

A parallel genetic algorithm of feature selection for analysis of complex system

V.V. Mokshin^{1,2}, I.R. Saifudinov², L.M. Sharnin², M.V. Trusfus², P.I. Tutubalin²

¹Siemens Engineering Center, Automation and Control System Laboratory, Karl Marx Str. 10, Kazan, Russia, 420111

²Kazan National Research Technical University named after A.N. Tupolev, Karl Marx Str. 10, Kazan, Russia, 420111

Abstract. The paper shows an approach of important features selection characterizing the evolution of the complex system. A parallel genetic algorithm is proposed to solve this problem. As a result of the proposed approach, the search for the best number of parallel evolutionary paths for the feature selection is carried out. The effectiveness of the proposed approach is demonstrated on the basis of the data analysis of production enterprise functioning. This paper also shows comparison results of parallel genetic algorithm with other algorithms of feature selection by standard deviation, Fisher criterion and multiple determination coefficient.

Keywords: complex systems, feature selection, parallel genetic algorithm, regression model.

1. Introduction

Most of the complex systems such as social systems, computer vision systems, companies are characterized by many input and output features. The Black box conception is used for an analyses such complex systems. Programs introduced in [1] could be used in modeling such systems. Developing methods of forming mathematical models is one of important cybernetic tasks [2, 3]. Moreover, the main task is selecting the most important features characterizing the complex system that will be used for constructing result functions like

$$y_i = f_j(x_1, x_2, \dots, x_M), \quad j = \overline{1, K}, \quad (1)$$

where K is the number of dependent result output features of the system functioning, M is the number of input features affecting to the system functioning, $(M + K)$ is the total number of features.

The well-known iteration methods of forward selection [4], backward selection [5] and correlation methods [6] have disadvantages of including to the model (1) are not important features. The ratio of the standard error to the mean value is significantly increased by using the test values of the input features x_i , $i = \overline{1, M}$ for equation (1) [3, 7]. Genetic algorithms for feature selection also could include to equation (1) not important features. It could happen due to the small value of mutation probability [3, 7, 8, 9].

2. Method of forming analysis model of complex system

The method for important features selection is shown in this paper. Some rather short parallel evolution algorithms of feature selection are started in this case. The frequency of occurrence

of each feature is determined by all parallel evolutionary paths. The introduced algorithm uses early stop like in neural networks [8]. It prevents the convergence of the algorithm of feature selection and reduces the ratio of the standard error to the mean values using input validation features $x_i, i = \overline{1, M}$ for equation (1) [3, 10]. If a feature is really important, the frequency of appearance will be high in all or in most of parallel evolution paths. As a result, average frequency of appearance of input features will be high only for the really important features. In the process of feature selection is defined the optimal number of parallel evolution paths.

We have used recursively self-organization of regression to form model of analysis of complex systems based on the group method of data handling in this paper. It is shown on figure 1.

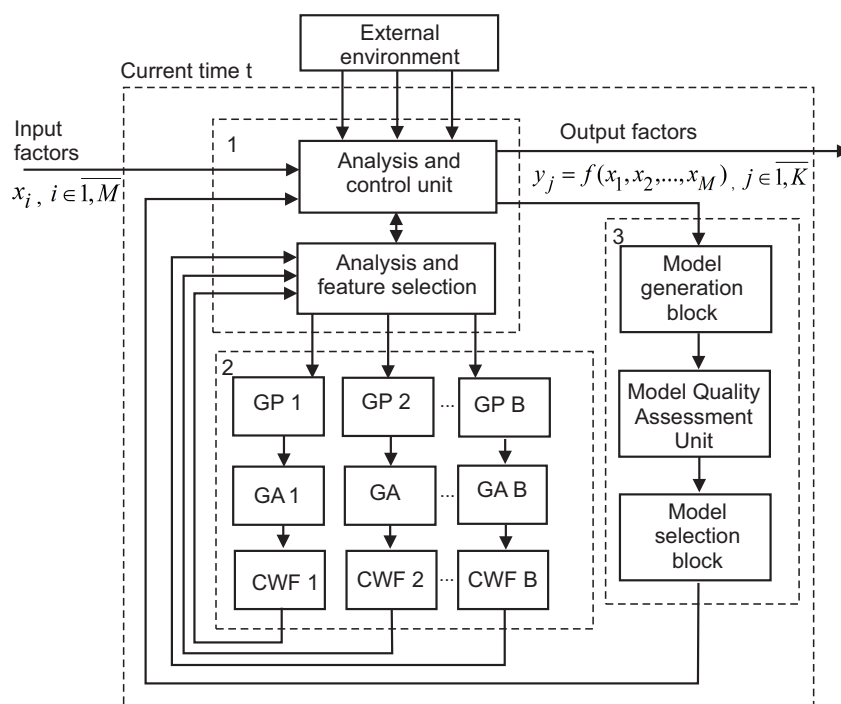


Figure 1. Block diagram of the regression equation formation taking into account changes of the system during the time (GA is a genetic algorithm with the standard genetic operators, GP is a generation of population, CWF is a calculating weights of features, B is a number of parallel evolution paths).

After that the orderly method of preference by similarity to ideal solution for decision making and solving optimization task are used [11, 12, 13, 14, 15].

The significant features selection and construction of regression equation based on generalized polynomial of Kolmogorov-Gabor equation (2) are shown on figure 1.

$$\begin{aligned}
 y = & a_0 + \sum_{i=1}^M a_i x_i + \sum_{i=1}^M \sum_{j=i}^M a_{ij} x_i x_j \\
 & + \sum_{i=1}^M \sum_{j=i}^M \sum_{k=j}^M a_{ijk} x_i x_j x_k + \dots,
 \end{aligned}
 \tag{2}$$

where weights a are decided from the following equation:

$$\mathbf{A}_j = (\mathbf{X}^T \cdot \mathbf{X})^{-1} (\mathbf{X}^T \cdot \mathbf{Y}_j), \quad j = \overline{1, K},$$

where \mathbf{A}_j is the matrix of a weights of j function (2), \mathbf{X} is the matrix of values of polynomial members (2), \mathbf{Y}_j is the matrix of j result feature y values.

As shown on figure 1 the first stage in the first block is collected the information about input features $x_i, \quad i = \overline{1, M}$ and output features $y_j, \quad j = \overline{1, K}$. Selection of significant features is realized in the first and second blocks of forming regression equation figure 1. The first and second blocks are stage of adjustment and selection of significant features. Finding the optimal number of parallel evolution ways B is realized at the stage of adjustment. Further significant features selection is implemented by the parallel genetic algorithm (PGA) with the given number of evolution ways B . After significant features selection is started, the algorithm of forming structure of regression equation 2 considering changes of the system during the time as it is shown at third stage, figure 1. and consist of models generation, models selection, analysis and control. The generation of models will be continued while the minimum value of the criterion of regularity not reached equation 3.

$$\Delta^2 = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} [y - \hat{y}], \quad (3)$$

where $n_{test} = n - n_{train}$ is the number of time observations of validation features, y is the real value of result feature in the i th time interval of the n_{test} test sample, \hat{y} is a result value in the is calculated value of result feature in the i th time interval of the n_{test} test sample in according with equation (2), n_{train} is the number of time observations of training features, n is the total number of time observations of features.

The algorithm of significant features selection and the algorithm of formation multi dimensional additive sequence equation (2) is shown in the item 1.1 and in the item 1.2.

2.1. Parallel genetic algorithm of significant feature selection and determining the optimal number of evolution paths

The idea of significant features selection includes multidimensional analysis data about features. Instead of long evolution path, some of the short parallel evolution paths B [12] are taken in consideration here. One of the reasons of such solution is to find global result. One long path of evolution can converge on a local solution. That is why for each evolution path $b = \overline{1, B}$ is running its own evolution path of selection significant features with number of generation N and number of population m . Let's say $P(b, t)$ is t -th generation of population on the b -th of the evolution path, where $t = \overline{1, N}, b = \overline{1, B}$. As a result it is needed to determine frequency of appearance of the i -th input feature during all parallel evolution paths $b = \overline{1, B}$.

The algorithm of significant features selection with B number of parallel evolution paths are shown further:

Input data:

\mathbf{A} is the matrix $n \times 1$ (time observations one of the result indicators $j = \overline{1, K}$ of the system functioning);

\mathbf{X} is the matrix $n \times M$ (time observations of input features $x_i, i = \overline{1, M}$);

n is a number of time observations of features;

m is a size of population;

N is a number of generations;

B is a number of parallel evolution paths;

v_t is a probability of mutation of t generation individuals (as default $v_t = \frac{1}{M}$);

Step 1: Sets the number of parallel evolutionary paths B .

Step 2: Let's say $P(b, N)$ is the genetic algorithm result with parameters $(\mathbf{Y}, \mathbf{X}, m, N, v_t)$.

As a fitness function in a genetic algorithm is used the generalized cross-validation criterion [3] equation (4).

$$F(w) = - \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i(w)}{1 - \text{tr}(\mathbf{H})/n} \right)^2 = - \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i(w)}{1 - (q + 1)/n} \right)^2, \quad (4)$$

where $\mathbf{H} = \mathbf{X}_w (\mathbf{X}_w^T \mathbf{X}_w)^{-1} \mathbf{X}_w^T$, $\text{tr}(\mathbf{H})$ is the elements of the main diagonal of matrix \mathbf{H} accordingly $i = \overline{1, n}$, n is the number of time observations of features, q is the number of input features used for describing result features $j = \overline{1, K}$, w is an individual formed combinations of input features consist of binary note [8]. Here 1 corresponds to the selected feature, 0 corresponds to ignore the feature.

Step 4: Is calculated density $r(i, b)$ of each input features x_i , $i = \overline{1, M}$ on the b -th of genetic paths characterized the frequency of occurrence equation (5),

$$r(i, b) = \frac{1}{m} \sum_{p=1}^m \omega_p(i), \quad (5)$$

where $\omega_p \in P(b, N)$, $p = \overline{1, m}$, $0 < r(i, b) < 1$.

Step 5: Determining the frequency of occurrence \bar{r}_i of the i input feature concerning to all $b = \overline{1, B}$ parallel evolutionary paths equation (6),

$$\bar{r}_i = \frac{1}{B} \sum_{b=1}^B r_i(i, b). \quad (6)$$

Step 6: Sorting the frequency of appearance \bar{r}_i in matrix of frequency in decreasing order and determining maximum of the distance d between frequency \bar{r}_i and \bar{r}_{i+1} , $i = \overline{1, M - 1}$,

Step 7: Calculating the distance d_{max} . It describes the maximum distance for selecting features between the ordered in descending order of frequencies of appearance of the features \bar{r}_i and \bar{r}_{i+1} , $i = \overline{1, M - 1}$,

$$d_{max} = \alpha \left(\frac{1}{4B} \right)^{\frac{1}{2}}, \quad (7)$$

where $\alpha = 1,645$. This parameter is described in [3]

Step 8: If the condition is true:

$$d \geq d_{max}, \quad (8)$$

then follow to step 9, else all features will be significant.

Step 9: Among the sorted weights of the features of the matrix \bar{r}_i are selected features for the next research whose weights are above the maximum distance d . Selected features are included to the array $\overline{R}(i)$, where $i = \overline{1, M}$ (figure 1). If feature is included to the model then in the corresponding position of the array $\overline{R}(i)$ will be written 1, else 0.

Output data:

$\overline{R}(i)$ is the array of selected features with the number of evolution paths B .

The main idea of searching the optimal number of the evolution paths B include the following. The algorithm of selecting significant features is run U times for each number of parallel evolution paths $b = \overline{1, B}$. As a result, is determined a three dimensional array of frequencies of appearance

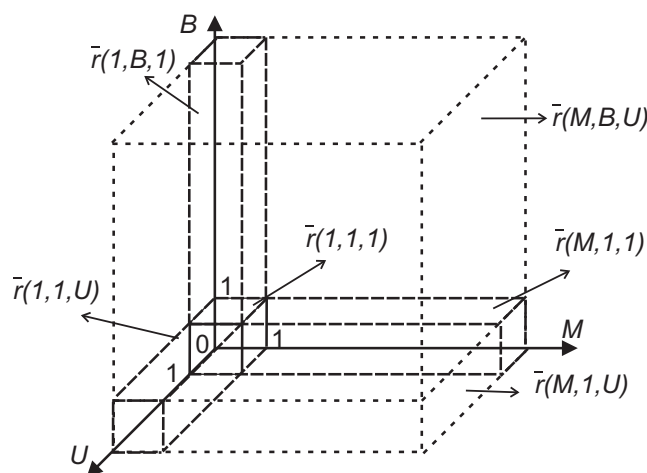


Figure 2. A schematic representation of a three dimensional array of frequency appearance of input features $x_i, i = \overline{1, M}$ for result feature y .

features $\bar{r}(M, B, U)$ for each input feature $x_i, i = \overline{1, M}$ with each number of parallel evolution paths $b = \overline{1, B}$ given the re-runs of the algorithm of features selection (figure 2).

For example, if $B = 1$ and $U = 1$ is generated array of frequencies $\bar{r}(M, 1, 1)$. If the algorithm of feature selection will be restarted U times then will be formed array of the features $\bar{r}(M, 1, U)$. For each of input feature $x_i, i = \overline{1, M}$ from each of selection $\bar{r}(i, 1, U)$ is determined $\Delta r = r_{max} - r_{min}$ and is formed the array of $\Delta \bar{r}$. After that the Δr_{min} and Δr_{max} is determined from the array $\Delta \bar{r}$. Similar reasoning is applied for each number of parallel evolutionary paths of $b = \overline{1, B}$.

The best number of parallel evolutionary paths B will be considered as a number in which the scattering of features will be the smallest.

$$\Delta r_{max} - \Delta r_{min} \rightarrow \min. \tag{9}$$

In the case when a set of selected features becomes stable, the increase the number of parallel evolutionary paths is terminated.

Selected features and the number of parallel evolutionary paths are passed to an analyzer and control (figure 2). In this case a selection of features is stopped and is started the algorithm of formation of structure of multivariate regression equation to reflect the changes in the system over time.

2.2. Formation a structure of regression equation

The formation of regression equation (2) using the method of group account of arguments can be described as genetic algorithm using the standard genetic operators. The description of population $P_2(t_2)$ with number of generation t_2 is possible to represent as a matrix form \mathbf{G} . It is a possible degrees of input features. \mathbf{A} is a matrix of coefficients of regression equations of candidates. And \mathbf{W} is a matrix of coordinates. Elements of matrix \mathbf{W} indicate the numbers of rows of the matrix \mathbf{G} . In this case, each individual $\xi \in P_2(t_2)$ with $i = \overline{1, m_2}$ and $j = \overline{1, l}$ of population size m_2 contains information about the structure of the regression equation (2) for $y_j, j = \overline{1, K}$, through the number of possible members l of the regression equation. It is shown in figure 3.

In this case a_{ij} is the coefficients of regression equations that form the matrix A . The values of w_{ij} are the elements of the matrix W and Δ^2 is a criterion of regularity (3). Thus, the

a_{i0}	a_{i1}	a_{i2}	a_{i3}	\dots	a_{il}	Δ^2
w_{i0}	w_{i1}	w_{i2}	w_{i3}	\dots	w_{il}	

Figure 3. Schematic representation of individuals $\xi_i, i = \overline{1, m_2}$ and contains information about the structure of some of the regression equation for the result feature y .

model-contender or a regression equation equation (10) given the matrices $\mathbf{G}, \mathbf{A}, \mathbf{W}$ can be represented in a formula as follows:

$$y_j = \sum_{i=0}^{l-1} a_{ij} \prod_{q=0}^{m-1} x_{q+1}^{g_{w_{j-1}iq}}, j = \overline{1, K}, \tag{10}$$

where m is number of input features, $q = \overline{0, m-1}$, $g_{w_{j-1}iq}$ are elements of the matrix \mathbf{G} , l is the number of members of the regression equation defined as equation (11)

$$l_M^p = 1 + \sum_{r=1}^p \sum_{i=1}^M i^{r-1}, \tag{11}$$

where M is a number of input features, p is a power of polynomial of equation (2).

For example, for some of the result feature y are generated regression equations from the two input features x_1 and x_2 . It's power not exceeding 1. These equations will have the following form:

$$\begin{aligned} y_1 &= a_{10} + a_{11}x_1 + a_{12}x_2, \\ y_2 &= a_{20} + a_{21}x_1 + a_{22}x_2 + a_{23}x_1x_2, \\ y_3 &= a_{30} + a_{31}x_2 + a_{32}x_1x_2, \\ y_4 &= a_{40} + a_{41}x_1 + a_{42}x_1x_2. \end{aligned}$$

$$\mathbf{G} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}, \mathbf{A} = \begin{pmatrix} a_{10} & a_{11} & a_{12} & 0 \\ a_{20} & a_{21} & a_{22} & a_{23} \\ a_{30} & a_{31} & a_{32} & 0 \\ a_{40} & a_{41} & a_{42} & 0 \end{pmatrix},$$

$$\mathbf{W} = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 0 & 1 & 2 & 3 \\ 0 & 2 & 3 & 0 \\ 0 & 1 & 3 & 0 \end{pmatrix},$$

and the regression equation with $y = y_1$ can be written in generalized form:

$$y_j = \sum_{i=0}^3 a_{1j} \prod_{q=0}^1 x_{q+1}^{g_{w_{j-1}iq}}, j = \overline{1, K},$$

Let's consider the algorithm of structure regression equation formation:

Input data:

\mathbf{Y}_j is a matrix of $n \times 1$ dimension (the time observations of one of the $j = \overline{1, K}$ of output features).

\mathbf{X} is a matrix of $n \times M$ dimension (the time observations of input features $i = \overline{1, M}$);

$R(i, B)$ is a matrix of selected features with B number of parallel evolution paths and number of input features $i = \overline{1, M}$;

m_2 is a size of population;

v_{t2} is a probability of mutation of individuals of generation t .

Step 1: The division of selection time observations of input features into two parts: training data n_{train} and test data n_{test} i.e. $n = n_{train} + n_{test}$. Generation the models of equation (10) based on training data n_{train} .

Step 2: Forming the initial population of size m_2 .

Step 3: Do Steps 3 - 5 while the value of regularity criterion Δ^2 is decreasing.

Step 4: For each model of equation (10) is calculated the criterion of regularity and selected the best models based on the test data.

Step 5: Implementation of genetic operators of equation (5).

Output data:

Matrices $\mathbf{G}, \mathbf{A}, \mathbf{W}$ and number of selected model with the least regularity criterion Δ^2 .

After the system obtained the best models based on equation (10) for the result features $y_j, j = \overline{1, K}$, optimization task can be solved.

2.3. An example of practical use of the proposed method

Lets consider the application of the proposed method for the formation the model of industrial enterprise functioning. Initially there is a quarterly set of input features $x_i, i = \overline{1, M}$ (borrowed funds, management costs, etc.) and output features $y_j, j = \overline{1, K}$ (for example, profit from sales) for some years of functioning of research object. For result feature y_3 the selection of important input features was conducted and the regression equation was constructed based on equation (10).

The results of changing the variance of the weights $\bar{r}_{min}, \bar{r}_{max}$ and $\Delta\bar{r}$ based on equation (5), (6) of input features are shown on figure 4.

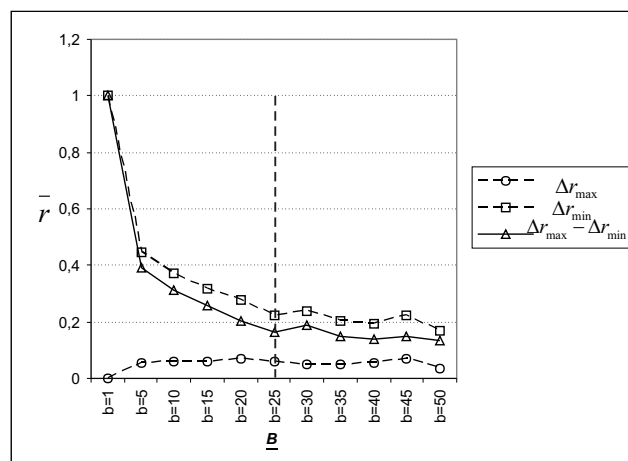


Figure 4. The variance changing of the weights of input features $x_i, i = \overline{1, 23}$ for result feature y_3 in repeated runs of the feature selection algorithm with each number of parallel evolution paths B . At each $b = 1, B$ selected the smallest and largest deviations in the weights of features $x_i, i = \overline{1, 23}$.

The result of the significant features selection determined the best number of parallel evolutionary paths $B = 25$, wherein the weights of selected features were stable. Deviations of the features weights (9) in repeated runs parallel genetic algorithm for the selection of important input features (see section 2.1) are the smallest. A further increase in the number of parallel evolutionary paths does not give significant changes of $\Delta\bar{r}$. After determining the best number of parallel evolutionary paths, selection of significant features is performed (see section 2.1.). The results of which are shown in figure 5.

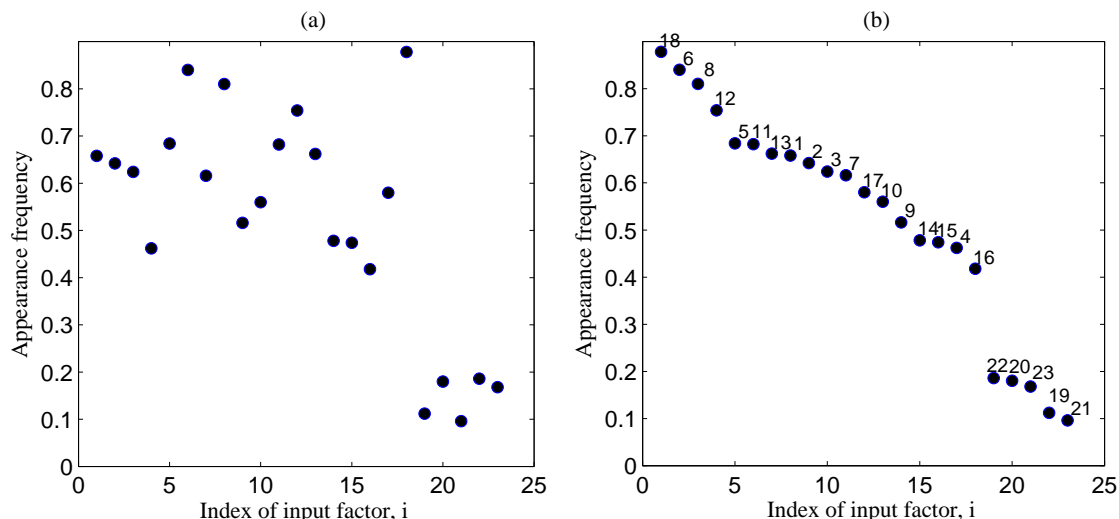


Figure 5. Significance \bar{r} of input features $x_i, i = \overline{1, M}$ for describing result feature y_3 a) in numerical order of the input features, b) sorted in descending order of \bar{r}_i . Number of evolution paths $B = 25$.

3. Conclusions

Given the algorithm for the selection of important features (see section 2.1.) and figure 5 the optimal features are shown with $x_i, i = \overline{1, 18}$. That features have high level of significance value \bar{r} (i.e. frequency of feature appearance). It means that selected features could be used for constructing regression equation of result feature y_3 as shown on equation (10). For it is used the algorithm of forming structure of regression equation (see section 2.2.).

$$\begin{aligned}
 y_3 = & -44296.34 + 5.96x_6 + 4.69 \cdot 10^{-3}x_1x_{13} + \\
 & + 8.73 \cdot 10^{-6}x_3x_5 - 5.05 \cdot 10^{-5}x_4x_9 + \\
 & + 5.32 \cdot 10^{-5}x_4x_{14} + 3.83 \cdot 10^{-5}x_6x_{18} - \\
 & - 7.93 \cdot 10^{-4}x_7x_{15} + 3.84 \cdot 10^{-4}x_8x_{14} + \\
 & + 1.62 \cdot 10^{-4}x_8x_{18} - 1.65 \cdot 10^{-3}x_{13}x_{14}.
 \end{aligned}
 \tag{12}$$

Assessment of the quality of regression equation of y_3 based on the forward and backward selection methods [4, 5], genetic algorithm of feature selection, correlation matrices and introduced method are shown on the table 1. The feature selection for the structure regression equation formation of equation (10) using the investigated method is conducted above (see

Table 1. Assessment of the regression equation quality for y_3 .

	S'_{std}	S'_{std}/\bar{y}	F	R^2
AIC	24300.46	0.04	133.79	0.990
	143471.91	0.49	10.23	0.857
BIC	28971.79	0.12	133.49	0.983
	138921.80	0.47	15.56	0.845
Correlation method	29021.63	0.12	62.61	0.992
GA	164191.84	0.56	5.09	0.864
	47452.56	0.19	32.21	0.967
	265847.88	0.91	1.48	0.504
PGA	26230.18	0.11	115.95	0.989
	51560.58	0.21	43.24	0.962

figure 5). All time series of features are separated into educational, checking and test, i.e. $n = n_{ed} + n_{check} + n_{test}$.

The table shows the results at the stage of formation and testing (forecasting) regression, using $n_{ed} + n_{test}$ and $n_{ed} + n_{test} + n_{test}$ time intervals, respectively.

Table 2. A slightly more complex table.

	Parameter 1.	p	Parameter 2.	p	Parameter 3.	p
F3	1.143	0.285	0.286	0.593	0.286	0.593
Fz	1.143	0.285	0.067	0.796	0.067	0.796

Table 3. Number features of the regression equation for y_3 .

	Q	Selected important features
AIC	10	$x_{17}, x_3, x_{11}, x_2, x_8, x_6, x_{12}, x_{15}, x_7, x_5$
BIC	7	$x_3, x_5, x_6, x_7, x_8, x_{11}, x_{15}$
Correlation method	15	$x_2, x_3, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{12}, x_{13}, x_{15}, x_{16}, x_{17}, x_{18}$
GA	14	$x_{21}, x_2x_{20}, x_3x_{14}, x_4x_{14}, x_5x_{11}, x_5x_{14}, x_6x_9, x_{13}x_{15}, x_{13}x_{23}, x_{15}, x_{15}x_{20}$
PGA	12	$x_6, x_1x_{13}, x_3x_5, x_4x_9, x_4x_{14}, x_6x_{18}, x_7x_{15}, x_8x_{14}, x_8x_{18}, x_{13}x_{14}$

The ratio of standard error to the mean S_{std}/\bar{y} at the stage of formation of the regression equation for y_3 based on described method in the paper is 0.11. And at the stage of predicting this value increased to 0.21. This increase is small relative to the results obtained using the methods of forward selection, backward elimination features, as well as using a genetic algorithm and a correlation matrix.

Given the high values of coefficients of multiple determination R^2 and Fisher criteria F at the prediction step using the proposed method, we can conclude about the high quality of the regression model generated by introduced method at the paper (see table 1). Regression equation generated by parallel genetic algorithm is shown on equation (12).

Finally, the proposed method of constructing models of complex systems based on multivariate nonlinear regression model using the methods of group accounting of arguments, a parallel genetic algorithm for the selection of important features. The proposed method allows to obtain higher quality aggregate input features. Application of parallel genetic algorithm and artificial intelligence methods allows you to generate the regression equation, allowing to make a qualitative forecast of development of complex systems.

4. References

- [1] I. Yakimov, A. Kirpichnikov, V. Mokshin, Z. Yakhina, R. Gainullin. *The comparison of structured modeling and simulation modeling of queueing systems* // Communications in Computer and Information Science, CCIS. - 2017. - Vol. 800. - P. 256-267.
- [2] Y.N. Matveev. *Basics of system theory and system analysis* // Tver State Technical University, Tver, 2007.
- [3] D.T. Larose. *Data mining methods and models* // John Wiley and Sons Inc., 2006.
- [4] H. Akaike. *Information Theory and an Extension of the Maximum Likelihood Principle* // Second International Symposium on Information Theory, 1973. - P. 267-281.
- [5] L. Breiman. *Bagging Predictors* // Machine Learning. - 1996. - Vol. 24 - P. 123-140.
- [6] M. Hall. *CFast Correlation Based Filter (FCBF) with a different search strategy: Computer and Information Sciences* // ISCIS'08. 23rd International Symposium on. IEEE, 2008.
- [7] Mu Zhu, Hugh A. Chipman. *Darwinian evolution in parallel universes: a parallel genetic algorithm for variable selection* // Technometrics. - 2006. - Vol. 48(4). - P. 491-502.
- [8] E. Cantu-Paz. *Efficient and Accurate Parallel Genetic Algorithms* // Massachusetts: Kluwer Academic Publishers, 2000. - P. 432.
- [9] V.M. Kureichik, V.V. Kureichik, L.A. Gladkov. *Theory and practise of evolution modeling* - M.: 2003.
- [10] I.I. Eliseeva, M.M. Uzbashev *General Theory of Statistics* // Finance and statistics, 1995. - P. 368.
- [11] R.U. Tzarev. Modification of the method of preference ordered by similarity to ideal solution for multi-task decision making // Information technology. - 2007. - Vol. 7(19). -P. 23.
- [12] H.R. Malada, A.G. Ivakhnenko *Inductive Learning Algorithms for Complex System Modeling* // Florida. CRC Press, 1994. - 368 p.
- [13] V.V. Mokshin. Parallel genetic algorithm of significant features selection influencing to the evolution of a complex system // Herald of Kaz. State Tech. Univ. - 2009. - Vol. 3. - P. 89-93.
- [14] V.V. Mokshin, I.M. Yakimov, R.M. Yulmetyev, A.V. Mokshin. Recursive-regression self organization of analysis and control models of complex systems // The non-linear world, 2009. - Vol. 1.- P. 48-63.
- [15] P.I. Tutubalin, V.V. Mokshin. *The Evaluation of the cryptographic strength of asymmetric encryption algorithms* // Second Russia and Pacific Conference on Computer Technology and Applications (RPC), IEEE, 2017. - P.180-183. DOI: 10.1109/RPC.2017.8168094).