

Czesław Domański  
Robert Kubacki

## APPLICATION OF THE DIFFERENTIAL EVOLUTION ALGORITHM TO GROUP A BANK'S INDIVIDUAL CLIENTS

### Abstract

*Objective:* The aim of the article is to present the results of grouping individual clients of a bank with the differential evolution algorithm.

*Research Design & Methods:* The research offers conclusions based on analysis of the bank's customer base and deductive and inductive reasoning.

*Findings:* The results of the authors' research show that the differential evolution algorithm correctly groups bank customers and can be used for this purpose.

*Implications/Recommendations:* The differential evolution algorithm is an alternative to the commonly used k-means algorithm. The algorithm generates several competing solutions in one iteration. It enables independence from starting vectors and greater effectiveness in searching for an optimal solution. The differential evolution algorithm was itself enriched with a variable that allows the optimal number of clusters to be selected. Each iteration contained proposed solutions (chromosomes) that were evaluated by the target function built on the CS measure proposed by Chou.

*Contribution:* The article presents the application of the differential evolution algorithm to group a bank's clients.

Czesław Domański, University of Lodz, Faculty of Economics and Sociology, Department of Statistical Methods, POW 3/5, 90-255 Łódź, Poland, e-mail: [czesdoman@uni.lodz.pl](mailto:czesdoman@uni.lodz.pl), ORCID: <https://orcid.org/0000-0001-6144-6231>.

Robert Kubacki, University of Lodz, Faculty of Economics and Sociology, Department of Statistical Methods, POW 3/5, 90-255 Łódź, Poland, e-mail: [robertkubacki@o2.pl](mailto:robertkubacki@o2.pl), ORCID: <https://orcid.org/0000-0003-0591-9529>.

This is an open access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 License (CC BY-NC-ND 4.0); <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

**Keywords:** CS measure, differential evolution algorithm, clustering methods, banking.  
**JEL Classification:** C38, M31, G21.

## 1. Introduction

In today's world more and more companies are having problems with the effective management of available data. The gap between the amount of data that is generated and stored and the degree of companies' understanding is constantly widening. According to a survey conducted by IBM among representatives of the largest banks, over 40% of them have problems with an excess of information and the lack of appropriate tools for analysing it (Giridhar *et al.* 2011).

Grouping methods are effective in describing populations. Many authors have studied these methods (Everitt *et al.* 2011, Feoktistov & Janaqi 2004, Gan, Ma & Wu 2007).

Classical grouping algorithms have two major disadvantages:

1. They easily fall into local optima in multidimensional spaces that have multimodal objective functions,
2. The efficiency of searching for a solution depends very much on the start vectors.

Grouping methods are known as methods without a supervisor, while most traditional algorithms require a priori knowledge of the number of clusters, which means that this is not a method without the interference of an outsider. On the other hand, in many practical applications it is impossible to provide even an approximate number of groups for an unknown data set.

The limitations of classical grouping methods, including the k-means algorithm, led the researchers to search for new, more effective grouping methods. One of the development trajectories for grouping algorithms has been to treat them as an optimization problem. Over time, the paradigm of evolutionary computation – the relationship between optimization and biological evolution – has evolved. Evolutionary calculations use the power of natural selection and allow the computing power of computers to be used for automatic optimization (Das, Abraham & Konar 2009).

The first part of the article presents selected issues related to the differential evolution algorithm. Next, the assumptions of the study and the steps taken to prepare the data are described. Selected results for the study are then outlined. In the final part of the article, the authors indicate possible directions for further research in this area.

## 2. Differential Evolution Algorithm – Selected Issues

The differential evolution algorithm is part of heuristic methods because the goal of optimization is not to find the exact equation describing the studied phenomenon, but to search the available space for solutions. These solutions are constructed using random elements. What is more, in one iteration of the algorithm several competing solutions are created. Subsequent solutions are created using similarities to the evolutionary mechanisms occurring in nature. These are the ones that, according to the defined objective function, are the best. The characteristic feature of the differential evolution algorithm is that solutions are created on the basis of real variable vectors, not vectors coded to zero-one sequences (Das, Mullick & Suganthan 2016).

Since 1995 the differential evolution algorithm (Storn 1995, Storn & Price 1997) has drawn the attention of optimization practitioners due to the degree of resistance, the speed of convergence and the accuracy of solutions for real optimization problems. The differential evolution algorithm has defeated many algorithms, such as genetic algorithms, evolutionary strategies and memetic algorithms (Das, Mullick & Suganthan 2016).

Suppose we have a set of objects  $Np$  vectors, each of which has  $D$  dimensions. In addition, we mark  $P_X$  as the current population of solutions to the optimization problem, which was created as an initial solution or at any subsequent stage of the algorithm's operation (Das, Abraham & Konar 2009).

$$P_{X,g} = (X_{i,g}), \quad i = 0, 1, \dots, Np - 1, \quad g = 0, 1, \dots, g_{\max}, \quad (1)$$

$$X_{i,g} = (x_{j,i,g}), \quad j = 0, 1, \dots, D - 1. \quad (2)$$

Index  $g = 0, 1, \dots, g_{\max}$  denotes the generation to which the vector belongs. Each vector is assigned to the corresponding population index  $i = 0, 1, \dots, Np - 1$ . The dimensions of the vector are marked by  $j = 0, 1, \dots, D - 1$ .

The differential evolution algorithm generates mutant vectors in the next step, which will be marked as follows (Das, Abraham & Konar 2009):

$$P_{V,g} = (V_{i,g}), \quad i = 0, 1, \dots, Np - 1, \quad g = 0, 1, \dots, g_{\max}, \quad (3)$$

$$V_{i,g} = (v_{j,i,g}), \quad j = 0, 1, \dots, D - 1. \quad (4)$$

However, the vectors after crossover will be marked as follows:

$$P_{U,g} = (U_{i,g}), \quad i = 0, 1, \dots, Np - 1, \quad g = 0, 1, \dots, g_{\max}, \quad (5)$$

$$U_{i,g} = (u_{j,i,g}), \quad j=0, 1, \dots, D-1. \quad (6)$$

The first stage, i.e. setting the initial vectors, consists in generating starting vectors. The initial parameters (for  $g = 0$ ) are set within limits that correspond to a range that is acceptable for the intended solution. Therefore, if  $j$ -th the search task parameter has ranges marked as  $x_{\min,j}$  and  $x_{\max,j}$  and  $rand_{i,j}(0, 1)$  means  $j$ -th realizations of a uniform distribution from the range from 0 to 1 for  $i$ -th vector, then  $j$ -th component  $i$ -th population element can be determined as (Das, Abraham & Konar 2009):

$$x_{i,j}(0) = x_{\min,j} + rand_{i,j}(0,1) \cdot (x_{\max,j} - x_{\min,j}). \quad (7)$$

The differential evolution algorithm searches for the global optimum in  $D$ -dimensional continuous hyperspace. It starts with a randomly selected population  $Np$   $D$ -dimensional values of parameter vectors. Each vector, also known as a genome/chromosome, is a proposed solution in a multidimensional optimization issue. The next generations of solutions in the differential evolution are marked as  $g = 0, 1, 2, \dots, g, g + 1$ .

The vector parameters may change with the appearance of new generations, therefore the notation for which it will be accepted, for which  $i$ -th population vector for the current generation over time ( $g = g$ ) as (Das, Abraham & Konar 2009):

$$\vec{X}_i(g) = [x_{i,1}(g), x_{i,2}(g), \dots, x_{i,D}(g)]^T, \quad (8)$$

where  $i = 1, 2, \dots, Np$ .

Mutation means a sudden change in the characteristics of the chromosome gene. In the context of evolutionary computation, a mutation means a change or disorder of a random component. Most evolutionary algorithms simulate the effect of mutations through the additivity of the component generated with a given probability distribution. In the differential evolution algorithm, a uniform distribution of the vector of the form differences was used (Das, Abraham & Konar 2009):

$$\Delta \vec{X}_{r2,r3} = (\vec{X}_{r2} - \vec{X}_{r3}). \quad (9)$$

In the differential evolution algorithm, the mutation creates a successor vector  $\vec{V}_i(g)$  for changing the population element  $\vec{X}_i(g)$  in every generation or iteration of the algorithm.

To create a vector  $\vec{V}_i(t)$  for each  $i$ -th element of the current population, the other three disjoint vectors  $\vec{X}_{r_1}(g)$ ,  $\vec{X}_{r_2}(g)$ ,  $\vec{X}_{r_3}(g)$  are randomly

selected from the current population. Indexes  $r_1^i, r_2^i, r_3^i$  are mutually exclusive integers selected from a range  $[1, Np]$ , which are also different from the index and the base vector. Indexes are generated randomly for each mutated vector. Then, the difference of any two of the three vectors is scaled by the number  $F$  and added to the third vector. In this way, we get a vector  $\vec{V}_i(g)$  expressed as (Das, Abraham & Konar 2009):

$$\vec{V}_i(g) = \vec{X}_{r_1^i}(g) + F \cdot (\vec{X}_{r_2^i}(g) - \vec{X}_{r_3^i}(g)). \quad (10)$$

The mutation scheme shows different ways of differentiating the proposed solutions.

The crossover operation is used to increase the diversity of the population of solutions. Crossing takes place after generating a donor vector through a mutation. The algorithms of the differential evolution family use two intersection schemes – exponential and binomial (zero-one). The donor vector lists the components with the target vector  $\vec{X}_i(g)$  to create a trial vector:

$$\vec{U}_i(g) = [u_{i,1}(g), u_{i,2}(g), \dots, u_{i,D}(g)]^T. \quad (11)$$

In exponential crossover, we first select a random integer  $n$  from range  $[0, D - 1]$ . The drawn number is the starting point for the target vector from which the components are crossed with the donor vector. An integer  $L$  is also selected from range  $[1, D]$ .  $L$  indicates the number of components in which the donor vector is involved. After selection  $n$  and  $L$ , the trial vector takes the form (Das, Abraham & Konar 2009):

$$u_{i,j}(g) = \begin{cases} v_{i,j}(g) & \text{for } j = \langle n \rangle_D, \langle n+1 \rangle_D, \dots, \langle n+L-1 \rangle_D \\ x_{i,j}(g) & \text{for other } j \in [0, D-1] \end{cases} \quad (12)$$

where the intervals denote the module modulo function  $D$ . Integer  $L$  is drawn from the sequence  $[1, 2, \dots, D]$  according to the following pseudocode:

```

L = 0;
Do
{
L = L + 1;
} while (rand(0, 1) < CR) AND (L < D));

```

As a result, the probability  $(L \geq \nu) = (CR)^{\nu-1}$  for any  $\nu > 0$ . The crossover rate ( $CR$ ) is a parameter the same as  $F$ . For each donor vector, a new set  $n$  and  $L$  must be drawn as described above.

On the other hand, binomial crossover is carried out for each  $D$  variable each time, when the number selected from 0 to 1 is less than or equal to the value  $CR$ . In this case, the number of parameters inherited from the donor has a very similar distribution to the binomial one. This scheme can be represented in the following way (Das, Abraham & Konar 2009):

$$u_{i,j,g} = \begin{cases} v_{i,j,g}, & \text{if } (rand_{i,j}(0,1) \leq CR \text{ or } j = j_{rand}) \\ x_{i,j,g} & \text{otherwise} \end{cases} \quad (13)$$

where  $rand_{i,j}(0,1) \in [0,1]$  is a randomly drawn number that is generated for every  $j$ -th of the  $i$ -th parameter of the vector.  $j_{rand} \in [1, 2, \dots, D]$  is a randomly selected index that ensures that  $\vec{U}_{i,g}$  contains at least one component from the vector  $\vec{V}_{i,g}$ .

This is determined once for each vector in a given generation.  $CR$  is an estimate of true probability  $p_{CR}$  the event that the component of the sample vector will be inherited from the parent. It may also happen that in the two-dimensional search space, three possible test vectors can be the result of one-dimensional mating of the mutant/donor vector  $\vec{V}_i(g)$  with the target vector  $\vec{X}_i(g)$ . Trial vectors:

- a)  $\vec{U}_i(g) = \vec{V}_i(g)$  both components  $\vec{U}_i(g)$  inherited from the vector  $\vec{V}_i(g)$ ,
- b)  $\vec{U}'_i(g) = \vec{V}_i(g)$  one component ( $j = 1$ ) comes from vector  $\vec{V}_i(g)$ , the second ( $j = 2$ ) from vector  $X_i(t)$ ,
- c)  $\vec{U}''_i(g) = \vec{V}_i(g)$  one component ( $j = 1$ ) comes from vector  $X_i(g)$ , the second ( $j = 2$ ) from vector  $\vec{V}_i(g)$ .

The last stage of the differential evolution algorithm is selection, i.e. the choice between the vector  $\vec{X}_i(g)$  and a newly designated test vector  $\vec{U}_i(g)$ . The decision which of the two vectors will survive in the next generation  $g + 1$  depends on the value of the matching function. If the values of the matching function for the sample vector are better than the value of the target vector, the existing vector is replaced with the new vector (Das, Abraham & Konar 2009).

$$\vec{X}_i(g+1) = \begin{cases} \vec{U}_i(g) & \text{for } f(\vec{U}_i(g)) \leq f(\vec{X}_i(g)) \\ \vec{X}_i(g) & \text{for } f(\vec{U}_i(g)) > f(\vec{X}_i(g)) \end{cases} \quad (14)$$

where  $f(\vec{X})$  is a minimized function. The selection process consists in selecting one of two variants. The adjustment of population members improves in subsequent generations or remains unchanged, but never deteriorates.

The *CS* (Candidate Solution) measure proposed by Chou (Chou, Su & Lai 2004) is an objective function in this study. Group centroids are determined as the average vectors belonging to a given cluster

$$\bar{m}_i = \frac{1}{N_i} \sum_{\bar{Z}_j \in C_i} \bar{Z}_j. \tag{15}$$

The distance between two points  $\bar{Z}_p$  and  $\bar{Z}_y$  is marked as  $d(\bar{Z}_p, \bar{Z}_y)$ . Then the *CS* measure can be defined as:

$$\begin{aligned} CS(k) &= \frac{\frac{1}{k} \sum_{i=1}^k \left[ \frac{1}{|C_i|} \sum_{\bar{Z}_y \in C_i} \max\{d(\bar{Z}_p, \bar{Z}_y)\} \right]}{\frac{1}{k} \sum_{i=1}^k \left[ \min_{j \in k, j \neq i} d(\bar{m}_i, \bar{m}_j) \right]} = \\ &= \frac{\sum_{i=1}^k \left[ \frac{1}{|C_i|} \sum_{\bar{Z}_y \in C_i} \max\{d(\bar{Z}_p, \bar{Z}_y)\} \right]}{\sum_{i=1}^k \left[ \min_{j \in k, j \neq i} d(\bar{m}_i, \bar{m}_j) \right]}. \end{aligned} \tag{16}$$

The measure is a function of the ratio of the amount of intra-group dispersion and the separation between groups. The *CS* measure is more effective at clusters with different density and/or different sizes than other measures.

### 3. Design of the Study

A database of a commercial bank’s clients from Europe was used for the study. It was limited to that part of the population for which the actions taken will translate in the maximum way into business benefits. In particular, clients meet the following criteria: individual clients with active products as on 1 March 2017, aged between 18 and 75 years, who are not bank employees, with positive marketing consent, and without delays in the repayment of loan products<sup>1</sup>.

As for the variables used for the study, the choice was not accidental. The variables selected for the study can be evaluated for each customer regardless of whether they have deposit, credit or investment products. By pre-processing data it was possible to eliminate outliers from the studied population. Due to the strong right-side skewness of the variables, a transformation was made by adding the constant 0.001, and then their

<sup>1</sup> The authors are not permitted to disclose the exact name of the bank which supported the data for this study.

logarithmisation. The resulting distributions of variables are thus more symmetrical.

The final set of variables used in the study is presented below:

- ZM1 (DEPOZYTY) – total funds on accounts and deposits in thousand PLN,
- ZM2 (INWESTYCJE) – total funds in investment products in thousand PLN,
- ZM3 (LUDNOSC) – number of inhabitants, based on the city from the correspondence address and data published by Statistics Poland,
- ZM4 (KREDYTY) – amount of bank loans taken in thousand PLN,
- ZM5 (SALDO\_BIK) – balance for repayment on credit products outside the bank, based on inquiries from BIK in thousand PLN,
- ZM6 (AVG\_TRN\_INCOMING\_ALL\_3M) – average monthly income on customer accounts in the last 3 months in thousand PLN,
- ZM7 (AVG\_TRN\_INCOMING\_CLEAN\_3M) – cleaned average monthly income on customer accounts in the last 3 months in thousand PLN,
- ZM8 (AVG\_TRN\_OUTGOING\_ALL\_3M) – average monthly outflows from customer accounts in the last 3 months in thousand PLN,
- ZM9 (AVG\_TRN\_OUTGOING\_CLEAN\_3M) – cleaned monthly average outflows from customer accounts in the last 3 months in thousand PLN,
- ZM10 (AVG\_TRN\_OUT\_DEBIT\_3M) – average monthly transaction amount on the debit card from the last 3 months in thousand PLN,
- ZM11 (AVG\_TRN\_OUT\_CREDIT\_3M) – average monthly amount of credit card transactions from the last 3 months in thousand PLN,
- ZM12 (WIEK\_LATA) – customer age in years,
- ZM13 (STAZ\_LATA) – customer experience in years.

Table 1 shows the constants used in the algorithm.

Table 1. Constants Used in the Study

Constant	Value	Description of Constant
<i>LZ</i>	13	number of variables describing the client
<i>LC</i>	13	number of chromosomes
<i>LK</i>	15	maximum number of clusters
<i>SA</i>	0.2	constant activation of the vector
<i>F</i>	0.7	mutation operator
<i>Iterations</i>	15	number of iterations
<i>CR</i>	1	crossover rate

Source: authors' own calculations.



For the purpose of optimizing the number of centroids, a dimensional matrix  $MR_{c,k,z}$  is created, where  $c$  means the number of chromosomes,  $k$  means the number of clusters, and  $z$  means the number of variables. The number of variables is increased by 1. An additional variable is used to store information on whether the cluster is active or inactive in the given iteration (Das, Abraham & Konar 2008). The values for individual matrix elements are generated according to formula (7). An additional variable indicating focus activation is determined on the basis of the following rule: if the randomly generated number from the range 0 to 1 is smaller than the activation constant (SA), then the variable takes the value 0, otherwise it takes the value 1.

#### 4. Results of Empirical Analyses

The smallest value of the CS function in the fifteenth iteration was obtained for chromosome number 3. This solution was chosen as the optimal solution.

Table 2 presents the characteristics of chromosome 3, which divided the surveyed population of the bank's clients into 9 groups (the maximum number of groups into which the population could be divided was 15).

Table 2. Numbers and Share of Groups for Chromosome 3

Group	Number of Clients	% of Total
8	92,109	45.71
4	44,545	22.11
6	29,047	14.41
3	20,003	9.93
5	5,476	2.72
14	3,582	1.78
1	2,839	1.41
15	2,075	1.03
12	1,832	0.91
Sum	201,508	100

Source: authors' own calculations.

The results of the grouping in Table 2 indicate that the distinguished groups are characterized by nonequal distribution of the number of clients in groups. Group 8 is more selective and includes 45.71% of clients, group 4

Table 3. Average Values of Variables ZM1–ZM13 for Clusters Obtained by the Differential Evolution Algorithm

Cluster	ZM1	ZM2	ZM3	ZM4	ZM5	ZM6	ZM7	ZM8	ZM9	ZM10	ZM11	ZM12	ZM13
8	8	0	452	137	88	5	4	5	3	0	0	43	5
4	2	5	273	5	1	5	4	5	4	0	0	43	6
6	20	10	500	182	146	22	17	22	17	1	0	42	6
3	28	4	453	7	45	1	1	1	0	0	0	47	6
5	60	61	627	325	0	26	20	27	19	0	1	41	6
14	113	97	758	414	144	113	84	106	73	2	1	42	6
1	20	67	473	223	114	4	3	4	2	0	0	43	6
15	24	48	321	264	10	10	8	7	5	0	0	41	5
12	0	2	131	7	117	15	11	18	15	0	0	41	3

Source: authors' own calculations.

contains 22.11% of clients, and group 6, the third group – 14.41% of clients. The three mentioned groups account for over 80% of the surveyed population.

More detailed characteristics of the distinguished groups of clients are presented in Table 3, which contains the average values of features in individual groups. The data presented in Table 3 indicate that individual groups differ from each other. Thanks to knowing the average values for particular groups, it is possible to indicate groups of transactionally active customers (groups 14, 5 and 6) and customers who use accounts less frequently (groups 3, 8, 4 and 1). The most affluent group of customers with very high means is without a doubt group 14.

Thanks to the use of the differential evolution algorithm to group the bank's clients, we can, in a relatively short period of time, get information on how many natural groups exist. Moreover, the number of groups is calculated by the algorithm, not imposed in advance. The algorithm evaluated and compared the obtained results for other candidate solutions in subsequent iterations, recognizing, according to the values of the objective function, that the optimal division of this group of customers contains 9 clusters.

## 5. Conclusions

The differential evolution algorithm is a promising approach to optimization because it generates a whole set of solutions that can be easily adapted to carry out the optimization again. The fact that a set of solutions is retained, and not just the best solution, allows faster adaptation to new conditions using the previously made calculations. It is resistant in terms of the choice of parameters as well as the regularity in which it finds the global optimum. The algorithm is a direct search solution method, versatile enough to solve problems whose objective function lacks the analytical description needed to determine the gradient. The algorithm is also very simple to use and modify.

Evolutionary algorithms, in particular the differential evolution algorithm, do well with continuous variables when grouping clients. Customers from particular groups can be synthetically described by the mean vector for variables used in clustering. Customers with the same basket of products, but differing in the level of individual variables, can be effectively separated.

The results of the study show that the differential evolution algorithm can be successfully applied to group retail banks' clients. Further research might be conducted with an extended list of variables i.e. those with a wider

window of observation (maximum balance in a deposit product in the last 6 months, maximum balance in an investment product in the last 6 months). It would also be advisable to analyse exclusively deposit clients or credit clients with specific variables calculated and populated for those groups.

## Bibliography

- Chou, C. H., Su, M. C. and Lai, E. (2004) "A New Cluster Validity Measure and Its Application to Image Compression". *Pattern Analysis and Applications* 7 (2): 205–20, <https://doi.org/10.1007/s10044-004-0218-1>.
- Das, S., Abraham, A. and Konar, A. (2008) "Automatic Clustering Using an Improved Differential Evolution Algorithm". *IEEE Transactions on Systems, Man, and Cybernetics – Part A: Systems and Humans* 38 (1): 218–37, <https://doi.org/10.1109/TSMCA.2007.909595>.
- Das, S., Abraham, A. and Konar, A. (2009) *Metaheuristic Clustering*, vol. 178. Springer.
- Das, S., Mullick, S. S. and Suganthan, P. N. (2016) "Recent Advances in Differential Evolution – an Updated Survey". *Swarm and Evolutionary Computation* 27: 1–30, <https://doi.org/10.1016/j.swevo.2016.01.004>.
- Everitt, B. S., Landau, S., Leese, M. and Stahl, D. (2011) *Cluster Analysis: Wiley Series in Probability and Statistics*. Wiley.
- Feoktistov, V. and Janaqi, S. (2004) "New Strategies in Differential Evolution" in *Adaptive Computing in Design and Manufacture VI*. London: Springer, pp. 335–46.
- Gan, G., Ma, C. and Wu, J. (2007) *Data Clustering: Theory, Algorithms, and Applications*, vol. 20. Siam.
- Giridhar, S., Notestein, D., Ramamurthy, S. and Wagle, L. (2011) "Od złożoności do orientacji na klienta". *Executive Report, IBM Global Business Services*.
- Storn, R. (1995) "Differential Evolution – a Simple and Efficient Adaptive Scheme for Global Optimization over Continuous Spaces". *Technical Report* 11. International Computer Science Institute.
- Storn, R. and Price, K. (1997) "Differential Evolution – a Simple and Efficient Heuristic for Global Optimization over Continuous Spaces". *Journal of Global Optimization* 11 (4): 341–59.