

Multi-dimensional scaling

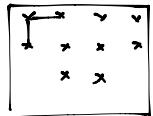
Distance matrix

$$D = (d_{ij}) \text{ where } d_{ij} = d(x_i, x_j) \quad i, j = 1, \dots, n$$

$n \times n$ symmetric

not necessarily observed

- want to represent (possibly unobserved) data in low dimensions



distances in 2D plot close to those in distance matrix

$$\text{Observe } X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \quad n \times p$$

Euclidean distances

$$d_{ij} = d(x_i, x_j) = [(x_i - x_j)^T (x_i - x_j)]^{1/2}$$

$$\text{Recover } \{d_{ij}\} \text{ from } B = XX^T = \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \cdots & b_{nn} \end{pmatrix}$$

$$d_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij} \quad n \times p$$

$$\text{SVD: } X = U \Lambda^{1/2} V^T$$

$$XX^T = U \Lambda U^T$$

$$U = \begin{pmatrix} u_1 & \cdots & u_p \end{pmatrix}$$

- optimal k dimensional representation of points $(\lambda_1^{1/2} u_1, \lambda_2^{1/2} u_2, \dots, \lambda_k^{1/2} u_k)$

Now the real problem: Given distance (similarity) matrix, how to get low dimensional representation for data?

Idea: Use $\{d_{ij}\}$ to construct matrix B .

$$b_{ij} = -\frac{1}{2} (d_{ij}^2 - d_{ii} - d_{jj} + d_{..}^2)$$

$$\text{where } d_{ii} = \frac{1}{n} \sum_{j=1}^n d_{ij}^2$$

$$d_{..}^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2$$

\tilde{B} is symmetric and so we can write

$$B = U \Lambda U^T \quad . \quad \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

ordered

- if all eigenvalues $\lambda_1, \dots, \lambda_n$ are non-negative then eigenvalues
 $\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_n}$ gives an idea of how good the k dim representation is.

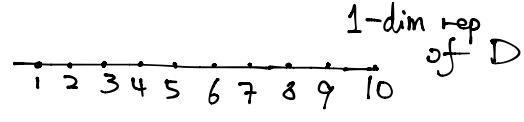
K=2 Plot $\lambda_1^{1/2} u_1$ vs $\lambda_2^{1/2} u_2$

Example

Start with very simple distance matrix

$$n=10 \quad D = (d_{ij}) \quad d_{ij} = |i-j|$$

$$= \begin{pmatrix} 0 & 1 & \cdots & 9 \\ 1 & 0 & \cdots & 8 \\ \vdots & \vdots & \ddots & \vdots \\ 9 & \cdots & 0 \end{pmatrix}$$



What happens to MDS using D? R cmdscale

- intuition suggests that $\lambda_1 > 0, \lambda_2 = \dots = \lambda_{10} = 0$

Output: $\lambda_1 = 82.5, \lambda_2 = \dots = \lambda_{10} = 0$

$$\lambda_1^{\frac{1}{2}} \underline{u}_1 = \begin{pmatrix} -4.5 \\ -3.5 \\ \vdots \\ 4.5 \end{pmatrix}$$

Now add noise to D

$$D_i = D + N$$

symmetric elements 0+1

- order of distances remaining unchanged for each row (& column.)

cmdscale : $\lambda_1 = 92.5, \lambda_2 = 5.87, \lambda_3 = 2.54$
 $\lambda_4 = 1.69, \lambda_5 = 1.35, \lambda_6, \dots, \lambda_{10} < 6$

$$\lambda_1^{\frac{1}{2}} \underline{u}_1 = \begin{pmatrix} -4.79 \\ -3.63 \\ -2.64 \\ -1.66 \\ -0.68 \\ 0.66 \\ 1.71 \\ 2.62 \\ 3.64 \\ 4.76 \end{pmatrix}$$

$$\lambda_2^{\frac{1}{2}} \underline{u}_2 = \begin{pmatrix} 0.96 \\ 0.58 \\ 0.18 \\ -0.62 \\ -1.16 \\ -1.26 \\ -0.18 \\ -1.06 \\ 0.71 \\ 0.85 \end{pmatrix}$$

Independent Components Analysis ICA

- Cousin of PCA

- Theory of PCA. X random vector with covariance matrix C , mean vector μ . Write $X = \mu + A Y$ where the components of $Y (y_1, \dots, y_p)$ are uncorrelated. $\text{Cov}(y_i, y_j) = 0$ with mean 0 & variance 1.

But Y is not uniquely determined:

$\text{Cov}(Y) = I \Rightarrow$ if Q is orthogonal then $\text{Cov}(QY) = I$

Look at C and choose A based on eigenvalues and eigenvectors of C .
⇒ not invariant to choice of scale.

Data analysis: Use data to estimate C (or R) and use these to estimate PC loadings

ICA: X mean vector μ , covariance matrix C

Model: Write $X = \mu + A Y$ where the components of $Y (y_1, \dots, y_p)$ are mutually independent

with $E(y_i) = 0, \text{Var}(y_i) = 1$.

mixing matrix

Model is invariant to linear transformations

$$\underline{X}' = \underline{\alpha} + B\underline{X} = \underline{\alpha} + B\underline{\mu} + BA\underline{X}$$

↑ mixing matrix

Key assumption: \underline{X} is not multivariate normal and at most one component of \underline{X} is normal.

- non-unique of mixing matrix A.

If \underline{X} is multivariate normal then ICA=PCA

$$\text{ICA: } \underline{X} = \underline{\mu} + A\underline{Y} \rightarrow \text{indep components}$$

↓
mixing matrix

Given data $\underline{x}_1, \dots, \underline{x}_n$, how to estimate A?

Two steps

① "Prewhtening": Transformed centred data $\underline{x}_1 - \bar{\underline{x}}, \dots, \underline{x}_n - \bar{\underline{x}}$ to make variables uncorrelated (with variance 1)
 $\underline{\underline{z}}_i = L(\underline{x}_i - \bar{\underline{x}})$ with $\frac{1}{n-1} \sum_{i=1}^n \underline{\underline{z}}_i \underline{\underline{z}}_i^T = I$ eg. PCs.
 ↓
pxp

② Component extraction

Find one-dimensional projection $\underline{w}_1, \dots, \underline{w}_p$ with $\underline{w}_i^T \underline{w}_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$
 and define $W = \begin{pmatrix} \underline{w}_1^T \\ \vdots \\ \underline{w}_p^T \end{pmatrix}$

want components
to be independent → $\underline{y}_i = W \underline{\underline{z}}_i = WL(\underline{x}_i - \bar{\underline{x}})$
 $\underline{x}_i = \bar{\underline{x}} + (WL)^{-1} \underline{y}_i = \bar{\underline{x}} + L^{-1} W^{-1} \underline{y}_i$

How to do steps ① + ②:

Look Fast ICA algorithm:

① can be done using PCA:

$$\hat{C} = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T = V \Lambda V^T$$

$$\underline{\underline{z}}_i = \Lambda^{-\frac{1}{2}} V^T (\underline{x}_i - \bar{\underline{x}})$$

Easy to verify $\frac{1}{n-1} \sum_{i=1}^n \underline{\underline{z}}_i \underline{\underline{z}}_i^T = I$

② Find ~~not~~ independent components using prewhitened data.

- find projections of prewhitened data that are as non-normal as possible
 \Rightarrow projection pursuit

Need indices of non-normality for a given projection $\underline{w} \rightarrow J(\underline{w})$

Example ① $H(x) = \ln[\cosh(x)] = \ln\left[\frac{\exp(x) + \exp(-x)}{2}\right]$

$$H'(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

$$H''(x) = 1 - [H'(x)]^2$$

Define $J(w) = \frac{1}{n} \sum_{i=1}^n H(w^T z_i)$ Kurtosis

② Kurtosis: Maximize $J(w) = \left| \frac{1}{n} \sum_{i=1}^n (w^T z_i)^4 - 3 \right|$ for normal dist'n

Why try to maximize non-normality?

Look at ICA model

$$\begin{aligned} \underline{x} &= \mu + A\underline{y} \\ \underline{y} &= A^{-1}(\underline{x} - \mu) \end{aligned}$$

$$\text{Cov}(\underline{y}) = I \quad \text{and} \quad \text{Cov}(Q\underline{y}) = I \quad \text{for any orthogonal matrix } Q = \begin{pmatrix} q_{11} & \cdots & q_{1p} \\ \vdots & \ddots & \vdots \\ q_{p1} & \cdots & q_{pp} \end{pmatrix}$$

$$\underline{s} = Q\underline{y} = \begin{pmatrix} s_1 \\ \vdots \\ s_p \end{pmatrix}$$

What do we know?

① s_1, \dots, s_p are uncorrelated with mean 0 & variance 1

② $s_i = \underbrace{q_{i1} y_1 + q_{i2} y_2 + \dots + q_{ip} y_p}_{\substack{\text{addition of } k=2 \\ \text{independent r.v.s}}} \text{ for } i \in [1, p]$

Central Limit Thm: "Sums of independent r.v.s are closer to normal than each of the summands."

s_1, \dots, s_p closer to normal than y_1, \dots, y_p
 \Rightarrow maximizing non-normality of projections is reasonable.