

APPLICATIONS OF CONVEX ANALYSIS TO MULTIDIMENSIONAL SCALING

JAN DE LEEUW

ABSTRACT. In this paper we discuss the convergence of an algorithm for metric and nonmetric multidimensional scaling that is very similar to the C -matrix algorithm of Guttman. The paper improves some earlier results in two respects. In the first place the analysis is extended to cover general Minkovski metrics, in the second place a more elementary proof of convergence based on results of Robert is presented.

This paper was originally presented at the *European Meeting of Statisticians*, 6–11 September, 1976, in Grenoble, France. It was published previously in J.R. Barra, F. Brodeau, G. Romier and B. Van Cutsem (eds), *Recent Developments in Statistics*, Amsterdam, North Holland Publishing Company, 1977, 133–146.

1. INTRODUCTION

In multidimensional scaling (MDS) problems the data consist of m nonnegative square matrices $\Delta_1, \dots, \Delta_m$ of order n , whose elements are interpreted as measures of *dissimilarity* between the n objects o_1, \dots, o_n measured at m replications r_1, \dots, r_m . Thus δ_{ijk} is the dissimilarity between objects o_i and o_j at replication r . In a psychological context the objects are often called *stimuli* and the replications are defined by the dissimilarity judgments of

2000 *Mathematics Subject Classification.* 62H25.

Key words and phrases. nonmetric scaling, multidimensional scaling, convex analysis, majorization.

different *subjects*. Moreover we assume that m nonnegative square matrices W_1, \dots, W_m of order n are given, whose elements are interpreted as *weights*, i.e. w_{ijk} indicates the relative importance or precision of measurement δ_{ijk} .

Multidimensional scaling techniques represent objects o_1, \dots, o_n as *points* x_1, \dots, x_n in a metric space $\langle \Omega, d \rangle$ in such a way that the *distances* $d(x_i, x_j)$ are approximately equal to the dissimilarities δ_{ijk} . We sometimes write d_{ij} for $d(x_i, x_j)$.

In this paper we study representations of $O = \{o_1, \dots, o_n\}$ in the space of all p -tuples of real numbers, in which the metric is defined by a norm $\|\bullet\|$. Thus $d_{ij} = \|x_i - x_j\|$. A representation of O is then an $n \times p$ matrix X , with row i representing o_i . We also define the notation $d_{ij}(X)$ for the distance between x_i and x_j .

The *loss function* we use in this paper to evaluate the badness-of-fit of a particular representation X is

$$\sigma(X) = \sum_{k=1}^m \sum_{i=1}^n \sum_{j=1}^n w_{ijk} (\delta_{ijk} - d_{ij}(X))^2$$

Clearly $\sigma(X) \geq 0$, and $\sigma(X) = 0$ if and only if $d_{ij}(X) = \delta_{ijk}$ for all i, j, k with $w_{ijk} > 0$. If $w_{ijk} = 0$ then the value of $\sigma(X)$ does not depend on δ_{ijk} . This provides us with a simple device for handling *missing data*: if the observation corresponding with the triple i, j, k is missing, then we can choose δ_{ijk} arbitrarily, and set $w_{ijk} = 0$.

The first and most basic MDS problem we study in this paper is the minimization of $\sigma(X)$ over all $n \times p$ *configuration matrices* X . This is usually called *metric* MDS, to distinguish it from the more general *nonmetric* problem in which the δ_{ijk} are only partially known. Or, more precisely, in the MDS problem as we have defined it so far the δ_{ijk} have to be either

completely known or completely unknown (missing), nonmetric MDS deals with various kinds of intermediate cases. In a later section of the paper we discuss straightforward extensions of our techniques that cover nonmetric MDS.

2. PREVIOUS WORK

Algorithms for the minimization of $\sigma(X)$ have been proposed earlier by Kruskal [1964a,b] and Guttman [1968]. In fact both Kruskal and Guttman propose algorithms to solve the more general nonmetric scaling problems. In this general nonmetric case there are substantial differences between the two approaches, but if we specialize them to the metric MDS problem they become very similar. A detailed discussion and comparison of the algorithms and the corresponding computer programs is available in Lingoes and Roskam [1973]. We only discuss the main ideas, and the major differences between the two approaches in the metric case.

Kruskal proposes a gradient method of the form

$$X \leftarrow X - \alpha \nabla \sigma(X),$$

where $\nabla \sigma(X)$ is the *gradient* of σ at X , i.e. the $n \times p$ matrix of partial derivatives, and where $\alpha > 0$ is a *step-size*. Guttman on the other hand shows that the stationary equation $\nabla \sigma(X) = 0$ can be rewritten in the form $X - C(X)X = 0$, where $C(X)$ is a square symmetric matrix valued function of X . He proposes the iterative process

$$X \leftarrow C(X)X.$$

By substituting Guttman's formula for the gradient in Kruskal's algorithm we find

$$X \leftarrow X - \alpha(X - C(X)X) = (1 - \alpha)X + \alpha C(X)X.$$

Thus Guttman's algorithm is a special case of Kruskal's with a constant step-size α equal to one. And Kruskal's algorithm can be interpreted as an over- or underrelaxed version of Guttman's algorithm. Interesting geometrical and mechanical interpretations of these algorithms have been discussed by Kruskal and Hart [1966], McGee [1966], and Gleason [1967].

There are two problems with the Kruskal-Guttman approach that specifically interest us. In the first place the distance function $d_{ij}(X)$ is typically not differentiable at all configurations X with $x_i = x_j$. This implies that a gradient method cannot be applied to $\sigma(X)$ without further specifications, it implies that the usual convergence theorems for gradient methods are invalid, and it also implies that local minimum points of $\sigma(X)$ need not satisfy the stationary equations. The second problem is that it has not been shown, for either Kruskal's "heuristic" or for Guttman's "constant" step-size procedure, that the resulting algorithms are indeed convergent. Kruskal [1969, 1971] has proved some partial results, and Guttman [1968] and Lingoes and Roskam [1973] have some heuristic arguments and some empirical results, but there is no complete convergence proof.

Until recently these problems have been ignored, or they have been "dissolved" by transforming the model and, through the model, the loss function. ALSCAL, for example, defines the loss on the squared distances and squared dissimilarities Takane et al. [1977]. Classical metric scaling methods apply both squaring and double centering to the dissimilarities, and then

define the loss on the scalar products Torgerson [1958]. These transformations do make the loss functions better behaved in some respects, but they do not really solve the problems with the Kruskal-Guttman approach, they merely transform them away. Moreover using transformations seems less direct, and does not generalize to other distance functions than the usual Euclidean one.

In this paper we derive a simple algorithm for directly minimizing $\sigma(X)$. that can easily proved to be convergent. Although the derivation of the algorithm does not use differentiation or stationary equations, it turns out that the algorithm is identical to Guttman's C -matrix method. One (modest) interpretation of the main result of this paper is that it provides a convergence proof for Guttman's algorithm. Another interpretation is that we show that the C -matrix method should not be interpreted as a gradient method. It is more natural to view it as a minimization method based on an analysis of the convexity properties of the distance function. In fact it may very well be better not to interpret Kruskal's algorithm as a gradient method, but as a relaxed version of the C -matrix method. This interpretation makes it possible, for example, to construct interesting optimal step-size procedures that do not use heuristic arguments and several arbitrary parameters.

3. PROBLEM REDUCTION: PARTITIONING

In this section we reduce the MDS problem to a more simple form by partitioning the loss function into additive components. For this purpose we define

$$\begin{aligned}
\underline{w}_{ij} &= \frac{1}{m} \sum_{k=1}^m w_{ijk}, \\
\underline{\delta}_{ij} &= \sum_{k=1}^m w_{ijk} \delta_{ijk} / m \underline{w}_{ij}, \\
w_{ij} &= \frac{1}{2} (\underline{w}_{ij} + \underline{w}_{ji}), \\
\delta_{ij} &= \frac{\underline{w}_{ij} \underline{\delta}_{ij} + \underline{w}_{ji} \underline{\delta}_{ji}}{\underline{w}_{ij} + \underline{w}_{ji}} \text{ for } i \neq j, \\
w_{ii} &= \delta_{ii} = 0.
\end{aligned}$$

As is customary in the analysis of variance we collect the components of the partitioning in a table,

SOURCE	LOSS COMPONENT
Proper Loss	$\sum_{i=1}^n \sum_{j=1}^n w_{ij} (\delta_{ij} - d_{ij}(X))^2.$
Symmetry	$\sum_{i \neq j}^n \{ \underline{w}_{ij} \underline{\delta}_{ij}^2 - w_{ij} \delta_{ij}^2 \}.$
Hollowness	$\sum_{i=1}^n \underline{w}_{ii} \underline{\delta}_{ii}^2.$
Individual differences	$\sum_{k=1}^m \sum_{i=1}^n \sum_{j=1}^n w_{ijk} (\delta_{ijk} - \underline{\delta}_{ij})^2.$
Total loss	$\sum_{k=1}^m \sum_{i=1}^n \sum_{j=1}^n w_{ijk} (\delta_{ijk} - d_{ij}(X))^2$

It is obvious that we minimize the total loss if we minimize the proper loss, and that the proper loss is more simple. In fact in defining the proper loss we can suppose without loss of generality that both the weights and the dissimilarities are symmetric and hollow. The only assumption we make about the δ_{ij} is that they are nonnegative. The weights are also assumed to be nonnegative, but we make the additional nondegeneracy assumption of *irreducibility*: we suppose that there is no partitioning of $\{1, 2, \dots, n\}$ such that $w_{ij} = 0$ whenever i and j belong to different members of the partition. Again this assumption causes no real loss of generality, because

if all between-subset weights are zero the MDS problem separates into a number of smaller problems corresponding with each of the subsets.

4. PROBLEM MANIPULATION: USE OF HOMOGENEITY

As we have shown in the previous section the metric MDS problem can be reformulated without loss of generality as the minimization of

$$\sigma(X) = \sum_{1 \leq i < j \leq n} \sum w_{ij} (\delta_{ij} - d_{ij}(X))^2$$

over the $n \times p$ configuration matrices X . In this section we study a closely related maximization problem, that is in some respects more simple. For the discussion of this alternative problem we need the following definitions.

$$\begin{aligned} \rho(X) &= \sum_{1 \leq i < j \leq n} \sum w_{ij} \delta_{ij} d_{ij}(X), \\ \eta^2(X) &= \sum_{1 \leq i < j \leq n} \sum w_{ij} d_{ij}^2(X), \\ \eta_\delta^2 &= \sum_{1 \leq i < j \leq n} \sum w_{ij} \delta_{ij}^2, \end{aligned}$$

and

$$\lambda(X) = \rho(X) / \eta(X) \eta_\delta.$$

Theorem 4.1. *For all X we have $0 \leq \lambda(X) \leq 1$. Moreover $\lambda(X) = 1$ if and only if the dissimilarities δ_{ij} and the distances $d_{ij}(X)$ for which $w_{ij} \neq 0$ are proportional.*

Proof. This follows directly from the Cauchy-Schwartz inequality applied to $\rho(X)$. □

Theorem 4.2. **a:** Suppose \hat{X} minimizes $\sigma(X)$. Then \hat{X} also maximizes $\lambda(X)$.
 b: Suppose \hat{X} maximizes $\lambda(X)$. Then $\{\rho(\hat{X})/\eta^2(\hat{X})\}\hat{X}$ minimizes $\sigma(X)$.

Proof. Because $d_{ij}(\beta X) = \beta d_{ij}(X)$ for all X and all $\beta \geq 0$ we can reformulate the MDS problem as the minimization of

$$\sum_{1 \leq i < j \leq n} \sum w_{ij} (\delta_{ij} - \beta d_{ij}(X))^2$$

over the $n \times p$ matrices X and over all $\beta \geq 0$. The minimum over β for fixed X is attained at

$$\hat{\beta} = \rho(X)/\eta^2(X),$$

and the value at the minimum is $\eta_\delta^2(1 - \lambda^2(X))$. The theorem follows from these computations. \square

It follows from Theorem 4.2 that we can solve the metric MDS problem by finding the configuration matrix that maximizes $\lambda(X)$.

5. THE EUCLIDEAN CASE

Suppose $d_{ij}(X)$ is Euclidean, i.e.

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

In this case it is convenient to derive some matrix expressions for $\rho(X)$ and $\eta(X)$. Define the matrix valued function $B(X)$ by

$$\begin{aligned} b_{ij}(X) &= -w_{ij}\delta_{ij}s_{ij}(X) \text{ if } i \neq j, \\ b_{ii}(X) &= \sum w_{ij}\delta_{ij}s_{ij}(X). \end{aligned}$$

Here

$$s_{ij}(X) = \begin{cases} d_{ij}^{-1}(X) & \text{if } d_{ij}(X) \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

We also define the matrix V by

$$\begin{aligned} v_{ij} &= -w_{ij} \text{ if } i \neq j, \\ v_{ii} &= \sum w_{ij}. \end{aligned}$$

Both $B(X)$ and V are real symmetric matrices with non-positive off-diagonal and non-negative diagonal elements, whose rows and columns sum to zero. By a familiar matrix theorem they are consequently both positive semi-definite of rank not exceeding $n-1$. Because V is irreducible by assumption we have in fact $\text{rank}(V) = n-1$, and the null space of V is the set of all vectors with constant elements. (Taussky [1949]; also Varga [1962, sections 1.4 and 1.5]). If e is the n -vector with all elements equal to one, then the Moore-Penrose inverse of V is simply

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

The following results can be verified easily.

Theorem 5.1. **a:** $\rho(X) = \text{tr } X' B(X) X$.

b: $\eta^2(X) = \text{tr } X' V X$.

We also define, for all pairs of configuration matrices,

$$\mu(X, Y) = \text{tr } X' B(Y) Y.$$

Theorem 5.2. $\mu(X, Y) \leq \rho(X)$ for all X, Y .

Proof. The Cauchy-Schwartz inequality implies

$$d_{ij}(X) \geq s_{ij}(Y) \sum_{s=1}^p (x_{is} - x_{js})(y_{is} - y_{js}).$$

If we multiply both sides with $w_{ij}\delta_{ij}$, sum all inequalities, and simplify, we find the inequality stated in the theorem. \square

Using the notation developed in this section we can now define the B -matrix algorithm for Euclidean metric multidimensional scaling as the recursion $X^{k+1} = V^+ B(X^k) X^k$. The only difference between Guttman's C -matrix and our B -matrix is due to the fact that we have removed the homogeneity from the problem in section 4, this makes B more simple than C .

Theorem 5.3. **a:** *The three sequences $\rho(X^k)$, $\eta(X^k)$, and $\lambda(X^k)$ are bounded and increasing. The limits are ρ_∞ , $\eta_\infty = \sqrt{\rho_\infty}$, and $\lambda_\infty = \rho_\infty / \eta_\infty$.*

b: *The sequence X^k has convergent subsequences. If X_∞ is the limit of a convergent subsequence, then $\lambda(X_\infty) = \lambda_\infty$. Moreover X_∞ is a fixed point, i.e. $X_\infty = V^+ B(X_\infty) X_\infty$, and if λ is differentiable at X_∞ then $\nabla \lambda(X_\infty) = 0$.*

c: $\|X^{k+1} - X^k\| \rightarrow 0$.

Proof. From the Cauchy-Schwartz inequality

$$\rho(X^k) = \mathbf{tr} (X^k)' V X^{k+1} \leq \eta(X^k) \eta(X^{k+1}).$$

From Theorem 5.2

$$\rho(X^{k+1}) \geq \mathbf{tr} (X^{k+1})' B(X^k) X^k = \mathbf{tr} (X^{k+1})' V X^{k+1} = \eta^2(X^{k+1}).$$

If we combine these inequalities we obtain

$$\eta(X^k) \leq \frac{\rho(X^k)}{\eta(X^k)} \leq \eta(X^{k+1}),$$

and

$$\rho(X^k) \leq \frac{\eta(X^k)}{\eta(X^{k+1})} \rho(X^{k+1}) \leq \rho(X^{k+1}).$$

Because $\rho(X^k) \leq \eta_\delta^2$ and $\eta(X^k) \leq \eta_\delta$ it follows that both $\rho(X^k)$ and $\eta(X^k)$ are convergent increasing sequences, with limits, say, ρ_∞ and η_∞ . Because $\lambda(X_k) \leq 1$ it follows that $\lambda(X^k)$ is another convergent increasing sequence with limit λ_∞ . Moreover $\rho_\infty = \eta_\infty^2$, and $\lambda_\infty = \rho_\infty / \eta_\delta \eta_\infty$. It also follows that the sequence X^k lies in the compact set $\eta(X) \leq \eta_\delta$ and has convergent subsequences. There is equality in the basic chain of inequalities if and only if X is a fixed point. This implies that subsequential limits are fixed points. Finally

$$\begin{aligned} \text{tr}(X^{k+1} - X^k)' V(X^{k+1} - X^k) &= \\ &= \eta^2(X^{k+1}) + \eta^2(X^k) - 2\rho(X^k) \rightarrow 2(\eta_\infty^2 - \rho_\infty) = 0, \end{aligned}$$

which implies part (c) of the theorem. \square

A very similar result appears in Robert [1967]. The general convergence theorems of Zangwill [1969] and Meyer [1976] are also relevant. Observe that Theorem 5.3 does not say that X^k converges. This follows only if we make some rather arbitrary additional assumptions, for example that there is only a finite number of fixed points, or that one of the subsequential limits is an isolated fixed point. If X^k does *not* converge, it follows that the set of limit points is a continuum (a result due to Ostrowski, cf Daniel [1971, section 6.3]).

Theorem 5.3 is our basic convergence theorem for metric Euclidean MDS. It is quite satisfactory, and it has been proved by very elementary methods. In fact the proof only uses some elementary properties of sequences, and the Cauchy-Schwartz inequality.

6. NON-EUCLIDEAN METRICS

Our method can be generalized to general Minkovski metrics, with the metric defined by a *gauge* ϕ . We shall use some elementary facts about gauges without proof. The proofs follow easily from the beautiful introduction to gauges and norms in Rockafellar [1970, chapter 15]. Introductions to Minkovski geometry are given in Busemann and Kelly [1953] and Busemann [1955].

A *gauge* is a function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying

$$\mathbf{G1:} \quad \phi(x) \geq 0.$$

$$\mathbf{G2:} \quad \phi(x) = 0 \text{ if and only if } x = 0.$$

$$\mathbf{G3:} \quad \phi(\mu x) = \mu \phi(x) \text{ for all } \mu \geq 0.$$

$$\mathbf{G4:} \quad \phi(x + y) \leq \phi(x) + \phi(y).$$

A gauge is a *norm* if we can replace **G3** by the stronger

$$\mathbf{G5:} \quad \phi(\mu x) = |\mu| \phi(x) \text{ for all } \mu.$$

A gauge defines a distance function by the rule

$$d_{ij}(X) = \phi(x_i - x_j).$$

Unless the gauge is a norm this distance is not necessarily symmetric. With some minor modifications our results are also valid if we replace **G2** by the weaker

$$\mathbf{G6:} \quad \phi(0) = 0.$$

In fact most of the results remain valid if we only assume **G3** and **G4**, i.e. for all homogeneous convex functions. Thus gauges are relevant for

our problem because they can be used to construct very general distance functions. But they are even more relevant because of the following result.

Theorem 6.1. *Both $\rho(X)$ and $\eta(X)$ are gauges on the space of all $n \times p$ configuration matrices.*

Proof. For both $\rho(X)$ and $\eta(X)$ property **G1** is obvious. For $\eta(X)$ property **G2** follows from irreducibility, for $\rho(X)$ property **G6** is obvious, We could assume **G2**, but we never need it. Properties **G3** and **G4** follows from the fact that each $d_{ij}(X)$ is convex and homogeneous on the space of configuration matrices. \square

The theorem shows that the metric MDS problem reduces to the maximization of a ratio of two gauges. Problems of that type have been studied by Robert [1967], Boyd [1974], Tao [1975, 1976]. Before we discuss their results and apply them to our problem we state some of the elementary facts about gauges.

First define the polar of a gauge as the function $\phi^\circ : \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$\phi^\circ(y) = \max_x \frac{\langle x, y \rangle}{\phi(x)}.$$

Here $\langle \bullet, \bullet \rangle$ denotes inner product.

- Fact 6.2.**
- a:** *The polar of a gauge is a gauge, the polar of a norm is a norm.*
 - b:** *The polar of the polar of a gauge ϕ is the gauge ϕ .*
 - c:** *The Euclidean norm $\sqrt{\langle x, x \rangle}$ is its own polar.*
 - d:** (Holder's inequality). *If ϕ and ϕ° are polar gauges, then $\langle x, y \rangle \leq \phi(x)\phi^\circ(y)$ for all $x, y \in \mathbf{R}^n$.*

We can study the conditions for equality in Holder's inequality by introducing subdifferentials. Remember that a *subgradient* of a function ϕ at a point

x is a vector y such that $\phi(z) \geq \phi(x) + \langle y, z - x \rangle$ for all $z \in \mathbb{R}^n$. The set of all subgradients of ϕ at a point x is the *subdifferential* of ϕ at x , and is written as $\partial\phi(x)$. Thus for each x the symbol $\partial\phi(x)$ stands for a subset of \mathbb{R}^n , possibly empty. Again we mention some facts about subdifferentials, without proof. The proofs can be found in part V of Rockafellar [1970].

- Fact 6.3.**
- a:** *If ϕ is a finite convex function, then for each $x \in \mathbb{R}^n$ the set $\partial\phi(x)$ is nonempty, convex, and compact.*
 - b:** *If ϕ is differentiable at x with gradient $\nabla\phi(x)$, then $\partial\phi(x) = \{\nabla\phi(x)\}$.*
 - c:** *The map $x \rightarrow \partial\phi(x)$ is closed, i.e. if $x_i \rightarrow x_\infty$, if $y_i \rightarrow y_\infty$, and if for each i also $y_i \in \partial\phi(x_i)$, then $y_\infty \in \partial\phi(x_\infty)$.*

By combining the results of Fact 6.2 and Fact 6.3 we find the following results.

- Fact 6.4.**
- a:** *Suppose ϕ is a gauge. Then $y \in \partial\phi(x)$ if and only if $\phi(z) \geq \langle y, z \rangle$ for all $z \in \mathbb{R}^n$, and $\phi(x) = \langle y, x \rangle$.*
 - b:** *Suppose ϕ and ϕ° are polar gauges. Then $x \in \partial\phi(y)$ if and only if $\langle x, y \rangle = \phi^\circ(y)$ and $\phi(x) = 1$. Moreover $y \in \partial\phi(x)$ if and only if $\langle x, y \rangle = \phi(x)$ and $\phi^\circ(y) = 1$.*
 - c:** *Suppose ϕ and ϕ° are polar gauges. Then $\langle x, y \rangle = \phi(x)\phi^\circ(y)$ if and only if $x \in \phi(x)\partial\phi(y)$ if and only if $y \in \phi^\circ(y)\partial\phi(x)$.*

Now consider the problem of maximizing the ratio

$$\lambda(x) = \frac{\phi(x)}{\psi(x)}.$$

with both ϕ and ψ gauges. From the definitions of gauges and their polars we obtain the following result.

Theorem 6.5. *Maximizing $\lambda(x)$ over \mathbb{R}^n is equivalent to maximizing*

$$\xi(x, y) = \frac{\langle x, y \rangle}{\psi(x)\phi^\circ(y)}$$

over $\mathbb{R}^n \otimes \mathbb{R}^n$, and this is equivalent to minimizing

$$\lambda^\circ(x) = \frac{\phi^\circ(x)}{\psi^\circ(x)}$$

over \mathbb{R}^n .

By using Holder's inequality we can derive the following necessary conditions for an extreme value.

Theorem 6.6. **a:** *If \hat{x}, \hat{y} maximizes $\xi(x, y)$ then $\hat{y} \in \phi^\circ(\hat{y})\partial\phi(\hat{x})$ and $\hat{x} \in \psi(\hat{x})\partial\psi(\hat{y})$.*
b: *If \hat{x} maximizes $\lambda(x)$ or minimizes $\lambda^\circ(x)$ then $\hat{x} \in \psi(\hat{x})\partial\psi^\circ(\partial\phi(\hat{x}))$.*

For maximizing $\lambda(x)$ the following algorithm was proposed by Robert (1967).

We start with x^0 such that $\psi(x^0) = 1$. Then define $y^k \in \partial\phi(x^k)$ and $x^{k+1} \in \partial\psi^\circ(y^k)$.

Theorem 6.7. **a:** *The sequence $\lambda(x^k)$ is increasing and convergent. The sequence $\lambda^\circ(y^k)$ is decreasing and convergent. Both sequences converge to the same limit $\underline{\lambda}$.*
b: *All accumulation points of (x^k, y^k) correspond with the same function value $\underline{\lambda} = \underline{\xi}$. Moreover all accumulation points satisfy the necessary conditions of Theorem 6.6.*

Proof. From our facts about gauges

$$\phi^\circ(y^k) = 1,$$

$$\psi(x^{k+1}) = 1,$$

$$\langle x^k, y^k \rangle = \phi(x^k),$$

$$\langle x^{k+1}, y^k \rangle = \psi^\circ(y^k).$$

By applying Holder's inequality

$$\phi(x^k) = \langle x^k, y^k \rangle \leq \psi^\circ(y^k)\psi(x^k) = \psi^\circ(y^k),$$

$$\psi^\circ(y^k) = \langle x^{k+1}, y^k \rangle \leq \phi(x^{k+1})\phi^\circ(y^k) = \phi(x^{k+1}),$$

and thus

$$\lambda(x^k) \leq 1/\lambda^\circ(y^k) \leq \lambda(x^{k+1}),$$

which implies part (a). Because the subdifferentials are closed and the iterations remain in a compact set we can apply the general convergence theorems of Zangwill [1969] to get (b). \square

If we compare 6.7 and 5.3 we see that 6.7 has no part (c), and is consequently weaker than 5.3. It is possible to prove that $\|x^{k+1} - x^k\| \rightarrow 0$ in this more general context too, but we need additional assumptions. One of the more natural ones is that ϕ or ψ° or both are differentiable at all stationary points, other possibilities are discussed by Meyer [1976]. A far more important difference between the algorithms of sections 5 and 6 is that in most cases the function ψ° cannot be computed in closed form. The same thing is true for the subdifferential $\partial\psi^\circ$. This means that we must compute x^{k+1} by maximizing $\langle x, y^k \rangle$ over $\{x \mid \psi(x) = 1\}$. This is a convex programming problem which cannot be solved in a finite number of steps in general.

Consequently we need a version of Theorem 6.7 in which this convex programming problem is truncated after a finite number of steps. Zangwill's general convergence theory shows how this truncating should be done.

The simplifications in section 5 are possible because the gauge ψ is ellipsoidal in this case. If both ψ and ϕ are ellipsoidal, then Robert has pointed out that the algorithm reduces to the power method for solving a generalized eigenvalue problem. Compare also Tao [1976]. It is of some interest that Guttman already pointed out that his C -matrix method for MDS looked like a sort of generalized power method. The analysis in this paper shows what the exact relationships are.

7. NON-METRIC MDS

In the simplest forms of nonmetric MDS we must minimize

$$\tau(X, \Delta) = \frac{\sum_{i=1}^n \sum_{j=1}^n w_{ij} (\delta_{ij} - d_{ij}(X))^2}{\sum_{i=1}^n \sum_{j=1}^n w_{ij} d_{ij}(X)^2}$$

over all $n \times p$ configuration matrices X and over all $n \times n$ disparity matrices Δ . The disparity matrices must be chosen from a known convex cone Γ , the metric MDS problem is the special case in which Γ is a ray, the additive constant problem is the special case in which Γ is a two-dimensional subspace. We briefly indicate the modifications needed to apply our ideas to nonmetric MDS in this simple form. More complicated partitioned loss functions, with more complicated normalizations, will be discussed in subsequent publications.

By using the homogeneity of the distance function as in section 4 we can show that the nonmetric MDS problem is equivalent to the maximization of

$$\lambda(X, \Delta) = \frac{\rho(X, \Delta)}{\eta(X)\eta(\Delta)},$$

with

$$\rho(X, \Delta) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} \delta_{ij} d_{ij}(X),$$

$$\eta(\Delta) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} \delta_{ij}^2$$

and $\eta(X)$ as before. If we define

$$\rho(X) = \max_{\Delta \in \Gamma} \frac{\rho(X, \Delta)}{\eta(\Delta)}$$

then $\rho(X)$ is a homogeneous convex function, in fact a gauge. Thus we have a problem of the familiar form, a ratio of gauges must be maximized, and the algorithm of section 6 can be applied.

If the distance function is Euclidean, the analysis of section 5 can be used. The only difference with metric MDS is in the definition of $\rho(X)$, in the nonmetric case we have to compute the optimum Δ for given X in order to compute $\rho(X)$. We can compute the optimum $\hat{\Delta}(X)$ as the unique minimizer of

$$\sum_{i=1}^n \sum_{j=1}^n w_{ij} (\delta_{ij} - d_{ij}(X))^2$$

over the cone Γ . After solving this regression problem we can normalize the solution such that $\eta(\hat{\Delta}(X)) = 1$, but this is not strictly necessary. The B -matrix algorithm for nonmetric Euclidean MDS is defined as the recursion

$$X^{k+1} = V^+ B(X^k, \hat{\Delta}(X^k)) X^k,$$

with $B(X^k, \hat{\Delta}(X^k))$ defined as $B(X)$, but with $\hat{\delta}_{ij}(X)$ substituted for δ_{ij} . The same inequalities and equations can be derived as in 5.1, 5.2, and the proof of 5.3.

Theorem 7.1. *Parts a,b,c of Theorem 5.3 are also true for the nonmetric Euclidean B-matrix algorithm.*

In the nonmetric case the differences between our B -matrix method and Guttman's C -matrix method are larger than in the metric case. One important reason is that Guttman uses rank images, while our convexity approach forces us to use monotone regression estimates of the δ_{ij} . I have not been able to find a rigorously defined optimization problem in which rank images can be used. This does not mean, of course, that we cannot use rank images in the earlier iterations of an MDS algorithm. In the earlier iterations we can do anything we please. As in TORSCA we can use the semi-nonmetric Young-Householder process, or as in MINISSA we can use rank images. We only have to switch to monotone regression and the B -matrix algorithm if things are getting out of hand (if the loss starts to increase, for example).

REFERENCES

- D.W. Boyd. The Power Method for ℓ^p Norms. *Linear Algebra and Its Applications*, 9:95–101, 1974.
- H. Busemann. *The Geometry of Geodesics*. Academic Press, New York, 1955.
- H. Busemann and P. Kelly. *Projective Geometry and Projective Metrics*. Academic Press, New York, 1953.
- J.W. Daniel. *The Approximate Minimization of Functionals*. Prentice Hall, Englewood Cliffs, 1971.
- T.C. Gleason. A General Model for Nonmetric Multidimensional Scaling. Report 3, Michigan Mathematical Psychology Program, 1967.
- L. Guttman. A General Nonmetric Technique for Fitting the Smallest Coordinate Space for a Configuration of Points. *Psychometrika*, 33:469–506, 1968.
- J. B. Kruskal. Multidimensional Scaling by Optimizing Goodness of Fit to a Nonmetric Hypothesis. *Psychometrika*, 29:1–27, 1964a.

- J.B. Kruskal. Nonmetric Multidimensional Scaling: a Numerical Method. *Psychometrika*, 29:115–129, 1964b.
- J.B. Kruskal. A New Convergence Condition for Methods of Ascent. Unpublished memo, Bell Telephone Laboratories, 1969.
- J.B. Kruskal. Monotone Regression: Continuity and Differentiability Properties. *Psychometrika*, 36:57–62, 1971.
- J.B. Kruskal and R.E. Hart. A Geometric Interpretation of Diagnostic Data from a Digital Machine. *Bell System Technical Journal*, 45:1299–1338, 1966.
- J.C. Lingoes and E.E. Roskam. A Mathematical and Empirical Analysis of Two Multidimensional Scaling Algorithms. *Psychometrika*, 38: Monograph Supplement, 1973.
- V.E. McGee. The Multidimensional Scaling of “Elastic” Distances. *British Journal of Mathematical and Statistical Psychology*, 19:181–196, 1966.
- R. R. Meyer. Sufficient Conditions for the Convergence of Monotonic Mathematical Programming Algorithms. *Journal of Computer and System Sciences*, 12:108–121, 1976.
- F. Robert. Calcul du Rapport Maximal de Deux Normes sur \mathbb{R}^n . *Revue Française d’Automatique, d’Informatique Et De Recherche Operationelle*, 1: 97–118, 1967.
- R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, New Jersey, 1970.
- Y. Takane, F.W. Young, and J. De Leeuw. Nonmetric Individual Differences in Multidimensional Scaling: An Alternating Least Squares Method with Optimal Scaling Features. *Psychometrika*, 42:7–67, 1977.
- Pham Dinh Tao. Calcul du Maximum d’une Forme Quadratique Définie Positive sur la Boule Unité de ψ_∞ . Seminaire d’Analyse Numérique 247, Université Scientifique et Médicale de Grenoble, 1976.

- Pham Dinh Tao. Eléments Homoduaux d'une Matrice A Relatifs a un Couple de Normes (ϕ, ψ) . Applications au Calcul de $S_{\phi\psi}(A)$. Seminaire d'Analyse Numérique 236, Université Scientifique et Médicale de Grenoble, 1975.
- O. Taussky. A Recurring Theorem on Determinants. *American Mathematical Monthly*, 56:672–676, 1949.
- W.S Torgerson. *Theory and Methods of Scaling*. Wiley, New York, 1958.
- R.S. Varga. *Matrix Iterative Analysis*. Prentice Hall, Englewood Cliffs, 1962.
- W. I. Zangwill. *Nonlinear Programming: a Unified Approach*. Prentice-Hall, Englewood-Cliffs, N.J., 1969.