Dario Villamaina

École normale supérieure & Institut Philippe Meyer

Random matrices and correlation data analysis



Motivation: Analysis of sequences





Motivation: Analysis of sequences



Motivation: Analysis of sequences



Sites in contact are associated with high localized eigenvectors of the correlation matrix

Population Covariance

Infinite statistics

C (<u>One</u> matrix)





Population Covariance *vs* **Empirical Covariance**

Infinite statistics





Finite statistics







by finite sampling effects!

Random Matrix Model





Efficiency parameter $r = \frac{N}{M}$

Perfect sampling limit

$$\hat{C} \to C$$

$$r \to 0$$

Toy model : Finite rank correlation matrix

$$C \propto \begin{pmatrix} \gamma & 0 \\ 1 & 0 \\ 0 & \ddots \\ 0 & \ddots \\ 1 \end{pmatrix} \qquad \qquad \boldsymbol{\xi}_1 \propto \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

continuous limit $\rho_{c}(\lambda) = (1 - \alpha)\delta(\lambda - 1) + \alpha\delta(\lambda - \gamma)$

Toy model : Finite rank correlation matrix

continuous limit $\rho_{c}(\lambda) = (1 - \alpha)\delta(\lambda - 1) + \alpha\delta(\lambda - \gamma)$



The transition in the eigenvalue spectrum



Noise undressing of empirical correlation matrices



Noise undressing of empirical correlation matrices

L Laloux, P Cizeau, JP Bouchaud, M Potters - Physical Review Letters, (1999) (Measured) (True) $(\hat{\lambda}, \boldsymbol{\hat{\xi}}_i)$ $(\lambda, \boldsymbol{\xi}_i)$ 4 Real data p(N) Null model 2 (uncorrelated variables) 0 3 1 λ (Only noise ?) Signal **Bulk eigenvalues**

What about eigenvectors ?

Estimation of true eigenvectors from the empirical ones

Scalar product between <u>the largest</u> true and empirical eigenvectors



Estimation of true eigenvectors from the empirical ones



Estimation of true eigenvectors from the empirical ones

Scalar product between the largest true and bulk eigenvectors





Summary on empirical eigenvectors



 $\left\langle |\boldsymbol{\xi}_1 \cdot \hat{\boldsymbol{\xi}}_1|^2 \right\rangle = \mathcal{O}(1) \qquad \left\langle |\boldsymbol{\xi}_1 \cdot \hat{\boldsymbol{\xi}}_{m>1}|^2 \right\rangle = \mathcal{O}(N^{-1})$

Two questions:

Is it possible to <u>calculate analytically</u> the scalar products ? Is <u>there information</u> stored in bulk empirical eigenvectors ?

Part 1: How to calculate the scalar products?

$$g(\lambda, \hat{\lambda}) = \lim_{N \to \infty} \frac{1}{N} \operatorname{Tr} \left[\left(\lambda - C \right)^{-1} \cdot \left\langle \left(\hat{\lambda} - \hat{C} \right)^{-1} \right\rangle \right]$$

Projection on the eigenvectors

$$\left(\lambda - C\right)_{ij}^{-1} = \sum_{m} (\lambda - \lambda_m)^{-1} \xi_{m,i} \xi_{m,j}$$
$$\left(\hat{\lambda} - \hat{C}\right)_{ij}^{-1} = \sum_{m} (\hat{\lambda} - \hat{\lambda}_m)^{-1} \hat{\xi}_{m,i} \hat{\xi}_{m,j}$$

Continuous limit of the spectrum

$$\hat{\rho}(\hat{\mu}) = \frac{1}{N} \sum_{m} \left\langle \delta(\hat{\mu} - \hat{\lambda}_m) \right\rangle$$

Part 1: How to calculate the scalar products?

$$g(\lambda, \hat{\lambda}) = \lim_{N \to \infty} \frac{1}{N} \operatorname{Tr} \left[\left(\lambda - C \right)^{-1} \cdot \left\langle \left(\hat{\lambda} - \hat{C} \right)^{-1} \right\rangle \right]$$

Projection on the eigenvectors

$$\left(\lambda - C\right)_{ij}^{-1} = \sum_{m} (\lambda - \lambda_m)^{-1} \xi_{m,i} \xi_{m,j}$$
$$\left(\hat{\lambda} - \hat{C}\right)_{ij}^{-1} = \sum_{m} (\hat{\lambda} - \hat{\lambda}_m)^{-1} \hat{\xi}_{m,i} \hat{\xi}_{m,j}$$

Continuous limit of the spectrum

$$\hat{\rho}(\hat{\mu}) = \frac{1}{N} \sum_{m} \left\langle \delta(\hat{\mu} - \hat{\lambda}_m) \right\rangle$$

$$W^{2}(\lambda,\hat{\lambda}) = \frac{g(\lambda + i\epsilon, \hat{\lambda} + i\epsilon) - g(\lambda + i\epsilon, \hat{\lambda} - i\epsilon)}{2\pi \rho(\lambda) \,\hat{\rho}(\hat{\lambda})}$$

An extension of the Edwards-Jones technique SF Edwards and RC Jones Journal of Physics A: Mathematical and General, 9(10), 1595 (1976)

$$\left\langle \operatorname{Tr}((\lambda - C)^{-1}(\hat{\lambda} - \hat{C})^{-1}) \right\rangle = \left\langle \int \mathcal{D}Y \mathcal{D}Z (Y \cdot Z)^2 \exp(-\frac{1}{2}(Y(\lambda - C)Y^{\dagger} + Z(\hat{\lambda} - \hat{C})Z^{\dagger})) \right\rangle$$

$$\langle \dots \rangle \propto \int \prod_{i,t} dx_{i,t} \exp\left(-\frac{1}{2}x_{i,t}C_{ij}^{-1}x_{j,t}\right) \qquad \qquad \hat{C}_{ij} = \frac{1}{M}\sum_{t} x_{i,t}x_{j,t}$$

An extension of the Edwards-Jones technique SF Edwards and RC Jones Journal of Physics A: Mathematical and General, 9(10), 1595 (1976)

$$g(\lambda,\hat{\lambda}) = \frac{2}{N} \int d\phi \, d\hat{\phi} \frac{\partial}{\partial \eta} \Big|_{\eta=0} e^{-NS_{\eta}(\phi,\hat{\phi})}$$

$$S_{\eta}(\phi,\hat{\phi}) = \frac{1}{2} \left(\int d\mu \rho(\mu) \ln((\lambda-\mu)(\hat{\lambda}-\mu\,\hat{\phi}) - \eta) + \frac{1}{r} \ln(1-r\phi) + \phi\,\hat{\phi} \right)$$

Equations for the densities

O. <u>Ledoit</u> and S. <u>Péché</u>, Probability Theory and Related Fields 151.1-2 (2011). JW <u>Silverstein</u> Journal of Multivariate Analysis, 55(2), (1995)

$$g(\lambda, \hat{\lambda}) = \int d\mu \frac{\rho(\mu)}{(\lambda - \mu)(\hat{\lambda} - \mu \, \hat{\phi}^*)}$$

Saddle point equation
$$\frac{\hat{\phi}^* - 1}{r\hat{\phi}^*} = \int d\mu \frac{\rho(\mu)}{(\hat{\lambda} - \hat{\mu}\hat{\phi}^*)}$$



$$\hat{\rho}(\hat{\lambda}) = \frac{1}{\pi} \operatorname{Im} \int d\mu \frac{\rho(\mu)}{(\hat{\lambda} - \mu \, \hat{\phi}^*)}$$

Equations for the densities

O. <u>Ledoit</u> and S. <u>Péché</u>, Probability Theory and Related Fields 151.1-2 (2011). JW <u>Silverstein</u> Journal of Multivariate Analysis, 55(2), (1995)

$$g(\lambda, \hat{\lambda}) = \int d\mu \frac{\rho(\mu)}{(\lambda - \mu)(\hat{\lambda} - \mu \, \hat{\phi}^*)}$$

Saddle point equation
$$\frac{\hat{\phi}^* - 1}{r\hat{\phi}^*} = \int d\mu \frac{\rho(\mu)}{(\hat{\lambda} - \hat{\mu}\hat{\phi}^*)}$$



$$\hat{\rho}(\hat{\lambda}) = \frac{1}{\pi} \operatorname{Im} \int d\mu \frac{\rho(\mu)}{(\hat{\lambda} - \mu \, \hat{\phi}^*)}$$





Estimation of the principal component

Inverse problem : How to estimate ξ_1 ?

$$\boldsymbol{\xi}_1^{(est)} = w_1 \hat{\boldsymbol{\xi}_1} + \sum_m^N w_m \hat{\boldsymbol{\xi}}_m$$

Estimation of the principal component

Inverse problem : How to estimate ξ_1 ?

$$\boldsymbol{\xi}_{1}^{(est)} = w_{1}\hat{\boldsymbol{\xi}_{1}} + \sum_{m}^{N} w_{m}\hat{\boldsymbol{\xi}}_{m}$$

The ws are Gaussian variables (fairly verified for large N)

$$P\left(\boldsymbol{\xi}_{1}^{(est)}|\{\hat{\xi}\}\right)$$

"standard approach"

$$\left\langle \boldsymbol{\xi}_{1}^{(est)} \right\rangle = \left\langle w_{1} \right\rangle \boldsymbol{\hat{\xi}}_{1}$$

Main limitation: reconstruction impossible for weak signals

$$\gamma < \gamma_c \qquad \longrightarrow \quad \langle w_1 \rangle = 0$$

One question left:

Is it possible to <u>calculate analytically</u> the scalar products ?

YES !!

Is <u>there information</u> stored in bulk empirical eigenvectors ?

The mutual information

$$I(x;y) = \left\langle \log\left(\frac{P(x,y)}{P(x)P(y)}\right) \right\rangle$$

It measures how much knowing X reduces uncertainty about Y

The mutual information

It measures how much knowing X reduces uncertainty about Y



The mutual information

It measures how much knowing X reduces uncertainty about Y



Three questions:

Is it possible to calculate analytically the scalar products ?

YES !!

Is there information stored in bulk empirical eigenvectors ? YES !!

Is it possible to exploit this information to improve the <u>noise undressing</u> procedure ?

discrete approximation : $w_m \to \sqrt{\langle w_m^2 \rangle} \sigma_m \equiv \overline{w}_m \sigma_m$ $\boldsymbol{\xi}_1^{(est)}(\boldsymbol{\sigma}) = \overline{w}_1 \hat{\boldsymbol{\xi}}_1 + \sum_m^N \overline{w}_m \hat{\boldsymbol{\xi}}_m \sigma_m$ (±1)

discrete approximation :
$$w_m \to \sqrt{\langle w_m^2 \rangle} \sigma_m \equiv \overline{w}_m \sigma_m$$

$$\boldsymbol{\xi}_1^{(est)}(\boldsymbol{\sigma}) = \overline{w}_1 \hat{\boldsymbol{\xi}}_1 + \sum_m^N \overline{w}_m \hat{\boldsymbol{\xi}}_m \sigma_m$$
(±1)

The estimation problem is mapped onto an Ising model $E(\boldsymbol{\sigma}) = -\sum_{i} \left(\xi_{1,i}^{(est)}(\boldsymbol{\sigma})\right)^4$

<u>Estimate</u>: finding $\overline{\sigma}$ for which $E(\overline{\sigma})$ attains its minimal value







Conclusions and Perspectives

R. <u>Monasson</u> and D. <u>Villamaina</u> (arXiv:1503.00287)

General formalism for the calculation of the scalar products between the empirical eigenvectors and the true ones of a correlation matrix

 $\operatorname{Tr}\left[\left(\lambda-C\right)^{-1}\cdot\langle\left(\hat{\lambda}-\hat{C}\right)^{-1}\rangle\right]$

Large deviations ?

Importance of minor components for the reconstruction of

eigenvectors



Analogy between spin models and estimations of localized eigenvectors



<u>Smarter algorithms ?</u>

Perspectives: analogy with the Hopfield model

Hopfield $m^{\mu} = \sum_{j} \xi_{j}^{\mu} \sigma_{j}$ $H_{opf}(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{\mu} (m^{\mu})^{2}$

Eigenvector estimation

$$m^i = \sum_m \overline{w}_m \hat{\xi}_{m,i} \sigma_m$$

$$H(\boldsymbol{\sigma}) = -\sum_{i} (m^{i})^{4}$$

Perspectives: analogy with the Hopfield model



Perspectives: analogy with the Hopfield model



References related to this talk

The transition in eigenvalue spectrum

SF <u>Edwards</u> and RC <u>Jones</u>, Journal of Physics A: Mathematical and General 9, 1595 (1976) J <u>Baik</u>, G <u>Ben Arous</u> and <u>S Péché</u>, Annals of Probability, 1643 (2005)

Data analysis in finance and bioinformatics

L <u>Laloux</u>, P <u>Cizeau</u>, JP <u>Bouchaud</u>, M <u>Potters</u>, Physical Review Letters 83, 1467 (1999) S. <u>Cocco</u>, R. <u>Monasson</u>, M. <u>Weigt</u>, PLoS computational biology 9, e1003176 (2013)

Eigenvalues and eigenvectors in empirical correlation matrices

O. Ledoit and S. Péché, Probability Theory and Related Fields 151, 233 (2011)

JW Silverstein, Journal of Multivariate Analysis 55, 331 (1995)

R. <u>Monasson</u> and D. <u>Villamaina</u> Estimating the principal component of correlation matrices from all their eigenmodes arXiv:1503.00287