

Majorizing Stress Formula Two

Jan de Leeuw - University of California Los Angeles

Started April 19 2024, Version of April 26, 2024

Abstract

Modifications of the smacof algorithm for multidimensional scaling are proposed that provide a convergent majorization algorithm for Kruskal's stress formula two.

Contents

1	Introduction	3
2	Problem	4
3	Notation	5
4	Majorization	6
5	Derivatives	9
6	Examples	10
6.1	Ekman	10
6.2	De Gruijter	12
7	Appendix: Code	14
7.1	stress2.R	14
	References	16

Note: This is a working paper which will be expanded/updated frequently. All suggestions for improvement are welcome.

1 Introduction

The loss function minimized in the current non-metric and non-linear R implementations of the smacof programs for MDS (De Leeuw and Mair (2009), Mair, Groenen, and De Leeuw (2022)) is Kruskal’s original *normalized stress* (Kruskal (1964a), Kruskal (1964b)). It is defined as

$$\sigma_1(X) := \frac{\sum \sum w_{ij}(\delta_{ij} - d_{ij}(X))^2}{\sum \sum w_{ij}d_{ij}^2(X)}. \quad (1)$$

In equation (1) we assume throughout that dissimilarities δ_{ij} and weights w_{ij} are non-negative, and, without loss of generality, that the weights add up to one. The double summation is over all pairs of indices (i, j) with $i > j$, i.e, over the elements below the diagonal of the matrices Δ , W , and $D(X)$.

In Kruskal (1965) a different loss function was used in the context of using monotone transformations when fitting a linear model. In MDS this loss function is

$$\sigma_2(X) := \frac{\sum \sum w_{ij}(\delta_{ij} - d_{ij}(X))^2}{\sum \sum w_{ij}(d_{ij}(X) - \bar{d}(X))^2}, \quad (2)$$

where

$$\bar{d}(X) = \sum \sum w_{ij}d_{ij}(X). \quad (3)$$

In Kruskal and Carroll (1969), in the section written by Kruskal (p. 652), we see

In several of my scaling programs, I refer to these expressions as “stress formula one” and “stress formula two”, respectively. Historically, stress formula one was the only badness-of-fit function used for some time. Stress formula two has been used more recently and I now tend to recommend it.

Another early adopter (Roskam (1968), p. 34) says

While the original formula is adequate for completely ordered B-data, we found it is not adequate with completely ordered A-data.

The distinction between A-data and B-data comes from Coombs (1964). For B-data the δ_{ij} are dissimilarities between pairs of elements of a single set, while for A-data they are dissimilarities between two different sets, a row-set and a column-set. Moreover both Kruskal and Roskam found that having the variance of the distances in the denominator of stress has major advantages for conditional A-data, in which only comparisons of dissimilarities with in the same row are meaningful.

In this paper we will extend the theory and algorithm of smacof to stress formula two. We emphasize that normalized loss functions such as σ_1 and σ_2 should be only used in non-linear or non-metric MDS problems. In metric MDS problems raw stress, without any normalization, can be used.

2 Problem

We want to minimize Kruskal's σ_2 from (2) over the $n \times p$ configuration matrices X .

It is convenient to have some notation for the numerator and denominator of the two stress formulas.

$$\sigma_R(X) := \sum \sum w_{ij} (\delta_{ij} - d_{ij}(X))^2, \quad (4a)$$

$$\eta_1^2(X) := \sum \sum w_{ij} d_{ij}^2(X), \quad (4b)$$

$$\eta_2^2(X) := \sum \sum w_{ij} (d_{ij}(X) - \bar{d}(X))^2, \quad (4c)$$

Kruskal terms σ_R from definition (4a) the *raw stress*.

There have not been any systematic comparisons of the two stress formulas, and the solutions they lead to, that I am aware of. Kruskal (in Kruskal and Carroll (1969), p. 652) says

For any given configuration, of course, stress formula two yields a substantially larger value than stress formula one, perhaps twice as large in many cases. However, in typical multidimensional scaling applications, minimizing stress formula two typically yields very similar configurations to minimizing stress formula one.

We can get some idea about the difference in scale of the two loss functions from the results

$$\frac{\sigma_1(X)}{\sigma_2(X)} = \frac{\eta_2^2(X)}{\eta_1^2(X)} \geq \min_X \frac{\eta_2^2(X)}{\eta_1^2(X)} \quad (5)$$

De Leeuw and Stoop (1984) show that in the one-dimensional case with $p = 1$ and with all w_{ij} equal, this implies

$$\sigma_1(X) \geq \frac{1}{3} \frac{n-2}{n} \sigma_2(X). \quad (6)$$

Thus in this special case σ_1 is three to nine times as large as σ_2 . In general the bound in equation (6) depends on the weights, on the dimensionality p , and on the order n of the problem.

As a qualitative statement, supported to some extent by the computations of De Leeuw and Stoop (1984), we can say that minimizing σ_1 will tend to give optimal configurations in which distances have less variance than those in configurations that minimize σ_2 . One thing is for sure, however. If X is a regular simplex in $n - 1$ dimensions then σ_2 is not even defined. Or, to put it differently, if all δ_{ij} are equal the minimum of σ_2 in $n - 1$ dimensions does not exist.

3 Notation

Now for some notation. As in standard MDS theory (De Leeuw (1977), De Leeuw and Heiser (1977), De Leeuw (1988)) we use the matrices

$$A_{ij} := (e_i - e_j)(e_i - e_j)', \quad (7)$$

where e_i are unit vectors with element i equal to one and the other $n - 1$ elements equal to zero. Thus A_{ij} has elements (i, i) and (j, j) equal to $+1$, elements (i, j) and (j, i) equal to -1 , and all other elements equal to zero. The usefulness of the A_{ij} in MDS derives mainly from the formula

$$d_{ij}^2(X) = \text{tr } X' A_{ij} X. \quad (8)$$

Using the A_{ij} we now define other matrices, also standard in MDS,

$$V := \sum \sum w_{ij} A_{ij}, \quad (9a)$$

$$B(X) := \sum \sum w_{ij} \frac{\delta_{ij}}{d_{ij}^2(X)} A_{ij}. \quad (9b)$$

Note that B is a matrix-valued function, not a single matrix. For completeness also define

$$\eta^2(\Delta) := \sum \sum w_{ij} \delta_{ij}^2. \quad (10)$$

Specifically because we are dealing with σ_2 we also need the non-standard definition

$$M(X) := \bar{d}(X) \sum \sum \frac{w_{ij}}{d_{ij}^2(X)} A_{ij}. \quad (11)$$

In both definitions (9b) and (11) the summation is over pairs (i, j) with $d_{ij}^2(X) > 0$. Of course we can also omit all pairs from the summation for which $w_{ij} = 0$.

4 Majorization

In this section we construct a convergent majorization algorithm (De Leeuw (1994)) (also known as an MM algorithm, Lange (2016)) to minimize σ_2 .

The first step is to turn the minimization of a ratio of two functions into the iterative minimization of a difference of the two functions. This is a classical trick in fractional programming, usually attributed to Dinkelbach (1967). Define

$$\omega(X, Y) := \sum \sum w_{ij}(\delta_{ij} - d_{ij}(X))^2 - \sigma(Y) \{ \sum \sum w_{ij}(d_{ij}(X) - \bar{d}(X))^2 \} \quad (12)$$

Lemma 4.1. *If $\omega(X, Y) < \omega(Y, Y) = 0$ then $\sigma(X) < \sigma(Y)$.*

Proof. This is embarassingly simple. Direct substitution shows $\omega(X, X) = 0$ for all X . Also $\omega(X, Y) < 0$ if and only if

$$\sum \sum w_{ij}(\delta_{ij} - d_{ij}(X))^2 < \sigma(Y) \{ \sum \sum w_{ij}(d_{ij}(X) - \bar{d}(X))^2 \} \quad (13)$$

Dividing both sides by $\{ \sum \sum w_{ij}(d_{ij}(X) - \bar{d}(X))^2 \}$ shows that $\sigma(X) < \sigma(Y)$. \square

It follows from lemma 4.1 that if we are in iteration k , with tentative solution $X^{(k)}$, then finding any $X^{(k+1)}$ such that $\omega(X^{(k+1)}, X^{(k)}) < 0$ will decrease stress. We will accomplish this in our algorithm by performing one or more majorization steps decreasing $\omega(X, X^{(k)})$.

We should note that as a general strategy we cannot use finding $X^{(k+1)}$ by minimizing $\omega(X, X^{(k)})$ over X . If the minimum exists this will work, but in general $\omega(\bullet, X^{(k)})$ may be unbounded below, and the minimum may not exist. This is easily seen from the example $f(x) = x'Ax/x'x$ for which Dinkelbach's maneuver gives $g(x, y) = x'Ax - f(y)x'x$. The minimum of g over x is zero if $f(y)$ is equal to $\lambda_{\min}(A)$, the smallest eigenvalue of A , which is actually the minimum of f . If $f(y) > \lambda_{\min}(A)$ the minimum does not exist (the infimum is $-\infty$). We can ignore the case $f(y) < \lambda_{\min}(A)$ because that is impossible. But if $f(y) > \lambda_{\min}(A)$ any x with $x'x = 1$ other than the eigenvector corresponding with the minimum eigenvalue satisfies $g(x, y) < 0$ and thus $f(x) < f(y)$.

Back to σ_2 . From definitions (9a), (9b), (10), and (11)

$$\omega(X, Y) = \eta^2(\Delta) + (1 - \sigma(Y)) \text{tr } X'VX - 2 \text{tr } X'B(X)X + \text{tr } X'M(X)X \quad (14)$$

Lemma 4.2. *For all X and Y*

$$\text{tr } X'B(X)X \geq \text{tr } X'B(Y)Y, \quad (15)$$

with equality if $X = Y$.

Proof. By Cauchy-Schwartz

$$d_{ij}(X) \geq \frac{1}{d_{ij}(Y)} \text{tr } X'A_{ij}Y \quad (16)$$

Multiplying both sides by $w_{ij}\delta_{ij}$ and summing proves the lemma. \square

Lemma 4.3. *For all X and Y*

$$\text{tr } X' M(X) X \leq \text{tr } X' M(Y) X, \quad (17)$$

with equality if $X = Y$.

Proof. Start with the trivial result

$$\sum \sum w_{ij} d_{ij}(X) = \sum \sum \frac{w_{ij}}{d_{ij}(Y)} d_{ij}(X) d_{ij}(Y). \quad (18)$$

By Cauchy-Schwartz

$$\bar{d}(X) \leq \sqrt{\sum \sum \frac{w_{ij}}{d_{ij}(Y)} d_{ij}^2(X)} \sqrt{\sum \sum \frac{w_{ij}}{d_{ij}(Y)} d_{ij}^2(Y)} \quad (19)$$

Squaring both sides proves the lemma. \square

We are now ready for the main result.

Theorem 4.1. *Suppose $\sigma_2(X^{(0)}) \leq 1$. The update*

$$X^{(k+1)} = \{(1 - \sigma_2(X^{(k)}))V + \sigma_2(X^{(k)})M(X^{(k)})\}^+ B(X^{(k)})X^{(k)} \quad (20)$$

defines a convergent majorization algorithm.

Proof. Using the definitions in equations (9a), (9b), (10), and (11) define

$$\xi(X, Y) := \eta^2(\Delta) + (1 - \sigma(Y))\text{tr } X' V X - 2\text{tr } X' B(Y)Y + \sigma(Y)\text{tr } X' M(Y)X. \quad (21)$$

From lemmas 4.2 and 4.3 $\omega(X, Y) \leq \xi(X, Y)$ with equality if $X = Y$. In particular

$$\omega(X^{(k+1)}, X^{(k)}) \leq \xi(X^{(k+1)}, X^{(k)}). \quad (22a)$$

The update $X^{(k+1)}$ minimizes $\xi(X, X^{(k)})$ and thus

$$\xi(X^{(k+1)}, X^{(k)}) \leq \xi(X^{(k)}, X^{(k)}) = \omega(X^{(k)}, X^{(k)}). \quad (22b)$$

Combining equations (22a) and (22b), and using lemma 4.1, shows that also $\sigma_2(X^{(k+1)}) \leq \sigma_2(X^{(k)})$. \square

In order to make our proof work we had to guarantee that for all k

$$(1 - \sigma_2(X^{(k)}))V + \sigma_2(X^{(k)})M(X^{(k)}) \succeq 0, \quad (23)$$

because otherwise the minimum of $\xi(\bullet, X^{(k)})$ does not exist. If $\sigma_2(X^{(k)}) \leq 1$ the matrix in inequality (23) is a convex combination of two positive semi-definite matrices, and is thus positive semi-definite. And because of theorem 4.1 it is sufficient to assume that $\sigma_2(X^{(0)}) \leq 1$, because

subsequent $X^{(k)}$ will have σ_2 values smaller than the value for $X^{(0)}$. Thus we need to start our majorization algorithm with a sufficiently good initial estimate of X . A random start may not work.

From the practical point of view the condition $\sigma_2(X^{(0)}) \leq 1$ is not really restrictive. As the introduction of this paper says, in metric MDS we do not use σ_2 . But even in metric MDS the Torgerson initial estimate usually takes σ_2 well below one. In non-linear or non-metric scaling the δ_{ij} are optimal transformations or quantifications. If the optimum transformation is better than the optimal constant transformation the condition $\sigma_2(X^{(k)}) \leq 1$ is automatically satisfied for all k . And even if the optimum transformation is the constant transformation we still have $\sigma_2(X^{(k)}) = 1$ and inequality (23) is satisfied.

If for some reason you want to proceed if the matrix in (23) is not positive semi-definite, then it suffices to choose any Y with

$$\text{tr } Y' \{((1 - \sigma_2(X^{(k)}))V + \sigma_2(X^{(k)})M(X^{(k)}))\} Y \geq 0 \quad (24)$$

and to minimize $\xi(X^{(k)} + \alpha Y, X^{(k)})$ over α .

5 Derivatives

The derivatives of σ_2 are

$$\mathcal{D}\sigma_2(X) = \frac{\mathcal{D}\sigma_R(X) - \sigma_2(X)\mathcal{D}\eta_2^2(X)}{\eta_2^2(X)} \quad (25)$$

Now

$$\mathcal{D}\sigma_R(X) = -2 \sum \sum w_{ij}(\delta_{ij} - d_{ij}(X))\mathcal{D}d_{ij}(X), \quad (26a)$$

$$\mathcal{D}\eta_2^2(X) = 2 \sum \sum w_{ij}\mathcal{D}d_{ij}^2(X) - 2\bar{d}(X) \sum \sum w_{ij}\mathcal{D}d_{ij}(X), \quad (26b)$$

and

$$\mathcal{D}d_{ij}(X) = \frac{1}{d_{ij}(X)}A_{ij}X. \quad (27)$$

And thus, using definitions (9a), (9b), and (11)

$$\mathcal{D}\sigma_R(X) = 2(V - B(X))X, \quad (28a)$$

$$\mathcal{D}\eta_2^2(X) = 2(V - M(X))X. \quad (28b)$$

It follows that $\mathcal{D}\sigma_2(X) = 0$ if and only if

$$X = \{(1 - \sigma_2(X))V + \sigma_2(X)M(X)\}^+ B(X)X. \quad (29)$$

We can summarize the results of our computations in this section.

Theorem 5.1. *X is a fixed point of the majorization iterations (20) if and only if $\mathcal{D}\sigma_2(X) = 0$.*

6 Examples

Although we mentioned in the introduction that it is unusual to use σ_2 in metric MDS problems we will nevertheless give some metric examples to illustrate the algorithm. In both examples we start with the Torgerson initial solution which takes the initial σ_2 way below one.

6.1 Ekman

Our first example are the obligatory color data from Ekman (1954). The stress2 program produces the following sequence of σ_2 values and converges in 28 iterations.

```
## itel  1 sold  0.1577255150 snw  0.1321216983
## itel  2 sold  0.1321216983 snw  0.1207395499
## itel  3 sold  0.1207395499 snw  0.1156260670
## itel  4 sold  0.1156260670 snw  0.1135043532
## itel  5 sold  0.1135043532 snw  0.1126543441
## itel  6 sold  0.1126543441 snw  0.1123159800
## itel  7 sold  0.1123159800 snw  0.1121798388
## itel  8 sold  0.1121798388 snw  0.1121239038
## itel  9 sold  0.1121239038 snw  0.1121002964
## itel 10 sold  0.1121002964 snw  0.1120900307
## itel 11 sold  0.1120900307 snw  0.1120854276
## itel 12 sold  0.1120854276 snw  0.1120833009
## itel 13 sold  0.1120833009 snw  0.1120822904
## itel 14 sold  0.1120822904 snw  0.1120817979
## itel 15 sold  0.1120817979 snw  0.1120815523
## itel 16 sold  0.1120815523 snw  0.1120814273
## itel 17 sold  0.1120814273 snw  0.1120813627
## itel 18 sold  0.1120813627 snw  0.1120813287
## itel 19 sold  0.1120813287 snw  0.1120813107
## itel 20 sold  0.1120813107 snw  0.1120813010
## itel 21 sold  0.1120813010 snw  0.1120812957
## itel 22 sold  0.1120812957 snw  0.1120812929
## itel 23 sold  0.1120812929 snw  0.1120812913
## itel 24 sold  0.1120812913 snw  0.1120812904
## itel 25 sold  0.1120812904 snw  0.1120812899
## itel 26 sold  0.1120812899 snw  0.1120812897
## itel 27 sold  0.1120812897 snw  0.1120812895
## itel 28 sold  0.1120812895 snw  0.1120812894
```

The optimum configuration is in figure 1, which can be compared with the solution minimizing raw stress (which is identical up to a scale factor with the solution minimizing stress formula one) in figure 2. The raw stress solution reaches stress formula one equal to 0.5278528 in 32 iterations. The two optimal configurations are virtually identical.

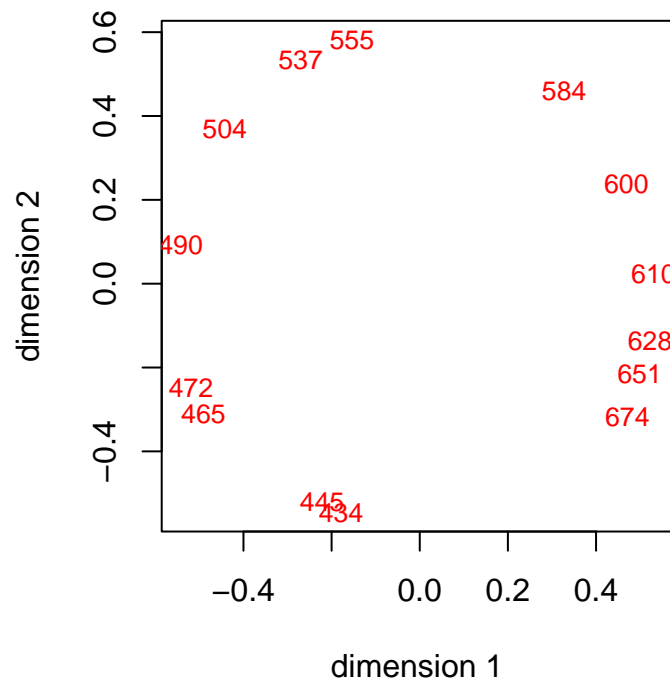


Figure 1: Ekman Metric Stress 2 Solution

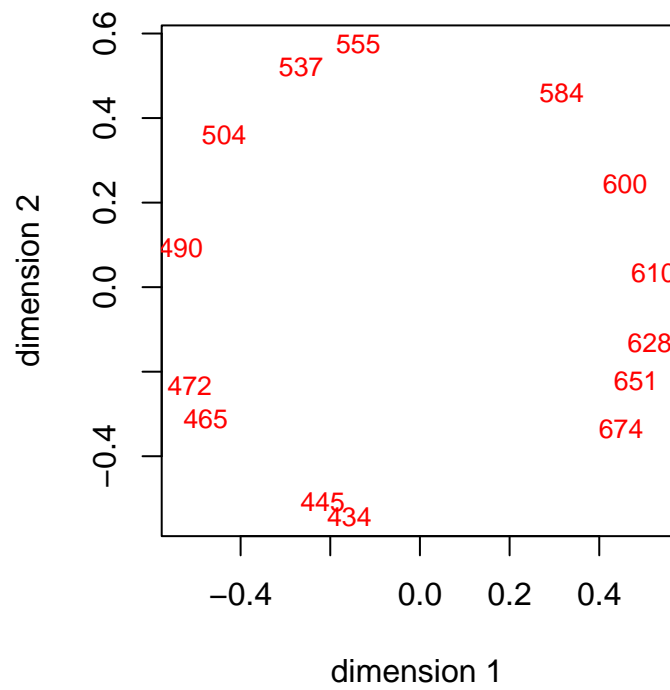


Figure 2: Ekman Metric Raw Stress Solution

6.2 De Gruijter

The Ekman data have an excellent fit in two dimensions and the optimum configuration is extremely stable over variations in MDS methods. The data from De Gruijter (1967) on the similarities between nine Dutch political parties in 1966 have a worse fit, and less stability.

The solution minimizing σ^2 has a loss of 0.3482919 and uses 230 iterations. Minimizing raw stress finds stress 9.4408856 and uses 244 iterations. The optimal configurations in figures 3 and 4 are similar, but definitely not the same. Specifically the position of D66 (a “pragmatic” party, ideologically neither left nor right, established only in 1966, i.e. in the year of the De Gruijter study) differs a lot between solutions.

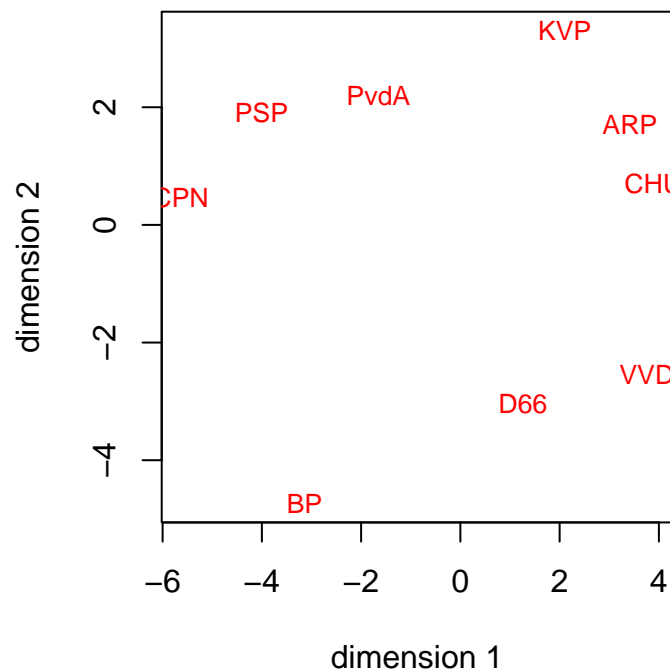


Figure 3: Gruijter Metric Stress 2 Solution

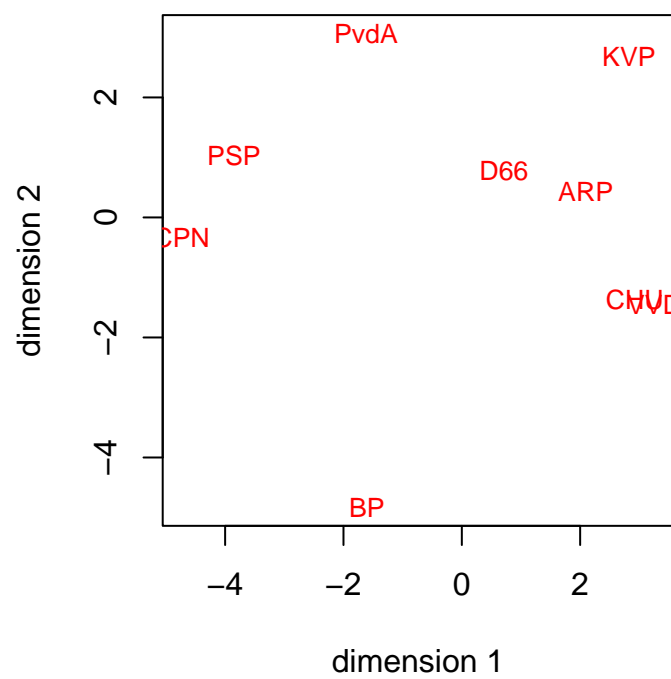


Figure 4: Gruijter Metric Raw Stress Solution

7 Appendix: Code

7.1 stress2.R

```
stress2 <-  
  function(delta,  
    wmat = 1 - diag(nrow(delta)),  
    ndim = 2,  
    itmax = 1000,  
    eps = 1e-10,  
    verbose = TRUE) {  
  itel <- 1  
  n <- nrow(delta)  
  wmat <- wmat / sum(wmat)  
  vmat <- -wmat  
  diag(vmat) <- -rowSums(vmat)  
  xold <- torgerson(delta, ndim)  
  dold <- as.matrix(dist(xold))  
  enum <- sum(wmat * delta * dold)  
  eden <- sum(wmat * dold ^ 2)  
  lbda <- enum / eden  
  dold <- lbda * dold  
  xold <- lbda * xold  
  aold <- sum(wmat * dold)  
  sold <- sum(wmat * (delta - dold) ^ 2) / sum(wmat * (dold - aold) ^ 2)  
  repeat {  
    mmat <- -aold * wmat / (dold + diag(n))  
    diag(mmat) <- -rowSums(mmat)  
    bmat <- -wmat * delta / (dold + diag(n))  
    diag(bmat) <- -rowSums(bmat)  
    umat <- ((1 - sold) * vmat) + (sold * mmat)  
    uinv <- solve(umat + 1/n) - 1/n  
    xnew <- uinv %*% bmat %*% xold  
    dnew <- as.matrix(dist(xnew))  
    anew <- sum(wmat * dnew)  
    snew <- sum(wmat * (delta - dnew) ^ 2) / sum(wmat * (dnew - anew) ^ 2)  
    if (verbose) {  
      cat(  
        "itel ",  
        formatC(itel, format = "d"),  
        "sold ",  
        formatC(sold, digits = 10, format = "f"),  
        "snew ",  
        formatC(snew, digits = 10, format = "f"),  
      )  
    }  
  }  
}
```

```

        "\n"
    )
}
if ((itel == itmax) || ((sold - snew) < eps)) {
    break
}
sold <- snew
dold <- dnew
xold <- xnew
aold <- anew
itel <- itel + 1
}
return(list(
    x = xnew,
    s = snew,
    d = dnew,
    b = bmat,
    m = mmat,
    w = wmat,
    a = anew,
    u = umat,
    itel = itel
))
}

torgerson <- function(delta, ndim) {
    dd <- delta ^ 2
    rd <- apply(dd, 1, mean)
    rr <- mean(dd)
    cc <- -.5 * (dd - outer(rd, rd, "+") + rr)
    ec <- eigen(cc)
    xx <- ec$vectors[, 1:ndim] %*% diag(sqrt(ec$values[1:ndim]))
    return(xx)
}

```

References

- Coombs, C. H. 1964. *A Theory of Data*. Wiley.
- De Gruijter, D. N. M. 1967. “The Cognitive Structure of Dutch Political Parties in 1966.” Report E019-67. Psychological Institute, University of Leiden.
- De Leeuw, J. 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.
- . 1988. “Convergence of the Majorization Method for Multidimensional Scaling.” *Journal of Classification* 5: 163–80.
- . 1994. “Block Relaxation Algorithms in Statistics.” In *Information Systems and Data Analysis*, edited by H. H. Bock, W. Lenski, and M. M. Richter, 308–24. Berlin: Springer Verlag. <https://jansweb.netlify.app/publication/deleeuw-c-94-c/deleeuw-c-94-c.pdf>.
- De Leeuw, J., and W. J. Heiser. 1977. “Convergence of Correction Matrix Algorithms for Multidimensional Scaling.” In *Geometric Representations of Relational Data*, edited by J. C. Lingoes, 735–53. Ann Arbor, Michigan: Mathesis Press.
- 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>.
- De Leeuw, J., and I. Stoop. 1984. “Upper Bounds for Kruskal’s Stress.” *Psychometrika* 49: 391–402.
- Dinkelbach, W. 1967. “On Nonlinear Fractional Programming.” *Management Science* 13: 492–98.
- Ekman, G. 1954. “Dimensions of Color Vision.” *Journal of Psychology* 38: 467–74.
- 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.
- . 1965. “Analysis of Factorial Experiments by Estimating Monotone Transformations of the Data.” *Journal of the Royal Statistical Society B27*: 251–63.
- 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.
- Lange, K. 2016. *MM Optimization Algorithms*. SIAM.
- 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>.
- Roskam, E. E. 1968. “Metric Analysis of Ordinal Data in Psychology.” PhD thesis, University of Leiden.