# Investigation of Intelligent Technologies for Formation Forecasting Models

**Elena Skakalina[1]***

*Poltava National Technical Yuri Kondratyuk University, Ukraine*
*Corresponding Author E-Mail:Elenaskakalina2501@Gmail.Com*

## Abstract

Actually much attention is paid to the development of new intelligent information technologies for solving forecasting problems in different subject areas. The goal of solving the problem of forecasting dynamic indicators is in most cases to increase the effectiveness of making managerial decisions in conditions of uncertainty for complex distributed systems, which include economic entities. The modern global business environment dynamically forms new markets, which in turn require the use of new innovative technologies, without which it is impossible to have a competitive efficient economy in general and successful business groups in particular. In paper the research of intellectual information technologies of construction of predictive models on the basis of modified adaptive prediction methods is carried out: a neuro-network group method of data handling and a hybrid genetic algorithm with fuzzy predictive block with the purpose of justification of their use for different subject areas. Exactly these technologies are relevant and promising for improving the accuracy of forecasts.

*Keywords: forecasting, genetic algorithm , information technologies, neuro-network group method of data handling, theory of fuzzy sets*

## 1 Introduction

The need for the transition of any national economy to sustainable development is impossible without the creation of conditions for the sustainable development of regional socio-economic systems and organizational and technical systems that in general belong to the class of complex distributed systems (CDS). CDS as management objects are characterized by multifactority, significant nonlinearity of functional dependencies between factors, active influence on the control system. Examples of CDS are oil and gas producing organizations, industrial and agricultural holdings, army units in combat situation, corporate information systems. The processes of creation, operation and management of geographically distributed systems are complex and expensive. To improve the efficiency of managing such CDS, Decision Support Systems (DSS) are widely used, with mathematical identification models used in the mathematical support. As one of the approaches to constructing mathematical models of complex systems of arbitrary nature, the group method of data handling (GMDH), proposed by academician A. Ivakhnenko [1]. The expediency of applying the GMDH for constructing the CDS models is determined by its ability to provide a perfectly acceptable forecast error in the conditions of the multifactority of the CDS as a managed object and the limited training sample. It should be borne in mind that when making managerial decisions based on forecasting the development of the situation, there is a weak mathematical formalization of some socio-economic processes, which is reflected in a limited amount of statistical data. It is possible to minimize these problems on the basis of applying the neuro-fuzzy method of group accounting of arguments, constructing hybrid prediction algorithms with the inclusion of genetic algorithms in the structure of information technologies, and using the theory of fuzzy sets to solve prediction problems.

## 2 Analysis of existing approaches to the construction of predictive models

Numerical and fractal analysis of time series, methods of partial differential equations, multifactor regression analysis, neural networks, genetic algorithms, fuzzy sets can be called quite widespread methods of mathematical modeling [2-4]. The formal formulation of the problem of constructing a predictive model (PM) in the most general case can be represented thus:

$$Y := F(X) + \Delta , \qquad (1)$$

where Y is the output value of the PM; $X = [x_1, \ldots x_n]$ - a vector of factors affecting the output value of the PM; F is a mapping operator or a function; $\Delta$ - deviation, characterizing the error of the PM.

It should be noted that at the moment the objects of forecasting have changed. If earlier the main studies were focused on the development of large complex programs, studying the dynamics of the development of large industries, then fundamentally new forecast objects appeared, to which the CDS belong. Among the main properties of CDS are:

- ✓ the presence of subsystems, connected by certain criteria;
- ✓ relationships and connections between subsystems are always stronger than links with the external environment;
- ✓ the principles of combining subsystems can be contradictory.

As part of the study of such CDS, the range of predicted indicators has significantly expanded, the forecasting problems themselves have changed. Often the forecasts are scenario-based, the

principle "what happens if ..." is a mandatory preliminary stage and justification for large international investment programs.

The accuracy of the forecasting result depends on such factors as the amount of historical information gathered for the forecast (value n), the required magnitude of the forecasting horizon, the presence of external factors represented in the form of concomitant time series [5] and affecting the predicted value, the presence of distortions (abnormal emissions and passes). A valid choice of the used forecasting choice is one of the determining factors for obtaining a reliable forecast.

Analysis of the use of existing prediction methods in different subject areas has shown that the share of methods based on artificial intelligence technologies is steadily growing. The classical methods of forecasting include the method of least squares, multiple linear regression, multiple nonlinear regression. Methods of artificial intelligence include neural network methods, evolutionary modeling, methods of self-organization, fuzzy logical inference. The disadvantage of most of the classical methods of forecasting is the impossibility of taking into account the change in the environment in which the process is realized, therefore the methods of artificial intelligence that can be integrated with other approaches, both classical and evolutionary, are actively developing. They differ in the ability to parallelize and adaptive capabilities. In addition, the development of hardware and software provides more and more powerful computing resources, on which it becomes possible to implement complex intelligent forecasting algorithms. In other words, simultaneously complicating the development of information technology makes it possible to solve increasingly complex forecasting problems.

## 2.1. The Paper Should Have the Following Structure

As the initial data there is a matrix $A=(X_1, X_2,…,X_n,Y)$, where $X_i$, i=1,n, and Y are column vectors of dimension m, $X_i$ are input factors, and Y is the resulting parameter. It is necessary to identify the dependence:

$$Y = F(X_1, X_2, …, X_n) \tag{2}$$

the Kolmogorov-Gabor polynomial.

When the order of the polynomial is increased, the accuracy of the approximation of the function F (x) to it increases the accuracy of the function approaching it first, and then decreases. When the maximum accuracy is attained, the process of complicating the polynomial ends. The array of experimental points is divided into at least two subsets: the training $N_1$ and the verification $N_2$. On the set $N_1$, the least-squares method defines the coefficients vectors $A_{ij} = \{A_{0tj}, A_{tj}, ..., A_j\}$ of each particular description. In the most general case, the GMDH solves the multidimensional optimization problem of the model by means of the procedure of enumeration:

$$g = \arg \min CR(g), \quad CR(g)= f(P,S,\eta,T,V) \tag{3}$$

where $g \in G$ is the set of models under consideration, CR is the external quality criterion for the model, P is the number of sets of variables, S is the complexity of the model, ŋ is the variance of the interference, T is the number of transformations in the sample of data, and V is the number of kinds of the reference function. The basic reference function corresponds to P = S and the problem has the form of a simple one-dimensional CR (g) = f (S). GMDH allows simultaneously obtaining optimal model structure and dependence of output parameters on the most significant initial parameters. GMDH successfully solves problems of long-term and short-term forecasting, approximation of multifactor processes, clustering of data samples, self-organization of multi-row neural networks with active neurons, recognition by probabilistic search algorithms.

## 2.2. Artificial Intelligence Tools

Methods of the theory of fuzzy sets together with neural networks and methods of evolutionary modeling belong to the paradigm of "Soft Computing" [6,7].

The essence of fuzzy logic (FL), proposed by L. Zadeh, boils down to the following points:
• it uses linguistic variables (instead of ordinary numeric variables) or in addition to them;
• Simple relations between variables are described using fuzzy statements;
• complex relationships are determined by fuzzy algorithms.

Let's briefly consider the essence of FL, which is a convenient way of mapping the input space to the output (Figure 1).
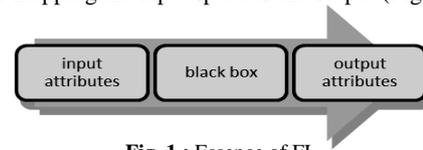


**Fig. 1 :** Essence of FL

As a black box, in principle, there can be different systems, in particular: neural networks, expert systems, differential equations, etc. In our case, the role of the black box is performed by FL. At the heart of FL, mapping input space to output, the mechanism of this mapping is used as a set of rules of the form: if - then (if - that). Example: if the water is hot, then you need to turn off the hot water tap. All rules are evaluated in parallel, their order is not important. Before building a system described by the rules, it is necessary to define all members that will be used in the system, and adjectives to describe them (for example, water can be cold, hot, warm, etc.).

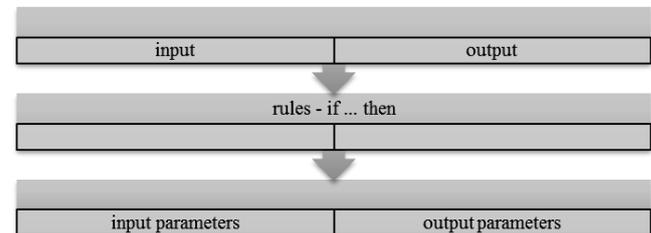The diagram in Figure 2, explains the operation of the system with FL.



**Fig. 2 :** Operation of the FL system

Using this diagram, you can define FL as a method that interprets the values of the output vector and, based on a set of rules, assigns values to the output vector.

We indicate the reasons for choosing the systems with FL:
- conceptually easier to understand;
- Flexible system and resistant to inaccurate input data;
- can simulate nonlinear functions of arbitrary complexity;
- the experience of experts is taken into account;
- based on the natural language of human communication.

A fuzzy set (FS) is defined as a set without clear, defined boundaries. It can only contain elements with a partial degree of membership. To better understand the difference between clear and fuzzy sets, let's compare them at the level of definitions.

Let $E$ be a universal set, $x$ - an element of the set $E$, and $R$ - a property. A clear subset $A$ of the universal set $E$ whose elements satisfy the property $R$ is defined as the set of ordered pairs $A = \{\mu A (x) / x\}$, where $\mu A (x)$ is the characteristic membership function (or simply membership function) taking values in some set $M$ (for example, $M = [0,1]$).

The membership function (MF) indicates the degree (or level) of the membership of the element x in the subset of $A$. The fuzzy sub-set differs from the usual one in that for the elements x of $E$ there is no unique "yes-no" answer to the property of $R$.

This concept is used in many areas. But there are situations in which this concept will lack flexibility.

Genetic algorithms are very popular methods of solving optimization problems. They are based on the genetic processes of biological organisms: biological populations develop over several generations, subject to the laws of natural selection and the "survival of the fittest" principle, discovered by Charles Darwin.

The genetic algorithm is based on the use of evolutionary principles to find the optimal solution. The very idea itself looks rather intriguing and curious to implement it, and numerous positive results only foment interest on the part of researchers. Genetic algorithms in various forms have applied to many scientific and technical problems. Genetic algorithms were used to create other computing structures, for example, automata or sorting networks. In machine learning, they were used in designing neural networks or controlling robots. They have also been used to model development in various subject areas, including biological (ecology, immunology and population genetics), social (such as economics and political systems), and cognitive systems.

A simple GA randomly generates an initial population of structures. The GA operation is an iterative process that continues until the specified number of generations or some other stopping criterion is satisfied. At each GA generation, selection is made proportionally to fitness, a single-point crossover and a mutation. First, proportional selection assigns to each structure the probability $Ps(i)$ equal to the ratio of its fitness to the total fitness of the population:

$$Ps(i) = \frac{f(i)}{\sum_{i=1}^{n} f(i)} \qquad (4)$$

Then all the n individuals are selected (with replacement) for further genetic processing, according to the *Ps (i)* value. The simplest proportional selection - roulette-wheel selection - selects individuals using *n* "launches" of the roulette. The roulette wheel contains one sector for each member of the population. The size of the *i*-th sector is proportional to the corresponding value *Ps (i)*. With such selection, members of a population with a higher fitness are more likely to be selected more often than individuals with low fitness.

After selection, *n* selected individuals undergo a crossover (sometimes called recombination) with a given probability *Pc*. *N* lines are randomly divided into *n / 2* pairs. For each pair with probability *Pc*, a crossover can be used. Accordingly, with a probability of *1-Pc*, the crossover does not occur, and unmodified individuals pass to the mutation stage. If a crossover occurs, the resulting offspring replace their parents and move on to the mutation.

The single-point crossover works as follows. First, one of the l-1 break points is randomly selected. (The break point is the section between adjacent bits in the row.) Both parent structures are split into two segments at this point. Then, the corresponding segments of different parents are stuck together and two genotypes of descendants are obtained.

For example, suppose one parent consists of 10 zeros, and the other - of 10 units. Suppose that the point 3 is chosen from the 9 possible points of discontinuity. The parents and their descendants are shown below.

*Crossover*
*Parent 1* 0000000000 000 ~ 0000000 -> 111 ~ 0000000 1110000000
*Descendant 1*
*Parent 2* 1111111111 111 ~ 1111111 -> 000 ~ 1111111 0001111111
*Descendant 2*

After the crossover stage ends, the mutation operators are executed. In each line that undergoes a mutation, each bit with the probability Pm changes to the opposite one. The population obtained after the mutation records over the old one and this completes the cycle of one generation. Subsequent generations are treated in the same way: selection, crossover and mutation.

Genetic algorithms have several advantages, for example, such as:

• GA do not require any information about the behavior of the function (for example, differentiability and continuity);

• gaps existing on the surface of the response have a negligible effect on the overall optimization efficiency;

• GA relatively stable to hit in local optima;

• GA are suitable for solving large-scale optimization problems;

• GA can be used for a wide range of tasks;

• GAs are easy to implement;

• GA can be used in tasks with a changing environment.

At the same time, there are a number of difficulties in the practical use of GA, namely:

• using GA it is problematic to find an exact global optimum;

• GA is ineffective in the case of optimizing a function that requires a long calculation time;

• GA is not easy to simulate to find all solutions of the problem;

• it is not possible to find optimal coding of parameters for all tasks;

• in multi-extremal problems, the GA faces a set of attractors;

• GA is difficult to apply for isolated functions. Isolation ("search for needles in a haystack") is a problem for any optimization method, since the function does not provide any information suggesting in which area to search for a maximum;

• additional noise (noise) spreads the fitness values of the chim, so often even good small-order shims do not pass the selection, which slows down the search for a GA solution;

• for some functions, small-scale shims lead the population to a local optimum. Such a characteristic of a function is called deception.

## 3 Selection of Input Data

The genetic algorithm for selection of input data - *Genetic Algorithm Input Selection* of *ST Neural Networks* package implements an elegant automated approach to the selection of meaningful input data. We can consider it an "intellectual" form of trial and error. This algorithm examines a large number of combinations of input variables using probabilistic and generalized regression neural networks. Networks of these types are chosen because for them the total learning / evaluation time is very short, and also because they suffer very much from the presence of unnecessary input variables (and therefore are a good means of detecting them). Each possible variant of a set of input variables can be represented as a bitmask. A zero in the corresponding position means that this input variable is not included in the input set, unit is what is included. Thus, the mask is a string of bits - one for each possible input variable - and the *Genetic Algorithm Input Selection* optimizes this bit string.

The algorithm follows a set of such masking strings, evaluating each of them for a control error (if monitoring observations were specified, if not, then a training error is used instead). According to the values of the error, the best masks variants are selected, which are combined with each other with the help of artificial genetic operations: crossing and mutation (random changes of individual bits). As initial data, data on 60 countries and 47 economic indicators that determine the level of economic development of the country were considered. It is necessary to select from 45 the minimum number of the most significant variables (those that most influence the definition of the economic level). Launch the *Neural Networks* module and go to the *Advanced* tab. Choose the tool - *Dimension reduction*. Algorithm checks only those variables that are designated as input variables.

All economic indicators are continuous and input, variable "*The level of economic development*" - input and categorical. Click *OK*. On the Quick tab, select the *Genetic Algorithm*. On the *Genetic algorithm* tab, set the *algorithm settings*. On the *Interactive* tab, select *Only the last (best) set of characteristics*. We leave the remaining settings unchanged and click *OK*.

After training, we have a table of the results of the work of the genetic algorithm. The last line shows a set of variables,

which is the best from the point of view of the procedure of the genetic algorithm.

As a result of the operation of reducing the dimension of the genetic algorithm, it can be noted that the selection errors were minimal (0.000909), which indicates the effectiveness of this method.

In the course of the selection, the following 12 variables out of 45 were noted as the most significant:

1. Gross capital formation ( $ )
2. Gross national expenditure ( $ )
3. GDP per capita  ( $ )
4. GNP per capita   ( $ )
5. Total reserves (including gold ( $ )
6. Import of goods and services ( $ )
7. Inflation, consumer prices (annual%)
8. Foreign direct investment, net inflow (  $)
9. Final consumption expenditure, etc. ( $ )
10. Adjusted savings: consumption of fixed capital  ( $ )
11. Net income from abroad  ( $ )
12. Export of goods and services  ( $ )

# 4 Forecasting

To solve the problem of forecasting economic indicators, we will consider the process of developing a fuzzy model of a hybrid network. The essence of this problem is that, knowing the dynamics of the change in the indicator for a fixed time interval, predict its value at a certain point in time in the future.

Traditionally, various models and technical analysis based on the use of various indicators are used to solve this problem. At the same time, the presence of implicit trends in the dynamics of changes in economic indicators makes it possible to apply the model of adaptive neuro-fuzzy networks.

The task of forecasting economic indicators, typical in practice of analysis and forecasting, will be implemented in two development environments of MATLAB Fuzzy Logic Toolbox and GMDH Shell for comparison with another method of forecasting - the method of group accounting of arguments. We will analyze the generated fuzzy inference system and the GMDH model, and give a comparative description of the results obtained.

The initial data are indicators of the economic sector of the United States of America between 1963 and 2014, provided by the World Bank website [8].

## 4.1 Implementing in the MATLAB Fuzzy Logic Toolbox Environment.

First of all, it is necessary to build a fuzzy model of a hybrid network. Suppose that the fuzzy model of the hybrid network will contain 4 input variables. In this case, the first input variable will correspond to the value of GDP for the current year, the second - to the previous year, i.e. for a year *(i - 1),* where i denotes the current year. Then the third input variable will correspond to the value at *(i - 2)* year, and the fourth - to *(i - 3)* year.

The relevant training data can be summarized in a separate table 1.

**Table 1:** Fuzzy GDP model for 1963-2010. (trillion US dollars)

| year | Input1 | Input2 | Input3 | Input4 | Output |
|---|---|---|---|---|---|
| 1966 | 3827,5271 | 3573,9412 | 3374,5152 | 3243,8431 | 4146,3166 |
| 1967 | 4146,3166 | 3827,5271 | 3573,9412 | 3374,5152 | 4336,4266 |
| 1968 | 4336,4266 | 4146,3166 | 3827,5271 | 3573,9412 | 4695,9234 |
| 1969 | 4695,9234 | 4336,4266 | 4146,3166 | 3827,5271 | 5032,1447 |
| 1970 | 5032,1447 | 4695,9234 | 4336,4266 | 4146,3166 | 5246,9617 |
| 1971 | 5246,9617 | 5032,1447 | 4695,9234 | 4336,4266 | 5623,5884 |
| 1972 | 5623,5884 | 5246,9617 | 5032,1447 | 4695,9234 | 6109,6924 |
| 1973 | 6109,6924 | 5623,5884 | 5246,9617 | 5032,1447 | 6741,1011 |
| 1974 | 6741,1011 | 6109,6924 | 5623,5884 | 5246,9617 | 7242,3242 |
| 1975 | 7242,3242 | 6741,1011 | 6109,6924 | 5623,5884 | 7819,9590 |
| 1976 | 7819,9590 | 7242,3242 | 6741,1011 | 6109,6924 | 8611,4615 |
| 1977 | 8611,4615 | 7819,9590 | 7242,3242 | 6741,1011 | 9471,5287 |
| 1978 | 9471,5287 | 8611,4615 | 7819,9590 | 7242,3242 | 10587,4160 |
| 1979 | 10587,4160 | 9471,5287 | 8611,4615 | 7819,9590 | 11695,3634 |
| 1980 | 11695,3634 | 10587,4160 | 9471,5287 | 8611,4615 | 12597,6455 |
| 1981 | 12597,6455 | 11695,3634 | 10587,4160 | 9471,5287 | 13993,3585 |
| 1982 | 13993,3585 | 12597,6455 | 11695,3634 | 10587,4160 | 14439,0151 |
| 1983 | 14439,0151 | 13993,3585 | 12597,6455 | 11695,3634 | 15561,2681 |
| 1984 | 15561,2681 | 14439,0151 | 13993,3585 | 12597,6455 | 17134,3157 |
| 1985 | 17134,3157 | 15561,2681 | 14439,0151 | 13993,3585 | 18269,2793 |
| 1986 | 18269,2793 | 17134,3157 | 15561,2681 | 14439,0151 | 19115,2403 |
| 1987 | 19115,2403 | 18269,2793 | 17134,3157 | 15561,2681 | 20100,7887 |
| 1988 | 20100,7887 | 19115,2403 | 18269,2793 | 17134,3157 | 21483,1145 |
| 1989 | 21483,1145 | 20100,7887 | 19115,2403 | 18269,2793 | 22922,4655 |
| 1990 | 22922,4655 | 21483,1145 | 20100,7887 | 19115,2403 | 23954,5234 |
| 1991 | 23954,5234 | 22922,4655 | 21483,1145 | 20100,7887 | 24404,9948 |
| 1992 | 24404,9948 | 23954,5234 | 22922,4655 | 21483,1145 | 25492,9556 |
| 1993 | 25492,9556 | 24404,9948 | 23954,5234 | 22922,4655 | 26464,7833 |
| 1994 | 26464,7833 | 25492,9556 | 24404,9948 | 23954,5234 | 27776,8065 |
| 1995 | 27776,8065 | 26464,7833 | 25492,9556 | 24404,9948 | 28782,3252 |
| 1996 | 28782,3252 | 27776,8065 | 26464,7833 | 25492,9556 | 30068,2272 |
| 1997 | 30068,2272 | 28782,3252 | 27776,8065 | 26464,7833 | 31572,6352 |
| 1998 | 31572,6352 | 30068,2272 | 28782,3252 | 27776,8065 | 32949,3138 |
| 1999 | 32949,3138 | 31572,6352 | 30068,2272 | 28782,3252 | 34620,8429 |
| 2000 | 34620,8429 | 32949,3138 | 31572,6352 | 30068,2272 | 36449,9295 |
| 2001 | 36449,9295 | 34620,8429 | 32949,3138 | 31572,6352 | 37273,5339 |
| 2002 | 37273,5339 | 36449,9295 | 34620,8429 | 32949,3138 | 38165,9892 |
| 2003 | 38165,9892 | 37273,5339 | 36449,9295 | 34620,8429 | 39677,3018 |
| 2004 | 39677,3018 | 38165,9892 | 37273,5339 | 36449,9295 | 41921,7141 |
| 2005 | 41921,7141 | 39677,3018 | 38165,9892 | 37273,5339 | 44307,8326 |
| 2006 | 44307,8326 | 41921,7141 | 39677,3018 | 38165,9892 | 46437,1073 |

| 2007 | 46437,1073 | 44307,8326 | 41921,7141 | 39677,3018 | 48061,4215 |
| 2008 | 48061,4215 | 46437,1073 | 44307,8326 | 41921,7141 | 48401,4865 |
| 2009 | 48401,4865 | 48061,4215 | 46437,1073 | 44307,8326 | 47001,4282 |
| 2010 | 47001,4282 | 48401,4865 | 48061,4215 | 46437,1073 | 48377,3938 |
| 2011 | 48377,3938 | 47001,4282 | 48401,4865 | 48061,4215 | 49803,4929 |
| 2012 | 49803,4929 | 48377,3938 | 47001,4282 | 48401,4865 | 51495,8748 |
| 2013 | 51495,8748 | 49803,4929 | 48377,3938 | 47001,4282 | 53041,9814 |

The volume of the training sample thus obtained is equal to 45, which corresponds to the dynamics of GDP per capita from 1963 to 2011. Thus, the data for 2012-2014 are not included in the training sample and will be used to verify the adequacy of the constructed fuzzy model.

Save the training sample with the extension .**dat,** run the *ANFIS*-editor using the command **anfisedit** from the command line MATLAB and load the created model (Figure 3).
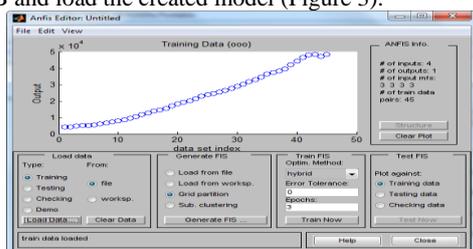


**Fig. 3:** Downloaded sample in Anfis editor

Generate a system of fuzzy output of Sugeno type by pressing the **Generate FIS** ... button. In the window that appears, for each input variable, set the **gbellmf** type membership functions. For the output variable, set the membership function to **constant.** To learn the hybrid network, we'll select the **hybrid** method with the error level 0 and the number of cycles 400. Let's start learning the hybrid network. After learning the hybrid network, you can look at the structure of the constructed fuzzy model. In this system, 81 rules were obtained (Figure 4).
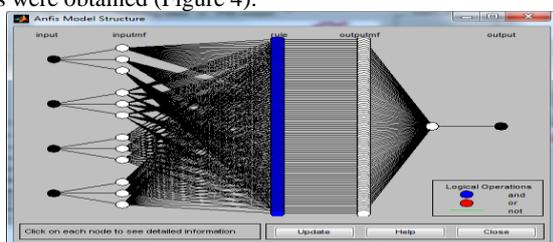


**Fig. 4:** Structure of fuzzy model

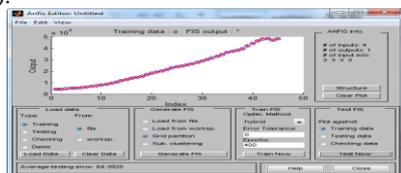We test the obtained system of fuzzy inference on the training set of data (Figure 5).



**Fig. 5 :** Comparative chart of test results and input sample

Now let's check the accuracy of the constructed system on those data that are not included in the training sample. Suppose that the current year is 2011, and we will make a forecast for GDP per capita for 2014. To do this, we use the **evalfis** command-line function, since the graphical tools of the **Fuzzy Logic Toolbox** package give too much error.

Define the incoming values (GDP per capita for 2011, 2010, 2009 and 2008)

```
>> input = [48377.3938 47001.4282 48401.4865 48061.4215]
input =
  1.0e + 004 *
  4.8377 4.7001 4.8401 4.8061
```

We connect the **fis**-structure

```
>> fis = readfis ('VVP')
fis =
    name: 'VVP'
    type: 'sugeno'
  andMethod: 'prod'
  orMethod: 'probor'
  defuzzMethod: 'wtaver'
  impMethod: 'prod'
  aggMethod: 'sum'
    input: [1x4 struct]
    output: [1x1 struct]
    rule: [1x81 struct]
```

We calculate the predicted value using the **evalfis** function

```
>> znach = evalfis (input, fis)
> In evalfis at 76
znach =
  4.9754e + 004
```

The value obtained is $ 49754. The present value is $ 49803. The difference is $ 49.40.

## 4.2 Stages Of Implementation Of Intellectual

### Technology

Information technology of FFM construction has the following structure:

Step 1 - the initial set of economic indicators is formed (the dimension of the set is 45).

Step 2 - using the genetic algorithm, the dimension of the initial set is reduced from 45 indicators to 12, selected by GA as the most significant. Selection errors were minimal (0.000909).

Step 3 - with the help of the apparatus of fuzzy sets, the FNNM of each indicator is constructed. After the generation of the *Sugeno* type fuzzy output system, a hybrid network is trained with a specified error level of 0 and a number of cycles of 400. In the FNNM, 81 rules are obtained.

Step 4-Test the FNNM on the training set of data. The test result showed the deviation of the actual US GDP from the model (FNNM) in the range (0,099376% - 0,18479%), which confirms the correctness of the proposed information technology.

Step 5 - with the help of FNNM we compute the forecast values of the 12 macroeconomic indicators of the US economy selected by step 2 up to 2020. The results of the computational experiment showed that the range of deviations of the model values of the parameter from the real ones lies in the range from 0.064326% to 9.2%.

Lowering the dimension of the initial set of macroeconomic indicators was carried out using the Genetic Algorithm Input Selection software package of ST Neural Networks. The formation of the FNNM was carried out using the software of the MATLAB Fuzzy Logic Toolbox, using the adaptive neural-fuzzy inference (ANFIS) system of the *Sugeno* type.

## 4. Conclusion

GMDH has the advantage of small sample data due to the choice of model complexity, optimally takes into account informative data. The efficiency of the method has been repeatedly confirmed by the decision of many specific tasks from the areas of ecology, economics, hydrometeorology, etc. GMDH is well-known and very actively developing. The basics of the theory of structural identification of models with the minimum dispersion of prediction error are developed. Effective apparatus of this theory is the method of critical variances, allowing the first analytically solve

the actual problem: a comparative analysis of the criteria of structural identification, planning experiments, analysis methods, etc. properties, both with limited sample, and in asymptotytsi. This examines the conditions of choosing the optimal structure model based on the variance (level) noise, the length of the sample input (experimental plan) and object parameters, and a close relationship between them. By means of this theory it is established that GMDH is a method of constructing models with a minimum dispersion of prediction error, and its comparison with other methods is carried out. In the last decade, interest in GMDH has been growing worldwide, which can be explained, in addition to the known method efficiency, as well as the growing popularity of artificial neural network technology. The fact is that the structure of the GMDH can be interpreted as a neural network, the originality of which consists in the self-organization of both its structure and parameters. It turns out that the obvious advantages of GMDH include automatic generation of network structure, simplicity and performance tuning options and the ability to "roll" customized network directly in explicit mathematical expression. At the same time, there are many successful GMDH applications for solving a wide range of Data Mining tasks: classification, prediction, identification of complex systems, finding empirical dependencies, clustering, etc. Therefore, the use of GMDH algorithms at both the "lower" and "upper" levels of management of complex distributed systems is considered relevant.

## References

[1] Ivakhnenko A.G., Ivakhnenko G.A. The Review of Problems Solvable by Algorithms of the Group Method of Data Handling. *Published in Pattern Recognition and Image Analysis*, Vol. 5, No. 4, 1995, pp.527-535.

[2] Peters E. *Fractal Market Analysis. Applying Chaos Theory to Investment&Economics.* / E Peters – J. Wiley&Sons, Inc. – New York, 1994.

[3] Chen S.M. Forecasting enrolments based on fuzzy time series. – *Fuzzy sets Systems*, 1996, vol. 81, №3, p.p. 311−319.

[4] Goldberg, D. E. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, Reading, Mass., 1989.

[5] Korablev, N.M. Parallel immune algorithm of short-term forecasting based on model of clonal selection / N.M. Korablev, G.S. Ivaschenko // *Radio Electronics, Computer Science, Control:* scientific journal. – 2014. – № 2(31). – Pp. 73-78. ISSN 1607-3274.

[6] R. Bellman and L. Zadeh, *Decision-making in vague terms, Questions analysis and decision-making procedures*. Moscow, Russia: Mir, 1976, p. 172–215.

[7] L.A. Zadeh, J. Kacprzyk (Eds.) *Fuzzy logic for the management of uncertainty*. Wiley, New York, 1992.

[8] / <http://data.worldbank.org/country/united-states / Data. United States> [lastaccessdon24thApril2018].