# Cluster Analysis in International Comparisons

*György Szilágyi*[1]

**Abstract:** This article investigates a possible application of distance measures and cluster analysis for international comparisons of macroeconomic structures. Its focus is the economic interpretation of various magnitudes which can be obtained by clustering. Theoretical and methodological issues are discussed along with a numerical illustration, covering twenty European countries and using eleven variables.

The main use of distance measures and cluster analysis is (i) the automatic process of grouping countries according to structural similarities, and (ii) the identification of the correspondence of the various magnitudes of the variables.

**Key words:** Economic statistics; cluster analysis; Euclidian distance; international comparisons.

## 1. Introduction

The subject of this study is the use of distance measures and cluster analysis in macroeconomic international comparisons. It is not our intention to develop new methodologies or modifications of existing methods, which are described in the literature; for instance, Dubes and Jain (1979), Gordon (1981), Hartigan (1975). Our basic concern is the economic interpretation of the various indicators obtained by this technique and the economic analysis based on them, rather than the technique itself. Consequently, the different distance measures and clustering strategies are not discussed, but the most common and simplest measures are used (Euclidian distance, centroid method).

Cluster analysis is a tool for comparing the economic structures of various countries. It classifies the countries according to their economic (and social) characteristics, which

are described by a set of variables such as production or consumption of selected goods and services, the stock of various producers' and consumers' goods, etc. A prerequisite to the use of cluster analysis in international comparisons is the availability of a set of variables which are internationally comparable for a large group of countries. Recently, international official statistics has produced a large selection of data which meet these requirements: the United Nations International Comparison Project (ICP). Its European counterpart, the European Comparison Programme (ECP) of the UN/ECE, produced the data set which is used in this article (United Nations 1988). This programme covered 20 European countries with 1985 as the reference year. The basic objective of the ICP and ECP was to compare Gross Domestic Product (GDP) per capita in real terms, together with the various subaggregates like consumption, investment, purchasing power parities, etc. United Nations (1988) provides a rich

[1] Hungarian Central Statistical Office, Keleti Károly u. 5–7, H-1024 Budapest, Hungary.

analysis of magnitudes, structure, inter-relationship of various components, effects of income level, and prices on consumption patterns, etc. As, however, statistical multi-variate analysis was not used in the ECP, it should be of some interest to learn what these methods can contribute to international comparisons.

## 2. A Summary Outline of the Method

Let us consider $m$ countries and $n$ variables, which characterize the different aspects of the countries' economies. To eliminate any effect of differences in measurement units, we standardize each variable so that its mean equals 0 and its variance 1. Let $z_{ji}$ be the standard value of variable $i$ in country $j$. The data are thus arranged in a matrix $\mathbf{Z} = [z_{ji}]$, of the size of $(m \times n)$. With the help of the vectors of any two countries ($h$ and $j$) in this matrix, a distance measure $d_{hj}$ can be constructed. To eliminate the effect of the different number of variables used in the different calculations, the average Euclidian distance, i.e., the usual Euclidian distance divided by the square root of the number of variables ($n$), is used in this study. The distance equals 0 if and only if the value of each variable is the same in country $h$ as in country $j$ or if $h = j$ (as the distance of a country from itself equals 0).

In international comparisons, the distance measure can be considered *the quantitative measure of dissimilarity between the economic structures of two countries*. The smaller the distance, the greater the similarity in the economic structures of the two countries. The set of distance measures constitutes the distance matrix $\mathbf{D} = [d_{hj}]$ with size $(m \times m)$. This is the basis for all further operations of cluster analysis – in our case hierarchical classification, agglomerative procedure, and the centroid method. This means that $c_{gi}$, the cluster value of the $i$th

variable in cluster $g$, will be the average value of variable $i$ for the countries contained in cluster $g$.

The $g$th cluster is therefore defined by the transposed vector

$$\mathbf{c}_g = [c_{g1}, \ldots, c_{gi}, \ldots, c_{gn}]'. \tag{1}$$

The distance between a cluster and a country, or between two clusters, are defined again by the average Euclidian distance. In this technique the ordering of the clusters goes from individual countries to larger and larger groups of countries, so that those countries or country groups enter the next cluster, which are separated by the smallest distance.

## 3. Theoretical Considerations Regarding the Selection of Variables

Cluster analysis is one method that supplies new information on and special insight into the economies of countries by the use and combination of the selected variables. Consequently, the content, results, interpretation, and explanatory power of the comparison depend largely on the selection of the variables.

Every variable describes the economy from some aspect – essential or less essential, – in an elementary, partial, or complex manner. The number of possible sets of economically meaningful variables is extremely large, even if we consider the limits of statistical availability. However, any study can involve only a limited number of variables.

One of the principles generally accepted for the selection is diversification: The variables should characterize the phenomenon from as many aspects as possible. For example, in a macroeconomic comparison, the greatest possible number of sectors, or components of economic development or components of living standards should be represented by the variables. Correlation

among the variables has an outstanding and multi-faceted importance in the selection. From this point of view that set of variables is optimal, in which

    i. the correlation among the selected variables is small and;

    ii. the correlation between the selected and not selected variables is large.

The second criterion is important in that it insures maximum information content from the selected set of variables.

The first requirement does not only reduce multicollinearity (which, generally, is not of great importance in cluster analysis) but avoids overemphasis of one or another sector or aspect of the economy to the omission or underemphasis of others. The second requirement implies the necessity of selecting variables which are characteristic of the given field. For instance, when one or two variables from the field of transportation are selected, efforts should be made to select ones that characterize transportation in the best possible manner, which is tantamount to saying that they should be in the closest possible relationship with the greatest number of aspects of transportation (quantity transported, distance, speed, comfort, reliability, technical level, etc.). One of the difficulties of selecting variables is precisely that we usually have little a priori information, and so we are often compelled to make assumptions. Consequently the selection is necessarily somewhat arbitrary. But the effect of this arbitrariness can at least be partially checked in a later stage of the analysis (see Section 6).

## 4. A Numerical Example

We now illustrate our thesis with a numerical example. As mentioned before, all of our data (except one) are taken from the same source: The European Comparison Programme (ECP) of the United Nations

Table 1. *Countries and real per capita GDP*

| Country name | Country symbol | Real per capita GDP in percentage of Austria |
|---|---|---|
| Luxembourg | L | 123.1 |
| Norway | N | 123.1 |
| Sweden | S | 112.8 |
| Denmark | DK | 112.3 |
| FR of Germany | D | 111.7 |
| France | F | 104.9 |
| Finland | SF | 103.6 |
| Netherlands | NL | 103.3 |
| United Kingdom | UK | 100.1 |
| Austria | A | 100.0 |
| Italy | I | 99.4 |
| Belgium | B | 97.9 |
| Spain | E | 69.7 |
| Ireland | IRL | 61.9 |
| Greece | GR | 53.8 |
| Portugal | P | 51.1 |
| Hungary | H | 47.2 |
| Yugoslavia | YU | 44.2 |
| Poland | PL | 37.1 |
| Turkey | TR | 31.2 |

Economic Commission for Europe (UN 1988). This programme covered 20 European countries with 1985 as reference year. Table 1 lists the countries together with their real per capita GDP compared to Austria which was used as the base country.

According to the principle outlined above, the basic criterion in the selection of variables was to provide as rich a description of the economic structure as possible. Eleven variables were selected, 10 of which from the ECP. (The abbreviations used in the tables appear in parentheses.)

    1. Area of country in $km^2$ (AREA).

    2. Mid-year population (POP).

Both variables 1 and 2 refer to country size, but from different points of view. Experience of comparative economics has shown that country size affects economic

characteristics. Economists (for example, Kuznets 1965) treat population as one of the basic factors of economic development (even if one may argue that size variables are fixed magnitudes which cannot be manipulated by economic policy). One may further object to the use of two variables. This procedure is, however, justified by the fact that the effects of the two variables of size are different and the correlation between them is small (0.583).

3. Per capita GDP in real terms (GDP). All countries' GDP is expressed in a common currency in accordance with the comparison programme's guidelines. This calculation was made with the help of purchasing power parities derived from actual price ratios, without using the official exchange rates. GDP in real terms is the most important indicator of general economic development.

4. Exchange rate deviation (XRATE) is defined as the ratio of official exchange rate to actual purchasing power parity. This is the variable that indicates the relative over- or undervaluation of a national currency.

5. Share of government expenditure to GDP (GOVERN). This variable indicates the extent of the government's role in the economy. This variable was calculated in terms of national currencies rather than in real terms, because the purchasing power parities obtained for government expenditures are less reliable than those obtained for other sectors of the economy.

6. Share of gross fixed capital formation of GDP (GFC). This variable reflects the extent of investments in a given year. Unlike GOVERN, it was calculated in terms of the common currency. (ECP provides special parity for capital formation, different from the overall parity for GDP.)

7. Share of net purchases abroad to total consumption of the population (NETABR). This variable was included to reflect the special effects of foreign travel. Net purchases abroad show the extent to which the resident population (within the country or abroad) spends its income within its own territory or abroad (in a foreign country).

8. Per capita consumption of medical goods and services (MEDICAL). Even though, as often pointed out, such measures do not reflect the state of a country's health, it can be assumed that medical expenditures are closely correlated with the level of health.

9. Price level of machinery and equipment relative to the price level of total gross capital formation (MACHIN). It is generally accepted that this ratio indicates the overall technical level of a given economy. The lower the relative prices of machines and equipment, the higher the assumed technical level.

10. Per capita consumption of animal protein (ANIMALP). This term is generally used to reflect the level of nutrition. The following items of the ECP classification have been added up in real terms: meat, fish, milk, cheese, and eggs.

11. Share of agriculture in economic activity (AGRIC). This is the only variable in the data set which was taken from sources other than the ECP. It is, however, a necessary characteristic of a macroeconomic structure.

Of course this set of variables is somewhat arbitrary, as would be any other set. It can be argued, for example, that important aspects of economic structure like housing or environment are missing, or final use rather than production is reflected, etc. Obviously, any result and any statement holds only for this group of countries, and only in the space generated by these variables.

The set of the distance measures (as

*Table 2. Distance measures between countries*

| Country | L | N | S | DK | D | F | SF | NL | UK | A |
|---|---|---|---|---|---|---|---|---|---|---|
| L | 0.000 | 1.201 | 1.200 | 0.597 | 1.228 | 1.417 | 1.141 | 1.076 | 1.308 | 0.599 |
| N | | 0.000 | 0.704 | 0.883 | 0.954 | 1.272 | 0.531 | 1.019 | 1.438 | 1.201 |
| S | | | 0.000 | 0.964 | 0.896 | 1.185 | 0.958 | 0.972 | 1.228 | 1.242 |
| DK | | | | 0.000 | 0.954 | 1.273 | 0.912 | 0.602 | 1.003 | 0.744 |
| D | | | | | 0.000 | 0.864 | 1.077 | 0.932 | 0.803 | 1.212 |
| F | | | | | | 0.000 | 1.386 | 1.305 | 1.055 | 1.393 |
| SF | | | | | | | 0.000 | 1.192 | 1.498 | 0.955 |
| NL | | | | | | | | 0.000 | 0.799 | 1.122 |
| UK | | | | | | | | | 0.000 | 1.274 |
| A | | | | | | | | | | 0.000 |

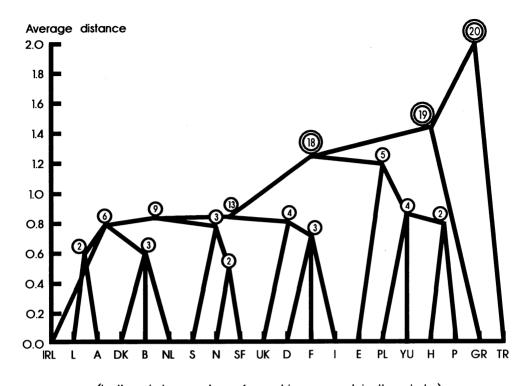| Country | I | B | E | IRL | GR | P | H | YU | PL | TR |
|---|---|---|---|---|---|---|---|---|---|---|
| L | 1.089 | 0.879 | 1.522 | 1.009 | 1.850 | 1.830 | 1.466 | 1.675 | 1.745 | 2.626 |
| N | 1.293 | 1.069 | 1.908 | 1.167 | 2.191 | 1.907 | 1.611 | 1.830 | 1.640 | 2.481 |
| S | 1.279 | 1.035 | 1.807 | 1.339 | 2.261 | 1.988 | 1.615 | 1.716 | 1.746 | 2.431 |
| DK | 1.002 | 0.595 | 1.477 | 0.710 | 1.670 | 1.533 | 1.271 | 1.397 | 1.476 | 2.392 |
| D | 0.760 | 0.991 | 1.575 | 1.252 | 2.066 | 1.829 | 1.527 | 1.595 | 1.501 | 2.277 |
| F | 0.725 | 1.093 | 1.328 | 1.498 | 1.897 | 1.993 | 1.812 | 1.724 | 1.793 | 2.370 |
| SF | 1.247 | 1.230 | 1.704 | 0.911 | 2.021 | 1.739 | 1.405 | 1.700 | 1.320 | 2.254 |
| NL | 1.235 | 0.710 | 1.694 | 0.993 | 1.675 | 1.589 | 1.316 | 1.249 | 1.533 | 2.380 |
| UK | 0.943 | 1.028 | 1.313 | 1.246 | 1.628 | 1.632 | 1.440 | 1.178 | 1.489 | 2.156 |
| A | 1.111 | 1.053 | 1.343 | 0.795 | 1.596 | 1.671 | 1.293 | 1.490 | 1.445 | 2.401 |
| I | 0.000 | 0.960 | 1.019 | 1.136 | 1.762 | 1.628 | 1.441 | 1.517 | 1.447 | 2.135 |
| B | | 0.000 | 1.458 | 0.980 | 1.699 | 1.531 | 1.370 | 1.488 | 1.708 | 2.504 |
| E | | | 0.000 | 1.280 | 1.310 | 1.338 | 1.396 | 1.316 | 1.424 | 1.806 |
| IRL | | | | 0.000 | 1.362 | 1.119 | 0.886 | 1.131 | 0.983 | 2.019 |
| GR | | | | | 0.000 | 1.374 | 1.522 | 1.145 | 1.585 | 2.128 |
| P | | | | | | 0.000 | 0.811 | 1.092 | 1.143 | 1.588 |
| H | | | | | | | 0.000 | 0.832 | 0.816 | 1.525 |
| YU | | | | | | | | 0.000 | 0.959 | 1.432 |
| PL | | | | | | | | | 0.000 | 1.213 |
| TR | | | | | | | | | | 0.000 |

defined in Section 2) between the 20 countries is shown in Table 2.

These figures reflect the similarity or dissimilarity between the economic structures of any two countries. The smaller the distance, the more similar the economic structure. For example, the economic structures of Norway and Finland are very similar (0.531), and that of Luxembourg and Turkey very different (2.626). We also see that the structure of France differs more from that of Spain than from that of Italy (1.328 and 0.725 respectively), etc.

## 5.  The Use of Distance Measures

### 5.1.  The clustering procedure

Using structural distance measures entails clustering; a stepwise arrangement of countries into clusters provides insight into the behaviour of these countries. This kind of investigation may reveal, for example, whether characteristics such as geographical neighbourhood, common social structure, membership in an economic integration, etc., affect structural similarities. As already

(In the circles: number of countries merged in the cluster)

*Fig. 1.   Clustering procedure*

mentioned, in this study the centroid method is used to define clusters. According to this method the two countries with the smallest distance (Norway and Finland), constitute the first cluster. In all further steps this pair of countries is treated as a single unit with the averages of the standard values of the two countries. The result of the complete clustering procedure is shown in Figure 1.

The circles in this dendogram indicate the number of countries that were merged to form the cluster. The positions of the clusters reflect the distances between clusters (or countries). The distance is measured on the vertical axis. To extract economic information from this picture, the dendogram can be read from both the bottom and the top. Countries with highly similar structures (i.e., with the smallest distances), like Norway and Finland or Austria and Luxem-

bourg, are merged at the beginning of the process and appear close to the horizontal axis. Reading the diagram from the top, it can be seen that the structure of Turkey differs drastically from the other 19 countries. And within the 19, there is another country, Greece, which is substantially distant from the remaining 18. Then these 18 countries can be subdivided into two groups, one embracing 13 Western European countries, the other Spain, Portugal and the Eastern European countries. From the top down to this level of clustering the overwhelming role of per capita GDP is obvious. Turkey, with the lowest level forms a one-country cluster, considerably far from the other 19. Then the group of 13 includes countries with higher per capita GDP (except Ireland) than the remaining 7 (i.e., the group of 5 plus Greece and Turkey).

Below this level of clustering, however, other factors than GDP determine which countries merge into given clusters. For example, the cluster of 13 is subdivided into two subclusters, one with 9 and the other with 4 countries. The effect of *size* (especially in terms of population) is obvious, as all the countries in the cluster on the right are large.

## 5.2. International correspondence analysis

A by-product of distance measures and cluster analysis is the analysis of correspondence between various magnitudes of the selected variables (in the sense of UNRISD (1970)). The value of two or more variables within the same cluster can be interpreted as magnitudes which correspond to each other on an international scale. The greater the clusters are in terms of the number of countries included, the more reliable the correspondence measures.

In order to illustrate this, the following three clusters are selected from our example: (i) nine small-size Western European countries, (ii) four large-size Western European countries (UK, FRG, France, Italy), (iii) five countries of Eastern and South-Western Europe (omitted from this part of the analysis are Greece and Turkey, each of which forms a one-country cluster). Table 3 presents the average standard value of each variable for these three country-groups.

Some conclusions can be drawn from

*Table 3.   Average standard values of the variables*
(in bracket the standard deviations)

| Variables | Cluster (i) | Cluster (ii) | Cluster (iii) |
|---|---|---|---|
| Number of countries | 9 | 4 | 5 |
| AREA | − 0.451 | 0.448 | 0.034 |
|  | (0.785) | (0.616) | (0.763) |
| POP | − 0.793 | 1.607 | 0.017 |
|  | (0.183) | (0.102) | (0.585) |
| GDP | 0.656 | 0.649 | − 1.146 |
|  | (0.574) | (0.163) | (0.362) |
| XRATE | − 0.698 | − 0.496 | 1.108 |
|  | (0.230) | (0.171) | (0.589) |
| GOVERN | − 0.3515 | − 0.021 | 0.086 |
|  | (0.554) | (0.689) | (0.662) |
| GFC | 0.582 | − 0.251 | − 0.535 |
|  | (0.958) | (0.623) | (0.917) |
| NETABR | 0.246 | 0.121 | − 0.340 |
|  | (1.073) | (0.502) | (1.023) |
| MEDICAL | 0.532 | 0.675 | − 0.918 |
|  | (0.735) | (0.615) | (0.358) |
| MACHIN | − 0.555 | − 0.607 | 0.925 |
|  | (0.583) | (0.239) | (0.923) |
| ANIMALP | 0.155 | 0.824 | − 0.656 |
|  | (0.548) | (1.024) | (0.875) |
| AGRIC | − 0.573 | − 0.661 | 0.734 |
|  | (0.269) | (0.216) | (0.462) |

Table 3 even at first glance. As regards size, the third group lies between the "small size" and the "large size" cluster. The behaviour of the two variables indicating country size is similar, but differences in population are more significant than differences in area. The average per capita GDP in the two Western clusters are very close to each other and substantially higher than in the third group.

However, caution should be exercised in analysing these figures. The individual clusters are, in general, sufficiently homogeneous, but this homogeneity does not hold automatically for all individual variables. For this reason the standard deviation of the standard values in the cluster is shown in parentheses in Table 3. In all cases where the deviation is large, we refrain from making analytical statements. As the standard deviation of the standard values equals 1, variables with cluster deviation close to 1 (and especially more than 1) cannot be analysed. This applies in the case of Gross Fixed Capital Formation, Net Purchase Abroad and Animal Protein. Comparing the third group with the two others, it can be stated that the official exchange rate substantially underestimates the purchasing power of the currencies of low income countries. Low income is coupled with a high share of government expenditures, a low level of medical consumption, a high relative price of machinery and equipment (i.e., relatively low technical level) and a high share of agriculture. On the other hand, small size (second group) seems to be coupled with more favourable exchange rates, lower share of government expenditures, etc.

*Table 4.   Correlation coefficients of the variables*

|  |  | AREA | POP | !<br>GDP | !<br>XRATE | GOVERN | GFC | NETABR |
|---|---|---|---|---|---|---|---|---|
| AREA |  | 1.00 | 0.58 | − 0.20 | 0.25 | − 0.19 | − 0.04 | 0.09 |
| POP |  |  | 1.00 | − 0.01 | 0.16 | 0.01 | − 0.03 | − 0.02 |
| GDP | ! |  |  | 1.00 | − 0.87 | − 0.32 | 0.38 | 0.19 |
| XRATE | ! |  |  |  | 1.00 | 0.12 | − 0.46 | − 0.13 |
| GOVERN |  |  |  |  |  | 1.00 | − 0.49 | − 0.36 |
| GFC |  |  |  |  |  |  | 1.00 | 0.21 |
| NETABR |  |  |  |  |  |  |  | 1.00 |

|  |  | !<br>MEDICAL | !<br>MACHIN | ANIMALP | !<br>AGRIC | First factor |
|---|---|---|---|---|---|---|
| AREA |  | − 0.06 | 0.19 | − 0.13 | 0.44 | 0.29 |
| POP |  | − 0.04 | − 0.00 | 0.09 | 0.15 | 0.17 |
| GDP | ! | 0.84 | − 0.78 | 0.58 | − 0.87 | − 0.96 |
| XRATE | ! | − 0.75 | 0.68 | − 0.59 | 0.86 | 0.92 |
| GOVERN |  | − 0.30 | 0.21 | 0.05 | 0.20 | 0.32 |
| GFC |  | 0.30 | − 0.27 | 0.03 | − 0.22 | − 0.45 |
| NETABR |  | 0.39 | − 0.22 | − 0.32 | − 0.09 | − 0.23 |
| MEDICAL | ! | 1.00 | − 0.77 | 0.46 | − 0.79 | − 0.89 |
| MACHIN | ! |  | 1.00 | − 0.26 | 0.69 | 0.81 |
| ANIMALP |  |  |  | 1.00 | − 0.63 | − 0.58 |
| AGRIC | ! |  |  |  | 1.00 | 0.93 |

! Indicates the core variables.

### 6. Checking the Selection of the Variables

As previously emphasized, all the analytical statements made in the preceding sections depend largely on the original set of variables. Other sets might result in more or less different conclusions. In Section 3, when discussing the general principle of the selection, two problems were identified. The first problem was the assumed effect of the variables *not* selected. There is a large number of possible variables and their effects cannot be measured unless the whole procedure is repeated several times using different sets. Such a series of repetitions was not conducted in the present study. Second, the composition of the selected set may lead to a bias. This is the question of weighting the variables. If two or more variables are influenced by the same components of economic structure, the selected set overemphasizes these components.

To check this effect, consider Table 4, the correlation matrix of the variables.

*Table 5. Modified (unweighted) distances*

| Country | L | N | S | DK | D | F | SF | NL | UK | A |
|---|---|---|---|---|---|---|---|---|---|---|
| L | 0.000 | 1.366 | 1.661 | 0.931 | 1.681 | 1.615 | 1.383 | 1.269 | 1.573 | 1.373 |
| N | | 0.000 | 1.346 | 1.072 | 1.459 | 1.219 | 0.605 | 0.996 | 1.504 | 1.580 |
| S | | | 0.000 | 1.365 | 1.521 | 1.518 | 1.350 | 1.281 | 1.588 | 1.430 |
| DK | | | | 0.000 | 1.337 | 1.586 | 1.147 | 0.845 | 1.306 | 1.442 |
| D | | | | | 0.000 | 1.488 | 1.430 | 1.257 | 1.382 | 1.246 |
| F | | | | | | 0.000 | 1.324 | 1.353 | 1.545 | 1.692 |
| SF | | | | | | | 0.000 | 1.183 | 1.482 | 1.294 |
| NL | | | | | | | | 0.000 | 0.928 | 1.496 |
| UK | | | | | | | | | 0.000 | 1.656 |
| A | | | | | | | | | | 0.000 |

| Country | I | B | E | IRL | GR | P | H | YU | PL | TR |
|---|---|---|---|---|---|---|---|---|---|---|
| L | 1.172 | 1.462 | 1.520 | 1.549 | 1.686 | 1.742 | 1.412 | 1.369 | 1.610 | 1.625 |
| N | 1.353 | 1.400 | 1.617 | 1.414 | 1.505 | 1.408 | 1.499 | 1.451 | 1.388 | 1.494 |
| S | 1.668 | 1.248 | 1.445 | 1.445 | 1.844 | 1.715 | 1.666 | 1.240 | 1.616 | 1.576 |
| DK | 1.045 | 1.098 | 1.416 | 1.069 | 1.423 | 1.376 | 1.640 | 1.167 | 1.307 | 1.407 |
| D | 1.119 | 1.315 | 1.760 | 1.573 | 1.685 | 1.524 | 1.678 | 1.508 | 1.323 | 1.567 |
| F | 1.117 | 1.330 | 1.431 | 1.598 | 1.487 | 1.785 | 1.470 | 1.386 | 1.500 | 1.703 |
| SF | 1.345 | 1.536 | 1.418 | 1.169 | 1.495 | 1.451 | 1.457 | 1.459 | 1.073 | 1.528 |
| NL | 1.307 | 1.159 | 1.589 | 1.265 | 1.327 | 1.362 | 1.336 | 0.887 | 1.274 | 1.563 |
| UK | 1.442 | 1.723 | 1.375 | 1.586 | 1.707 | 1.592 | 1.634 | 1.306 | 1.431 | 1.777 |
| A | 1.533 | 1.557 | 1.520 | 1.548 | 1.464 | 1.681 | 1.514 | 1.477 | 1.463 | 1.759 |
| I | 0.000 | 1.117 | 1.318 | 1.254 | 1.609 | 1.530 | 1.515 | 1.382 | 1.202 | 1.437 |
| B | | 0.000 | 1.565 | 1.170 | 1.522 | 1.352 | 1.477 | 1.212 | 1.524 | 1.690 |
| E | | | 0.000 | 1.378 | 1.637 | 1.468 | 1.696 | 1.499 | 1.543 | 1.629 |
| IRL | | | | 0.000 | 1.561 | 1.539 | 1.580 | 1.220 | 0.849 | 1.758 |
| GR | | | | | 0.000 | 1.618 | 1.815 | 1.284 | 1.587 | 1.711 |
| P | | | | | | 0.000 | 1.436 | 1.640 | 1.618 | 1.628 |
| H | | | | | | | 0.000 | 1.247 | 1.362 | 1.729 |
| YU | | | | | | | | 0.000 | 1.135 | 1.287 |
| PL | | | | | | | | | 0.000 | 1.447 |
| TR | | | | | | | | | | 0.000 |

There exists a core of five variables (indicated by "!" in the table) which are closely correlated (positively or negatively) with each other, and especially with the per capita GDP. These are: the exchange rate deviation, the medical consumption, the relative prices of machinery and equipment, and the share of agriculture. Is there any way to check the effect of this phenomenon? The most trivial way, i.e., the removal of GDP alone, does not change the picture. It would be more promising to exclude all the correlating variables (except one), but this solution would reduce the number of variables and impoverish the analysis. In order to keep the original data set, a slightly more sophisticated procedure is necessary.

It is often convenient to do cluster analysis not only on the original data, but also on their first *few principal components*. The last column of the correlation matrix (Table 4) shows the patterns of the first principal component (first factor). As can be seen, the correlations are very high with the GDP and GDP-related variables (mentioned above). And the combined weight (not explicit in the table) of this first factor is also important: 0.44 of the variance of the total.

If cluster analysis is done on the first few principal components and these components are in standardized form, then the weights of the various principal components are removed. The removal of weights means that the group of GDP-related variables, or any other overemphasized component is "deprived" of its overwhelming effect. The new set of distance measures can be labeled as "unweighted," in contrast to the original ("weighted") distances. The unweighted distance measures appear in Table 5.

When the values in Table 2 and Table 5 are compared, the results obtained by the original (weighted) measures appear to be justified in many respects. For example, the structures of Norway and Finland are the

most similar, in both cases. Both weighted and unweighted distances show the structures of Turkey and Greece to differ the most from the structures of the other countries (in terms of the average of the distance measures). In other cases, structural similarities disappear when the weights are removed, for example, France and the Federal Republic of Germany. (This means that the low value of their weighted distance is mainly due to the closeness of their per capita GDP.)

On the other hand, the general picture obtained by the unweighted figures are less characteristic than that of the weighted distances, because the differences in the various distance measures are substantially smaller. Whereas the average distance of the two sets are close to each other, the variance of the unweighted measures is only about half that of the weighted measures:

Table 6. *Comparison of weighted and unweighted measures*

|  | Distances | |
|---|---|---|
|  | Weighted | Unweighted |
| Average | 1.316 | 1.365 |
| Standard deviation | 0.518 | 0.371 |
| Smallest distance | 0.531 (N-SF) | 0.605 (N-SF) |
| Largest distance | 2.626 (L-TR) | 1.777 (UK-TR) |

The conclusion to be drawn from this is that unweighted measures might be useful for checking the reliability of the clustering procedure and for providing a supplementary picture, but they do not constitute an analytical tool, per se.

## 7.  Some Concluding Remarks

The use of cluster analysis in international comparisons provides some insight into structural differences and similarities of the

various national economies. Economic relationships and similarities missed by other types of international comparisons can be uncovered by cluster analysis. The example presented in this article has illustrated the general procedure and some possible conclusions, but is in no way exhaustive. Further research – with different (more sophisticated) techniques in defining distances and calculating clusters, with different sets of countries and variables – might improve our understanding of the methods and the economic interpretation of the results obtained.

## 8. References

Dubes, R. and Jain, A.K. (1979). Validity Studies in Clustering Methodologies. Pattern Recognition, 11, 235–254.

Gordon, A. (1981). Classification: Methods for the Exploratory Analysis of Multivariate Data. London: Chapman and Hall.

Hartigan, J.A. (1975). Clustering Algorithms. New York: John Wiley and Sons.

Kuznets, S. (1965). Toward a Theory of Economic Growth. Economic Growth and Structure. Selective Essays. New York: Yale University Press, 1–81.

United Nations (1988). International Comparison of Gross Domestic Product in Europe 1985. Report on the European Comparison Programme. Statistical Standards and Studies – No. 41, United Nations, New York, 1988.

United Nations Research Institute for Social Development (UNRISD) (1970). Content and Measurement of Socio-Economic Development: An Empirical Enquiry. Geneva: United Nations.