

UL = MULTIDIMENSIONAL SCALING (MDS)

"Multidimensional scaling (MDS) is a method that represents measurements of similarity (or dissimilarity) among pairs of objects as distances among points of a low-dimensional multidimensional space". (Borg & Groenen (2005))

Instead of the usual observations $x_1, \dots, x_n \in \mathbb{R}^d$,

the data consists of measurements of similarity/dissimilarity among pairs of objects: d_{ij} (e.g. $d_{ij} = d(x_i, x_j)$)
 object i object j $d = \text{distance in a metric space}$

Usually, 2D, for visualization purposes.

- An object can be
- a person
 - a country
 - an attribute
 - a stimulus ... / ...

• Goal: to represent these objects in a 2D space

• Ex: Perception of Morse Signals (Rothkopf (1957))

There are 36 Morse signals (26 for the letters + 10 for the numbers). The task is to judge if two morse signals are the same or different. The conducted experiment consisted of 598 subjects (unfamiliar with the Morse code), who were shown 351 pairs. Each pair was presented in two orders. For example, for the pair $(A, C) = (.-, -.-.)$, first $.-$ then $-.-.$, and also $-.-.$ then $.-$.

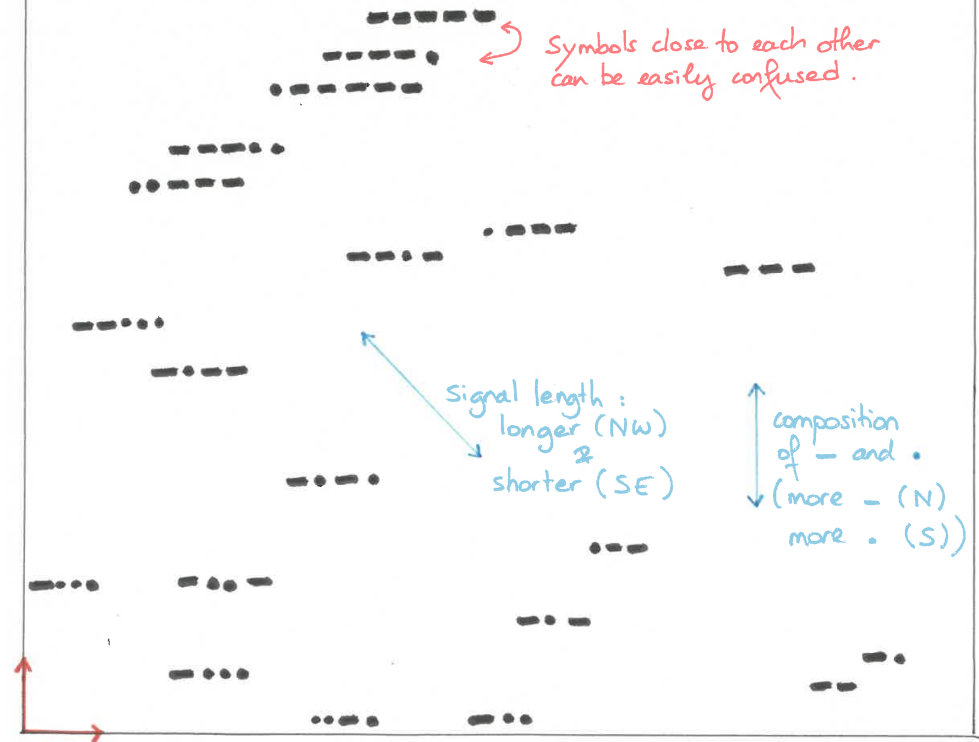
Which Morse codes are the most easy to be confused?

At the end of the experiment, we obtain a confusion matrix

| | A | B | ... | g | o |
|-----|-----|-----|-----|---|-----|
| A | 92 | 4 | ... | | 3 |
| B | 5 | 84 | ... | | 4 |
| ... | ... | ... | ... | | ... |
| g | | | | 9 | 3 |
| o | | | | | 94 |

matrix of similarities

Since distances are always symmetric, the matrix is first symmetrized, and the diagonal terms are set to 0 ($d(x_i, x_i) = 0 \forall i$). MDS would provide a visual representation of the confusion matrix. It would look something like this:



I - CLASSICAL MDS. (3)

- An $(n \times n)$ matrix $D_X = (\delta_{ij})_{i,j}$ is called a DISTANCE or AFFINITY MATRIX if
 - it is symmetric $\delta_{ij} = \delta_{ji}$
 - $\delta_{ii} = 0$
 - $\delta_{ij} > 0 \quad \forall i \neq j$.

Typically, δ_{ij} represents the distance between two observations x_i and $x_j \in \mathbb{R}^d$ a high dimensional space.

$$\delta_{ij} = d(x_i, x_j) = \|x_i - x_j\|$$

We assume here that $d =$ Euclidean distance, but it could be something else.

Given D_X , the task is to find n points $z_1, \dots, z_n \in \mathbb{R}^r$ ($r \ll d$) such that $D_Z := (d_{ij})_{i,j}$; $d_{ij} = d(z_i, z_j) = \|z_i - z_j\|$ is "similar" to D_X :

$$\min_{\substack{z_1, \dots, z_n \\ \in \mathbb{R}^r}} \sum_{i,j=1}^n (\delta_{ij} - d_{ij})^2 \quad \text{(MDS 1)}$$

$\delta_{ij} = \|x_i - x_j\|$ given z_i, z_j depend on z_i, z_j
 $d_{ij} = \|z_i - z_j\|$
MDS OPTIMIZATION PROBLEM.

Fact: Consider $D_X^{(2)} := (\delta_{ij}^2)_{i,j}$ = the matrix of square Euclidean distances. Then $-\frac{1}{2} J_n D_X^{(2)} J_n = X X^T$, where

$$X = \begin{pmatrix} x_{11} & x_{1n} \\ \vdots & \vdots \\ x_{n1} & x_{nn} \end{pmatrix} = \begin{pmatrix} -x_1^t \\ \vdots \\ -x_n^t \end{pmatrix} = \text{matrix of obs.}$$

(assume centered columns)

$$J_n = I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^t = \text{centering matrix} \quad \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

($1 \times n$)

→ First, note that for an $(n \times n)$ square matrix $A = (a_{ij})$, $A J_n$ removes the mean of each row of A , while $J_n A$ removes the mean of each column of A . (4)

$$J_n A = A - \frac{1}{n} \mathbf{1} \mathbf{1}^t A = \begin{pmatrix} \sum_{i=1}^n a_{i1} & \dots & \sum_{i=1}^n a_{in} \\ \vdots & \dots & \vdots \\ \sum_{i=1}^n a_{i1} & \dots & \sum_{i=1}^n a_{in} \end{pmatrix}$$

$$\Rightarrow J_n A = \begin{pmatrix} a_{11} - a_{.1} & \dots & a_{1n} - a_{.n} \\ \vdots & \dots & \vdots \\ a_{n1} - a_{.1} & \dots & a_{nn} - a_{.n} \end{pmatrix}$$

where columns are centered

where $a_{.j} = \frac{1}{n} \sum_{i=1}^n a_{ij} = \text{column } j \text{ mean}$

$$\text{likewise, } A J_n = \begin{pmatrix} a_{11} - a_{1.} & \dots & a_{1n} - a_{1.} \\ \vdots & \dots & \vdots \\ a_{n1} - a_{n.} & \dots & a_{nn} - a_{n.} \end{pmatrix}$$

where rows are centered

where $a_{i.} = \frac{1}{n} \sum_{j=1}^n a_{ij} = \text{row } i \text{ mean}$

In addition,

$$J_n A J_n = \begin{pmatrix} a_{11} - a_{1.} - a_{.1} + a_{..} & \dots & a_{1n} - a_{1.} - a_{.n} + a_{..} \\ \vdots & \dots & \vdots \\ a_{n1} - a_{n.} - a_{.1} + a_{..} & \dots & a_{nn} - a_{n.} - a_{.n} + a_{..} \end{pmatrix}$$

where $a_{..} = \frac{1}{n^2} \sum_{i,j=1}^n a_{ij}$

$$\Rightarrow (J_n A J_n)_{ij} = a_{ij} - a_{i.} - a_{.j} + a_{..}$$

(*)

→ Put $B := -\frac{1}{2} J_n D_x^{(2)} J_n$. (5)
 We show next that $B = X X^t$ (assuming X has centered columns)

Note that the (i, j) element of the matrix $-2B$ is $d_{ij}^2 - d_{i.}^2 - d_{.j}^2 + d_{..}^2$, where

$$\begin{cases} d_{i.}^2 = \frac{1}{n} \sum_{j=1}^n \|x_i - x_j\|^2 \\ d_{.j}^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - x_j\|^2 \\ d_{..}^2 = \frac{1}{n^2} \sum_{i,j=1}^n \|x_i - x_j\|^2 \end{cases}$$

(follows from (*) page 4)

In addition, $\|x_i - x_j\|^2 = \|x_i\|^2 + \|x_j\|^2 - 2 x_i^t x_j$. ($i, j = 1, \dots, n$)
 Assuming that X has centered columns: $\sum_{i=1}^n x_{ij} = 0 \quad \forall j$,

we see that

$$\begin{aligned} \sum_{i=1}^n x_i^t x_j &= \sum_{i=1}^n \sum_{k=1}^d x_{ik} x_{jk} \\ &= \sum_{k=1}^d x_{jk} \left(\sum_{i=1}^n x_{ik} \right) = 0 \end{aligned}$$

$= 0 \quad \forall k=1, \dots, d$

$$\Rightarrow \sum_{i=1}^n \|x_i - x_j\|^2 = \sum_{i=1}^n \|x_i\|^2 + n \|x_j\|^2 - 2 \sum_{i=1}^n x_i^t x_j$$

$$d_{.j}^2 = \frac{1}{n} \sum_{i=1}^n \|x_i\|^2 + \|x_j\|^2$$

Likewise, $\sum_{j=1}^n \|x_i - x_j\|^2 = \sum_{j=1}^n \|x_j\|^2 + n \|x_i\|^2$, and

$$d_{i.}^2 = \frac{1}{n} \sum_{j=1}^n \|x_j\|^2 + \|x_i\|^2$$

&

$$d_{..}^2 = \frac{2}{n} \sum_{i=1}^n \|x_i\|^2$$

$$\Rightarrow d_{ij}^2 - d_{i.}^2 - d_{.j}^2 + d_{..}^2 = \|x_i - x_j\|^2 - \|x_i\|^2 - \|x_j\|^2 = -2 x_i^t x_j \quad (6)$$

$$\Rightarrow -\frac{1}{2} (d_{ij}^2 - d_{i.}^2 - d_{.j}^2 + d_{..}^2) = x_i^t x_j \quad (i, j = 1, \dots, n)$$

$$= (X X^t)_{ij}$$

Provided X has centered columns, we only need to know $d_{ij} = \|x_i - x_j\|$ to calculate $\langle x_i, x_j \rangle$, we do not need to know x_i and x_j .

⇒ Instead of considering the original (MDS 1) (page 3) optimization problem, we may consider the matrix of square Euclidean distances $D_x^{(2)}$, apply the transform $B = -\frac{1}{2} J D_x^{(2)} J$, and search for low-dimensional representations z_1, \dots, z_n such that

$$\min_{\substack{z_1, \dots, z_n \\ \in \mathbb{R}^r}} \sum_{i,j=1}^n (\langle x_i, x_j \rangle - \langle z_i, z_j \rangle)^2$$

(MDS 2) aka CLASSICAL MDS.

Remarks:

- (i) $\langle x_i, x_j \rangle = (i, j)$ element of $B = -\frac{1}{2} J D_x^{(2)} J$
- (ii) $z_i \in \mathbb{R}^r$, $x_i \in \mathbb{R}^d$, with $r \ll d$ (usually $r=2$)
- (iii) Columns of $X = \begin{pmatrix} - & x_1^t & - \\ - & x_n^t & - \end{pmatrix}$ are assumed centered, so that $B = -\frac{1}{2} J D_x^{(2)} J$. Since the same linear transformation is applied to $D_z^{(2)}$, the solution to (MDS 2) must also be centered.

→ Solution to (MDS 2) optimization problem. (7)

Consider the SVD decomposition of $X = U \Lambda V^t$
(nxd) (nxd)(dxd) (dxd)

Then $XX^t = U \Lambda^2 U^t$, and the (i,j) entry of XX^t is $\langle x_i, x_j \rangle = (U \Lambda^2 U^t)_{ij}$

↑ term appearing in the optimization problem.

Then, the Eckart-Young theorem (p. 8 in UL: PCA) ensures that the best rank r approximation to XX^t (measured with respect to Frobenius norm) is given by

$U_r \Lambda_r^2 U_r^t$, where

→ Λ_r^2 = diagonal matrix retaining the largest r values of Λ^2

→ U_r = $(n \times r)$ matrix retaining the r columns of U associated with the r largest values of Λ^2 .

⇒ Solution Z to (MDS 2) satisfies $ZZ^t = U_r \Lambda_r^2 U_r^t = (U_r \Lambda_r)(U_r \Lambda_r)^t$

⇒ $Z = U_r \Lambda_r$
(nxr) (nxr) (rxr)
 = PCA solution.

Classical MDS & PCA return the same solution.

• Algorithm. (Classical MDS)

- Start with a distance matrix D_x
- Compute $B = -\frac{1}{2} J_n D_x J_n$
- Consider the eigenvalue / eigenvector decomposition of B
- Retain the largest r components

II. THE SMACOF ALGORITHM (8)

The SMACOF (Scaling by MAjorizing a COmplicated Function) approach searches to minimize the stress function

$$\sigma(Z) = \sum_{i < j} w_{ij} (\delta_{ij} - d_{ij})^2$$
 (MDS 3) De Leeuw & Heiser (1977)

$Z = \begin{pmatrix} -z_1^t \\ \vdots \\ -z_n^t \end{pmatrix}$
(nxr)
 or
 $Z = \begin{pmatrix} | & \dots & | \\ z_1 & \dots & z_r \\ | & \dots & | \end{pmatrix}$

weights (≥ 0)

given

$d_{ij} = \|z_i - z_j\|$

Weight w_{ij} measures the importance of the residual $(\delta_{ij} - d_{ij})^2$.
 Setting $w_{ij} = 0$ can be used to handle missing data.

Note that the overall scaling does not matter; scaling δ_{ij} by a factor α and coordinates z_1, \dots, z_n by the same factor, so that $d_{ij}(\alpha Z) = \alpha d_{ij}(Z)$, yields $\alpha^2 \sigma(Z)$.
 De Leeuw fixes $\sum_{i < j} w_{ij} \delta_{ij}^2$ to some fixed constant, in order to standardize the solution. He chooses $\sum_{i < j} w_{ij} \delta_{ij}^2 = \frac{n(n-1)}{2}$ so that disparities are on average dispersed around 1.

⇒
$$\sigma(Z) = \underbrace{\sum_{i < j} w_{ij} \delta_{ij}^2}_{\frac{n(n-1)}{2}} + \underbrace{\sum_{i < j} w_{ij} d_{ij}^2}_{\eta^2(Z)} - 2 \underbrace{\sum_{i < j} w_{ij} \delta_{ij} d_{ij}}_{e(Z)}$$

• Expressions for $\eta^2(Z)$ and $e(Z)$: First, note that

$$d_{ij}^2 = \sum_{l=1}^r (z_{il} - z_{jl})^2 = \sum_{l=1}^r (z_l^t (e_i - e_j))^2$$

$$\begin{aligned} \dots d_{ij}^2 &= \sum_{l=1}^r z_l^t (e_i - e_j)(e_i - e_j)^t z_l \\ &= \text{Tr} \left\{ Z^t (e_i - e_j)(e_i - e_j)^t Z \right\} \\ &= \text{Tr} \left\{ Z^t A_{ij} Z \right\}, \end{aligned} \quad (9)$$

where we defined $A_{ij} := (e_i - e_j)(e_i - e_j)^t = \begin{pmatrix} 1 & & & \\ & -1 & & \\ & & & \\ & & & 1 \end{pmatrix} \leftarrow \begin{matrix} i \\ j \end{matrix}$

It follows that

$$\begin{aligned} \rho^2(Z) &= \sum_{i < j} w_{ij} d_{ij}^2 = \sum_{i < j} w_{ij} \text{Tr} \left\{ Z^t A_{ij} Z \right\} \\ &= \text{Tr} \left\{ Z^t \left(\sum_{i < j} w_{ij} A_{ij} \right) Z \right\} \\ &= \text{Tr} \left\{ Z^t V Z \right\}, \end{aligned}$$

where $V := \sum_{i < j} w_{ij} A_{ij}$.

When all $w_{ij} = 1$, V simplifies to $V = nI - \underline{1}\underline{1}^t$.

Assuming all $d_{ij} > 0$,

$$\begin{aligned} \rho(Z) &= \sum_{i < j} w_{ij} \delta_{ij} d_{ij} = \sum_{i < j} w_{ij} \frac{\delta_{ij}}{d_{ij}} d_{ij}^2 \\ &= \sum_{i < j} w_{ij} \delta_{ij} d_{ij}^{-1} \text{Tr} \left\{ Z^t A_{ij} Z \right\} \\ &= \text{Tr} \left\{ Z^t \left(\sum_{i < j} w_{ij} \delta_{ij} d_{ij}^{-1} A_{ij} \right) Z \right\} \\ &= \text{Tr} \left\{ Z^t B(Z) Z \right\} \end{aligned}$$

If $d_{ij} = 0$, set $s_{ij}(Z) = 0$

$$\Rightarrow \rho(Z) = \text{Tr} \left\{ Z^t B(Z) Z \right\},$$

where $B(Z) := \sum_{i < j} w_{ij} s_{ij}(Z) A_{ij}$

$$s_{ij}(Z) = \begin{cases} \delta_{ij} d_{ij}^{-1} & \text{if } d_{ij} > 0 \\ 0 & \text{if } d_{ij} = 0 \end{cases}$$

In summary,

$$\sigma(Z) = \frac{n(n-1)}{2} + \text{Tr} \left\{ Z^t V Z \right\} - 2 \text{Tr} \left\{ Z^t B(Z) Z \right\}. \quad (10)$$

To minimize this function, we use the SMA COF approach.

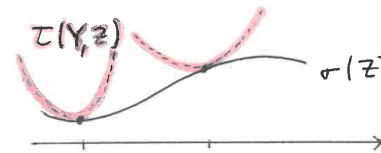
Idea: Instead of minimizing $\sigma(Z)$ directly, minimize a majorizing function $\tau(Y, Z)$ which satisfies:

• $\tau(Y, Z)$ is easier to minimize (e.g. quad.)

• $\sigma(Z) \leq \tau(Y, Z)$

• $\sigma(Z) = \tau(Z, Z)$

Y is called a supporting point



In our setting, consider

$$\tau(Y, Z) := \frac{n(n-1)}{2} + \text{Tr} \left\{ Z^t V Z \right\} - 2 \text{Tr} \left\{ Z^t B(Y) Y \right\} \quad (*)$$

↳ Quadratic in Z

↳ Clearly τ and σ coincide at $Y = Z$

↳ $\sigma(Z) \leq \tau(Y, Z)$; a consequence of the CS inequality.

Indeed, we want to show that $\text{Tr} \left\{ Z^t B(Y) Y \right\} \leq \text{Tr} \left\{ Z^t B(Z) Z \right\}$

Assuming all $d_{ij}(Y) > 0$,

$$\text{Tr} \left\{ Z^t B(Y) Y \right\}$$

$$= \sum_{i < j} w_{ij} \frac{\delta_{ij}}{d_{ij}(Y)} \sum_{l=1}^d (z_{il} - z_{jl})(y_{il} - y_{jl})$$

$$\langle z_i - z_j, y_i - y_j \rangle = \sum_{l=1}^d (z_{il} - z_{jl})(y_{il} - y_{jl}) \quad (11)$$

CS $\hookrightarrow \leq d_{ij}(z) d_{ij}(Y)$

$$\Rightarrow \text{Tr} \{ z^t B(Y) Y \} \leq \sum_{i < j} w_{ij} \frac{\delta_{ij}}{d_{ij}(Y)} d_{ij}(z) d_{ij}(Y)$$

$$= \text{Tr} \{ z^t B(z) z \}$$

Note that if $d_{ij}(Y) = 0$, then $y_i = y_j$ and $\text{Tr} \{ z^t B(Y) Y \} = 0$ & the inequality is still true. ■

• A general iterative majorization algorithm:

- (i) Start with $Y_0 = z_0$
- (ii) For $k=0, 1, 2, \dots$
 - Find z_{k+1} s.t. $\tau(Y_k, z_{k+1}) \leq \tau(Y_k, Y_k)$
 - Set $Y_{k+1} = z_{k+1}$

And indeed, $\sigma(z_{k+1}) \leq \tau(z_k, z_{k+1}) \leq \tau(z_k, z_k) = \sigma(z_k)$

• It remains to compute the gradient of $\tau(Y, z)$ with respect to z :

$$\nabla \tau(Y, z) = 2Vz - 2B(Y)Y = 0$$

\hookrightarrow Need to solve $Vz = B(Y)Y$. (**)

$$z = V^+ B(Y)Y,$$

where V^+ is the Moore-Penrose inverse of V .

SMACOF Starting with z_0 , iterate $z_{k+1} = V^+ B(z_k) z_k$

• Remark: SMACOF & Gradient Descent. (12)

The stress function is $\sigma(z) = \text{Tr} \{ z^t V z \} - 2 \text{Tr} \{ z^t B(z) z \} + \text{constant}$.

Its gradient is $\nabla \sigma(z) = 2Vz - 2B(z)z$

\nearrow Not "that" obvious.

Consider

$$\frac{d}{dz_{kl}} \left\{ \sum_{i < j} w_{ij} \delta_{ij} \left(\sum_{m=1}^r (z_{im} - z_{jm})^2 \right)^{1/2} \right\}$$

$1 \leq k \leq n$
 $1 \leq l \leq r$ \nearrow either i or j must be equal to k , otherwise the derivative vanishes.

2 cases \rightarrow $i=k < j$
 \searrow $i < j=k$.

$$\bullet \underline{i=k < j} = \sum_{k=i < j} w_{kj} \delta_{kj} \frac{1}{2} z \frac{(z_{kl} - z_{jl})}{d_{kj}(z)}$$

$$\bullet \underline{i < j=k} = \sum_{i < j=k} w_{ik} \delta_{ik} \frac{1}{2} z - \frac{(z_{il} - z_{kl})}{d_{ik}(z)}$$

\uparrow
Compare these terms with the (k, l) entry of $B(z)z$:

$$B(z)z = \sum_{i < j} w_{ij} \delta_{ij} d_{ij}^{-1} \boxed{A_{ij} z}$$

they are the same!

$$i=k < j \rightarrow \begin{pmatrix} (z_{i1} - z_{j1}) \dots (z_{il} - z_{jl}) \dots \dots \\ \vdots \\ (z_{j1} - z_{i1}) \dots (z_{jl} - z_{il}) \dots \dots \end{pmatrix}$$

\downarrow l

⇒ Gradient descent algorithm for stress minimization: (13)

$$\begin{aligned} z_{kh} &= z_k - \eta_k \nabla \sigma(z_k) \\ &= z_k - 2\eta_k \{ V z_k - B(z_k) z_k \} \end{aligned}$$

↑
Compare this iteration with the SMA-COF algorithm

$$\begin{aligned} z_{kh} &= V^T B(z_k) z_k \\ &= (z_{kh} - z_k) + V^T B(z_k) z_k \\ &= z_{kh} - \frac{1}{2} V^T \{ 2 V z_k - 2 B(z_k) z_k \} \\ &= \text{Gradient descent algorithm with constant step size.} \end{aligned}$$

x Remark: Constrained MDS

From the expression of $\tau(Y, Z)$ [(*) p. 10] and $\bar{Z} := V^T B(Y) Y$ [(**) p. 11], we have

$$\begin{aligned} \tau(Y, Z) &= \text{Tr} \{ z^t V z \} - 2 \text{Tr} \{ z^t B(Y) Y \} + \text{constant} \\ &= \text{Tr} \{ z^t V z \} - \text{Tr} \{ z^t V \bar{Z} \} + \text{constant} \\ &= \text{Tr} \{ (z - \bar{Z})^t V (z - \bar{Z}) \} - \text{Tr} \{ \bar{Z}^t V \bar{Z} \} + \text{constant} \end{aligned}$$

↑
The minimization of τ handles constraints on Z .

E.g. Linear constraints: $Z = HW$ → interpretable solution in terms of external variables.
observed covariates ← weights

$$\Rightarrow W_{\text{update}} = (H^t V H)^{-1} H^t V \bar{Z}$$

References:

- [1] P.J.F. Groenen & M. van de Velden.
Multidimensional Scaling by Majorization: A Review
Journal of Statistical Software, vol 73, issue 8 (2016)
- [2] P.J.F. Groenen & I. Borg.
The Past, Present, and Future of Multidimensional
Scaling.
Econometric Institute Report EI 2013-07.