

Multidimensional scaling

In this section, we introduce mapping methods on visual and geometric presentation of multivariate data, with the focus of maintaining the concepts of similarity and ranking among data points. These methods can be viewed as a type of unsupervised learning. Traditionally they belong to the ordination methods.

1 The concepts of Multidimensional Scaling

Associations among subjects or interdependence between variables are often quantified by pairwise relations, such as distance, similarity, and correlations. **Multidimensional scaling** (MDS) is a set of visualization techniques using pairwise information to map high dimensional data to a lower, usually two dimensional space, such that more "similar" objects are closer in the low dimensional configuration. The geometric representation provides insight into the relationships among objects of interest. The MDS process can be described as

Pairwise similarity or dissimilarity in high dimensions \implies A configuration, a map, in low dimensions

MDS Data

The data are $n(n-1)/2$ pairs of pairwise dissimilarities (or distances) of n items.

Presumably the n items were from a space of high dimension $p \geq 3$.

The $n(n-1)/2$ dissimilarity measures, denoted as d_{ji} , are often displayed as the lower triangular part of an $n \times n$ matrix, in which the (i, j) th entry is the dissimilarity or distance between item i and item j .

$$\begin{bmatrix} 0 & & & & & \\ d_{21} & 0 & & & & \\ \vdots & \vdots & & & & \\ d_{i1} & d_{i2} & \cdots & 0 & & \\ \vdots & \vdots & \vdots & \vdots & & \\ d_{n1} & d_{n2} & \cdots & d_{nj} & \cdots & 0 \end{bmatrix}$$

Quite often, the pairwise measures are given in the following misleading format:

$$\begin{bmatrix} 0 & 0 & \cdots & 0 & \cdots & 0 \\ d_{21} & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{i1} & d_{i2} & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{n1} & d_{n2} & \cdots & d_{nj} & \cdots & 0 \end{bmatrix}$$

If we assume the dissimilarity measure is symmetric with $d_{ij} = d_{ji}$, as dissimilarity measures most often are, it is suitable to convert the given $n(n-1)$ values to a symmetric $n \times n$ **dissimilarity matrix**, which has diagonal elements = 0.

$$\begin{bmatrix} 0 & d_{12} & \cdots & d_{1j} & \cdots & d_{1n} \\ d_{21} & 0 & \cdots & d_{2j} & \cdots & d_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{i1} & d_{i2} & \cdots & d_{ij} & \cdots & d_{ip} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{n1} & d_{n2} & \cdots & d_{nj} & \cdots & 0 \end{bmatrix}$$

It is not uncommon that the data are $n(n-1)/2$ measurements of pairwise similarities of n items, instead of dissimilarities. A **similarity matrix** should have the largest value on the diagonal, often = 1, indicating that an item is the most similar to itself.

$$\begin{bmatrix} 1 & s_{12} & \cdots & s_{1j} & \cdots & s_{1n} \\ s_{21} & 1 & \cdots & s_{2j} & \cdots & s_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ s_{i1} & s_{i2} & \cdots & s_{ij} & \cdots & s_{ip} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ s_{n1} & s_{n2} & \cdots & s_{nj} & \cdots & 1 \end{bmatrix}$$

A dissimilarity measure can be created based on the given similarity measure, and vice versa.

MDS objective

A $q < p$ dimensional map of data maintaining as much as possible the original dissimilarity relations. Usually $q \leq 3$, preferably $q \leq 2$.

Remarks on MDS

- The data format used by MDS is different from the datasets used by many other multivariate methods, which typically apply to a dataset with n observations of p -variate vectors.

MDS uses so called proximity data, which consists of similarity (or dissimilarity) information for pairs of objects. If the dataset is of n objects (n observations) of p -variate vectors, a similarity or dissimilarity matrix of dimensions $n \times n$ needs to be created. This matrix is the input data of MDS.

- Similar to PCA and FA, MDS achieves dimension reduction. However PCA and FA are linear methods with respect to the original data, while MDS is a non-linear procedure. The low dimensional configuration achieved by MDS aims to maintain similarity, at the loss of linearity if necessary.

Illustrative examples

- Example 1

Four observations A, B, C, D are in a high dimensional space $\mathbb{R}^p, p > 1$. Pairwise similarities of the points form a similarity matrix:

$$\begin{matrix} & A & B & C & D \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} 1 & & & \\ 0.980 & 1 & & \\ 0.995 & 0.955 & 1 & \\ 0.875 & 0.955 & 0.82 & 1 \end{bmatrix} \end{matrix}$$

The similarity matrix gives an ordering of the pairwise similarities:

$$s_{AC} \geq s_{AB} \geq s_{BC} \geq s_{BD} \geq s_{AD} \geq s_{CD} \quad (1)$$

There are many ways to convert a similarity measure s_{ik} between objects i and k to a distance measure d_{ik} , we may assign d to be $1/s$, $1/(s+1)$, $\sqrt{1-s^2}$, etc. Here we choose (actually with theoretical motivation)

$$d = 10\sqrt{2(1-s)}$$

Then the distance or dissimilarity matrix of the four points is

$$\begin{matrix} & A & B & C & D \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} 0 & & & \\ 2 & 0 & & \\ 1 & 3 & 0 & \\ 5 & 3 & 6 & 0 \end{bmatrix} \end{matrix}$$

From the distance matrix we have

$$d_{AC} \leq d_{AB} \leq d_{BC} \leq d_{BD} \leq d_{AD} \leq d_{CD} \quad (2)$$

which, by construction, matches the ordering of the original similarity measures in the opposite direction perfectly.

A question in MDS is: Can we assign coordinates $\hat{A}, \hat{B}, \hat{C}, \hat{D}$ to the points in a low dimensional space such that the corresponding distances maintain the ordering in (2)? For this example, the answer is yes. In fact, we can even find a $q = 1$ dimensional representation by assigning one-dimensional coordinates as

$$\hat{A} = 1, \quad \hat{B} = 3, \quad \hat{C} = 0, \quad \hat{D} = 6.$$

The distances of the one-dimensional representation of the points satisfy

$$d_{\hat{A}\hat{C}} \leq d_{\hat{A}\hat{B}} \leq d_{\hat{B}\hat{C}} \leq d_{\hat{B}\hat{D}} \leq d_{\hat{A}\hat{D}} \leq d_{\hat{C}\hat{D}}$$

thus maintaining the distance ordering of the points in the original high dimensional space.

• Example 2

Three observations A, B, C are in a high dimensional space $\mathbb{R}^p, p > 2$, with dissimilarity matrix

$$\begin{array}{c} A \quad B \quad C \\ \begin{array}{c} A \\ B \\ C \end{array} \begin{bmatrix} 0 & & \\ 4 & 0 & \\ 5 & 3 & 0 \end{bmatrix} \end{array}$$

From the distance matrix we have

$$d_{BC} \leq d_{AB} \leq d_{AC}$$

Can we assign coordinates $\hat{A}, \hat{B}, \hat{C}$ in a low dimensional space such that the corresponding distances maintain the ordering above? Again the answer is yes. In fact, we can even find a $q = 1$ dimensional representation by assigning one-dimensional coordinates as

$$\hat{A} = 1, \quad 3.5 < \hat{B} < 6, \quad \hat{C} = 6.$$

The distances of the one-dimensional representation of the points satisfy $d_{\hat{B}\hat{C}} \leq d_{\hat{A}\hat{B}} \leq d_{\hat{A}\hat{C}}$, thus maintaining the distance ordering of the points in the original high dimensional space.

However the pairwise similarity values are not perfectly maintained. A “loss” of information has occurred in the reduced dimension representation.

A $q = 2$ dimensional representation of

$$\hat{A} = (1, 1), \quad \hat{B} = (1, 4), \quad \hat{C} = (4, 6).$$

not only satisfying $d_{\hat{B}\hat{C}} \leq d_{\hat{A}\hat{B}} \leq d_{\hat{A}\hat{C}}$, but also maintaining the original distance perfectly, an ideal case of lower dimension representation.

• Example 3

Four observations A, B, C, D are in a high dimensional space $\mathbb{R}^p, p > 3$, with dissimilarity matrix

$$\begin{array}{c} A \quad B \quad C \quad D \\ \begin{array}{c} A \\ B \\ C \\ D \end{array} \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 1 & 1 & 1 & 0 \end{bmatrix} \end{array}$$

A tetrahedron in dimension $q = 4$ gives a geometric representation with perfect match of the dissimilarity. No geometric representation with perfect match of the original dissimilarity exists in lower dimension spaces $\mathbb{R}, \mathbb{R}^2, \mathbb{R}^3$.

Types of multidimensional scaling

Generally, N items form $N(N-1)/2$ pairwise measures of dissimilarities or distance. Multidimensional Scaling aims to find a map, a graphical representation of the items in lower dimensions such that the inter-item proximities match the original similarities as closely as possible.

If only the rank orders of the $N(N-1)/2$ original similarities are used, the process is called **non-numeric** scaling. If the actual magnitudes of the original similarities are used, the process is called **metric** scaling. The original metric MDS with Euclidean distance is **classical** scaling, also known as **principal coordinates** or principal coordinates analysis.

Multidimensional scaling method

Denote the similarity between items i and j by s_{ij} . Assume there are no ties in similarities. Arrange the similarities in a strictly ascending order as

$$s_{i_1 k_1} < s_{i_2 k_2} < \dots < s_{i_M k_M}, \quad M = N(N-1)/2 \quad (3)$$

If instead the measure is in terms of dissimilarity or distance d_{ik} , then strictly decreasing order is used. Multidimensional scaling attempts to find a lower q -dimensional configuration of the N items such that the distance $d_{ik}^{(q)}$ between the pairs of items match the order in (3). A perfect match occurs if

$$d_{i_1 k_1}^{(q)} > d_{i_2 k_2}^{(q)} > \dots > d_{i_M k_M}^{(q)} \quad (4)$$

Perfect match can be achieved if q = the dimension of the original data space, using the same proximity metric. However the interest is on finding (4) or its approximation in $q = 2$ or 3 dimensional space.

2 Measure of goodness of fit in MDS

For a given dimension q , there may not exist a configuration of the points whose pairwise distances match the ordering of the original pairwise dissimilarities. **Stress** (Kruskal, 1964) is a measure of closeness between the original similarities and the fitted values of similarities.

$$Stress(q) = \left\{ \frac{\sum \sum_{i < k} \left(d_{ik}^{(q)} - \hat{d}_{ik}^{(q)} \right)^2}{\sum \sum_{i < k} \left(d_{ik}^{(q)} \right)^2} \right\}^{1/2} \in [0, 1]$$

where $\hat{d}_{ik}^{(q)}$'s are the fitted distances in a q dimensional configuration, $d_{ik}^{(q)}$'s are ideal q -dimensional pairwise numbers (with respect to a reference dissimilarity matrix) that correspond to the perfect match with the original similarities s_{ik} or their dissimilarity counterpart.

An alternative measure of closeness is **SStress** (Takane, 1977). Dropping the index q in $d_{ik}^{(q)}$,

$$SStress(q) = \left\{ \frac{\sum \sum_{i < k} (d_{ik}^2 - \hat{d}_{ik}^2)^2}{\sum \sum_{i < k} d_{ik}^4} \right\}^{1/2} \in [0, 1]$$

Remarks

- Typically the value of *Stress* or *SStress* less than 0.1 is considered a good representation of the objects by the points in the given low dimension configuration. In practice, the goodness of fit of *Stress* is based on the range of the values:

	Perfect	Excellent	Good	Fair	Poor
Stress value	0%	2.5%	5%	10%	20%

- In the formula of *Stress* or *SStress*, the ideal distance $d_{ik}^{(q)}$ between the i th and k th objects often uses the corresponding entry D_{ik} in the input $n \times n$ dissimilarity matrix D , presumably the distance of the two objects in a higher dimensional space.
- The two stress measures often give comparable results. The newer one, *SStress*, is preferred sometimes. In our example demo, *SStress* consistently yields smaller values than *Stress*.
- If Euclidean distance in the lower dimension configuration space \mathbb{R}^q is used for dissimilarity, then $\hat{d}_{ik}^{(q)}$ in *Stress* has the simple form $\hat{d}_{ik}^{(q)} = \|x_i - x_k\|$ where $x_i = (x_{i1}, \dots, x_{iq})$ represents the coordinates of the i th object under the q -dimensional configuration, $\|\cdot\|$ is the usual ℓ_2 Euclidean norm.
- The *Stress* measure can be viewed as a function on \mathbb{R}^{qn} . An optimization algorithm can be carried out by numerical methods such as gradient descent.

Algorithm outline

- Obtain the $N(N-1)/2$ similarities between distance pairs of items. Order the similarities as in (3). Typically, dissimilarities or distances are used instead of similarities.
- Using a trial configuration in q dimensions, determine the inter-item distances $d_{ik}^{(q)}$ and estimate $\hat{d}_{ik}^{(q)}$ to minimize the *Stress* or *SStress* measure.
- Using $\hat{d}_{ik}^{(q)}$, moving the points around to obtain an improved configuration.
- Choose the dimension q by evaluating the *Stress* or *SStress* measure.

3 General measure of closeness

Similarity, dissimilarity, Proximity, distance, distance metric

- Similarity measure reflects how close two objects are. The “closer” the two objects are to each other, the larger is their similarity value.
- Dissimilarity measure indicates how different two objects are. The farther the two objects are to each other, the larger is their dissimilarity value.
- Proximity can be either similarity or dissimilarity.
- Distance sometime loosely refers to a dissimilarity measure.
On the other hand, distance as a metric is a more strictly defined mathematical concept.
- Definition of metric

A function $d(x, y)$ is a (distance) metric or a distance function if the following holds.

- $d(x, y) = d(y, x)$ for any x, y . (Symmetry)
 - $d(x, y) \geq 0$
 - $d(x, y) = 0$ if and only if $x = y$.
 - $d(x, y) \leq d(x, w) + d(w, y)$ for all x, y, w . (Triangle inequality)
- All sensible dissimilarity measures fulfill conditions 1 and 2. Most but not all dissimilarities satisfy condition 3. However it is not so uncommon that a dissimilarity measure fails to satisfy condition 4, which is usually the harder condition to verify.

Example

If correlation $\rho(\mathbf{x}, \mathbf{y})$ is used for similarity measure, then $d(\mathbf{x}, \mathbf{y}) = \sqrt{2(1 - \rho)}$ satisfies 1,2, 4 but not the “only if” part of 3. Such measure is called a pseudometric.

- Proof.*
- By the definition of correlation, $\rho(\mathbf{x}, \mathbf{y}) = \rho(\mathbf{y}, \mathbf{x})$, hence $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}) \geq 0$, property 1 holds.
 - By the definition of correlation, $\rho \leq 1$, $2(1 - \rho) \geq 0$, thus $d(\mathbf{x}, \mathbf{y}) \geq 0$, property 2 holds.
 - By the definition of correlation, if $\mathbf{y} = c\mathbf{x}$ ($c \neq 0$) is a non-zero constant multiple of \mathbf{x} , then $\rho(\mathbf{x}, \mathbf{y}) = 1$ thus $d(\mathbf{x}, \mathbf{y}) = 0$. Consequently, $d(\mathbf{x}, \mathbf{x}) = 0$ is true, thus the “if” part in property 3 holds. However $d(\mathbf{x}, 2\mathbf{x}) = 0$ also, which means the “only if” part in property 3 is not true.
 - To use the triangular property of Euclidean norm, let's consider centered, normed vectors.
Let $\mathbf{x}_c = \mathbf{x} - \bar{\mathbf{x}}$, $\mathbf{y}_c = \mathbf{y} - \bar{\mathbf{y}} \in \mathbb{R}^n$, where $\bar{\mathbf{x}}$ represents an n -vector of constant component value \bar{x} .
Since correlation does not change when the variables subtract constants,

$$\rho(\mathbf{x}, \mathbf{y}) = \rho(\mathbf{x} - \bar{\mathbf{x}}, \mathbf{y} - \bar{\mathbf{y}}) = \rho(\mathbf{x}_c - \bar{\mathbf{x}}, \mathbf{y}_c - \bar{\mathbf{y}}) = \rho(\mathbf{x}_c, \mathbf{y}_c)$$

The Euclidean distance between \mathbf{x}_c and \mathbf{y}_c is

$$\|\mathbf{x}_c - \mathbf{y}_c\|^2 = (\mathbf{x}_c - \mathbf{y}_c)'(\mathbf{x}_c - \mathbf{y}_c) = \mathbf{x}_c' \mathbf{x}_c + \mathbf{y}_c' \mathbf{y}_c - 2\mathbf{x}_c' \mathbf{y}_c = \|\mathbf{x}_c\|^2 + \|\mathbf{y}_c\|^2 - 2\|\mathbf{x}_c\| \|\mathbf{y}_c\| \rho(\mathbf{x}, \mathbf{y})$$

Consider the normed vectors

$$\mathbf{x}^* = \frac{\mathbf{x}_c}{\|\mathbf{x}_c\|}, \quad \mathbf{y}^* = \frac{\mathbf{y}_c}{\|\mathbf{y}_c\|}$$

Since the normed vectors are constant multiple of the original vectors, we have $\rho(\mathbf{x}, \mathbf{y}) = \rho(\mathbf{x}_c, \mathbf{y}_c) = \rho(\mathbf{x}^*, \mathbf{y}^*)$.

$$\|\mathbf{x}^* - \mathbf{y}^*\|^2 = \|\mathbf{x}^*\|^2 + \|\mathbf{y}^*\|^2 - 2\|\mathbf{x}^*\| \|\mathbf{y}^*\| \rho(\mathbf{x}^*, \mathbf{y}^*) = 2(1 - \rho(\mathbf{x}^*, \mathbf{y}^*))$$

Therefore

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{2(1 - \rho(\mathbf{x}, \mathbf{y}))} = \|\mathbf{x}^* - \mathbf{y}^*\|.$$

That is, the distance $d(\mathbf{x}, \mathbf{y})$ is the Euclidean distance of the centered, normed vectors \mathbf{x}^* and \mathbf{y}^* , and Euclidean norm has triangular property. We have

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x}^* - \mathbf{y}^*\| \leq \|\mathbf{x}^* - \mathbf{w}^*\| + \|\mathbf{w}^* - \mathbf{y}^*\| = d(\mathbf{x}, \mathbf{w}) + d(\mathbf{w}, \mathbf{y}).$$

Therefore the distance measure $d(\mathbf{x}, \mathbf{y})$ satisfies the triangular inequality, thus property 4 holds. \square

4 Classical metric scaling

How do we find a low dimension approximation of a given distance matrix?

Given coordinations of points, we can find their pairwise distances. Metric multidimensional scaling (including classical multidimensional scaling) deals with the inverse: from pairwise distance, can we find a set of coordinates for the points?

Metric scaling is an algebraic reconstruction from dissimilarity to a configuration of points.

Construct coordinates from pairwise distance

If we knew the coordinates of the r th observation point to be (x_{r1}, \dots, x_{rp}) , we can obtain the Euclidean distance (or any other distance) between the r th and s th observations as d_{rs} ,

$$d_{rs}^2 = \sum_{j=1}^p (x_{rj} - x_{sj})^2$$

Classical scaling aims to construct a set of coordinates $X = [x_{rs}]_{n \times p}$ from the knowledge of $D = [d_{rs}^2]$, only.

The construction considers an intermediate matrix

$$B = [b_{rs}]_{n \times n} = XX', \quad b_{rs} = \sum_{j=1}^p x_{rj}x_{sj}.$$

Then pairwise distance squares can be expressed in terms of entries of matrix B .

$$d_{rs}^2 = b_{rr} + b_{ss} - 2b_{rs}. \quad (5)$$

Define

$$d_{r.}^2 = \frac{1}{n} \sum_{s=1}^n d_{rs}^2, \quad d_{.s}^2 = \frac{1}{n} \sum_{r=1}^n d_{rs}^2, \quad d_{..}^2 = \frac{1}{n^2} \sum_{s=1}^n \sum_{r=1}^n d_{rs}^2,$$

which are the row, column and overall averages of the square-distance matrix $D = [d_{rs}^2]$ respectively.

To recover or to solve for a set of coordinates in \mathbb{R}^p , the solutions are not unique even when $q = p$. Now assume X has the constraint that column mean zero.

$$\sum_{k=1}^n x_{kj} = 0, \quad j = 1, \dots, n, \\ \sum_{r=1}^n b_{rs} = 0$$

Then

Given pairwise distances d_{ij} , sum over (5),

$$nd_{r.}^2 = T + nb_{ss}, \quad nd_{.s}^2 = T + nb_{ss}, \quad n^2 d_{..}^2 = 2nT$$

where $T = \text{trace}(B) = \sum_{j=1}^p b_{jj}$. B entries can be recovered by the relation

$$b_{rs} = -\frac{1}{2} (d_{rs}^2 - d_{r.}^2 - d_{.s}^2 + d_{..}^2)$$

By the symmetry of matrix B , the Spectral Theorem gives an eigenvalue-eigenvector decomposition of B .

$$B = V\Lambda V' = (V\Lambda^{1/2})(V\Lambda^{1/2})',$$

where

$$V = [e_1 \dots e_n]_{n \times n}$$

is the eigenvector matrix of B ,

$$\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_n\}$$

has B 's eigenvalues on the diagonal,

$$Be_i = \lambda_i e_i, \quad i = 1, \dots, n.$$

By the positive semi-definiteness of symmetric matrix B , for $k = \text{rank}(B) \leq n$, the non-zero eigenvalues of B can be ordered as

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0.$$

By construction,

$$X = V\Lambda^{1/2} = [\sqrt{\lambda_1}e_1 \dots \sqrt{\lambda_\ell}e_\ell \dots \sqrt{\lambda_p}e_p]_{n \times p}$$

provides a set of coordinates for the n items in \mathbb{R}^p , with pairwise Euclidean distance matrix D .

A "best" set of coordinates $\{(x_{r1}, \dots, x_{r\ell}), r = 1, \dots, n\} \in \mathbb{R}^\ell$ as a ℓ -dimensional representation of the n data points can be defined by

$$X = [\sqrt{\lambda_1}e_1 \dots \sqrt{\lambda_\ell}e_\ell]_{n \times \ell}$$

Pairwise distance, usually Euclidean distance, can be obtained and compared with the original distance, presumably from a space with higher dimensions $\geq p$.

In particular, a "best" two-dimensional map of the n points have the coordinates

$$X = [\sqrt{\lambda_1}e_1 \quad \sqrt{\lambda_2}e_2]_{n \times 2}$$

which are used in MDS plots. The above shows the process of finding a set of coordinates from pairwise distance:

$$d_{rs}^2 \implies b_{rs} \implies x_{rs}$$

Remarks

- The reconstruction of the coordinates assumes or treated the given pairwise distances as Euclidean distance in the p dimensional space. Therefore the method is particularly appropriate when the dissimilarities are actually or at least approximately Euclidean distances.
- The classical scaling method is by singular value decomposition, or equivalently, using principal components, therefore sometimes classical scaling is also called "principal coordinates analysis".

Note: Relevant section in Johnson and Wichern: section 12.6.