

# Smacof at 50: A Manual

## Part 1: Smacof Notation and Theory

Jan de Leeuw - University of California Los Angeles

Started February 21 2024, Version of April 16, 2024

### **Abstract**

TBD

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Kruskal's Stress</b>	<b>4</b>
2.1	Metric MDS . . . . .	4
2.2	Non-metric MDS . . . . .	5
2.3	Normalization . . . . .	5
2.4	Some thoughts on ALS . . . . .	6
2.4.1	The Single-Phase approach . . . . .	7
2.4.2	The Two-Phase Approach . . . . .	7
<b>3</b>	<b>Smacof Notation and Terminology</b>	<b>9</b>
<b>4</b>	<b>Intermezzo: Explicit Normalization</b>	<b>11</b>
<b>5</b>	<b>Properties of Smacof Loss</b>	<b>12</b>
5.1	Derivatives . . . . .	12
5.1.1	Gradient . . . . .	12
5.1.2	Hessian . . . . .	13
5.2	Lagrangian . . . . .	14
5.2.1	Kuhn-Tucker Points . . . . .	14
	<b>References</b>	<b>15</b>

**Note:** This is a working manuscript which will be expanded/updated frequently. All suggestions for improvement are welcome. All Rmd, tex, html, pdf, R, and C files are in the public domain. Attribution will be appreciated, but is not required. The files can be found at <https://github.com/deleeuw/smacofCode>.

# 1 Introduction

In *Multidimensional Scaling (MDS)* the data consists of information about the similarity or dissimilarity between pairs of objects selected from a finite set  $\mathcal{O} = \{o_1, \dots, o_n\}$ .

In *metric MDS* we have numerical dissimilarity measures and we want to map the objects  $o_i$  into  $n$  points  $x_i$  of some metric space in such a way that the distances between the points approximate the dissimilarities between the objects. In *smacof*, our framework for MDS theory, algorithms, and computer programs, the metric space is  $\mathbb{R}^p$ , the space of all  $p$ -tuples of real numbers, and in the code documented in this manual we assume the distance is the usual Euclidean distance.

In *non-metric MDS* the information about the dissimilarities is incomplete. It is usually *ordinal*, i.e. it tells us in some way or another that some dissimilarities are larger or smaller than others. Somewhere between metric and non-metric MDS is MDS with *missing data*, in which some dissimilarities are known numbers while others are unknown. MDS with missing data is a form of *distance matrix completion* (Fang and O’Leary (2012)).

## 2 Kruskal's Stress

In the pioneering papers Kruskal (1964a) and Kruskal (1964b) the MDS problem was formulated for the first time as minimization of an explicit loss function, which measures the quality of the approximation of the dissimilarities by the distances.

### 2.1 Metric MDS

The loss function in least squares metric Euclidean MDS is called *raw stress* or *Kruskal's raw stress* and is defined as

$$\sigma_R(X) := \frac{1}{2} \sum w_{ij} (\delta_{ij} - d_{ij}(X))^2. \quad (1)$$

The symbol  $:=$  is used for definitions. In definition (1) the  $w_{ij}$  are known non-negative *weights*, the  $\delta_{ij}$  are the known non-negative *dissimilarities* between objects  $o_i$  and  $o_j$ , and the  $d_{ij}(X)$  are the *distances* between the corresponding points  $x_i$  and  $x_j$ . The summation is over all  $\binom{n}{2}$  pairs  $(i, j)$  with  $j > i$ , i.e. over elements below the diagonal of the matrices  $W$  and  $\Delta$ . The subscript  $R$  in  $\sigma_R$  stands for “raw”. From now on we use “metric MDS” to mean Least Squares Metric Euclidean MDS.

The  $n \times p$  matrix  $X$ , which has the coordinates  $x_i$  of the  $n$  points as its rows, is called the *configuration*, where  $p$  is the *dimension* of the Euclidean space in which we make the map. Thus

$$d_{ij}(X) = \sqrt{\sum_{s=1}^p (x_{is} - x_{js})^2}. \quad (2)$$

The metric MDS problem (of dimension  $p$ , for given  $W$  and  $\Delta$ ) is the minimization of (1) over the  $n \times p$  configurations  $X$ .

The weights  $w_{ij}$  can be used to quantify information about the precision or importance of the corresponding dissimilarities. Some of the weights may be zero, which can be used to code *missing data*. If all weights are positive we have *complete data*. If we have complete data, and all weights are equal to one, we have *unweighted* metric MDS. The pioneering papers by Shepard, Kruskal, and Guttman only consider the unweighted case. Weights were only introduced in MDS in De Leeuw (1977).

We assume throughout that the weights are *irreducible* (De Leeuw (1977)). This means there is no partitioning of the index set  $I_n := \{1, 2, \dots, n\}$  into subsets for which all between-subset weights are zero. A reducible metric MDS problems decomposes into a number of smaller independent metric MDS problems, so the irreducibility assumption causes no real loss of generality.

The fact that the summation in (1) is over all  $j < i$  indicates that the diagonal elements of  $\Delta$  are not used (they are assumed to be zero) and the elements above the diagonal are not used as well (they are assumed to be equal to the corresponding elements below the diagonal). The somewhat mysterious factor  $\frac{1}{2}$  in definition (1) is there because it simplifies some of the formulas in later sections of this paper.

## 2.2 Non-metric MDS

Kruskal was not really interested in metric MDS and the “raw” loss function (1). His papers are really about non-metric MDS, by which we mean least squares non-metric Euclidean MDS. Non-metric MDS differs from metric MDS because we have incomplete information about the dissimilarities. As we have seen, that if some dissimilarities are missing metric MDS can handle this by using zero weights. In some situations, however, we only know the rank order of the non-missing dissimilarities. We do not know, or we refuse to use, their actual numeric values. Or, to put it differently, even if we have numerical dissimilarities we are looking for a *transformation* of the non-missing dissimilarities, where the transformation is chosen from a set of admissible transformations (for instance from all linear or monotone transformations). If the dissimilarities are non-numerical, for example rank orders or partitionings, we choose from the set of admissible *quantifications*.

In non-metric MDS raw stress becomes

$$\sigma_R(X, \Delta) := \frac{1}{2} \sum w_{ij} (\delta_{ij} - d_{ij}(X))^2, \quad (3)$$

where  $\Delta$  varies over the quantified or transformed dissimilarities. In MDS parlance they are also called *pseudo-distances* or *disparities*. Loss function (3) must be minimized over both configurations and disparities, with the condition that the disparities  $\Delta$  are an admissible transformation or quantification of the data. In Kruskal’s non-metric MDS this means requiring monotonicity. In this paper we will consider various other choices for the set of admissible transformations. We will use the symbol  $\mathfrak{D}$  for the set of admissible transformations

The most familiar examples of  $\mathfrak{D}$  (linear, polynomial, splines, monotone) define convex cones with apex at the origin. This means that if  $\Delta \in \mathfrak{D}$  then so is  $\lambda\Delta$  for all  $\lambda \geq 0$ . But consequently minimizing (3) over all  $\Delta \in \mathfrak{D}$  and over all configurations has the trivial solution  $\Delta = 0$  and  $X = 0$ , corresponding with the global minimum  $\sigma(X, \Delta) = 0$ . We need additional constraints to rule out this trivial solution, and in non-metric MDS this is done by choosing a *normalization* that keeps the solution away from zero.

Kruskal’s original solution is to define *normalized stress* as

$$\sigma(X, \Delta) := \frac{\sum w_{ij} (\delta_{ij} - d_{ij}(X))^2}{\sum w_{ij} d_{ij}^2(X)}. \quad (4)$$

To be precise, in Kruskal’s formulation there are no weights, and he actually takes the square root of (4) to define *Kruskal’s stress*. The non-metric Euclidean MDS problem now is to minimize loss function (4) over all  $n \times p$  configurations  $X$  and all admissible disparities  $\Delta$ .

## 2.3 Normalization

Equation (4) is only one way to normalize raw stress. Some obvious alternatives are discussed in detail in Kruskal and Carroll (1969) and De Leeuw (1975). In the terminology of De Leeuw (1975) there are *explicit* and *implicit* normalizations.

In implicit normalization we minimize either

$$\sigma(X, \hat{D}) := \frac{\sum w_{ij}(\hat{d}_{ij} - d_{ij}(X))^2}{\sum w_{ij}\hat{d}_{ij}^2} \quad (5)$$

or

$$\sigma(X, \hat{D}) := \frac{\sum w_{ij}(\hat{d}_{ij} - d_{ij}(X))^2}{\sum w_{ij}d_{ij}^2(X)} \quad (6)$$

over  $X$  and  $\Delta \in \mathfrak{D}$ .

As we have seen, Kruskal (1964a) chooses definition (6) and calls the explicitly normalized loss function *normalized stress*. Note that we overload the symbol  $\sigma$  to denote any one of the least squares loss functions. It will always be clear from the text which  $\sigma$  we are talking about.

In explicit normalization we minimize the raw stress  $\sigma_R(X, \hat{D})$  from (3), but we add the explicit constraint

$$\sum w_{ij}d_{ij}^2(X) = 1, \quad (7)$$

or the constraint

$$\sum w_{ij}\hat{d}_{ij}^2 = 1. \quad (8)$$

Kruskal and Carroll (1969) and De Leeuw (2019) show that these four normalizations all lead to essentially the same solution for  $X$  and  $\hat{D}$ , up to scale factors dictated by the choice of the particular normalization. It is also possible to normalize both  $X$  and  $\hat{D}$ , either explicitly or implicitly, and again this will give the same solutions, suitably normalized. These invariance results assume the admissible transformations form a closed cone with apex at the origin, i.e. if  $\hat{D}$  is admissible and  $\lambda \geq 0$  then  $\lambda\hat{D}$  is admissible as well. The matrices of Euclidean distances  $D(X)$  form a similar closed cone as well. The non-metric MDS problem is to find an element of the  $\hat{D}$  cone  $\mathcal{D}$  and an element of the  $D(X)$  cone where the angle between the two is as small as possible.

In the R version of smacof (De Leeuw and Mair (2009), Mair, Groenen, and De Leeuw (2022)) we use explicit normalization (8). This is supported by the result, also due to De Leeuw (1975), that projection on the intersection of the cone of disparities and the sphere defined by (8) is equivalent to first projecting on the cone and then normalizing the projection (see also Bauschke, Bui, and Wang (2018)).

In the version of non-metric MDS discussed in this manual we need more flexibility. For algorithmic reasons that may become clear later on, we will go with the original (4), i.e. with the implicitly normalized Kruskal's stress. For the final results the choice between normalizations should not make a difference, but the iterative computations will be different for the different choices.

## 2.4 Some thoughts on ALS

I will take this opportunity to clear up some misunderstandings and confusions that have haunted the early development of non-metric MDS.

### 2.4.1 The Single-Phase approach

In Kruskal (1964a) defines

$$\sigma(X) := \min_{\hat{D} \in \mathfrak{D}} \sigma(\hat{D}, X) = \sigma(X, \hat{D}(X)), \quad (9)$$

where  $\sigma(\hat{D}, X)$  is defined by (6). The minimum in (9) is over admissible transformations. In definition (9)

$$\hat{D}(X) := \operatorname{argmin}_{\hat{D} \in \mathfrak{D}} \sigma(X, \hat{D}). \quad (10)$$

Normalized stress defined by (9) is now a function of  $X$  only. Under some conditions, which are true in Kruskal's definition of non-metric MDS, there is a simple relation between the partials of (6) and those of (9).

$$\mathcal{D}\sigma(X) = \mathcal{D}_1\sigma(X, \hat{D}(X)), \quad (11)$$

where  $\mathcal{D}\sigma(X)$  are the derivatives of  $\sigma$  from (9) and  $\mathcal{D}_1\sigma(X, \hat{D}(X))$  are the partial derivatives of  $\sigma$  from (6) with respect to  $X$ . Thus the partials of  $\sigma$  from (9) can be computed by evaluating the partials of  $\sigma$  from (6) with respect to  $X$  at  $(X, \hat{D}(X))$ . This has created much confusion in the past. The non-metric MDS problem in Kruskal's original formulation is now to minimize  $\sigma$  from (9), which is a function of  $X$  alone.

Guttman (1968) calls this the *single-phase approach*. A variation of Kruskal's single-phase approach defines

$$\sigma(X) = \sum w_{ij} (d_{ij}^{\#}(X) - d_{ij}(X))^2, \quad (12)$$

where the  $d_{ij}^{\#}(X)$  are *Guttman's rank images*, i.e. the permutation of the  $d_{ij}(X)$  that makes them monotone with the  $\delta_{ij}$  (Guttman (1968)). Or, alternatively, define

$$\sigma(X) := \sum w_{ij} (d_{ij}^{\%}(X) - d_{ij}(X))^2, \quad (13)$$

where the  $\hat{d}_{ij}^{\%}(X)$  are *Shepard's rank images*, i.e. the permutation of the  $\delta_{ij}$  that makes them monotone with the  $d_{ij}(X)$  (Shepard (1962a), Shepard (1962b), De Leeuw (2017b)).

Minimizing the Shepard or Guttman single-phase loss functions is computationally more complicated than Kruskal's *monotone regression* approach, mostly because the rank-image transformations are not differentiable, and there is no analog of (11) and of the equivalence of the different implicit and explicit normalizations.

### 2.4.2 The Two-Phase Approach

The *two-phase approach* or *alternating least squares (ALS)* approach alternates minimization of  $\sigma(\hat{D}, X)$  over  $X$  for our current best estimate of  $\hat{D}$  with minimization of  $\sigma(\hat{D}, X)$  over  $\hat{D} \in \mathfrak{D}$  for our current best value of  $X$ . Thus an update from iteration  $k$  to iteration  $k + 1$  looks like

$$\hat{D}^{(k)} = \operatorname{argmin}_{\hat{D} \in \mathfrak{D}} \sigma(\hat{D}, X^{(k)}), \quad (14)$$

$$X^{(k+1)} = \operatorname{argmin}_X \sigma(\hat{D}^{(k)}, X). \quad (15)$$

This ALS approach to MDS was in the air since the early (unsuccessful) attempts around 1968 of Young and De Leeuw to combine Torgerson's classic metric MDS method with Kruskal's monotone regression transformation. All previous implementations of non-metric smacof use the two-phase approach, and we will do the same in this paper.

As formulated, however, there are some problems with the ALS algorithm. Step (14) is easy to carry out, using monotone regression. Step (15) means solving a metric scaling problem, which is an iterative process that requires an infinite number of iterations. Thus, in the usual implementations, step (14) is combined with one of more iterations of a convergent iterative procedure for metric MDS, such as smacof. If we take only one of these *inner iterations* the algorithm becomes indistinguishable from Kruskal's single-phase method. This has also created much confusion in the past.

In the usual implementations of the ALS approach we solve the first subproblem (14) exactly, while we take only a single step towards the solution for given  $\hat{D}$  in the second phase (15). If we have an infinite iterative procedure to compute the optimal  $\hat{D} \in \mathfrak{D}$  for given  $X$ , then a more balanced approach would be to take several inner iterations in the first phase and several inner iterations in the second phase. How many of each, nobody knows. In our current implementation of smacof we take several inner iteration steps in the first phase and a single inner iteration step in the second phase.



### 3 Smacof Notation and Terminology

We discuss some standard MDS notation, first introduced in De Leeuw (1977). This notation is useful for the second phase of the ALS algorithm, in which solve the metric MDS problem of we minimizing unnormalized  $\sigma(X, \hat{D})$  over  $X$  for fixed  $\hat{D}$ . We will discuss the first ALS phase later in the paper.

Start with the unit vectors  $e_i$  of length  $n$ . They have a non-zero element equal to one in position  $i$ , all other elements are zero. Think of the  $e_i$  as the columns of the identity matrix.

Using the  $e_i$  we define for all  $i \neq j$  the matrices

$$A_{ij} := (e_i - e_j)(e_i - e_j)'. \quad (16)$$

The  $A_{ij}$  are of order  $n$ , symmetric, doubly-centered, and of rank one. They have four non-zero elements. Elements  $(i, i)$  and  $(j, j)$  are equal to  $+1$ , elements  $(i, j)$  and  $(j, i)$  are  $-1$ .

The importance of  $A_{ij}$  in MDS comes from the equation

$$d_{ij}^2(X) = \text{tr } X' A_{ij} X. \quad (17)$$

In addition we use the fact that the  $A_{ij}$  form a basis for the  $\binom{n}{2}$ -dimensional linear space of all doubly-centered symmetric matrices.

Expanding the square in the definition of stress gives

$$\sigma(X) = \frac{1}{2} \left\{ \sum w_k \delta_k^2 - 2 \sum w_k \delta_k d_k(X) + \sum w_k d_k^2(X) \right\}. \quad (18)$$

It is convenient to have notation for the three separate components of stress from equation (18). Define

$$\eta_D^2 = \sum w_{ij} \hat{d}_{ij}^2, \quad (19)$$

$$\rho(X) = \sum w_{ij} \hat{d}_{ij} d_{ij}(X), \quad (20)$$

$$\eta^2(X) = \sum w_{ij} d_{ij}^2(X). \quad (21)$$

which lead to

$$\sigma(X) = \frac{1}{2} \left\{ \eta_D^2 - 2\rho(X) + \eta^2(X) \right\}. \quad (22)$$

We also need

$$\lambda(X) = \frac{\rho(X)}{\eta(X)}. \quad (23)$$

Using the  $A_{ij}$  makes it possible to give matrix expressions for  $\rho$  and  $\eta^2$ . First

$$\eta^2(X) = \text{tr } X' V X, \quad (24)$$

with

$$V := \sum w_{ij} A_{ij}. \quad (25)$$

In the same way

$$\rho(X) = \text{tr } X' B(X) X, \quad (26)$$

with

$$B(X) := \sum w_{ij} r_{ij}(X) A_{ij}, \quad (27)$$

with

$$r_{ij}(X) := \begin{cases} \frac{\delta_{ij}}{d_{ij}(X)} & \text{if } d_{ij}(X) > 0, \\ 0 & \text{if } d_{ij}(X) = 0. \end{cases} \quad (28)$$

Note that  $B$  is a function from the set of  $n \times p$  configurations into the set of symmetric doubly-centered matrices of order  $n$ . All matrices of the form  $\sum x_{ij} A_{ij}$ , where summation is over all pairs  $(i, j)$  with  $j < i$ , are symmetric and doubly-centered. They have  $-x_{ij}$  as off-diagonal elements while the diagonal elements  $(i, i)$  are  $\sum_{j=1}^n x_{ij}$ .

Because  $B(X)$  and  $V$  are non-negative linear combinations of the  $A_{ij}$  they are both positive semi-definite. Because  $W$  is assumed to be irreducible the matrix  $V$  has rank  $n - 1$ , with only vectors proportional to the vector  $e$  with all elements equal to one in its null-space (De Leeuw (1977)).

Summarizing the results so far we have

$$\sigma(X) = \frac{1}{2} \{ \eta_D^2 - \text{tr } X' B(X) X + \text{tr } X' V X \}. \quad (29)$$

Next we define the *Guttman transform* of a configuration  $X$ , for given  $W$  and  $\Delta$ , as

$$G(X) = V^+ B(X) X, \quad (30)$$

with  $V^+$  the Moore-Penrose inverse of  $V$ . In our computations we use

$$V^+ = (V + \frac{1}{n} ee')^{-1} - \frac{1}{n} ee' \quad (31)$$

Also note that in the unweighted case with complete data  $V = nJ$ , where  $J$  is the centering matrix  $I - \frac{1}{n} ee'$ , and thus  $V^+ = \frac{1}{n} J$ . The Guttman transform is then simply  $G(X) = n^{-1} B(X) X$ .

## 4 Intermezzo: Explicit Normalization

$$\sigma(X, \hat{D}) = \frac{1}{2} \frac{\sum w_{ij}(\hat{d}_{ij} - d_{ij}(X))^2}{\sum w_{ij}d_{ij}^2(X)}$$

Majorize

$$\sigma(X, \hat{D}) \leq \frac{1}{2} \frac{\eta^2(\hat{D}) - 2\text{tr } X'V\bar{Y} + \text{tr } X'VX}{\text{tr } X'VX} = \frac{\omega(X, Y)}{\eta^2(X)}$$

Stationary equations

$$\eta^2(X)(VX - VG(Y)) - \omega(X, Y)VX = V\{(\eta^2(X) - \omega(X, Y))X - \eta^2(X)\bar{Y}\}$$

So at a minimum  $X$  is proportional to  $\bar{Y}$  or  $X = \alpha\bar{Y}$  for some  $\alpha$ . For ... to be zero we must have

$$\alpha(\alpha^2\eta^2(\bar{Y}) - (\eta^2(\hat{D}) - 2\alpha\eta^2(\bar{Y}) + \alpha^2\eta^2(\bar{Y})) = \alpha^2\eta^2(\bar{Y})$$

which works out to be

$$\alpha = \frac{\eta^2(\hat{D})}{\eta^2(\bar{Y})}$$

$$\hat{X} = \frac{\eta^2(\hat{D})}{\eta^2(\bar{Y})} \bar{Y}$$

The minimum is equal to

$$\frac{-\frac{(\eta^2(\bar{Y}))^2}{\eta^2(\hat{D})} + \eta^2(\bar{Y})}{\eta^2(\bar{Y})} = 1 - \frac{\eta^2(\bar{Y})}{\eta^2(\hat{D})}$$

Use homogeneity of the Guttman transform.

More generally suppose we update with

$$X = \bar{Y} + \alpha(Y - \bar{Y})$$

Write

$$\omega(X, Y) = \eta^2(\hat{D}) + \text{tr } (X - \bar{Y})'V(X - \bar{Y}) - \eta^2(\bar{Y})$$

Thus if  $X(\alpha) = \bar{Y} + \alpha(Y - \bar{Y})$  we have

$$\omega(\alpha) = \eta^2(\hat{D}) + \alpha^2\text{tr } (Y - \bar{Y})'V(Y - \bar{Y}) - \eta^2(\bar{Y})$$

and

$$\eta^2(\alpha) = \eta^2(\bar{Y}) + 2\alpha\text{tr } (Y - \bar{Y})'V\bar{Y} + \alpha^2\text{tr } (Y - \bar{Y})'V(Y - \bar{Y})$$

$$\omega(Y, Y) = \eta^2(\hat{D}) + \text{tr } (Y - \bar{Y})'V(Y - \bar{Y}) - \eta^2(\bar{Y})$$

$$\frac{\omega(\alpha)}{\eta^2(\alpha)} \leq \sigma(Y)$$

## 5 Properties of Smacof Loss

### 5.1 Derivatives

The Euclidean distance function  $d_{ij}$  from ... is not differentiable at configurations  $X$  with  $x_i = x_j$ . If  $d_{ij}(X) > 0$  then

$$\mathcal{D}\sigma(X) = \frac{1}{d_{ij}(X)} A_{ij} X$$

If  $d_{ij}(X) = 0$  then

$$D_+ d_{ij}(X, Y) = \lim_{\epsilon \downarrow 0} \frac{d_{ij}(X + \epsilon Y) - d_{ij}(X)}{\epsilon} = d_{ij}(Y)$$

which is non-linear in  $Y$ , showing non-differentiability.

$$D_+ \sigma(X, Y) = \text{tr } Y'(V - B(X))X + \sum \{w_{ij} \delta_{ij} d_{ij}(Y) \mid d_{ij}(X) = 0\}$$

This form of the directional derivative is used by De Leeuw (1984) to show that two independent necessary conditions for a local minimum are  $(V - B(X))X = 0$  and  $d_{ij}(X) > 0$  for all  $(i, j)$  with  $w_{ij} \delta_{ij} > 0$ .

#### 5.1.1 Gradient

$$\mathcal{D}\sigma(X) = (V - B(X))X$$

At a stationary point  $B(X)X = VX$  or  $V^+ B(X)X = X$ . Thus a necessary condition for a local minimum is that  $V^+ B(X)$  has at least  $p$  eigenvalues equal to one. De Leeuw (2014) has shown that if  $V^+ B(X) \lesssim I$  then actually  $X$  is a global minimizer of stress.

$$\rho(X) = \sum w_{ij} \delta_{ij}(X)$$

$$\nabla d_{ij}(X) = \begin{bmatrix} 0 \\ \frac{x_i - x_j}{d_{ij}(X)} \\ 0 \\ -\frac{x_i - x_j}{d_{ij}(X)} \\ 0 \end{bmatrix}$$

$$\partial d_{ij}(X) = \left\{ \begin{bmatrix} 0 \\ y \\ 0 \\ -y \\ 0 \end{bmatrix} \mid y' y \leq 1 \right\}.$$

### 5.1.2 Hessian

The results on the Hessian of stress are largely unpublished. So we summarize them here in this manual, so they'll be even more unpublished.

$$H_{st}(X) := \sum w_{ij} \frac{\delta_{ij}}{d_{ij}(X)} \left\{ \frac{A_{ij} x_s x'_t A_{ij}}{d_{ij}^2(X)} \right\}$$

$$H_{st}(X) = \sum w_{ij} \frac{\delta_{ij}}{d_{ij}^3(X)} (x_{is} - x_{js})(x_{it} - x_{jt}) A_{ij}$$

$$\mathcal{D}_{st}\sigma(X) = \begin{cases} H_{st}(X) & \text{if } s \neq t, \\ V - B(X) + H_{st} & \text{if } s = t. \end{cases}$$

If  $I_p$  is the identity matrix of order  $p$ , and  $\otimes$  is the Kronecker product, then

$$\mathcal{D}^2\sigma(X) = I_p \otimes (V - B(X)) + H(X)$$

$$\sum_{s=1}^p \sum_{t=1}^p y'_s H_{st} y_t = \sum w_{ij} \frac{\delta_{ij}}{d_{ij}(X)} \left\{ \frac{(\text{tr } Y' A_{ij} X)^2}{d_{ij}^2(X)} \right\} \leq \sum w_{ij} \frac{\delta_{ij}}{d_{ij}(X)} \text{tr } Y' A_{ij} Y = \text{tr } Y' B(X) Y.$$

Thus

$$0 \preceq H \preceq I_p \otimes B(X),$$

and

$$I_p \otimes (V - B(X)) \preceq \mathcal{D}^2\sigma(X) \preceq I_p \otimes V$$

At a local minimum of  $\sigma$

$$0 \preceq \mathcal{D}^2\sigma(X) \preceq I_p \otimes V$$

In comparing the lower bounds on  $\mathcal{D}^2\sigma(X)$  in ... and ... De Leeuw (2014) shows that  $V - B(X) \succeq 0$  is sufficient for a *global* minimum of stress (but far from necessary).

Also

$$\sum_{t=1}^p H_{st} y_t = \sum w_{ij} \frac{\delta_{ij}}{d_{ij}(X)} \left\{ \frac{\text{tr } Y' A_{ij} X}{d_{ij}^2(X)} \right\} A_{ij} x_s$$

If  $Y = X$  then  $H(X)y = (I_p \otimes B(X))x$  and thus

$$\mathcal{D}^2\sigma(X)x = (I_p \otimes V)x.$$

In the unweighted case this means that  $X$  is an eigenvector of  $\mathcal{D}^2\sigma(X)$  with eigenvalue  $n$ . Inequalities ... show that this is actually the largest eigenvalue. Or  $(I_p \otimes V)^+ \mathcal{D}^2\sigma(X) \preceq I$ .

If  $Y = XT$  with  $T$  anti-symmetric then  $\text{tr } Y' A_{ij} X = 0$  then thus  $H(X)y = 0$ . Thus

$$\sum_{t=1}^p \mathcal{D}_{st}\sigma(X)y_t = (V - B(X))y_t$$

which is zero if  $\mathcal{D}\sigma(X)$  is zero. Thus at a stationary point of stress  $\mathcal{D}^\sigma(X)$  has  $\frac{1}{2}p(p-1)$  zero eigenvalues.

There are several ways to think of the Hessian. The simplest one (perhaps) is as an  $np \times np$  symmetric matrix (corresponding to

column-major R vector of length  $\frac{1}{2}np(np+1)$ ). This is what we would use for a straightforward version of Newton-Raphson.

It is more elegant, however, to think of  $H$  as a symmetric super-matrix of order  $p$ , with as elements  $n \times n$  matrices. And, for some purposes, such as the pseudo-confidence ellipsoids in De Leeuw (2017a), as a super-matrix of order  $n$  with as elements  $p \times p$  matrices. Both the super-matrix interpretations lead to four-dimensional arrays, the first a  $p \times p \times n \times n$  array, the second an  $n \times n \times p \times p$  array. The different interpretations lead to different ways to store the Hessian in memory, and to different ways to retrieve its elements. Of course we can write routines to transform from one interpretation to another.

## 5.2 Lagrangian

In our implementation of the smacof algorithm we minimize stress over configurations with  $\eta(X) = 1$ , or, equivalently,  $\sum w_{ij}d_{ij}^2(X) = 1$ . This means we do not look for  $X$  with  $\mathcal{D}\sigma(X) = (V - B(X))X = 0$ , but we look for solutions of

$$(V - B(X))X - \lambda VX = 0, \text{tr } X'VX = 1.$$

At the solution

$$\lambda = 1 - \rho(X)$$

and

$$X = \frac{\Gamma(X)}{\eta(\Gamma(X))}$$

Also it is necessary for a local minimum that

$$\Gamma(X) = \rho(X)X$$

Because the Guttman transform is homogeneous of degree zero this implies

$$\Gamma(\Gamma(X)) = \Gamma(X),$$

so although  $X$  is not a fixed point of the Guttman transform,  $\Gamma(X)$  is.

The second order necessary condition is that

$$H(X) \succeq I_p \otimes (\rho(X)V - B(X))$$

is positive

### 5.2.1 Kuhn-Tucker Points

bozo

## References

- Bauschke, H. H., M. N. Bui, and X. Wang. 2018. “Projecting onto the Intersection of a Cone and a Sphere.” *SIAM Journal on Optimization* 28: 2158–88.
- De Leeuw, J. 1975. “A Normalized Cone Regression Approach to Alternating Least Squares Algorithms.” Department of Data Theory FSW/RUL.
- . 1977. “Applications of Convex Analysis to Multidimensional Scaling.” In *Recent Developments in Statistics*, edited by J. R. Barra, F. Brodeau, G. Romier, and B. Van Cutsem, 133–45. Amsterdam, The Netherlands: North Holland Publishing Company.
- . 1984. “Differentiability of Kruskal’s Stress at a Local Minimum.” *Psychometrika* 49: 111–13.
- . 2014. “Bounding, and Sometimes Finding, the Global Minimum in Multidimensional Scaling.” UCLA Department of Statistics. <https://jansweb.netlify.app/publication/deleeuw-u-14-b/deleeuw-u-14-b.pdf>.
- . 2017a. “Pseudo Confidence Regions for MDS.” 2017. <https://jansweb.netlify.app/publication/deleeuw-e-17-q/deleeuw-e-17-q.pdf>.
- . 2017b. “Shepard Non-metric Multidimensional Scaling.” 2017. <https://jansweb.netlify.app/publication/deleeuw-e-17-e/deleeuw-e-17-e.pdf>.
- . 2019. “Normalized Cone Regression.” 2019. <https://jansweb.netlify.app/publication/deleeuw-e-19-d/deleeuw-e-19-d.pdf>.
- De Leeuw, J., and P. Mair. 2009. “Multidimensional Scaling Using Majorization: SMACOF in R.” *Journal of Statistical Software* 31 (3): 1–30. <https://www.jstatsoft.org/article/view/v031i03>.
- Fang, H., and D. P. O’Leary. 2012. “Euclidean Distance Matrix Completion Problems Euclidean Distance Matrix Completion Problems Euclidean Distance Matrix Completion Problems.” *Optimization Methods and Software* 27 (4-5): 695–717.
- Guttman, L. 1968. “A General Nonmetric Technique for Fitting the Smallest Coordinate Space for a Configuration of Points.” *Psychometrika* 33: 469–506.
- Kruskal, J. B. 1964a. “Multidimensional Scaling by Optimizing Goodness of Fit to a Nonmetric Hypothesis.” *Psychometrika* 29: 1–27.
- . 1964b. “Nonmetric Multidimensional Scaling: a Numerical Method.” *Psychometrika* 29: 115–29.
- Kruskal, J. B., and J. D. Carroll. 1969. “Geometrical Models and Badness of Fit Functions.” In *Multivariate Analysis, Volume II*, edited by P. R. Krishnaiah, 639–71. North Holland Publishing Company.
- Mair, P., P. J. F. Groenen, and J. De Leeuw. 2022. “More on Multidimensional Scaling in R: smacof Version 2.” *Journal of Statistical Software* 102 (10): 1–47. <https://www.jstatsoft.org/article/view/v102i10>.
- Shepard, R. N. 1962a. “The Analysis of Proximities: Multidimensional Scaling with an Unknown Distance Function. I.” *Psychometrika* 27: 125–40.
- . 1962b. “The Analysis of Proximities: Multidimensional Scaling with an Unknown Distance Function. II.” *Psychometrika* 27: 219–46.