c - 2c\lambda_i \operatorname{Re}[\check{m}_F(\lambda_i)]}{\lambda_i}$$

Conclusion

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation
- Gives location of sample eigenvectors relative to:

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation
- Gives location of sample eigenvectors relative to:
 - Population eigenvectors

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation
- Gives location of sample eigenvectors relative to:
 - Population eigenvectors
 - Population covariance matrix as a whole

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation
- Gives location of sample eigenvectors relative to:
 - Population eigenvectors
 - Population covariance matrix as a whole
 - Inverse of population covariance matrix

Conclusion

- Generalization of the Marčenko-Pastur (1967)/Silverstein (1995) Equation
- Gives location of sample eigenvectors relative to:
 - Population eigenvectors
 - Population covariance matrix as a whole
 - Inverse of population covariance matrix
- We do for sample eigenvectors what MP67/S95 did for sample eigenvalues

Directions for Future Research

Directions for Future Research

1. Construct *bona fide* nonlinear shrinkage estimator of the covariance matrix

Directions for Future Research

1. Construct *bona fide* nonlinear shrinkage estimator of the covariance matrix
2. Construct *bona fide* nonlinear shrinkage estimator of the inverse of the covariance matrix

Directions for Future Research

1. Construct *bona fide* nonlinear shrinkage estimator of the covariance matrix
2. Construct *bona fide* nonlinear shrinkage estimator of the inverse of the covariance matrix
3. Show that $N|u_i^* v_j|^2$ is even closer to

$$\frac{c\lambda_i\tau_j}{|\tau_j [1 - c - c\lambda_i\check{m}_F(\lambda_i)] - \lambda_i|^2}$$

than we have shown in this paper