

CLUSTER ANALYSIS

R. Capra, A. Lena, U. Santarelli and P. Vescovi

Cluster Analysis of a multivariate sample X:

$$X = \{ \underline{x}_1 \mid \underline{x}_1 \in R^L, i=1, \dots, N \}$$

is the partitioning of X into M homogeneous groups according to a minimization criterion E based on the distance d:

$$d : R^L \times R^L \rightarrow R^+$$

The algorithm elaborated has the following operational characteristics:

- Sequential organization of data.
- No interdistance matrix.
- Processing times linearly dependent on N.
- No effective limitation on N, L, M.
- Euclidean distance.

Moreover, in the Fortran implementation of the algorithm, it is possible to have:

- Standardization of the input data (in order to avoid scale effect).
- Weighting of the variables (in order to force a scale effect).

The minimization criterion used to measure the nonhomogeneity of a partition of X into groups X_1, X_2, \dots, X_M is given by the function:

$$E = \sum_{k=1}^M E_k$$

where;

$$E_k = \sum_{\underline{x}_i \in X_k} d^2(\underline{x}_i, \underline{b}_k)$$

and,

$$\underline{b}_k = \frac{1}{n_k} \sum_{\underline{x}_i \in X_k} \underline{x}_i$$

This function, in addition to its clear geometric significance - total sum of the moments of inertia of the systems of points in relation to its own barycenter \underline{b}_k - has useful computational properties.

It can be seen that if $\underline{x}_i \in X_k$, we have:

$$E_k^- = E_k(X_k) - E_k(X_k - \{\underline{x}_i\}) = \frac{n_k}{n_k - 1} d^2(\underline{x}_i, \underline{b}_k)$$

Likewise, if $\underline{x}_i \notin X_k$, we have:

$$E_k^+ = E_k(X_k \cup \{\underline{x}_i\}) - E_k(X_k) = \frac{n_k}{n_k + 1} d^2(\underline{x}_i, \underline{b}_k)$$

It follows that, for each partition, the ΔE variation of E , brought about by the shifting of one element \underline{x}_i from group X_k to group X_h :

$$\Delta E = E_h^+ - E_k^-$$

depends only on:

$$\underline{x}_i, \underline{b}_k, n_k, \underline{b}_h, n_h.$$

In this way it is possible to have in core only one observation at once.

A partition is considered stable when, for every \underline{x}_i of X , we have:

$$\Delta E > 0$$

In this case, we see that if $\underline{x}_i \in X_k$:

$$d(\underline{x}_i, \underline{b}_k) \leq \min d(\underline{x}_i, \underline{b}_h) \quad \text{for } h \neq k,$$

with \underline{b}_h' barycenter of $X_h \cup \{\underline{x}_i\}$.

Since it is not possible to guarantee that a stable partition is the best, because of the dependence of the classification on the order in

CLUSTER ANALYSIS - Continued

which the data are presented, the algorithm envisages the subsequent application of a C procedure so as to reach a solution independent of the initial sequence of the data while maintaining performance unchanged in relation to the E criterion.

The C procedure sorts the set $Y = \{Y_i | Y_i \in R^L, i=1, \dots, N\}$

into ng groups, utilizing the initial ng barycenters of B_0 :

$$B_0 = \{b_1^0, \dots, b_{ng}^0\},$$

as specified by an Isw switch, and placing in B,

$$B = \{b_1, \dots, b_{ng}\},$$

the barycenters of the classification is obtained.

This gives:

$$B \leftarrow C(Y, ng, B_0, Isw).$$

C operates in the following way:

a) Make $B=B_0$ and $n_k=0, Y_k=\emptyset, k=1, \dots, ng$

b) If Isw=0 go to e)

c) Compute $d(Y_i, b_k^0) = \min_{h \neq k} d(Y_i, b_h^0)$

for every i

$$\text{make } Y_k = Y_k \cup \{Y_i\};$$

$$\text{make } n_k = n_k + 1$$

d) Compute $b_k = \frac{1}{n_k} \sum Y_i$

for every k

e) Compute

$$E_k^- = \frac{n_k}{n_k - 1} d^2(Y_i, b_k)$$

for every i

$$E_h^+ = \frac{n_h}{n_h + 1} d^2(Y_i, b_h)$$

for every i and for $h \neq k$

$$\text{If } E_k^- \leq \min_h E_h^+,$$

Y_i remains in Y_k

$$\text{If } E_{h_0}^+ = \min_h E_h^+$$

and $E_{h_0}^+ < E_k^-$ move Y_i

from Y_k to Y_{h_0}

and make:

CLUSTER ANALYSIS - Continued

$$Y_k = Y_k - \{Y_1\}; \quad b_k = (b_k n_k - Y_1) / (n_k - 1); \quad n_k = n_k - 1;$$

$$Y_{h_0} = Y_{h_0} \cup \{Y_1\}; \quad b_{h_0} = (b_{h_0} n_{h_0} + Y_1) / (n_{h_0} + 1); \quad n_{h_0} = n_{h_0} + 1;$$

- f) If any Y_1 has been moved (E is reduced) go to e); otherwise stop.

The algorithm works in the following way:

- a) Make $ng = \sqrt{MN}$, $B_0 = \{0, 0, \dots, 0\}$

- b) Compute $B' = C(X', ng, B_0, 0)$

where X' is possibly a subsample of X if N/M is considered excessive depending on economic considerations.

- c) Make $ng=M$.

Choose $\bar{B} \subset B'$ card $(\bar{B}) = ng$ in the following way:

- 1) $k = 1$; $\bar{B} = (b'_1)$ 2) $k = 2, \dots, ng$; compute:

$$d(b'_{i_0}, \bar{b}_{j_0}) = \max_{b'_i \in B' - \bar{B}} \min_{\bar{b}_j \in \bar{B}} d(b'_i, \bar{b}_j); \quad \bar{B} = \bar{B} \cup (b'_{i_0}).$$

This makes it possible to choose the sparsest subset of M barycenters in B' .

- d) $B'' = C(B', M, \bar{B}, 0)$

- e) $B = C(X, M, B'', 1)$

which performs the final classification on the entire X sample, utilizing B'' as seeds of the classification.

Because of the absence of specific tests to evaluate the performances of Cluster Analysis Algorithms, tests were based upon the comparison of the automatic solution and the subjective classification of known samples.

The similarity between the clustering results achieved by the two methods confirms the reliability of the algorithm, at least when it is possible to give a coherent Euclidean interpretation of the subjective similarity criterion.