

Th. 4.1. (Marchenko-Pastur, 1967)

Let $X_1, \dots, X_n \in \mathbb{R}^p$, i.i.d. n.vectors with $E(X_i) = 0$

and $\text{Cov}(X_i) = \sigma^2 I_p$. $X = (X_1, \dots, X_n) \in \mathbb{R}^{p \times n}$,

$$S_n = \frac{1}{n} XX^T = \frac{1}{n} \sum_{i=1}^n X_i X_i^T \in \mathbb{R}^{p \times p},$$

$\lambda_1, \dots, \lambda_p$ eigenvalues of S_n .

Let $p, n \rightarrow \infty$ such that $\frac{p}{n} \rightarrow \gamma \in (0, 1] (\gamma \neq 0)$.

Then the sample distribution of $\lambda_1, \dots, \lambda_p$ (the histogram) converges a.s. to the density

$$f_\gamma(u) = \frac{1}{2\pi\sigma^2 u \gamma} \sqrt{(b-u)(u-a)}, \quad a \leq u \leq b$$

$$\text{with } a(\gamma) = \sigma^2 (1 - \sqrt{\gamma})^2 \quad b(\gamma) = \sigma^2 (1 + \sqrt{\gamma})^2.$$

Remark: If $\gamma > 1$ there will be a mass point at zero.

Conclusion: even in the i.i.d. uncorrelated case there is a wide spectrum of eigenvalues.

What happens, if there is a low-dimensional structure in the data?

4.1.5. Spike models

Model assumptions : $X_1, \dots, X_n \in \mathbb{R}^p$, i.i.d.

$$\text{Cov}(X_i) = \Sigma = I_p + \beta v v^\top$$

for some $v \in \mathbb{R}^p$, $\|v\|=1$, $\beta \geq 0$

Interpretation:

$$X_i = U_i + \sqrt{\beta} V_i v$$

$U_i \sim N(0, I_p)$ noise

$V_i \sim N(0, 1)$ signal, independent from U_i ,
multiplied by a fixed vector $\sqrt{\beta} v$.

$$\text{Then } \text{Cov}(X_i) = \text{Cov}(U_i) + \beta \text{Var}(V_i) v v^\top = I_p + \beta v v^\top.$$

Example: $p=500$, $n=1000$, $v=e_1$, $\beta=1.5$

$$\Sigma = I_p + 1.5 e_1 e_1^\top$$

$$\lambda_{\max}(\Sigma) = 1 + 1.5 = 2.5$$

all other eigenvalues 1

Question: Is there a threshold for β above which we will one eigenvector popping out.

Th 4.2. (BBP transition, Baik, Ben Arous, Peche [2005])

Assume $X_1, \dots, X_n \in \mathbb{R}^p$ r.vectors, with $E(X_i) = 0$

$\text{Cov}(X_i) = I_p + \beta v v^\top$, $\beta \geq 0$, $v \in \mathbb{R}^p$, $\|v\|=1$.

$$S_n = \frac{1}{n} X X^\top. \quad n, p \rightarrow \infty, \quad \frac{p}{n} \rightarrow \gamma.$$

If $\beta \leq \sqrt{\gamma}$ then $\lambda_{\max}(S_n) \rightarrow (1 + \sqrt{\gamma})^2$.

$$\text{and } (v_{\max}^\top v)^2 \rightarrow 0$$

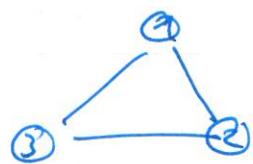
If $\beta > \sqrt{\gamma}$ then $\lambda_{\max}(S_n) \rightarrow (1 + \beta)(1 + \frac{\gamma}{\beta}) > (1 + \sqrt{\gamma})^2$.

$$\text{and } (v_{\max}^\top v)^2 \rightarrow \frac{1 - \gamma/\beta^2}{1 - \gamma/\beta} \downarrow$$

Proof. Bandera, references therein.

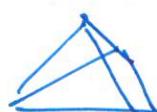
4.2. Multidimensional Scaling

Given pairwise distances d_{ij} between points.



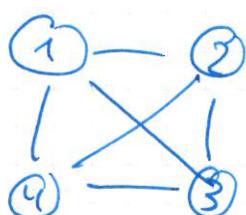
$$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$

Euclidean
embedding
in dim. 2 ?
Yes!



$$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

Euclidean
embedding
in dim. ?? no!
in dim 3? yes!



$$\begin{pmatrix} 0 & 1 & \sqrt{2} & 1 \\ 1 & 0 & 1 & \sqrt{2} \\ \sqrt{2} & 1 & 0 & 1 \\ 1 & \sqrt{2} & 1 & 0 \end{pmatrix}$$

Euclidean
embedding
in dim. 2? yes!

Given n objects O_1, \dots, O_n and pairwise dissimilarities δ_{ij} between object i and j .

Assume that $\delta_{ij} = \delta_{ji} \geq 0$ and $\delta_{ii} = 0$, $\forall i, j = 1, \dots, n$.

Define $\Delta = (\delta_{ij})_{1 \leq i, j \leq n}$ as the dissimilarity matrix.

and

$$\mathcal{M}_n = \left\{ \Delta = (\delta_{ij})_{1 \leq i, j \leq n} \mid \delta_{ij} = \delta_{ji} \geq 0, \delta_{ii} = 0 \quad \forall i, j \right\}$$

the set of dissimilarity matrices.

Objective: Find n points x_1, \dots, x_n in a Euclidean space, typically \mathbb{R}^k , such that the distances $\|x_i - x_j\|$ fit the dissimilarities δ_{ij} best.

Notation: $X = [x_1, \dots, x_n]^T \in \mathbb{R}^{n \times k}$

$$d_{ij}(X) = \|x_i - x_j\| \text{ distances}$$

$$D(X) = (d_{ij}(X))_{1 \leq i, j \leq n} \in \mathbb{R}^{n \times n}$$

$$\Delta^{(q)} = (\delta_{ij}^{(q)})_{i, j} \text{ and } D^{(q)}(X) = (d_{ij}^{(q)}(X))_{i, j}$$

Optimization problem: given $q \geq 1$

$$\min_{X \in \mathbb{R}^{n \times k}} \| \Delta^{(q)} - D(X) \| \quad (*)$$

typically $q=2$

Consider the case that $(*) = 0$,

4.2.1. Characterizing Euclidean Distance Matrices

$\Delta = (\delta_{ij}) \in \mathbb{R}^n$ is called Euclidean distance

matrix or it has a Euclidean embedding

in \mathbb{R}^k , if there are $x_1, \dots, x_n \in \mathbb{R}^k$

such that $\delta_{ij}^2 = \|x_i - x_j\|^2 \forall i, j$.

where $\|\cdot\|$ denotes the Euclidean norm $\|y\| = \sqrt{\sum_{i=1}^k y_i^2}$.

That means opt.-problem (*) has a solution with value 0.

Projection $E_n = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$ plays an important role.

$$E_n = \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{pmatrix}$$