# An Algorithm for Nonlinear Principal Components Analysis with B-splines by Means of Alternating Least Squares

H. Coolen, J. van Rijckevorsel, Rotterdam, and J. de Leeuw, Leiden

The algorithm described in this paper is suitable for nonmetric principal components analysis with B-splines, multiple correspondence analysis with B-splines as well as a mixture of the two. For the theory of these different forms of nonlinear principal components analysis and further literature see Gifi(1981), De Leeuw(1982) and the companion paper by Van Rijckevorsel(1982).

In nonlinear principal components analysis we have to minimize a least squares loss function $\sigma$:

$$\sigma = \Sigma_{j=1}^{m} \, \mathrm{tr} \, (X - G_j U_j)'(X - G_j U_j)$$

over $X$ and $U_j$ with normalization $X'X = I$. Here $G_j$ is a (pseudo)-indicator matrix with $k$(order of the spline) nonzero entries per row. Also define $D_j = G_j'G_j$. For single variables we require in addition that $U_j = a_j t_j'$, where $a_j'D_j a_j = 1$ and $a_j$ is restricted to lie in a given convex cone $K_j$. Minimizing the loss function $\sigma$ can be done by alternating least squares using direct iteration. Each iteration cycle consists of five or ten substeps; in each substep the loss $\sigma$ is minimized over one set of parameters for fixed values of the other set(s). Each iteration starts with $X^0$ and $U_j^0$ and gives updates $X^+$ and $U_j^+$ which conditionally minimize the loss $\sigma$(De Leeuw and Van Rijckevorsel (1980)). The main steps in the algorithm are:

$$\text{(1)} \qquad U_j^0 = D_j^{-1} G_j' X^0$$

$$\text{(2)} \qquad \sigma_j = p - \mathrm{tr} \, U_j^0{}'D_j U_j^0$$

$$\text{(3)} \qquad a_j^0 = P_j(U_j^0 t_j^0) \qquad \qquad , \text{ where } P_j \text{ projects on } K_j$$

$$\text{(4)} \qquad a_j^+ = a_j^0 (a_j^0{}'D_j a_j^0)^{-\frac{1}{2}}$$

$$\text{(5)} \qquad t_j^+ = a_j^+ U_j^0$$

$$\text{(6)} \qquad U_j^+ = a_j^+ t_j^+{}'$$

$$\text{(7)} \qquad \sigma_j = \sigma_j - t_j^+{}' t_j^+$$

$$\text{(8)} \qquad Z = \Sigma_{j=1}^{m} \, G_j U_j^+$$

(9)    Z in deviations from column means

(10)    $X^+ = Z(Z'Z)^{-\frac{1}{2}}$.

The first step of the algorithm consists of computing the category quantifications for every variable on p dimensions simultaneously. The startconfiguration $X^0$ is a normalized random configuration. $G_j$ is constructed rowwise by a subprogram that uses the recurrence relation for computing B-splines taking their small support into account(De Boor, 1978). $D_j$ is banded symmetric positive definite having bandwidth 2k-1. A Cholesky factorization of $D_j$, by the method of Gauss elimination adapted to the symmetry and bandedness of $D_j$, gives us $S_j$ which is lower triangular(De Boor, 1978). The system $D_j U_j^0 = S_j S_j' U_j^0 = G_j' X^0$ is solved first for $S_j' U_j^0$ by means of forward substitution and then for $U_j^0$ by backward substitution. In steps (3), (4) and (5) we are looking for the minimizing updates $a_j^+$ with $X^0$ and $t_j^0$ fixed, and for $t_j^+$ with $a_j^0$ and $X^0$ fixed. Finding $a_j^+$ is a cone regression problem, while computing $t_j^+$ amounts to an ordinary least squares problem; both problems have a unique solution(De Leeuw and Van Rijckevorsel, 1980). Finding $X^+$ which minimizes the loss for fixed $a_j^+$ and $t_j^+$ is an orthogonal procrustes problem that is solved here, with essentially the same results, by Gram-Schmidt because it is less expensive. Steps (1) through (7) are executed subsequently per variable. For multiple variables step (8) computes $G_j U_j^0$. In the case of multiple correspondence analysis, all variables being multiple, only steps (1), (2), (8), (9) and (10) are alternatingly computed. When the difference in loss between two successive iterations no longer exceeds a predetermined criterion, execution of the algorithm stops after step (2) or step (7). After this the solution is rotated to its principal components and the corresponding eigenvalues are computed. All other quantifications are recomputed with the rotated solution. For some results on convergence see De Leeuw and Van Rijckevorsel(1980).

In order to investigate practical examples we have tested an APL-version of the algorithm. The results can be found in Van Rijckevorsel (1982). A FORTRAN-version is available from Vakgroep M&T, FSCW, Erasmus University Rotterdam, Postbus 1738, Rotterdam, Netherlands.

### References

De Leeuw J. (1982), Nonlinear principal components analysis, Compstat 82, Physica, Wien.
Van Rijckevorsel J.L.A. (1982), Canonical analysis with B-splines, Compstat 82, Physica, Wien.
For all other references see Van Rijckevorsel(1982).

SPSS-X is the newest re
number of enhancements
the syntax rules, and a
allows for a variety of

1) GROUPED data
formats are read and co
standard multiple-recor
within a case may be in
some input records are

FILE TYPE    GROUPED    F
RECORD TYPE 1
DATA LIST  /MOHIRED YRH
RECORD TYPE 2
DATA LIST  /SALARY79 TO
        PROMO81 72
RECORD TYPE 3
DATA LIST  /JOBCAT 6 NA
END FILE TYPE

Commands start in colu
FORTRAN formats may be

DATA LIST  /JOBCAT NAME
    or as
DATA LIST  /JOBCAT 6 NA

The scratch variable (#
each record. For exampl
with the second DATA LIS
reported as errors.  If
duplicate record type or
a warning message be pri

2) MIXED data fil
of several different dat

FILE TYPE    MIXED    FIL
.        RECORD TYPE 21,22
.        DATA LIST  /SEX 5
.            RECORD TY
.            DATA LIST
END FILE TYPE

Note that the '.' (or a
the first line of a comm
variable is kept for lat