

1. Introduction

Cluster analysis as a quantitative methods of homogeneous group formation, can be easily and fruitfully applied to the regression modelling process. At the first stage of regression model building, especially in socio-economic researches, we have to define the aim and the scope of empirical investigations. The exact definition of the aim implies a concrete class of the model. The definition of the scope has a triple character. First, it is a definition of essential, second - spatial and third-time scope of empirical investigations. The definition of the scope of investigation involves many problems to be solved by statisticians.

The essential scope of the definition of regression modelling depends on the specification of endogenous and exogenous variables. The specification of the variables can be based on an adequate theory. On the other hand, there are some formal requirements with regard to variables.

It is the matter that in statistics and econometrics the variable selection problem is widely discussed. Among many normal methods which have been proposed there exists an extensive group of cluster analysis methods. The problem of variable selection is one of the fields where the cluster analysis methods and regression modelling meet.

Another problem is the construction of a model for some aggregates of variables. One possible approach to the problem is to build some models for synthetic variables. The creation of synthetic variable bases on application of cluster analysis methods, especially linear ordering methods. It is the second common level of cluster analysis and regression theory.

One of the basic conditions of correctness and effectiveness of regression modelling is construction of a model for homogeneous statistical population. Only in this type of population, relations among variables under consideration will have statistically synonymous and stable character.

Effective methods for obtaining such homogeneous populations are cluster analysis method. Problems of how to single out homogeneous statistical unit subsets is the third field where cluster analysis and regression modelling meet.

The fourth level is application of cluster analysis methods for verification and interpretation of the results of regression modelling.

Indices of taxonomic similarity could be applied to the measurement of similarity among structural parameters or stochastic structure parameters of the estimated models.

An alternative name of cluster analysis, more frequently used in socio-economic researches in Poland, is numerical taxonomy. We understand "cluster analysis" and "numerical taxonomy" or simply "taxonomy" as synonyms in this paper.

Econometrics is the most important application field of the regression theory in socio-economic researches. Then the regression model will be interpreted as an econometric model in this paper.

Econometric models built with the aid of

taxonomic methods are also called taxonomic models. In the following parts of the paper taxonomic methods of homogenization of statistical population, taxonomic procedures of explanatory variable selection, methods of synthetic variable construction and taxonomic methods of verification and interpretation of the results of econometric modelling are discussed.

2. Taxonomic methods of homogenization of statistical population

The homogeneity of statistical population is the basic condition if satisfactory results of regression modelling are to be obtained, such homogeneity can be understood as either spatial or time homogeneity. It leads to the application of cluster analysis methods for obtaining such homogeneous subsets of statistical units.

The regression function is equivalent to the approximation function and it is possible to utilize some methods from the approximation theory. The basic problem of approximation is to define the function which describes the relation among the variables considered, based on information from the sample:

$$(y_1, x_{11}, x_{21}, \dots, x_{m1}), (y_2, x_{12}, x_{22}, \dots, x_{m2}), \dots \\ \dots, (y_n, x_{1n}, x_{2n}, \dots, x_{mn}) \quad (1)$$

where:

$$j=1, 2, \dots, m - \text{number of variables,} \\ i=1, 2, \dots, n - \text{number of observation.}$$

In functions $h(x_1, x_2, \dots, x_m)$ are undistinguishable if:

$$|\Psi(x_1, x_2, \dots, x_m) - h(x_1, x_2, \dots, x_m)| < \Delta, \quad (2)$$

where:

$$x_1 \in [a_1, b_1], x_2 \in [a_2, b_2], \dots, x_m \in [a_m, b_m];$$

$$\Psi(x_1, x_2, \dots, x_m) \in R_k;$$

$$h(x_1, x_2, \dots, x_m) \in H_k; k \in \mathbb{N};$$

Δ - tolerance (admissible error of approximation),

R_k - a set approximated function (first type of regression),

H_k - a set approximating functions (second type of regression).

The stochastic approximation problem is to find a family of R_k - functions which satisfy the relation:

$$P\{|\Psi(x_1, x_2, \dots, x_m) - \Delta < Y < \Psi(x_1, x_2, \dots, x_m) + \Delta|\} \geq 1 - \alpha, \quad (3)$$

where:

$$1 - \alpha - \text{the likelihood of approximation.}$$

Relation (3) can be more generally described as:

$$P\{(Y, x_1, x_2, \dots, x_m) \in Q\} \geq 1 - \alpha, \quad (4)$$

where:

Q - space limited by $V_1(X_1, X_2, \dots, X_m)$ and $V_2(X_1, X_2, \dots, X_m)$ also by (a_1, a_2, \dots, a_m) and (b_1, b_2, \dots, b_m) .

In socio-economic researches it is difficult to analyse multidimensional distributions and for that reason it is convenient to define the spectrum and trace of distribution. The spectrum of distribution in the set spaces $Q\alpha_i$ in which for sequence $\alpha_1, \alpha_2, \dots, \alpha_i, \dots$, where $Q < \alpha_i < 1$,

The following relations are observed:

$$P[(Y, X_1, X_2, \dots, X_m) \in Q\alpha_i] = 1 - \alpha_i \quad (5)$$

and

$$f(y_i, x_{i1}, x_{i2}, \dots, x_{im}) > f(y_i, x_{i1}, x_{i2}, \dots, x_{im}), \quad (6)$$

where:

$$(Y_i, X_{i1}, X_{i2}, \dots, X_{im}) \in Q\alpha_i;$$

$$(Y_i, X_{i1}, X_{i2}, \dots, X_{im}) \in \bar{Q}\alpha_i;$$

$\bar{Q}\alpha_i$ is a complement of $Q\alpha_i$.

The trace of distribution is defined as:

$$P[(Y, X_1, X_2, \dots, X_m) \in Q] = 1 - \alpha, \quad (7)$$

where α is near 0.

For calculation of the distribution trace a variable can be utilized. Some information about the trace of distribution can also be obtained by application of cluster analysis methods.

Application of cluster analysis methods leads to homogeneous subsets of observations which could be understood as a distribution trace. It ought to be of ellipsoidal shape in multidimensional space if the regression model is a linear one. For separation of an ellipsoidal distribution trace from a set of data, the deviation method proposed by pluta (1986) can be employed. In the first stage a set of variables must be divided into two subsets:

- stimulant variables,
- destimulant variables.

Such a variable is recognized as a stimulant variable for which the greater the value the better the situation.

Destimulant variable is a reverse to stimulant.

Then destimulants into stimulants are transformed by the formula:

$$X_{ij} = -Y_{ij}, \quad (8)$$

where:

Y_{ij} - is i-observation of j variable which is recognized as a destimulant.

In the next stage the upper lower pole of the set of observation is defined. Coordinates of the upper and lower pole are maximal and minimal

values of variables respectively.

When the lower pole of the coordinates is shifted to the lower pole, it means that X_{ij} is transformed into u_{ij} according to the

formula:

$$u_{ij} = x_{ij} - x_{oj}, \quad (9)$$

where:

x_{oj} - are coordinates of the lower pole.

A line through the lower and upper pole is drawn and this is named the axis of the set. Deviations of observations from the axis of the set inform us about the shape of the distribution trace. Subsequently perpendicular projections of observations on the axis of the set are calculated. The coordinates of these projections are:

$$R_i = [X_{i1}, X_{i2}, \dots, X_{ij}, \dots, X_{im}], \quad (10)$$

where:

$$X_{ij} = x'_{oj} t_i, \quad (11)$$

x'_{oj} - coordinates of the upper pole,

$$t_i = \frac{\sum_{j=1}^m x'_{oj} u_{ij}}{\sum_{j=1}^m (x'_{oj})^2} \quad (12)$$

Then the measures M^* and W^* are defined and calculated as follows:

$$M^* = [(R_i - P)(R_i - P)^T]^{1/2} \quad (13)$$

$$W^* = [(P_i - R_i)(P_i - R_i)^T]^{1/2}, \quad (14)$$

where:

P - coordinates of the lower pole,

P_i - coordinates of the real data.

Values of M^* and W^* are presented on the graph, where M is the abscissa and W is the ordinate. This graph presents the distribution trace. If it is an ellipsoidal shape a linear regression function can be estimated, if it is not, ellipsoidal subsets must be separated by application of isomorphous subset procedure proposed by pluta (1986).

3. Taxonomic selection of variables

The most universal measure of taxonomic similarity among variables is the correlation coefficient. Among many taxonomic procedures of selection of variables, based on the correlation matrix the method proposed by pluta (1972) has interesting properties. Let:

$$R = [r_{jk}]; \quad j, k = 1, \dots, m, \quad (15)$$

be a correlation matrix among m preliminary proposed explanatory variables.

$$R_o = [r_{oj}]; \quad j = 1, \dots, m, \quad (16)$$

be a vector of correlation coefficients between

an endogenous variable and the proposed explanatory variables.

In the next stage, correlation coefficients are tested using the following statistics :

$$r^{**} = \left[\frac{t_{\alpha}^2}{t_{\alpha}^2 + n - 2} \right]^{1/2} \quad (17)$$

where:

n - number of observations,

t_{α} - critical value of the student's distribution for n - 2 degrees of freedom at level of significance .

These variables for which $r_{\bullet j} \leq r^{**}$ are eliminated from the vector R_{\bullet} .

The variable for which $r_{\bullet j}$ has the highest value is chosen as the first explanatory variable. The following chosen variables are those, which are those, which are successively maximally correlated with the endogenous variable and insignificantly correlated with the previously chosen explanatory variables . This method has a simple graph interpretation and it can be presented as a symmetric, full and not-oriented graph.

Another method based on the correlation matrix, which also has a graph interpretation has been proposed by Bartosiewicz (1976). On the base of the correlation matrix R a graph G is constructed the proposed explanatory variables are nodes of the graph and significant correlation coefficients between variables are arrows of the graph. Those variables with nodes of graph with maximal number of connection with other nodes and also variables represented by isolated nodes of the graph belong to the optimal vector of explanatory variables .

4. Methods of synthetic variable construction

Problems of aggregation are very essential in socio-economic research, especially in construction of forecasting models, input-output analysis, operation research and cluster analysis. An original method of aggregation in economics, i.e. method of synthetic variable construction has been proposed by Hellwing (1968). On the base of Hellwing's proposition many methods of synthetic variable construction have been explained. They are based on the following criteria of

- classification of synthetic variables :
- the way of allowance of stimulant and destimulant variable,
 - the way of defining of the point of reference of coordinates,
 - the way of normalization of variables,
 - analytic form of the aggregative function,
 - weighting system of importance of variables.

Among the methods of construction of synthetic variables there are procedures based either only on stimulants or only on destimulants or on both types of variables. The reference point of the coordinates could be chosen on the basis of expert's opinions, international comparisons and in a statistical way. A hypothetical point with maximal observed values of stimulants and minimal

observed values of destimulants can be adopted as a statistical pattern point.

There are following methods of normalization of variables :

- 1) rank method, which changes the observed values into their ranks,
- 2) quotient transformations,
- 3) standardizations,
- 4) unitarizations.

The following functions can be chosen as aggregative functions:

A - additive :

$$S_i^{(1)} = \sum_{j=1}^m \alpha_j x'_{ij} \quad , \quad (18)$$

$$S_i^{(2)} = \frac{1}{m} \sum_{j=1}^m \alpha_j x'_{ij} \quad , \quad (19)$$

$$S_i^{(3)} = \frac{\sum_{j=1}^m \alpha_j}{\sum_{j=1}^m \frac{\alpha_j}{x'_{ij}}} \quad , \quad (20)$$

where:

- α_j - weight of variable j ,
- x'_{ij} - normalized value of variable j .

B - multiplicative:

$$S_i^{(4)} = \prod_{j=1}^m (x'_{ij})^{\alpha_j} \quad (21)$$

$$S_i^{(5)} = \left[\prod_{j=1}^m (x'_{ij})^{\alpha_j} \right]^{\frac{1}{\sum_j \alpha_j}} \quad (22)$$

Most frequently every variable is assumed to carry the same weight .

Among many possibilities of synthetic variable construction two of them are most popular standardized value method and development pattern method . synthetic variable in the standardized value method can be defined by the formula :

$$S_i^{(sv)} = \frac{1}{m} \sum_{j=1}^m \frac{x_{ij} - \bar{x}_j}{S_j} \quad . \quad (23)$$

In the pattern method it can be defined by:

$$S_i^{(dp)} = \left[\frac{1}{m} \sum_{j=1}^m (x'_{ij} - x'_{\bullet j})^2 \right]^{1/2} \quad (24)$$

where:

- x'_{ij} - standardized value of x_{ij} ,
 $x'_{\bullet j}$ - standardized value of the pattern .

Another approach to aggregation of "simple" variables is formation of a homogeneous aggregated variable by application of Hellwig's stochastic dependence coefficient⁹. It is done by the formula :

$$S = \left[\frac{1 - \sum_{i,j} \min(p_{ij}, p_i, q_j)}{1 - \max(\sum_i p_i^2, \sum_j q_j^2)} \right]^{1/2} \quad (25)$$

where :

- $p_{ij} = P(X=x_i, Y=y_j)$;
 $p_i = P(X=x_i); i=1, \dots, r$;
 $q_j = P(Y=y_j); j=1, \dots, s$;
 r - number of rows in the contingency table ,
 s - number of columns ,
 $0 \leq S \leq 1$ (26)

For a variable, which is supposed to be an aggregate of m elementary variables, the matrix Δ of dependence can be calculated:

$$\Delta = \left[\delta_{ki} \right], \quad k, i = 1, \dots, m \quad (27)$$

Any cluster analysis method could be applied for obtaining homogeneous subsets of elements of matrix Δ . These subsets could be recognized as homogeneous aggregative variables .

REFERENCES

=====

- 1- A. Koutsoyiannis . (1983) Theory of econometrics (M. P.) .
- 2- Bartosiewicz . (1976) Econometrics (P. W. B)
- 3- G.C. Chow (1985) Econometrics (M.H) .
- 4- G.S. Maddala (1977) Econometrics (E.H) .
- 5- H.J. Brennan (1973) Preface to econometrics .
- 6- Norman Draper (1961) Applied Regression analysis (J.WO) .
- 7- R.Sokal & P.A. Sneath (1963) Numerical taxonomy (U S A) .

8- S. Wheelwright (1977) forecasting methods . (J.W.) .

9- Therl (1971) principles of Econometrics (J.W.) .