

What if our data were intervals?

Lina Oliveira & M. Rosário Oliveira

CAMGSD, CEMAT, and Department of Mathematics, Instituto Superior Técnico

MMAC Day, March 14, 2024

Joint with R. Girão Serrão and M.Rosário Oliveira

“Theoretical derivation of interval principal component analysis”,

Information Sciences, 2023

Conventional data/**Symbolic Data**:

- Each **object** is characterised by several **variables**
- Data is organized as an $n \times p$ matrix: rows correspond to **objects**, columns to **variables**

$$\begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}$$

$$\begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}$$

For example,

x_{11} = real number

x_{11} = interval $[a, b]$ of real numbers

Conventional data/Symbolic Data:

Data matrix

$$\begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}$$

random vector $\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \rightarrow \mathbb{R}^p$

Data matrix

$$\begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}$$

random vector $\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \rightarrow \mathbb{I}\mathbb{R}^p$

Random vector

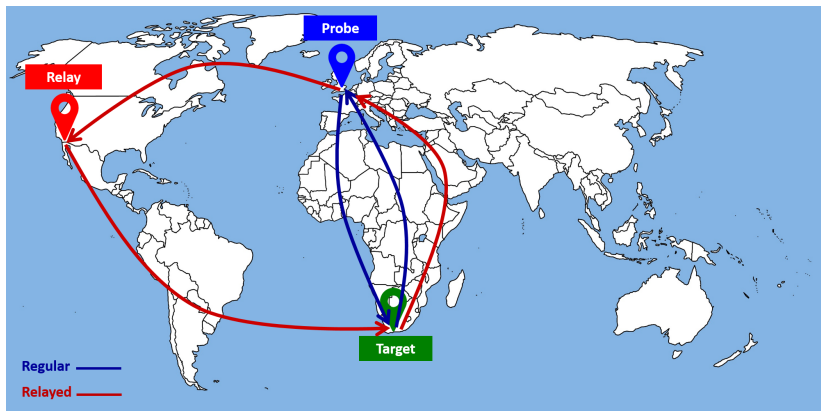
$[x_{11} \cdots x_{1p}]^T \in \mathbb{R}^p$ is a realisation of \mathbf{X}

Interval-valued random vector

$[x_{11} \cdots x_{1p}]^T \in \mathbb{I}\mathbb{R}^p$ is a realisation of \mathbf{X}

WHY INTERVALS?

Detecting Internet traffic redirection attacks



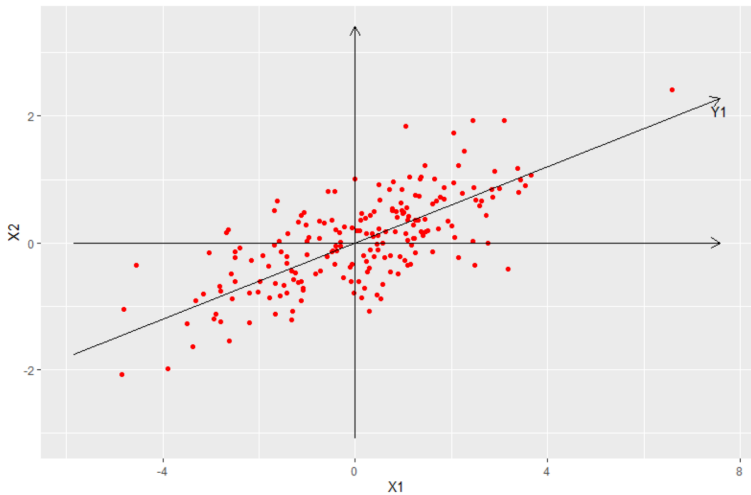
in A. Subtil, M.R. Oliveira, R. Valadas, A. Pacheco, and P. Salvador “Internet-Scale Traffic Redirection Attacks Using Latent Class Models”, *Intelligent Systems and Computing*, 2020

Conventional data:

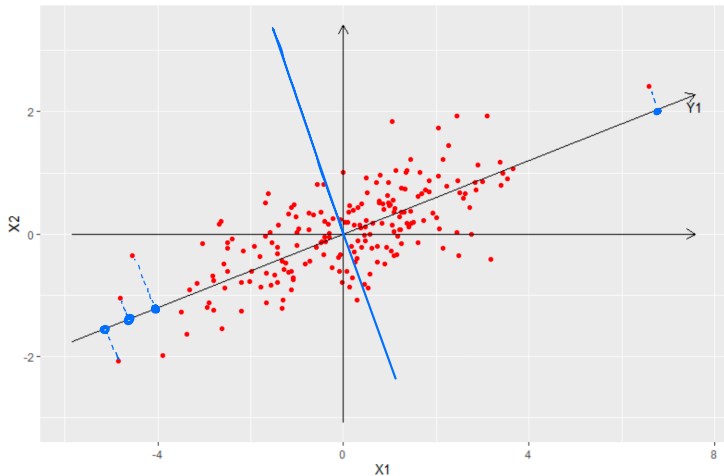
- Each **object** is characterised by several **variables**
- Data is organized as an $n \times p$ matrix: rows correspond to **objects**, columns to **variables**

$$\text{Data matrix } \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \quad \text{random vector } \mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \rightarrow \mathbb{R}^p$$

Example: $[x_{11} \ \cdots \ x_{1p}]^T \in \mathbb{R}^p$ is a realisation of \mathbf{X} .



Work with these blue data instead!



Principal Component Analysis (PCA)

covariance matrix $\Sigma = [\text{cov}(X_i, X_j)] \quad i, j = 1, \dots, p$

$\Sigma = \Sigma^T$ Σ is symmetric

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ \leftarrow eigenvalues of Σ

$\Gamma = [\gamma_1 | \gamma_2 | \dots | \gamma_p]$ \leftarrow eigenvectors $\Gamma \Gamma^T = I = \Gamma^T \Gamma$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \quad \leftarrow \text{eigenvalues of } \Sigma$$

$$\Gamma = [\gamma_1 | \gamma_2 | \dots | \gamma_p] \quad \leftarrow \text{eigenvectors in } \mathbb{R}^p$$

$$\Gamma^T \Sigma \Gamma = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \quad \Rightarrow \quad \gamma_i^T \Sigma \gamma_i = \lambda_i, \quad i = 1, \dots, p.$$

$$\gamma_1 = \arg \max_{\|\gamma\|=1} \gamma^T \Sigma \gamma$$

$$\begin{cases} \gamma_j = \arg \max_{\|\gamma\|=1} \gamma^T \Sigma \gamma \\ \gamma^T \gamma_i = 0, \quad i = 1, \dots, j-1, \quad j > 1 \end{cases}$$

principal components of $\mathbf{X} \in \mathbb{R}^p$

$$Y_1 = \gamma_1^T \mathbf{X}, \quad Y_2 = \gamma_2^T \mathbf{X}, \quad \dots, \quad Y_p = \gamma_p^T \mathbf{X}$$

Symbolic data:

How to reduce dimensionality/How to find an SPCA?

- Data is organized as an $n \times p$ matrix, where rows correspond to **objects**, columns to **interval-valued variables**:

$$\text{Data matrix } \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}, \quad \text{e.g., } x_{11} = [a_{11}, b_{11}]$$

$$\text{random vector } \mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \rightarrow \mathbb{IR}^p$$

Example: $[x_{11} \ \cdots \ x_{1p}]^T \in \mathbb{IR}^p$ is a realisation of X .

What do we need?

$\mathbb{I}\mathbb{R}^p$ Principal components of \mathbf{X}

$$Y_1 = \gamma_1^T \mathbf{X}, \quad Y_2 = \gamma_2^T \mathbf{X}, \quad \dots, \quad Y_p = \gamma_p^T \mathbf{X}$$

$$Y_1 = \gamma_1^T \mathbf{X} = \gamma_1^T \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} = \gamma_1^T \begin{bmatrix} [a_1, b_1] \\ \vdots \\ [a_p, b_p] \end{bmatrix}$$

- 1 Define $\mathbb{I}\mathbb{R}^p$
- 2 interval arithmetic, **linear combinations** of intervals
- 3 “orthogonality”

- Definition of \mathbb{IR}^p

Let

$$\mathbb{IR} = \{[a, b] : a, b \in \mathbb{R}, a \leq b\}$$

be the set of all real closed and bounded intervals,

$[a, b]$ seen as a point (C, R) in \mathbb{R}^2

$$C = \frac{a+b}{2} \quad \text{centre of } [a, b]$$

$$R = b - a \quad \text{range of } [a, b]$$

- Definition of \mathbb{IR}^p

Let

$$\mathbb{IR}^p = \mathbb{IR} \times \mathbb{IR} \times \dots \times \mathbb{IR}$$

be the cartesian product of p copies of \mathbb{IR} .

$$\mathbb{IR} \longleftrightarrow \mathbb{R}^2$$

Let $\mathbf{X} = (X_1, \dots, X_p)^T$ be an **interval-valued random vector**

$$X_1 \leftrightarrow (C_1, R_1), \quad \dots, \quad X_p \leftrightarrow (C_p, R_p)$$

$$\mathbf{C} = \begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ C_p \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_p \end{bmatrix}$$

① Definition of $\mathbb{IR}^p = \mathbb{IR} \times \mathbb{IR} \times \dots \times \mathbb{IR} \longleftrightarrow \mathbb{R}^{2p}$

Let $\mathbf{X} = (X_1, \dots, X_p)^T = (\mathbf{C}^T, \mathbf{R}^T)^T = \begin{bmatrix} \mathbf{C} \\ \mathbf{R} \end{bmatrix}$ be an **interval-valued random vector** with realisations in \mathbb{R}^{2p}

$$X_1 \leftrightarrow (C_1, R_1), \quad \dots, \quad X_p \leftrightarrow (C_p, R_p)$$

$$\mathbf{C} = \begin{bmatrix} C_1 \\ \vdots \\ C_p \end{bmatrix} = \frac{1}{2} \left(\begin{bmatrix} a_1 \\ \vdots \\ a_p \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_p \end{bmatrix} \right) = \frac{1}{2}(\mathbf{A} + \mathbf{B}) \quad \text{random vector of centres}$$

$$\mathbf{R} = \begin{bmatrix} R_1 \\ \vdots \\ R_p \end{bmatrix} = \mathbf{B} - \mathbf{A} \quad \text{random vector of ranges}$$

- **Moore's interval arithmetic**

For $X = [a, b] = (C, R)$, $Y = [c, d] = (C', R')$ in \mathbb{IR} , define

$$X + Y = (C + C', R + R') = \{x + y : x \in X, y \in Y\}$$

$$\alpha X = (\alpha C, |\alpha| R)$$

Linear combination

For $\mathbf{X} \in \mathbb{IR}^p$ and $\alpha \in \mathbb{R}^p$, define

$$\alpha^T \mathbf{X} = \sum_{i=1}^p \alpha_i X_i = (\alpha^T \mathbf{C}, |\alpha|^T \mathbf{R})^T = \left(\sum_{i=1}^p \alpha_i C_i, \sum_{i=1}^p |\alpha_i| R_i \right)$$

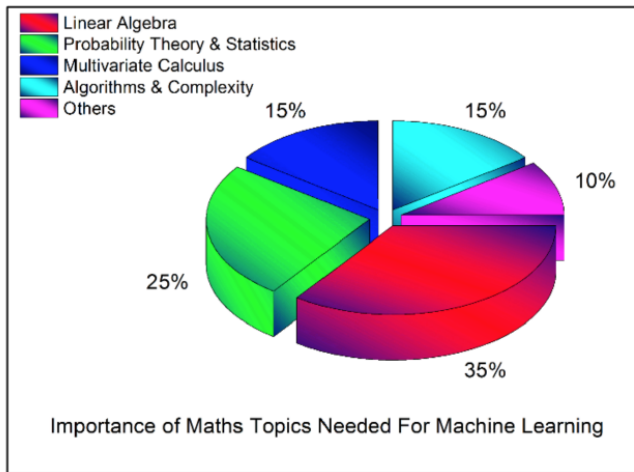
Theorem

The p *orthogonal-pcs* of \mathbf{X} are Y_1, \dots, Y_p with $Y_i = \gamma_i^T \mathbf{X}$, where $\gamma_i \in \mathbb{R}^p$ is

$$\gamma_i = \arg \max_{l=1, \dots, 2^p} \left(\gamma_{il}^T \Sigma_{CC} \gamma_{il} + \gamma_{il}^T \mathbf{M}_{RR} \gamma_{il} \right),$$

and γ_{il} , $l = 1, \dots, 2^p$, solves a subproblem:

$$\gamma_{il} = \begin{cases} \arg \max_{\gamma: \|\gamma\|=1} \left(\gamma^T \Sigma_{CC} \gamma + \gamma_{il}^T \mathbf{M}_{RR} \gamma \right) \\ \gamma^T \gamma_j = 0, j = 1, \dots, i-1 \text{ if } i > 1 \end{cases},$$



(by Wale Akinfaderin, <https://tinyurl.com/y6hc9sh8>)