

THE LINEAR NONMETRIC MODEL

Jan de Leeuw

June 1969

1. General formulation.

Let  $X$  be a finite subset of the set of all  $m$ -tuples of real numbers.  $X$  has  $n$  elements  $x_1, x_2, \dots, x_n$ . The coordinate values can be collected in a matrix  $X = \{x_{ij}\}$ ,  $i=1, \dots, n$ ;  $j=1, \dots, m$ . Suppose  $Y$  is another set of  $n$  elements, partially ordered by the relation  $\succsim$ . We are interested in the existence of real numbers  $v_1, \dots, v_m$  such that

$$\sum_{k=1}^m x_{ik} v_k \geq \sum_{k=1}^m x_{jk} v_k \text{ iff } v_i \succsim v_j. \tag{1}$$

To establish representation theorems the easiest way is to use the method of Tversky (1964) and Scott (1964). Each of the requirements (1) can be translated into a homogeneous linear inequality. If  $(v_1, v_m) \in \succsim$  then we require

$$\sum (x_{i1} - x_{j1}) v_1 \geq 0. \tag{2}$$

The representation theorems are the familiar theorems giving necessary and sufficient conditions for the existence of solutions of systems of linear inequalities. A very complete review of these theorems is given by Fan (1956). The uniqueness theorems, that state properties of the solution set, are most easily understood by using the geometric theory of polyhedral convex cones and their extreme faces.

The constructive approach to measurement theory, as practiced by people as Luce, Krantz, and Suppes, can also be used in this case. Let  $\oplus$  denote (componentwise) vector addition in  $X$ . The axioms are chosen in such a way that  $\langle X, \oplus, \succsim \rangle$  is a fully ordered Archimedean group, and by Hölder's theorem there exists a function  $f$  of  $X$  into the reals such that for all  $x, y \in X$

$$i) \quad f(x \oplus y) = f(x) + f(y), \tag{3}$$

$$ii) \quad f(x) \geq f(y) \text{ iff } x \succsim y. \tag{4}$$

It is well known that the only continuous solutions of the functional equation (3) are the linear functions. Continuity of  $f$  can be guaranteed by imposing (axiomatically) restrictions on the order topology induced

on  $X$  by  $\sum$ . We shall assume in the sequel

i) the columns of the matrix  $X$  are centered, i.e.

$$\sum_{i=1}^n x_{ij} = 0 \text{ for all } j=1, \dots, m.$$

ii)  $Y$  is a set of real numbers. Without loss of generality

$$\text{we may suppose } \sum y_i = 0.$$

iii) We assume that  $\text{rank}(X) = m$ , and consequently also that

$$n \geq m.$$

In later sections we shall show that all three assumptions can be avoided and are not essential for the algorithms.

2: Alternating least squares algorithm.

In order to find the best value of  $w$  for our problem we use a Krukal-type error theory. Define the loss function

$$S = \frac{\sum_{i=1}^n (z_i - \hat{z}_i)^2}{\sum_{i=1}^n (z_i - \bar{z})^2}, \quad (5)$$

with

$$z = Xw, \quad (6)$$

$\bar{z}$  is the mean value of the  $z_i$ , and  $\hat{z}$  is any vector monotone with  $y$ .

Throughout this paper  $j$  denotes a vector with all elements equal to unity. Because  $X$  is centered columnwise we have  $j'X = 0$  and  $\bar{z} = \frac{1}{n} j'Xw = 0$ . Thus

$$S = \frac{(z - \hat{z})'(z - \hat{z})}{z'z} = \frac{w'X'Xw - 2w'X'\hat{z} + \hat{z}'\hat{z}}{w'X'Xw}. \quad (7)$$

Consider  $\hat{z}$  as a fixed vector. Then  $S$  is minimized if we minimize its numerator under the condition that  $w'X'Xw$  is some constant value (say unity). The Lagrangian function is

$$\Phi(w) = w'X'Xw - 2w'X'\hat{z} + \hat{z}'\hat{z} - \mu(w'X'Xw - 1), \quad (8)$$

and

$$\frac{\partial \Phi}{\partial w'} = 2X'Xw - 2X'\hat{z} - 2\mu X'Xw. \quad (9)$$

We must solve the system of equations

$$(1 - \mu)X'Xw = X'\hat{z}, \quad (10)$$

$$w'X'Xw = 1, \quad (11)$$

for  $w$  and  $\mu$ . The solution is

$$w = (\hat{z}'X(X'X)^{-1}X'\hat{z})^{-\frac{1}{2}} (X'X)^{-1}X'\hat{z}, \quad (12)$$

$$\mu = 1 - (\hat{z}'X(X'X)^{-1}X'\hat{z})^{\frac{1}{2}}. \quad (13)$$

For a first approximation we take  $\hat{z} = y$ . Observe that the problem of minimizing  $S$  for fixed  $\hat{z}$  is closely related to the oblique Procrustus problem discussed by Browne (1967). There we require  $w'w = 1$  in stead of  $w'X'Xw = 1$  and, as a consequence, the stationary equations must be solved iteratively. The next step is to compute new values of  $z$ , and to find  $\hat{z}$  in such a way that  $S$  is minimized under the condition that  $\hat{z}$  remains monotone with  $y$ . This problem is solved by the amalgamation method used by Kruskal. In summary

- 1) set  $\hat{z}^{(0)} = y$ , and set iteration counter  $s = 1$  ;
- 2) find  $w^{(s)}$  by (12) with  $\hat{z} = \hat{z}^{(s-1)}$  ;
- 3) compute  $z^{(s)} = Xw^{(s)}$  ;
- 4) find  $\hat{z}^{(s)}$  by the amalgamation method, compute  $S^{(s)}$  ;
- 5) if  $s = s_{\max}$  then go to step 6 (see section 3), otherwise  $s = s + 1$  and goto step 2.

### 3: Newton-Rahson iterations.

The iterative scheme in the previous section has the obvious property that  $S^{(s+1)} \leq S^{(s)}$ . In fact in each iteration the value of  $S$  is diminished two times, in step (2) and in step (4). We have  $S^{(s+1)} = S^{(s)}$  iff  $w^{(s+1)} = w^{(s)}$  iff  $\hat{z}^{(s+1)} = \hat{z}^{(s)} = z^{(s)}$  iff  $S^{(s)} = 0$ . If  $S_{\min} > 0$  then the method does not converge in a finite number of steps. The rate of convergence is somewhat unsatisfactory (in general it can be compared to that of the optimal gradient method, which means that the method is very inefficient near the minimum). In our complete algorithm we use only a few ( $s_{\max}$ )

steps of this method. They are supposed to bring us close to the minimum, where  $S$  is dominated by the second order terms of its Taylor expansion, i.e. nearly quadratic. There the Newton-Raphson method becomes very efficient. It has the additional advantage that we obtain information about the distance from the minimum and about the curvature of the function. In order to develop the necessary formula's let (superscripts denoting the iteration number  $s$  are deleted)

$$e = X'z = X'Xw, \quad (14)$$

$$f = X'(z - \hat{z}), \quad (15)$$

$$S^{\#} = (z - \hat{z})'(z - \hat{z}), \quad (16)$$

$$T^{\#} = z'z. \quad (17)$$

Then

$$\frac{\partial S^{\#}}{\partial w_k} = 2f_k, \quad (18)$$

and

$$\frac{\partial T^{\#}}{\partial w_k} = 2e_k. \quad (19)$$

Thus

$$\frac{\partial S}{\partial w_k} = \frac{2}{T^{\#}} (f_k - S e_k). \quad (20)$$

These partial derivatives are collected in the gradient vector  $g$ . The second order derivatives are given by

$$\frac{\partial^2 S}{\partial w_k \partial w_l} = \frac{2}{T^{\#}} \left\{ (1 - S)c_{kl} - (g_l e_k + g_k e_l) \right\}, \quad (21)$$

with  $C = X'X$ . These elements are collected in the symmetric Hessian matrix  $V$ , and we assume that  $V$  is nonsingular. The Newton-Raphson scheme is simply

$$w^{(s+1)} = w^{(s)} - (V^{(s)})^{-1} g^{(s)}. \quad (22)$$

The estimated distance from the minimum is given by

$$S^{(s)} - S_{\min} \approx \frac{1}{2} g^{(s)' (V^{(s)})^{-1} g^{(s)}.$$

The quantity (say  $\rho$ ) on the right hand side of (23) is a good stopping criterion for the iterations. Observe that the form of the second order derivatives suggests another modified gradient method which will also be quite efficient near the minimum and does not require computation and inversion of  $V$ . Simply take

$$w^{(s+1)} = w^{(s)} - \left[ \frac{2(1-S^{(s)})}{(T^{\bar{x}})^{(s)}} \right]^{-1} C^{-1} g^{(s)}. \quad (24)$$

The weighting matrix  $C^{-1}$  must be computed anyway. This method is an analogon of the method of scoring for parameters in ML-estimation. To conclude the scheme from the previous section

- 6)  $t = 1, w^{(0)} = w^{(s)}, \hat{z}^{(0)} = z^{(s)}, S^{(0)} = S^{(s)}$  ;
- 7) compute  $g^{(t)}, V^{(t)}, (V^{(t)})^{-1}$ , and  $\rho^{(t)}$  ;
- 8) if  $\rho^{(t)} < \xi$  then stop ;
- 9) compute  $w^{(t)}$  by (22) ;
- 10) compute  $z^{(t)} = Xw^{(t)}$  ;
- 11) compute  $\hat{z}^{(t)}$  by the amalgamation method ;
- 12)  $t = t+1$ , and go to step 7 .

#### 4: Restrictive conditions.

The first restrictive condition we imposed in section 1 was that  $\sum_i x_{ij} = 0$  for all  $j$ . Suppose this is not true. Then there are real numbers  $\lambda_j$  ( $j=1, \dots, m$ ) such that

$$x_{ij} = \tilde{x}_{ij} + \lambda_j, \quad (25)$$

where  $\tilde{x}_{ij}$  is our original centered variate. Then, for a particular value of the vector  $w$ ,

$$z_i = \sum_k w_k x_{ik} = \sum_k w_k \tilde{x}_{ik} + \sum_k \lambda_k w_k = \tilde{z}_i + \lambda \text{ (say)}. \quad (26)$$

It follows from the nature of the amalgamation algorithm that

$$\hat{z}_i = \hat{\tilde{z}}_i + \lambda. \quad (27)$$

Moreover

$$\bar{z} = \frac{1}{n} \sum (\hat{\tilde{z}}_i + \lambda) = \bar{\tilde{z}} + \lambda = \lambda. \quad (28)$$

And thus

$$z_i - \hat{z}_i = \tilde{z}_i - \hat{\tilde{z}}_i, \quad (29)$$

and

$$z_i - \bar{z} = \tilde{z}_i - \bar{\tilde{z}} = \tilde{z}_i. \quad (30)$$

Consequently  $S = \tilde{S}$  for all vectors  $w$ . The minima of two identical functions are, of course, also identical.

The second restriction was that  $y$  must be a vector of real numbers. This restriction was used only in the computation of an initial estimate of  $w$ . It is easy to generalize the approach to an arbitrary weakly ordered set  $\langle Y, \succ_0 \rangle$ . We use the CDARDA technique (De Leeuw 1968) to obtain an initial estimate of  $w$ . Let

$$s_{ij} = \begin{cases} +1 & \text{if } y_i \succ_0 y_j, \\ 0 & \text{if } y_i \approx_0 y_j, \\ -1 & \text{if } y_i \prec_0 y_j. \end{cases} \quad (31)$$

Then we must maximize the coefficient

$$F = \frac{\sum_i \sum_j s_{ij} (z_i - z_j)}{\sum_i z_i^2}. \quad (32)$$

Simplifying the numerator (by using the fact that  $s_{ij} = -s_{ji}$ ) yields

$$\sum_i \sum_j s_{ij} (z_i - z_j) = 2 \sum_i z_i \sum_j s_{ij}. \quad (33)$$

It is easily seen that

$$\sum_j s_{ij} = 2r_i - (n + 1), \quad (34)$$

where  $r_i$  is the rank number of  $y_i$  in the order, and where ties have equal (averaged) rank numbers. Moreover

$$\sum_i \sum_j s_{ij} = 2 \sum_i r_i - n(n + 1) = 0. \quad (35)$$

Writing  $\tilde{r}_i$  for the centered rank numbers we obtain

$$F = \frac{2\tilde{r}'Xw}{w'X'Xw}. \quad (36)$$

The optimal solution for  $w$  is the same as in formula (12), with  $\hat{z} = \tilde{r}$ .

In the case of a partially ordered set the situation is less simple. We let  $s_{ij} = 0$  if the order relation between  $y_i$  and  $y_j$  is undetermined. Formula's (32) and (33) can still be used but the simple relation of  $\sum_i s_{ij}$  with rank numbers is not true. Define  $\tilde{s}_i = \sum_j s_{ij}$  and use  $\tilde{s}$  as  $\tilde{z}$  in (12). It is also possible to adapt the almagamation algorithm for partial orders. In the relevant step of the iterative program we have to solve the quadratic programming problem

$$\begin{cases} (z - \tilde{z})'(z - \tilde{z}) \text{ min !} & (37a) \\ Az \geq 0, & (37b) \end{cases}$$

where  $A$  is a matrix constructed in such a way that any  $\tilde{z}$  satisfying (37b) satisfies the order relations. Perhaps the problem (37) is most easily solved by the Gauss-Seidel type method for quadratic programming proposed by Hildreth and d'Esopo and discussed by Kunzi and Krelle (1962).

The third (and most important) restriction in section 1 was that  $c = X'X$  must be nonsingular. Observe that the only place in which we use this is section 2, where  $C^{-1}$  is computed (only once). Consider the normal equations

$$X'Xw = X'\tilde{z}. \tag{38}$$

These equations always have at least one solution. Proof: a system of linear equations  $Ax = b$  is consistent iff for all  $y$  for which  $y'A = 0$  it is also true that  $y'b = 0$ . Because  $X'X$  is positive semidefinite,  $y'X'Xy = 0$  implies  $y'X'X = 0$ . The converse is obvious. Thus:  $y'X'X = 0$  iff  $y'X'Xy = 0$  iff  $y'X' = 0$ , which implies that  $y'X'\tilde{z} = 0$ . Q.E.D.

It is well known that the solution of (38) for  $w$  is unique iff  $X'X$  is of full rank. If (38) is consistent, but  $\text{rank}(C) < m$  then there are more solutions. A general inverse (g-inverse) of an  $n \times m$  matrix  $A$  is defined as an  $m \times n$  matrix  $A^-$  such that for each consistent system  $Ax = b$  the vector  $x = A^-b$  is a solution. It is not difficult to see that  $A^-$  is the g-inverse of  $A$  iff  $AA^-A = A$  iff  $A^-AA^- = A^-$ . Because the



normal equations are consistent the vector

$$w = (X'X)^{-1}X'z \tag{39}$$

is a solution. The general solution is given by

$$w = C^{-1}X'z + (C^{-1}C - I)b, \tag{40}$$

with arbitrary  $b$ . In our procedure we have to scale the vector  $w$  afterwards in such a way that  $w'X'Xw = 1$ . Some practical ways to compute a  $g$ -inverse of  $C$  in the special case of normal equations are given by Rao (1965). Let the rank of  $X$  be  $r < m$ , then it is possible to construct an  $(m-r) \times m$  matrix  $H$  such that the matrix

$$\tilde{X} = \begin{bmatrix} X \\ \dots \\ H \end{bmatrix} \tag{41}$$

is of rank  $m$ . Then  $\tilde{C} = \tilde{X}'\tilde{X} = X'X + H'H$  possesses a true inverse, and this inverse  $\tilde{C}^{-1}$  is a  $g$ -inverse of  $C$ . Another method (which is especially useful if we also need the canonical form of  $C$  or the Eckart-Young decomposition of  $X$ ) is to write  $C$  as  $K\Lambda K'$ , with  $K'K = I$  and  $\Lambda$  diagonal. Let  $\Lambda^{-1}$  denote the diagonal matrix with elements  $\lambda_{ii}^{-1}$  if  $\lambda_{ii} > 0$ , and 0 otherwise. Then  $\Lambda^{-1}$  is a  $g$ -inverse of  $\Lambda$ , and  $K\Lambda^{-1}K'$  is a  $g$ -inverse of  $C$ . The final method we discuss is to delete the  $m - r$  dependent rows and columns of  $C$  (without loss of generality this may be taken to be the last ones) and to invert the resulting  $r \times r$  matrix  $C_r$ . Then the  $m \times m$  matrix

$$\begin{bmatrix} C_r^{-1} & \vdots & 0 \\ \dots & \dots & \dots \\ 0 & \vdots & 0 \end{bmatrix} \tag{42}$$

is a  $g$ -inverse of  $C$ .

For the similar assumption in the Newton-Raphson equations the situation is less simple. The Newton-Raphson method results from approximating a function by its second order Taylor expansion

$$f(x-h) \approx F(x,h) = f(x) - h'g_x + \frac{1}{2}h'V_x h. \tag{43}$$

A necessary condition for an extreme value of  $F$  is

$$V_x h = g_x. \tag{44}$$

This system is not necessary consistent. If it is, a solution is given by  $V_x^{-1} g_x$ . Otherwise we must go back to the method of section 2.

5: Psychometric Applications.

5.1: Nonmetric multiple linear regression.

This is the standard case, the analogon of ordinary multiple regression. The matrix X is measured in the ordinary way (at least on an interval scale), the 'dependent' variable y is measured on an ordinal scale.

5.2: Nonmetric polynomial regression analysis.

The matrix X is constructed by using as columns either powers of an independent variable x, or the orthogonal polynomial coefficients. Again y is measured on an ordinal scale. Observe that the polynomial of degree zero (a constant) is irrelevant and must not be included in the analysis. If we have two independent variables  $x_1$  and  $x_2$  the columns of X can be constructed by using the vectors  $x_1^2, x_2^2, x_1 x_2, x_1,$  and  $x_2$ . In this case we fit a quadratic function to the (ordinal) data. This approach generalizes easily to all 'separable' functions of p variables.

5.3: Nonmetric discriminant analysis.

In the metric case discriminant analysis can be considered as a degenerate form of multiple regression in which the 'dummy' dependent variable consists entirely of zeroes and ones. In the nonmetric case the same thing is true. In stead of requiring maximum distance between the two populations we require minimum overlap (and of course we use Kruskal's primary approach to ties in this case).

5.4: Nonmetric canonical discriminant analysis.

We adapt our algorithm of section 2 to the case where there are p groups. We require that the projections Xw of the groups are pairwise disjoint. To do this we form  $\binom{p}{2}$  vectors  $\hat{z}_{st}$  ( $s < t$ ), we let  $\hat{z} = \sum_s \sum_t \hat{z}_{st}$ , and we minimize

$$F = \sum_{s < t} \sum \frac{w'X'Xw - 2 w'X'z_{st} + z_{st}'z_{st}}{w'X'Xw} \quad (45)$$

This is equivalent (with the  $z_{st}$  fixed) to maximizing

$$F = \frac{w'X'z}{w'X'Xw}, \quad (46)$$

and no new aspects have to be considered.

6.5: Additive conjoint measurement.

In additive conjoint measurement the matrix X must be constructed in a rather special way. Consider the case of a two-factor design, with facets A and B. A has l elements (levels, structs), B has m - 1 elements. For the  $d_{ij} = (a_i, b_j)$  element of the data structure  $D = A \times B$  we construct a row of X as follows: the row  $x_k$  has m elements, both  $x_{ki}$  and  $x_{k,l+j}$  are equal to unity, all other elements of  $x_k$  are zero. It follows that

$$x_k \cdot w = w_i + w_{l+j}, \quad (47)$$

which is what we want. Generalization to more than two factors is obvious.

Because of the special structure of X the equations for the iterative program can be simplified. We develop them again for the two-factor case with no missing data. Facet A has k elements, B has l elements,  $k + l = m$ ,  $kl = n$ . The scale values are collected in  $w' = [u' : v']$ . Moreover

$$z_{ij} = u_i + v_j, \quad (48)$$

$z_{ij}$  is the corresponding monotone matrix,  $\sum_i \sum_j z_{ij} = 0$ . We must minimize

$$F = \frac{\sum \sum (z_{ij} - z_{ij})^2}{\sum \sum z_{ij}^2}, \quad (49)$$

which is equivalent to maximizing

$$F' = \frac{\sum \sum z_{ij} z_{ij}}{\sum \sum z_{ij}^2} = \frac{2 \sum \sum z_{ij} (u_i + v_j)}{2 \sum (u_i + v_j)^2}. \quad (50)$$

The derivatives of the Lagrangian function are ( $\lambda$  as Lagrange multiplier)

$$\frac{\partial \phi}{\partial u_s} = 2 \sum_j z_{sj} - 2 \lambda u_s - 2 \lambda \sum_j v_j, \quad (51a)$$

$$\frac{\partial \phi}{\partial v_t} = 2 \sum_i z_{it} - 2\lambda v_t - 2\lambda \sum_j u_s. \quad (51b)$$

Let  $\tilde{u}_s = \sum_j z_{sj}$  and  $\tilde{v}_t = \sum_i z_{it}$ , then  $\sum_s \tilde{u}_s = \sum_t \tilde{v}_t = \sum_i \sum_j z_{ij} = 0$ . It follows that, if  $\tilde{z}_{ij} = \tilde{u}_i + \tilde{v}_j$  and  $\gamma = \left( \sum_i \sum_j \tilde{z}_{ij}^2 \right)^{-\frac{1}{2}}$ , then  $u = \gamma \tilde{u}$ , and  $v = \gamma \tilde{v}$ , are a solution to (51). The theory in this section can easily be translated into a matrix formulation using the g-inverse (we add some rows to  $X$  to make up for the deficiency in rank). Moreover it can equally easily be generalized to cases with partial orders, zero-one responses, non-numerical responses, and multifactor designs. Quite similar simplifications are possible in the equations of the Newton-Raphson method.

7: References.

- Browne, M.                      Oblique Procrustus rotation.  
Psychometrika 1967.
- De Leeuw, J.                    Canonical discriminant analysis of relational data.  
Department of data theory, Leyden University, RN006-68
- Fan, K.                            Systems of linear inequatilities  
Annals of mathematics studies, no 38, 1956.
- Kunzi, H.P. &  
Krelle, W.                        Nichtlineare Programmierung  
Berlin, Göttingen, Heidelberg, Springer-Verlag, 1962.
- Scott, D.                         Measurement structures and linear inequalities.  
Journal of mathematical psychology, 1964.
- Tversky, A.N.                    Finite additive structures.  
Michigan mathematical psuchology program, 1964.

Appendix:

Upon further investigation the method proposed in formula (24) turns out to be somewhat of a hoax. The proposal was

$$w^{(s+1)} = w^{(s)} - \frac{(T^{\#})^{(s)}}{2(1 - S^{(s)})} (X'X)^{-1} g^{(s)}. \quad (52)$$

From (14), (15), and (20)

$$g^{(s)} = \frac{2}{(T^{\#})^{(s)}} \left\{ (1 - S^{(s)}) X'X w^{(s)} - X'z^{(s)} \right\}. \quad (53)$$

Substitution of (53) into (52) yields

$$w^{(s+1)} = \frac{1}{(1 - S^{(s)})} (X'X)^{-1} X'z^{(s)}, \quad (54)$$

and thus the method of (24) is identical to the method of section (2). This shows that the faster convergence of the Newton-Raphson method is due entirely to the terms  $g_{1k} e_k + g_{k1} e_1$  in the expression for  $v_{kl}$ . Nevertheless, if we are close to the minimum, these terms must necessarily be quite small.