

# The Fuzzy Group Method of Data Handling with Fuzzy Input Variables

Zaychenko Yu. (Ukraine)

NTUU “Kiev Polytechnic Institute”, Institute for Applied System Analysis, Peremogy avenue 37,  
03056, Kiev-56, Ukraine

zaych@i.com.ua

**Abstract.** *The problem of constructing forecasting models with incomplete and fuzzy input data is considered in this paper. For its solution Fuzzy Group Methods of Data Handling (FGMDH) with fuzzy inputs is suggested. The method enables to construct a forecasting fuzzy model using experimental data which are not distinct.*

*The method was implemented as software kit and experimental investigations of were carried out in the problem forecasting stock-prices at the Russian stock-exchange. The comparison of the suggested method with known methods: GMDH and fuzzy GMDH is also presented.*

## Keywords

Group method of Data Handling, fuzzy, economic indexes, forecasting

## 1 Introduction

The problem of forecasting models constructing using experimental data in terms of fuzziness, when input variables are not known exactly and determined as intervals of uncertainty is considered in this paper. The fuzzy group method of data handling is proposed to solve this problem. The theory of this method was suggested and researched in [1-5]. As it is well known, fuzzy GMDH allows to construct fuzzy models and has the following advantages:

1. The problem of optimal model finding is transformed to the problem of linear programming, which is always solvable;
2. There is interval regression model built as the result of method work out;
3. There is a possibility of adaptation of the obtained model.

The mathematical model of the problem mentioned above is built and fuzzy GMDH with fuzzy inputs is elaborated in the paper. The corresponding program, which uses the suggested algorithm, was developed. And the experimental researches and comparison of FGMDH with GMDH in the problem of stock prices forecasting were carried out and presented in this paper.

## 2 Math model of group method of data handling with fuzzy input data

Sometimes we face the problem of constructing the model using experimental data when the initial input data are fuzzy and are presented in the form of intervals. Such situation arises in the problem of macroeconomic indexes forecasting as the monthly values of indexes may be given in interval form. Let's consider a linear interval regression model:

$$Y = A_0 Z_0 + A_1 Z_1 + \dots + A_n Z_n, \quad (1)$$

where  $A_i$  are fuzzy numbers, which are described by threes of parameters  $A_i = (\underline{A}_i, \check{A}_i, \overline{A}_i)$ ,

where  $\check{A}_i$  – interval center,  $\overline{A}_i$  – upper border of the interval,  $\underline{A}_i$  – lower border of the interval,

and  $\underline{Z}_i$  – also fuzzy numbers, which are determined by parameters  $(\underline{Z}_i, \check{Z}_i, \overline{Z}_i)$ ,  $\underline{Z}_i$  – lower border,  $\check{Z}_i$  – center,  $\overline{Z}_i$  – upper border of a fuzzy number.

Then  $Y$  – output fuzzy number, which parameters are defined as follows (in accordance with L-R numbers multiplying formulas):

Center of interval:

$$\check{y} = \sum \check{A}_i * \check{Z}_i,$$

Deviation in the left part of the membership function:

$$\check{y} - \underline{y} = \sum (|\check{A}_i| * (\check{Z}_i - \underline{Z}_i) + (\check{A}_i - \underline{A}_i) * |\check{Z}_i|),$$

And the lower border of the interval:

$$\underline{y} = \sum (\check{A}_i * \check{Z}_i - |\check{A}_i| * (\check{Z}_i - \underline{Z}_i) - (\check{A}_i - \underline{A}_i) * |\check{Z}_i|),$$

The upper border of the interval

$$\overline{y} = \sum (|\check{A}_i| * (\check{Z}_i - \check{Z}_i) + |\check{Z}_i| * (\overline{A}_i - \check{A}_i) + \check{A}_i * \check{Z}_i).$$

For the interval model to be correct, the real value of input variable  $Y$  is needed to lay in the interval got by the method workflow.

So, the general requirements to estimation a linear interval model are to find such values of parameters  $(\underline{A}_i, \check{A}_i, \overline{A}_i)$  of fuzzy coefficients, which allow:

- Observed values  $y_k$  should lay in the estimation interval for  $Y_k$ ;
- The total width of estimation interval be minimal.

Input data for this task is  $Z_k = [Z_{ki}]_i$  – input training sample, and  $y_k$  are known output values,  $k = \overline{1, M}$ ,  $M$  is the number of observation points.

There are two cases of fuzzy membership functions used in this work:

- Triangular membership functions;
- Gaussian membership functions.

Quadratic partial descriptions were chosen:

$$f(x_i, x_j) = A_0 + A_1 x_i + A_2 x_j + A_3 x_i x_j + A_4 x_i^2 + A_5 x_j^2.$$

### 3 FGMDH with fuzzy input data for triangular membership function

The form of math model for triangular MF.

Let's consider the linear interval regression model:

$$Y = A_0 Z_0 + A_1 Z_1 + \dots + A_n Z_n,$$

The current task contains the case of symmetrical membership function for parameters  $A_i$ , so they can be described via a pair of parameters  $(a_i, c_i)$ .

$\underline{A}_i = a_i - c_i$ ,  $\overline{A}_i = a_i + c_i$ ,  $a_i$  is a center of an interval,  $c_i$  is the interval width,  $c_i \geq 0$ ,

$Z_i$  are also fuzzy numbers of triangular form, which are defined by parameters  $(\underline{Z}_i, \check{Z}_i, \overline{Z}_i)$ ,

$\underline{Z}_i$  – lower border,  $\check{Z}_i$  – center,  $\overline{Z}_i$  – upper border of fuzzy number.

Then  $Y$  is a fuzzy number, which parameters are defined as follows:

the center of the interval:

$$\check{y} = \sum a_i * \check{Z}_i,$$

Deviation in the left part of the membership function:

$$\check{y} - \underline{y} = \sum (a_i * (\check{Z}_i - \underline{Z}_i) + c_i |\check{Z}_i|),$$

Lower border of the interval:

$$\underline{y} = \sum (a_i * \underline{Z}_i - c_i |\check{Z}_i|)$$

Deviation in the right part of the membership function:

$$\bar{y} - \check{y} = \sum (a_i * (\bar{Z}_i - \check{Z}_i) + c_i |\check{Z}_i|) = \sum a_i \bar{Z}_i - a_i \check{Z}_i + c_i |\check{Z}_i|,$$

Upper border of the interval:

$$\bar{y} = \sum (a_i * \bar{Z}_i + c_i |\check{Z}_i|)$$

For the interval model to be correct, the real value of input variable Y should lay in the interval got by the method work.

It can be described in such a way:

$$\begin{cases} \sum (a_i * \underline{Z}_{ik} - c_i |\check{Z}_{ik}|) \leq y_k \\ \sum (a_i * \bar{Z}_{ki} + c_i |\check{Z}_{ik}|) \geq y_k, k = \overline{1, M} \end{cases}$$

Where  $Z_k = [Z_{ki}]_i$  is an input training sample,  $y_k$  are known output values,  $k = \overline{1, M}$ , M is a number of observation points.

So, the general requirements to estimation linear interval model are to find such values of parameters  $(a_i, c_i)$  of fuzzy coefficients, which enable:

- Observed values  $y_k$  lay in estimation interval for  $Y_k$ ;
- Total width of estimation interval be minimal.

These requirements can be redefined as a task of linear programming:

$$\min_{a_i, c_i} \sum_{k=1}^M (\sum (a_i * \bar{Z}_i + c_i |\check{Z}_i|) - \sum (a_i * \underline{Z}_i - c_i |\check{Z}_i|)), \quad (2)$$

under constraints:

$$\begin{cases} \sum (a_i * \underline{Z}_{ik} - c_i |\check{Z}_{ik}|) \leq y_k \\ \sum (a_i * \bar{Z}_{ki} + c_i |\check{Z}_{ik}|) \geq y_k, k = \overline{1, M} \end{cases} \quad (3)$$

Formalized problem formulation in case of triangular membership functions

Let's consider partial description

$$f(x_i, x_j) = A_0 + A_1 x_i + A_2 x_j + A_3 x_i x_j + A_4 x_i^2 + A_5 x_j^2. \quad (4)$$

Rewriting it in accordance with the model (1) needs such substitution:  $z_0 = 1$ ,  $z_1 = x_i$ ,

$$z_2 = x_j, z_3 = x_i x_j, z_4 = x_i^2, z_5 = x_j^2.$$

Then math model (2)-(3) will take the form

$$\begin{aligned} \min_{a_i, c_i} & (2Mc_0 + a_1 \sum_{k=1}^M (\bar{x}_{ik} - \underline{x}_{ik}) + 2c_1 \sum_{k=1}^M |\check{x}_{ik}| + a_2 \sum_{k=1}^M (\bar{x}_{jk} - \underline{x}_{jk}) + 2c_2 \sum_{k=1}^M |\check{x}_{jk}| + \\ & + a_3 \sum_{k=1}^M (|\check{x}_{ik}| (\bar{x}_{jk} - \underline{x}_{jk}) + |\check{x}_{jk}| (\bar{x}_{ik} - \underline{x}_{ik})) + 2c_3 \sum_{k=1}^M |\check{x}_{ik} \check{x}_{jk}| + 2a_4 \sum_{k=1}^M |\check{x}_{ik}| (\bar{x}_{ik} - \underline{x}_{ik}) + \\ & + 2c_4 \sum_{k=1}^M \check{x}_{ik}^2 + 2a_5 \sum_{k=1}^M |\check{x}_{jk}| (\bar{x}_{jk} - \underline{x}_{jk}) + 2c_5 \sum_{k=1}^M \check{x}_{jk}^2) \end{aligned} \quad (5)$$

with the following conditions:

$$\begin{aligned}
& a_0 + a_1 \underline{x}_{ik} + a_2 \underline{x}_{jk} + a_3 (-|\tilde{x}_{ik}|(\tilde{x}_{jk} - \underline{x}_{jk}) - |\tilde{x}_{jk}|(\tilde{x}_{ik} - \underline{x}_{ik}) + \tilde{x}_{ik} \tilde{x}_{jk}) + \\
& + a_4 (-2|\tilde{x}_{ik}|(\tilde{x}_{ik} - \underline{x}_{ik}) + \tilde{x}_{ik}^2) + a_5 (2|\tilde{x}_{jk}|(\tilde{x}_{jk} - \underline{x}_{jk}) + \tilde{x}_{jk}^2) - c_0 - c_1 |\tilde{x}_{ik}| - \\
& - c_2 |\tilde{x}_{jk}| - c_3 |\tilde{x}_{ik} \tilde{x}_{jk}| - c_4 \tilde{x}_{ik}^2 - c_5 \tilde{x}_{jk}^2 \leq y_k \tag{6}
\end{aligned}$$

$$\begin{aligned}
& a_0 + a_1 \bar{x}_{ik} + a_2 \bar{x}_{jk} + a_3 (|\tilde{x}_{ik}|(\bar{x}_{jk} - \tilde{x}_{jk}) + |\tilde{x}_{jk}|(\bar{x}_{ik} - \tilde{x}_{ik}) - \tilde{x}_{ik} \tilde{x}_{jk}) + a_4 (2|\tilde{x}_{ik}|(\bar{x}_{ik} - \\
& - \tilde{x}_{ik}) - \tilde{x}_{ik}^2) + a_5 (2|\tilde{x}_{jk}|(\bar{x}_{jk} - \tilde{x}_{jk}) - \tilde{x}_{jk}^2) + c_0 + c_1 |\tilde{x}_{ik}| + c_2 |\tilde{x}_{jk}| + c_3 |\tilde{x}_{ik} \tilde{x}_{jk}| + \\
& c_4 \tilde{x}_{ik}^2 + c_5 \tilde{x}_{jk}^2 \geq y_k
\end{aligned}$$

$$c_l \geq 0, \quad l = \overline{0,5}$$

As we can see, this is the linear programming problem, but there are still no limitations for non-negativity of variables  $a_i$ , so we need go to dual problem, introducing dual variables  $\{\delta_k\}$  and  $\{\delta_{k+M}\}$ .

## 4 THE DESCRIPTION OF FUZZY GMDH ALGORITHM

Let's present the brief description of the algorithm.

1. Choose of the general model type by which the dependence to be sought will be described.
2. Choose the external criterion of optimality (the criteria of regularity or non-biasedness)
3. Choice a general type of coordinate functions (type of partial description), for example, linear or quadratic one.
4. Divide the sample into training and test sub-samples.
5. Put zero values to the counter of number of models  $\mathbf{k}$  and to the counter of number of rows  $\mathbf{r}$ .
6. Generate a new partial model or the kind (5) at a training sample. Solve the LP problem (5), (6) using training sample and find the values of parameters of fuzzy coefficients.
7. Calculate using test sample the value of external criterion: non-biasedness ( $N_{CM}^{(r)}$ ) or regularity.
8.  $k = k + 1$ . If  $k \geq C \frac{2}{F}$ , then  $k = 0$ ,  $r = r + 1$ .
9. Calculate a mean value of criterion for the models of  $r$ -th iteration. Then go to step 6 if  $r = 1$ , otherwise, go to step 10.
10. If  $|N_{CM}^{(r)} - N_{CM}^{(r-1)}| \leq \mathcal{E}$  then go to step 11, otherwise select  $F$  best models and assigning  $\mathbf{r} = \mathbf{r} + 1$ ,  $\mathbf{k} = \mathbf{1}$ , go to step 6 and execute the next  $(r+1)$ -th iteration.
11. Out of  $F$  models of the previous row we select the best model using the external criterion of regularity.

The difference between GMDH and fuzzy GMDH lies in the following: in the GMDH we find out model using LSM method while in FGMDH we solve the LP problem of the form (5), (6) for finding out the model. Therefore we exclude the problem of ill-conditioned matrices occurring while using GMDH. The additional advantage of FGMDH is that as it was proven corresponding LP problem is always soluble.

### 4 Experimental results of FGMDH with fuzzy input data workflow in RTS index forecasting

For estimation of efficiency of the suggested FGMDH method with fuzzy inputs the corresponding software kit was elaborated and numerous experiments of financial markets forecasting were carried out. Some of them are presented below.

Forecasting of RTS index.

Experiment 1. RTS index forecasting (opening price)

In this experiment we used 5 fuzzy input variables, which represent stock prices of leading Russian energetic companies, which are included to the list of computations of RTS index:

LKOH – shares of “LUKOIL” joint-stock company,

EESR – shares of “РАО ЕЭС России” joint-stock company,

YUKO – shares of “ЮКОС” joint-stock company,

SNGSP – privileged shares of “Сургутнефтегаз” joint-stock company,

SNGS – common shares of “Сургутнефтегаз” joint-stock company.

Output variable is the RTS (opening price) index value of the same period (03.04.2006 – 18.05.2006).

Sample size – 32 values.

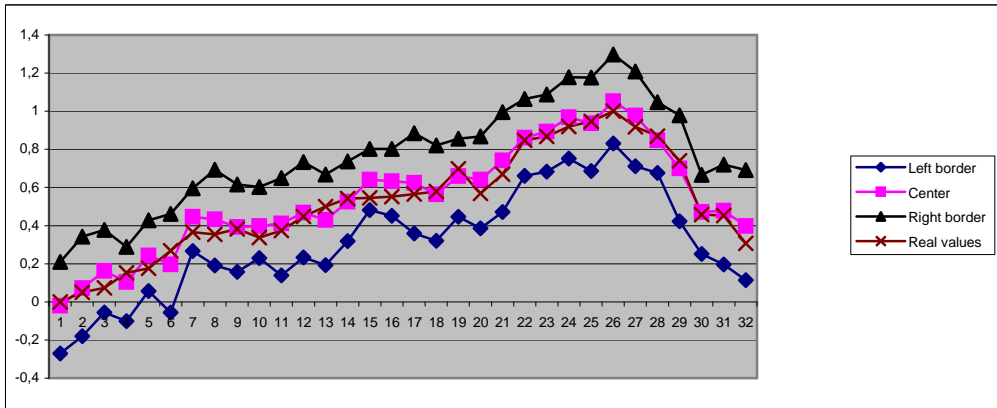
Training sample size – 18 values (optimal size of training sample for current experiment was determined experimentally by varying the size training sample).

The following results were obtained:

1. For triangular membership function

a) For normalized input data

Criterion value for current experiment were:  $MSE = 0.055557$

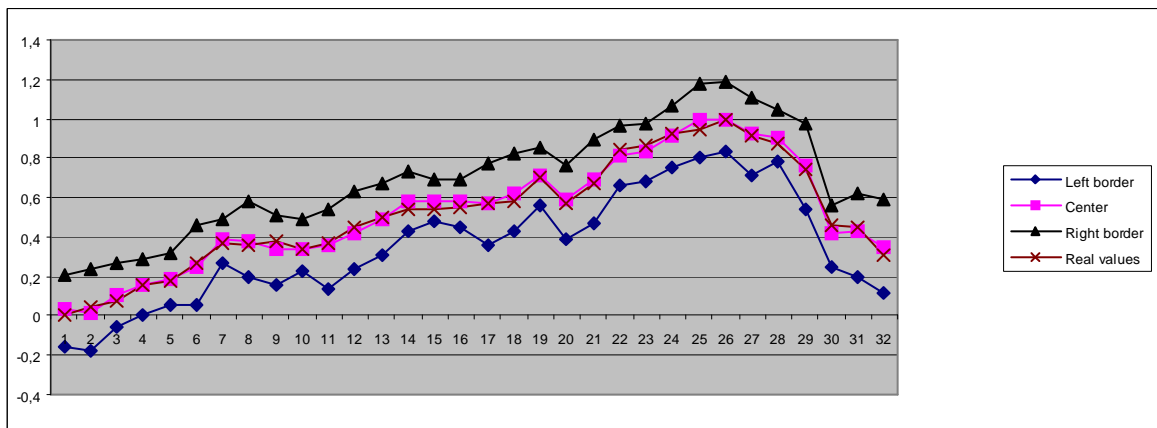


**Fig. 1.** Experiment 1 results for triangular membership function and normalized values of input variables

2. For the case of Gaussian membership function (optimal level is  $\alpha=0.8$ ) (see Fig. 2)

a) For normalized input data

Criterion values for this experiment were:  $MSE = 0.028013$



**Fig. 2.** Experiment 1 result for Gaussian MF and normalized input data

As we can see from the results of experiment 1, forecasting using triangular and Gaussian membership functions gives good results. Results of experiments with Gaussian MF are better than results of experiments with triangular MF.

**Tab.1.** For normalized data

	<b>Triangular MF</b>	<b>Gaussian MF</b>
MSE	0.055557	0.028013

**Tab.2.** For non-normalized data

	<b>Triangular MF</b>	<b>Gaussian MF</b>
MSE	18.48657	9.321461
MAPE	0.8%	0.4%

Experiment 2. Stock price forecasting

The following experiment uses stock prices of 4 leading energetic companies of Russia:

EESR – shares of “РАО ЕЭС России” joint-stock company,

YUKO – shares of “ЮКОС” joint-stock company,

SNGSP – privileged shares of “Сургутнефтегаз” joint-stock company,

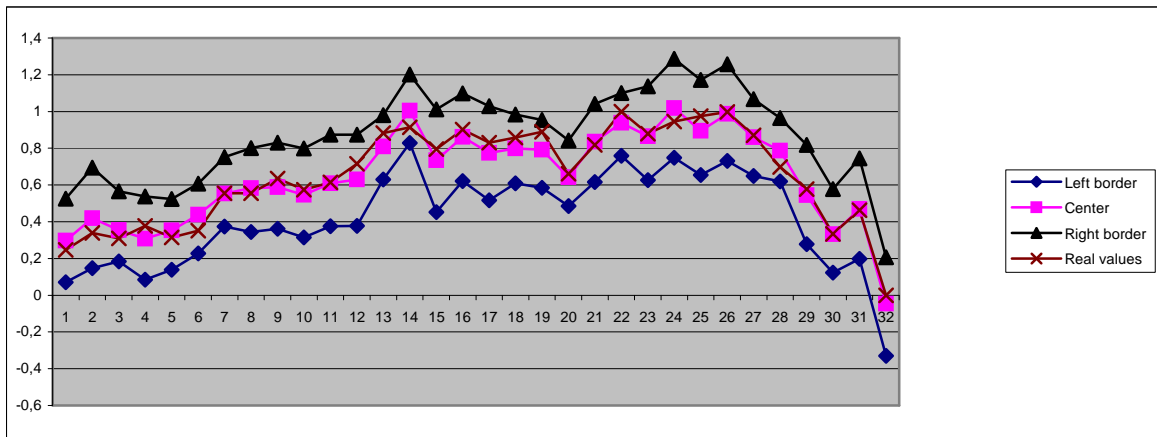
SNGS – ordinary shares of “Сургутнефтегаз” joint-stock company.

Stock prices of other company – “LUKOIL” joint-stock for the same period (03.04.2006 – 18.05.2006) were forecasted.

Sample size – 32 values. Training sample size – 17 values (optimal size of training sample for this experiment).

The following results were obtained:

1. For triangular membership function. For normalized input data: Criterion value: MSE=0.056481



**Fig.3.** Experiment 2 results for triangular MF and normalized values of input variables

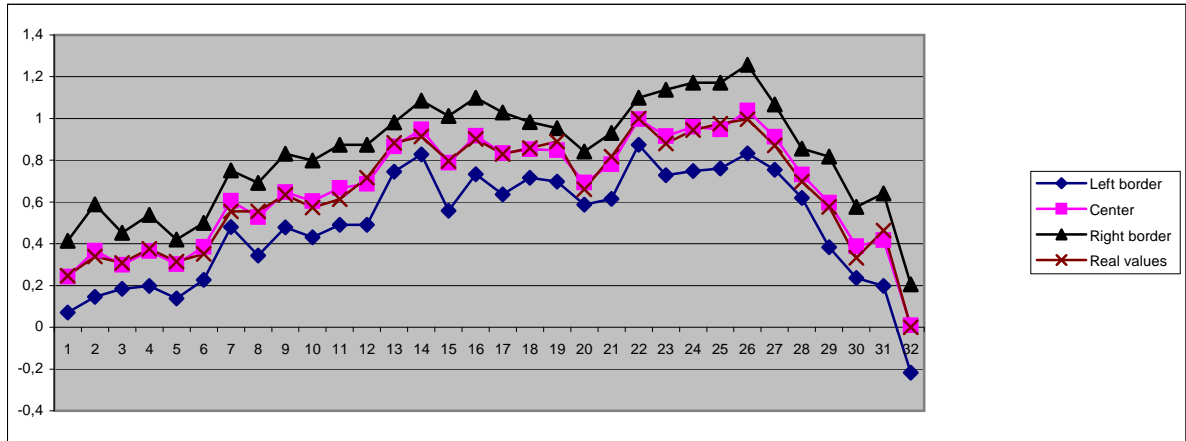
b) for non-normalized input data:

Criterion values: MSE = 0.914998; MAPE = 0.73%

2. For Gaussian membership function (optimal level of  $\alpha=0.9$ )

a) For normalized input data:

Criterion value for this experiment: MSE = 0.030464



**Fig. 4.** Experiment 2 results for Gaussian MF and normalized values of input variables

As we can see from the results of experiment 4, forecasting using triangular and Gaussian membership functions gives good results. The results of experiments with Gaussian MF are better than the results of experiments with triangular MF.

**Tab.3.** For normalized data

	<b>Triangular MF</b>	<b>Gaussian MF</b>
MSE	0.056481	0.030464

**Tab.4.** For non-normalized data

	<b>Triangular MF</b>	<b>Gaussian MF</b>
MSE	0.914998	0.493511
MAPE	0.73%	0.33%

## 5 The comparison of GMDH, FGMDH and FGMDH with fuzzy input data

In the next experiments the comparison of the suggested method FGMDH with fuzzy inputs with known methods: classical GMDH and Fuzzy GMDH was performed

Experiment 3. Forecasting of RTS index (opening price)

Current experiment contains 5 fuzzy input variables, which are the stock prices of leading Russian energetic companies included into the list of RTS index calculation:

Output variable is the value of RTS index (opening price) of the same period (03.04.2006 – 18.05.2006).

Sample size – 32 values.

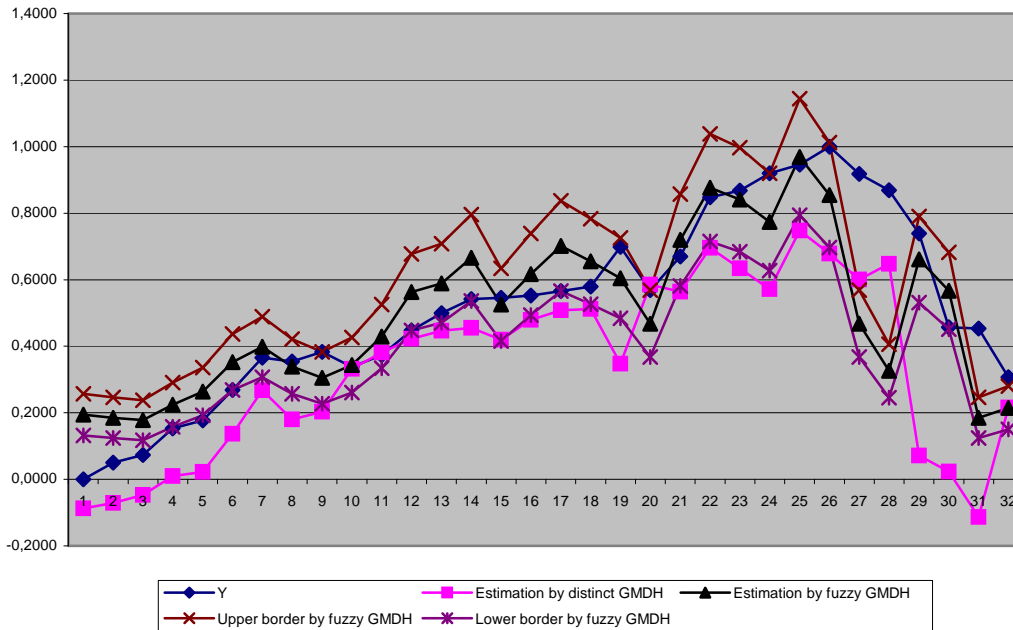
Training sample size – 18 values (optimal size of the training sample for current experiment).

The following results were obtained:

For normalized input when using Gaussian MF in group method of data handling with fuzzy input data: MSE for FGMDH with fuzzy inputs=0, 028013

For normalized values using GMDH and FGMDH:

MSE for GMDH = 0,1129737 MSE for FGMDH = 0,0536556



**Fig.5.** Experiment 3 results using GMDH and FGMDH

As the results of experiment 1 show, fuzzy group method of data handling with fuzzy input data gives more accurate result than FGMDH with triangular membership function or Gaussian membership function. In case of triangular MF FGMDH with fuzzy data gives a little worse than FGMDH with Gaussian MF.

**Tab. 5.** MSE comparison for different methods of experiment 3

	GMDH	FGMDH	FGMDH with fuzzy inputs, Triangular MF	FGMDH with fuzzy inputs, Gaussian MF
MSE	0,1129737	0,0536556	0,055557	0,028013

Experiment 4. RTS-2 index forecasting (opening price)

Sample size – 32 values.

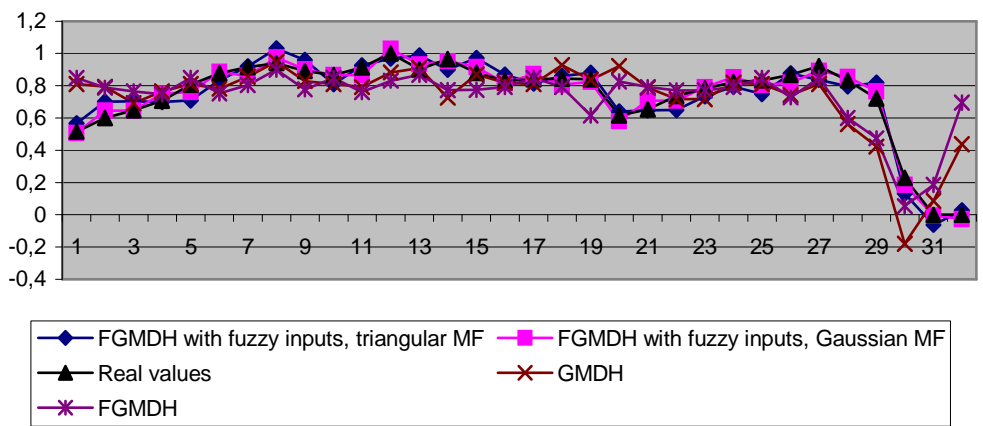
Training sample size – 19 values (optimal size of training sample for current experiment).

The following results were obtained, which are presented in Table 6 and on F ig. 6..

**Tab. 6.** MSE of different methods of experiment 4 comparison

	GMDH	FGMDH	FGMDH with fuzzy inputs, triangular MF	FGMDH with fuzzy inputs, Gaussian MF
MSE	0,051121	0,063035	0,061787	0,033097





**Fig.6.** GMDH, FGMDH (center of estimation), and FGMDH with fuzzy inputs (center of estimation) result comparison

As the results of the experiment 4 show, fuzzy group method of data handling with fuzzy input data gives the better result than GMDH and FGMDH in case of Gaussian membership functions. At the same time in this experiment GMDH gives the better results, than FGMDH and FGMDH with fuzzy input data in the case of triangular membership functions.

## 6 Conclusion

In this article new method of inductive modeling FGMDH with fuzzy inputs was suggested. This method represents the development of fuzzy GMDH when information is fuzzy and given in the form of uncertainty intervals. The mathematical model was constructed and corresponding algorithm was elaborated. The experimental results of application of the suggested method in the forecasting of market index and stock prices are presented and discussed. The comparison of the suggested method with classical GMDH and Fuzzy GMDH were performed and presented. The main advantages of the suggested method are following:

- It operates with fuzzy and uncertain input information and constructs the fuzzy model;
- The constructed model has the minimal possible total width and in this sense is optimal;
- For finding an optimal model we solve the corresponding linear programming problem which is always solvable for this task.

## References

- [1] Zaychenko Yu. “The Fuzzy Group Method of Data Handling and Its Application for Economical Processes Forecasting” - Scientific Inquiry, - Vol. 7, No.1, June, 2006 - p.83-96.
- [2] Zaychenko Yu. “Fuzzy method of inductive modeling in problems of macroeconomic indexes forecasting.” System researches and informational technologies, #3 of 2003, p. 25-45.
- [3] Zaychenko Yu. P., and Zayetz I.O. “The synthesis and adaptation of fuzzy forecasting model based of self-organization method. Science News of NTUU “KPI”, #2 of 2001.
- [4] Zaychenko Yu. P., Zayetz I.O., O.V. Kamotsky, O.V. Pavlyuk. Research of different kinds of membership functions of fuzzy forecasting models parameters in fuzzy group method of data handling. USiM, 2003, #2, p.56-67.
- [5] Zaychenko Yu. P. and Zayetz I.O. Comparative analysis of GMDH algorithms using different method of single-step adaptation of coefficients. The NTUU Herald.