

Time Series Analysis and Forecasting of Oilseeds Production in India – An Application of ARIMA and GMDH-Neural Network

ABSTRACT

Oilseeds have been the backbone of India's agricultural economy since long. Oilseed crops play the second most important role in Indian agricultural economy, next to food grains, in terms of area and production. Oilseeds production in India has increased with time, however, the increasing demand for edible oils necessitated the imports in large quantities, leading to a substantial drain of foreign exchange. The need for addressing this deficit motivated a systematic study of the oilseeds economy to formulate appropriate strategies to bridge the demand-supply gap. In this study, an effort is made to forecast oilseeds production by using Autoregressive Integrated Moving Average (ARIMA) model, which is the most widely used model for forecasting time series. One of the main drawbacks of this model is the presumption of linearity. The Group Method of Data Handling (GMDH) model has also been applied for forecasting the oilseeds production because it contains nonlinear patterns. Both ARIMA and GMDH are mathematical models well-known for time series forecasting. The results obtained by the GMDH are compared with the results of ARIMA model. The comparison of modeling results shows that the GMDH model perform better than the ARIMA model in terms of mean absolute error (MAE), mean absolute percentage error (MAPE), and root mean square error (RMSE). The experimental results of both models indicate that the GMDH model is a powerful tool to handle the time series data and it provides a promising technique in time series forecasting methods.

Keywords: Oilseeds, Forecasting, Autoregressive Integrated Moving Average, Group Method of Data Handling, Mean Absolute Percentage Error, Mean Absolute Error, Root Mean Square Error.

JEL Classification: Q10, C45, C53,

1. INTRODUCTION:

India is one among world's largest producers and consumers of vegetable oils. Oilseeds have been the backbone of India's agricultural economy since long. Indian vegetable oil economy is the fourth largest in the world, next to USA, China, and Brazil. The country's contribution is 7 percent of the global vegetable oils production with 14 per cent share in the area. Oilseed crops play the second most important role in the

33 Indian agricultural economy next to food grains in terms of area and production. The Indian climate is
34 suitable for the cultivation of oilseed crops; therefore, large varieties of oilseeds are cultivated here. The
35 major oilseeds cultivated in our country are Groundnut, Rapeseed and Mustard, Castor seed, Sesame,
36 Niger seed, Linseed, Safflower, Sunflower and Soybean. However, Groundnut, Rapeseed and Mustard,
37 Sesame, Soybean and Sunflower account for a major chunk of the output. At present, more than 27 million
38 hectares of land is under oilseeds cultivation. The area under oilseeds has been increasing over time and
39 the production has registered many fold increase; however, the productivity is still low as compared to the
40 other oilseed producing countries in the world. The low and fluctuating productivity is primarily because
41 cultivation of oilseed crops is mostly done on marginal lands, which are lacking in irrigation and low levels
42 of input is used here. To improve the situation of oilseeds in the country, Government of India has been
43 pursuing several development programs, such as Oilseed Growers Cooperative Project, National Oilseed
44 and Development Project, Technology Mission Oilseeds (TMO) and Integrated Scheme on Oilseeds,
45 Pulses, Oil Palm and Maize (ISOPOM) etc. The concerted efforts of these development programs register
46 significant improvement in annual growth of productivity and area under oilseed crops [2]. The combined
47 efforts have been reflected in oilseeds production. But the growth in the domestic production of oilseeds
48 has not been able to keep pace with the increase in demand in the country. As a result of which, India still
49 imports a significant proportion of its requirement of edible oil. Edible oil is the largest imported (30 percent)
50 commodity in India next only to petroleum products even though India had the world's second largest area
51 under oilseeds [23].

52 In this paper, an effort has been made to forecast oilseeds production for the next five years (2016-
53 17 to 2020-21). The model used for forecasting is an Autoregressive Integrated Moving Average (ARIMA)
54 model. As the model was introduced by Box and Jenkins in 1960, this model is also known as Box-Jenkins
55 model. The model is used for forecasting a single variable. Although it is used across various functional
56 areas, its application is very limited in agriculture, mainly because of unavailability of required data and
57 because agricultural output depends typically on monsoon and other factors [7]. The primary reason behind
58 choosing ARIMA model for forecasting is that it assumes non-zero autocorrelation between the successive
59 values of the time series data [12]. But ARIMA model can only capture linear feature of time series data
60 [18] to deal with non-linearity of time series data, Group Method of Data Handling (GMDH) has also been

61 used in our analysis for forecasting oilseeds production. This model was first used in 1968 by Prof. Alexey
62 G. Ivakhnenko [11].

63 **2. REVIEW OF LITERATURE:**

64 Padhan Purna Chandra (2012) [17], has applied ARIMA model on a 60years' time series data (from 1950
65 to 2010) to forecast annual productivity of selected agricultural product (34 different products). The validity
66 of the model is verified with various model selection criteria such as minimum of AIC (Akaike Information
67 Criteria) and lowest MAPE (Mean Absolute Percentage Error) values. Among the selected crops, tea
68 provides the lowest MAPE values, whereas cardamom provides lowest AIC values.

69 Kumar Manoj and Anand Madhu (2014) [12] forecasted sugarcane production in India by using ARIMA
70 model. The order of the best ARIMA model was found to be (2, 1, 0). They suggested that the forecast
71 results have shown the annual sugarcane production will grow in 2013, then there will be a sharp dip in
72 2014 and in subsequent years 2015 through 2017, it will continuously grow with an average growth rate of
73 approximately 3 percent year-on-year.

74 Arivarasi R and Ganesan Madhavi (2015) [6] have also used the ARIMA Model to forecast the area and
75 production of vegetables in the in the feeder zones (zone 1-Kancheepuram district & zone 2 -Thiruvallur
76 district) of Chennai city. The ARIMA (0, 1, 2) model is suitable for the cultivation area of the zone 2 and
77 ARIMA (2, 0, 1) model is suitable for zone 1. ARIMA (2, 0, 1) model is highly suitable for the vegetable
78 production in both the zones. The model performances are validated by comparing the regression co-
79 efficient values. While the model was used for forecasting for the period 2011-12 to 2014-15, decreasing
80 trend was found in cultivated area and production of vegetables in zone 1. However, in zone 2 increasing
81 trend was found in cultivated areas but decreasing trend was found for the vegetable production. Hence, it
82 can be concluded that if this situation remained the same for a long period, then the further cultivation of
83 vegetable crops will no longer be possible in both the zones.

84 Borkar Prema & Bodade V.M, (2017) [7] have applied the ARIMA model to forecast annual productivity of
85 selected pulse crops. Applying annual data from 1950-51 to 2014-15, forecasted values have been
86 obtained for another 5 years since 2016. The evaluation of forecasting of pulses production has been
87 carried out with Root Mean Squares Percentage Error (RMSPE), Mean Absolute Percentage Error (MAPE)
88 and Relative Mean Absolute Percentage Error (RMAPE).

89 Amanifardet. al., (2008a) [4] presented two meta-models based on the evolved group method of data
90 handling (GMDH) type neural networks for modeling of both pressure drop (ΔP) and Nusselt number (Nu).
91 It was shown that some interesting and important relationships like useful optimal design principles involved
92 in the performance of micro-channels can be discovered by Pareto based multi-objective optimization of the
93 obtained polynomial meta-models representing their heat transfer and flow characteristics. They concluded
94 that, such important optimal principles would not have been obtained without the use of both GMDH type
95 neural network modeling and the Pareto optimization approach.

96 Amanifardet. al., (2008b) [5] presented a quadratic model based upon some experimental results, using
97 evolved GMDH-type neural networks for modeling of the transient evolution of spiky stall cells in an axial
98 compressor. They concluded that the methodology applied in this work could sufficiently derive such
99 complex model of unstable flow of rotating stall based on experimental input–output data. The prediction
100 ability of such polynomial model has also been presented for some unforeseen data.

101 Ahmadiet. al., (2015) [3] proposed an intelligent approach to determine the output power and torque of a
102 Stirling heat engine. The approach employs the GMDH method to develop an accurate predictive tool for
103 determining output power and torque of a Stirling heat engine in manner that is inexpensive, fast and
104 precise. Consequently, based on the output results, the GMDH approach can help energy experts to design
105 Stirling heat engines with high levels of performance, reliability and robustness and with a low degree of
106 uncertainty.

107 Osman Dag and Ceylan Yozgatligil (2016) [16] in their study, the R package GMDH is presented to make
108 short term forecasting through GMDH-type neural network algorithms. The GMDH package has options to
109 use different transfer functions (sigmoid, radial basis, polynomial, and tangent functions) simultaneously or
110 separately. Data on cancer death rate in Pennsylvania from 1930 to 2000 are used to illustrate the features
111 of the GMDH package. The results based on ARIMA models and exponential smoothing methods are
112 included for comparison. GMDH algorithms show the same or even better performance than the other
113 methods.

114 **3. MATERIAL AND METHOD:**

115 The specific objective of the study is to attempt a short-term forecasting of the future oilseeds
116 production by using Autoregressive Integrated Moving Average (ARIMA) forecasting model and also

117 through Group Method of Data Handling (GMDH) - neural network which is an important model of time
118 series data (one the sub-model of Artificial Neuron Networks).

119 3.1 Data

120 The study used data of oilseeds production in India for the last 50 years, i.e., from 1966-67 to 2015-16
121 which have been collected from “Latest APY State Data”, uploaded by the Ministry of Agriculture and
122 Farmers Welfare, Govt. of India.

123 3.2 Autoregressive Integrated Moving Average (ARIMA)

124 The model used in this study is the autoregressive integrated moving average (ARIMA). The ARIMA is an
125 extrapolationⁱ method, which requires historical time series data of underlying variable.

126 The model in specific and general forms may be expressed as follows.

127 Let Y_t is a discrete time series variable which takes different values over a period of time. The
128 corresponding AR (p) model of Y_t series,

129 Which is the generalizations of autoregressive model, can be expressed as:

130 AR (p) Y_t

$$131 Y_t = \mu + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t \dots\dots\dots(1)$$

132 Where, Y_t is the response variables at time t,

133 $Y_{t-1}, Y_{t-2}, \dots, Y_{t-p}$ is the respective variables at different time with lags;

134 μ is the constant mean of the series. $\phi_1, \phi_2, \dots, \phi_p$ are the coefficients; and ε_t is the error factor. ε_t is a
135 white noise process, where $E(\varepsilon_t) = 0$, $\text{var}(\varepsilon_t) = \sigma^2 > 0$, $\text{cov}(\varepsilon_t, \varepsilon_{t-h}) = 0$, $t, h \neq 0$

136 Similarly, the MA (q) model which is again the generalization of moving average model may be specified
137 as:

$$138 \text{MA (q): } Y_t = \mu + \varepsilon_t - \delta_1 \varepsilon_{t-1} - \delta_2 \varepsilon_{t-2} - \dots - \delta_q \varepsilon_{t-q} \dots\dots\dots(2)$$

139 Where, μ is the constant mean of the series;

140 $\delta_1, \delta_2, \dots, \delta_q$ is the coefficients of the estimated error term; ε_t is the error term.

141 By combining both the models, we get the Autoregressive Moving Average or ARMA models, which has
142 general form as:

$$143 Y_t = \mu + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t - \delta_1 \varepsilon_{t-1} - \delta_2 \varepsilon_{t-2} - \dots - \delta_q \varepsilon_{t-q} \dots\dots\dots(3)$$

144 Box and Jenkins argue that a non-stationary series can be transformed either into a stationary or an
145 almost stationary series, if it is differenced an appropriate number of times. Thus, if we have a stochastic
146 process $\{Y_t, t = 0, \pm 1, \pm 2, \dots\}$ which is non-stationary and has a trend, we can find a positive integer 'd' such
147 that the transformed series $W_t = \nabla^d Y_t$ becomes stationary, ∇ being the difference operator, viz. $\nabla Y_t = Y_t - Y_{t-1}$,
148 $\nabla^2 Y_t = Y_t - 2Y_{t-1} + Y_{t-2}$ and so on. After the transformed into a stationary or to an almost stationary series, the
149 model transforms to ARIMA [9].

150 If Y_t is stationary at level or $I(0)$ or at first difference $I(1)$ or at second difference $I(2)$ determines the
151 order of integration. After the stationary of the series was attained, ACF (Auto Correlation Function) and
152 PACF (Partial Auto Correlation Function) of the stationary series are employed to select the order p and q
153 of the ARIMA model. The parameters were estimated using the non-linear least square method as
154 suggested by Box and Jenkins (1976). ε_t is a white noise process, where $E(\varepsilon_t) = 0$, $\text{var}(\varepsilon_t) = \sigma^2 > 0$, cov
155 $(\varepsilon_t, \varepsilon_{t-h}) = 0$, $t, h \neq 0$. Based on the model diagnostic tests and parsimony we obtained the best fitting
156 ARIMA model.

157 The complete procedure of model building and forecasting are fully described by Box and Jenkins
158 1976. In short, they have suggested four basic steps viz., (i) Identification of the model, (ii) Estimation of
159 parameters of the model, (iii) Diagnostic Checking of the model, and (iv) Forecasting. The details of the
160 estimation and forecasting process are discussed below.

161 **Identification:** The first step of applying Box-Jenkins forecasting model is to identify the appropriate order
162 of ARIMA (p, d, q) model. Identification of ARIMA model implies selection of order of AR(p), MA(q) and I(d).
163 The order of d is estimated through $I(1)$ or $I(2)$ process of unit root stationary tests. The model specification
164 and selection of order p and q involved plotting of autocorrelations functions (ACF) and partial
165 autocorrelations functions (PACF) or correlogram of variables at different lag length. If the PACF displays a
166 sharp cutoff while the ACF decays more slowly (i.e., has significant spikes at higher lags), we say that the
167 series displays an AR signature. However, if the ACF displays a sharp cutoff while the PACF decay more
168 slowly, we say that the series displays an MA signature [14]. The autocorrelation functions specify the order
169 of moving average process, q and partial autocorrelations function select the order of autoregressive
170 process p.

171 **Estimation of the model:** ARIMA models are fitted and accuracy of the model has tested based on
172 diagnostics statistics. Once the order of p, d, and q are identified, their statistical significance can be judged

173 by t-distribution. The next step is to specify appropriate regression model and estimate it. ARIMA models
174 are fitted and accuracy of the model was tested based on diagnostics statistics.

175 **Diagnostic checking:** Now a question may arise that how we know whether the identified model is
176 appropriate. One simple way to figure that out is by diagnostic checking the residual term obtained from
177 ARIMA model by applying the same ACF and PACF functions. First obtaining the ACF and PACF of
178 residual term up to certain lags of the estimated ARIMA model, and then checking whether the coefficients
179 are statistically significant or not. The best model was selected based on the following diagnostics,

180 (i) Low Akaike Information Criteria (AIC): AIC is estimated by $AIC = -2\log_e(L) + 2m$, where $m = p + q$
181 and L is the likelihood function.

182 (ii) Low Bayesian Information Criteria (BIC): The Bayesian information criterion is a criterion for model
183 selection among a finite set of models. It is based, in part, on the likelihood function, and it is closely related
184 to Akaike information criterion (AIC). Sometimes, Bayesian Information Criteria (BIC) is also used and
185 estimated by $BIC = -2\log_e(L) + \log_e(N) m$. Where N is number of observation and m is the number of
186 parameters.

187 (iii) The minimum Root Mean Square Error (RMSE) and Mean Absolute Percent Error (MAPE) are
188 used as a measure of accuracy of the models. $RMSE = \sqrt{\frac{\sum_{t=1}^n (X_{Actual,t} - X_{Forecast,t})^2}{n}}$ and MAPE

189 $= \frac{1}{n} \sum_{t=1}^n \left[\frac{|X_{Actual,t} - X_{Forecast,t}|}{X_{Forecast,t}} \right] \times 100$.

190 Where, $X_{Actual,t}$ and $X_{Forecast,t}$ are actual and forecast output at time t ,

191 (iv) These may also be judged by Ljung-Box Q (LBQ) statisticⁱⁱ under null hypothesis that
192 autocorrelation co-efficient up to lag k is equal to zero. LBQ is used to assess assumptions after fitting a
193 time series model (ARIMA), to ensure that the residuals are independent.

194 **Forecasting:** Once the first three steps of ARIMA model are over, then we can obtain the forecasted
195 values by estimating the appropriate model, which is free from problems. The forecasted values are
196 reported for a maximum of 5 years, as long-term forecasting might not be appropriate.

197 *The major drawback of ARIMA model is presumption of linearity, hence, no nonlinear patterns can*
198 *be recognized by ARIMA model. Sometimes, the time series often contain nonlinear components; under*
199 *such condition the ARIMA models are not adequate in modeling and forecasting [23]. To overcome this*
200 *difficulty, GMDH (Group Method of Data Handling) model has been successfully used. To deal with*

201 *uncertainty, linearity or nonlinearity of time series data in a wide range of disciplines GMDH is more*
 202 *effective.*

203 **3.3 Group Method of Data Handling (GMDH):**

204 GMDH is a family of inductive algorithms for computer-based mathematical modeling of multi-
 205 parametric datasets that features fully automatic structural and parametric optimization of models [26].
 206 GMDH is an original method for solving problems of structural and parametric identification under
 207 conditions of uncertainty [13]. It is an important model of time series data which is one sub-model of ANNⁱⁱⁱ
 208 (Artificial Neural Network). The main idea of the GMDH is to build an analytical function in a feed-forward
 209 network based on a quadratic node transfer function whose coefficients are obtained by using a regression
 210 technique. The GMDH is a self-organizing, unidirectional structure with multiple layers, each of which is
 211 composed of several neurons that have a similar structure. Weight is inserted inside each neuron as
 212 definite and constant values based on singular value decomposition method by solving normal equations
 213 [15].

214 The GMDH was introduced as a multivariate analysis method for modeling and identification of
 215 complex systems. In this model, the general connection between inputs and output variables can be
 216 expressed by a complicated polynomial series in the form of the Volterra series, known as the Kolmogorov-
 217 Gabor polynomial [11].

218
$$y_n = a_0 + \sum_{i=1}^n a_i x_i + \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j + \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk} x_i x_j x_k + \dots , \quad \dots\dots (4)$$

219 where{ $x_1, x_2, \dots, x_k, \dots$ } is the vector of input variables and { $a_0, a_i, a_{ij}, a_{ijk}, \dots$ } is the vector of
 220 coefficients of variables in the polynomial, n is the number of inputs, Y is a response variable, x_i and x_j are
 221 the lagged time series to be regressed. However, for most application the quadratic form are called as
 222 partial descriptions (PD) for only two variables is used in the form

223
$$y_n = G(x_i, x_j) = a_0 + a_1 x_i + a_2 x_j + a_3 x_i x_j + a_4 x_i^2 + a_5 x_j^2$$

224 to predict the output. The input variables are set to { $x_1, x_j, x_k, \dots, x_n$ } and output is set to { y }. The aim of
 225 the GMDH algorithm is to find a_i unknown coefficients of Volterra series. The coefficients (weights) a_i , for $i =$
 226 0, 1, 2, 3, 4, 5 are determined using the least square method for each pair of x_i and x_j input variables [10].

227 The GMDH algorithm considers all pairwise combinations of p lagged time series. Therefore, each
 228 combination enters each neuron. Using these two inputs, a model is constructed to estimate the desired

229 output. In other words, two input variables go in a neuron, one result goes out as an output. The structure
230 of the model is specified by the Ivakhnenko polynomial in equation 4 where $n = 2$. This specification
231 requires six coefficients in each model to be estimated [16].

232 The main function of GMDH is based on the forward propagation of signal through nodes of the net
233 similar to the principle used in classical neural nets. Every layer consists of simple nodes, each of which
234 performs its own polynomial transfer function and passes its output to nodes in the next layer. The
235 computation process comprises three basic steps [8]:

236 **Step 1:** Select input variables $\{x_1, x_2, x_k, \dots, x_n\}$ where n is the total number of input. The data are
237 separated into training and testing data sets. The training data set is used to construct a GMDH model and
238 the testing data set is used to evaluate the estimated GMDH model.

239 **Step 2:** Construct L numbers of new variables $Z = \{z_1, z_2, z_3, \dots, z_L\}$ in the training data set for all
240 independent variables and choose a PD of the GMDH. Conventional GMDH has been developed using
241 polynomial, PD of the following form

$$242 \quad z_l = G(x_i, x_j) = a_0 + a_1x_i + a_2x_j + a_3x_ix_j + a_4x_i^2 + a_5x_j^2 \text{ for } l=1,2,3,\dots, L.$$

243 where, $L = n(n-1)/2$

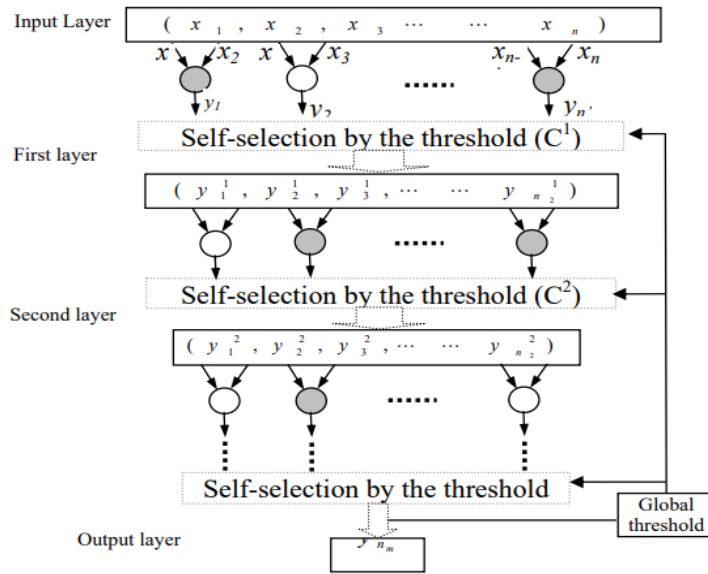
244 Select new variables as input of the next middle layer and truncate the subsequent computation.
245 With the identification of the optimal output of partial polynomials at each layer, the selection of new
246 variables enables the network to quickly converge to the target solution. Once the partial polynomial
247 equations at each unit are selected, the residual error in each layer is further checked to determine whether
248 the set of equations of the model should be further improved within the subsequent computation.

249 **Step 3:** Estimate the coefficient of the PD. The vectors of coefficients of the PDs are determined using the
250 least square method.

251 **Step 4:** Determine new input variables for the next layer. There are several specific selection criteria to
252 identify the input variables for the next layer. In our study, we used two criteria. The first criteria, the single
253 best neuron out of these L neurons, Z' identified according to the value of mean square error (MSE) of
254 testing dataset. In second criteria, eliminate the least effective variables, replace the column of $\{x_1, x_2, x_k, \dots, x_n\}$
255 by those column $\{z_1, z_2, z_3, \dots, z_L\}$ that best estimate the dependent variable y in the testing
256 dataset.

257 **Step 5:** Build the final model and compute the predicted value. The final prediction model can be obtained
 258 with selected variables in each layer and the coefficients of partial polynomials between the connected
 259 layers. Check the stopping criterion. The lowest value of selection criteria using GMDH model at each layer
 260 obtained during this iteration is compared with the smallest value obtained at the previous one.

261 The structure of the GMDH algorithm is illustrated in Figure 1. Those shadowed nodes in Fig 1 that
 262 have significant contribution to the output and are selected to be input in the next layer [25].



263

264 Figure (Fig.) 1: Structure of the GMDH algorithm.
 265

266 The GMDH algorithm is a system of layers in which there exist neurons. The number of neurons in a
 267 layer is defined by the number of input variables. To illustrate, assume that the number of input variables is
 268 equal to p ; since we include all pair-wise combinations of input variables, the number of neurons is equal to
 269 $h = {}^p C_2$ [16].

270 3.3.1 Time series prediction by GMDH

271 A classical method for time series forecasting problem, the number of input nodes of nonlinear
 272 model, such as the GMDH is equal to the number of lagged variables $(y_{t-1}, y_{t-2}, y_{t-3}, \dots, y_{t-p})$, where p is the
 273 number of chosen lagged. The outputs, y_t , the predicted value of a time series defined as

$$274 y_t = f(y_{t-1}, y_{t-2}, y_{t-3}, \dots, y_{t-p}),$$

275 However, there is no suggested systematic way to determine the optimum number of lagged p . The
 276 number of lagged p is chosen either in an adhoc basis or from traditional Box Jenkins methods. The lagged
 277 variables obtained from the Box-Jenkins analysis are the most important variables to be used as input

278 nodes in the input layer of the GMDH model [19]. In our study, a time series model is considered as
 279 nonlinear function of several past observations and random errors as follows:

280
$$y_t = f[(y_{t-1}, y_{t-2}, y_{t-3}, \dots, y_{t-p}), (a_{t-1}, a_{t-2}, a_{t-3}, \dots, a_{t-p})]$$

281 where f is a nonlinear function determined by the GMDH.

282 **3.3.2 Data structure of GMDH**

283 An illustration of time series data structure in GMDH algorithms is presented in Table 1. Since we
 284 have a time series data set with t time points and p inputs. We construct the model for the data with time
 285 lags, the number of observations presented under the subject column in the table is equal to $t-p$; and the
 286 number of inputs i.e, lagged time series, is p . In this table, the variable called y is put in the models as a
 287 response variable, and the rest of the variables are taken into models as lagged time series x_i , where $i =$
 288 $1, 2, \dots, p$. The notations in Table 1 are followed throughout this paper.

289 Table 1: An illustration of time series data structure in GMDH algorithms

Subjects	Y	X ₁	X ₂	X ₃		X _p
1	y _t	y _{t-1}	y _{t-2}	y _{t-3}		y _{t-p}
2	y _{t-1}	y _{t-2}	y _{t-3}	y _{t-4}		y _{t-p-1}
3	y _{t-2}	y _{t-3}	y _{t-4}	y _{t-5}		y _{t-p-2}
...						
t-p	y _{p+1}	y _p	y _{p-1}	y _{p-2}		y ₁

294
 295 A better model which explains the relation between response and lagged time series is captured via
 296 transfer functions.

297 **4. RESULTS AND DISCUSSION:**

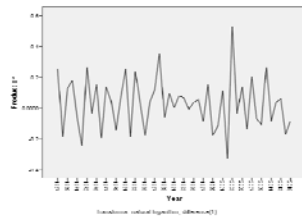
298 **4.1 ARIMA Model:**

299 The preliminary understating about the nature of data showed that there is no consistency in the production
 300 of oilseeds over the time period (Fig. 2). The variable shows increasing trend.

301



302 Fig 2. Time series plots of oilseeds



303 Fig 3. Plots of 1st difference

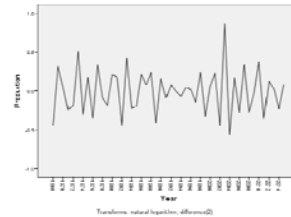


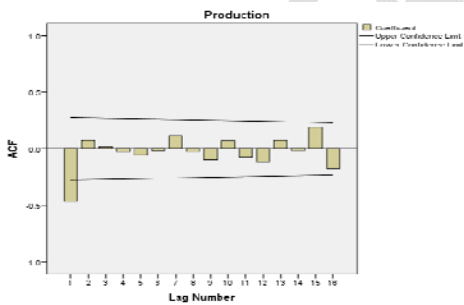
Fig 4. Plots of 2nd difference

304 **Identification:**

305 Identification of the model was concerned with deciding the appropriate values of (p, d, q). Auto regressive
306 and moving average terms are identified based on ACF and PACF values. The ACF helps in choosing the
307 appropriate values for ordering of moving average terms (MA) and PACF for those autoregressive terms
308 (AR).

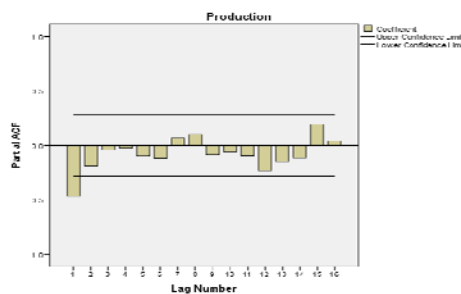
309 ARIMA model is generally applied for stationary time series data. Stationary vs. non-stationary can
310 check through correlogram or autocorrelation functions. If autocorrelation coefficients don't die out slowly,
311 then the series is probably non-stationary. The general procedure to convert a non-stationary series to a
312 stationary series is through first difference or second difference. In general, most of the variables are I (1)
313 i.e., first difference or I (2) i.e., second difference, thereby ARMA model is applied at I(1) or maybe I(2).
314 Both the first differences and the second difference time series data of production are given in Fig 3 and Fig
315 4, respectively. Comparing the figures, it has been observed that in the first figure, difference magnitude of
316 auto correlation is lower than that in the second difference data. Hence, we considered I(1) for making the
317 series stationary.

318 ACF and PACF of production of oilseeds are presented in Figs 5 and 6. Based on these figures, the
319 initial ARIMA model has been developed. It can be seen from Figs 5 and 6 that there is a slow decay in the
320 PACF, but it also has a cut-off only at lag1, suggesting AR (1). The ACF also has one significant spikes at
321 lag1. This pattern is typical to an MA process of orders 1.



322

323 Fig. 5: ACF of 1st differenced series by lag



324 Fig.6: PACF of 1st differenced series by lag

325 **Estimation of the model:** Once the orders of p, d, and q are identified, the next step is to specify
326 appropriate ARIMA model and estimate it. With the help of SPSS software, various orders of ARIMA model
327 has been estimated. After the identification process has completed, the number of possible models are
328 identified. According to identification process, the model has been identified as ARIMA (1, 1, 1). However,

328 the coefficient of AR (1) is not statistically significant. Hence in addition to ARIMA (1, 1, 1), the study also
 329 attempts to estimate ARIMA (1, 1, 0) and ARIMA (0, 1, 1) model. The results of ARIMA (1, 1, and 1),
 330 ARIMA (1, 1, 0) and ARIMA (0, 1, and 1) are summarized in Table 2.

331 Table2: Coefficients of estimated values of fitted ARIMA models

Sl.No	Variable	Model	Constant	AR(1)	MA(1)
1	Production	ARIMA(0,1,1)	0.028	-	0.596
		SE	0.009		0.126
		t- value	3.216		4.718
2	Production	ARIMA(1,1,0)	0.027	-0.479	-
		SE	0.015	0.128	-
		t- value	1.804	-3.743	-
3	Production	ARIMA(1,1,1)	0.028	-0.093	0.519
		SE	0.010	0.261	0.234
		t- value	2.901	-0.357	2.214

332 We proceeded to further statistically analyze these two possible models. The best model is selected based
 333 on the diagnostics checking.

334 **Diagnostic checking:** Now a question may arise that how we know whether the identified model is
 335 appropriate. After an estimation of the parameters, we test the adequacy of the model based on Box-Pierce
 336 (Q) and Ljung-Box (LB) statistics. The statistics is calculated from the ACF of residual term up to 16 lags of
 337 the estimated ARIMA model. We also check the statistical significance of the parameters. An adequate
 338 model does not always generate good forecasts. Further, we select the model having low Bayesian
 339 Information Criteria (BIC), lowest root means square error (RMSE), lowest mean absolute percent error
 340 (MAPE), and highest stationary R-Square and R-Square.

342 Comparing these three models, the ARIMA (0,1,1) model is found to be the best for oilseeds
 343 production. Only in this model, the estimated coefficient is statistically significant. LB and Q statistics of the
 344 model is also statistically significant. At the same time, RMSE, MAPE, MAE and BIC of ARIMA (0,1,1) have
 345 shown a value lower than that of ARIMA(1,1,0) and ARIMA(1,1,1) models. The summary of the estimates of
 346 ARIMA (0,1,1) models is given in Table 3.

347 Table3: RMSE, MAPE, BIC values and Q statistics of fitted ARIMA models

Sl. No	Variable	Model	RMSE	MAPE	MAE	BIC	Stationary R ²	R ²	Ljung Box Q Statistics	Df
1	Production	ARIMA (0,1,1)	2811.85	11.72	2008.66	16.04	0.26	0.87	13.03	17

348

349 Based on the parameter estimates in the Table 2 and model statistics presented in the table 3, the study
 350 chose the ARIMA (0,1,1) as the best model for the oilseeds production in the India. The model is thus given
 351 as:

$$\nabla Y_t = (1 - 0.596B)e_t$$

353 This model is a special case of ARIMA model, which is called an Integrated Moving Average Model.

354 **Forecasting:** Once the identification, estimation of the model and diagnostic checking steps of ARIMA
 355 model is over, then we can obtain forecasted values by estimating the appropriate model, which is free
 356 from problems. The forecasted values obtained from ARIMA model are reported in Table 4. The forecasted
 357 values are reported for a maximum 5 years as long-term forecasting might not be appropriate.

358 Table 4: Forecast values with ARIMA model

Model ARIMA (0,1,1)	Variable	Value	Years				
	Production (000 tonnes)		2016-17	2017-18	2018-19	2019-20	2020-21
	Forecast		30062	30987	31939	32922	33934
	Lower		22069	22181	22330	22510	22715
	Upper		40062	42195	44372	46601	48887

359
 360 In our study, ARIMA (0,1,1) is the best model for oilseeds production. Based on this model,
 361 forecasted values of oilseeds production will be 30062 thousand tonnes, 30987 thousand tonnes, 31939
 362 thousand tonnes, 32922 thousand tonnes and 33934 thousand tonnes during 2016-17, 2017-18, 2018-19,
 363 2019-20 and 2020-21, respectively. It is clear that oilseeds production will be slightly increasing over time.

364 4.2 GMDH Model

365 In this section we analyze the short-term forecasting results of oilseeds production through GMDH-
 366 type neural network algorithms^{iv} by using GMDH Shell software. GMDH-neural network selects the model
 367 of optimal complexity and such a selection depends on the form of external criterion realization. *K*-fold
 368 cross validation is one of such criteria. In our study, we used this *k* fold validation method. In this validation,
 369 original sample was randomly partitioned into *k* subsamples. A single subsample was taken as the
 370 validation data for testing model, and the other *k* – 1 sub-samples were used as training data. The cross-
 371 validation process was repeated *k* times using each of the *k* subsamples exactly once. The value of *k*
 372 obtained from the *K* folds can produce a single estimation. The advantage of this method over repeated
 373 random sub-sampling is that all observations are used for both training and validation, and each
 374 observation is used for validation exactly once. The experiment was carried out using RMSE validation

375 criterion [13]. Therefore, the optimal time series forecasting model was selected by minimum value of
 376 RMSE, calculated for the testing sample. This validation criterion defines model selection criterion for both
 377 the core algorithm^v and variables ranking^{vi} (Solver, GMDH shell documents). In our time series analysis
 378 under GMDH-neural network model, based on k- cross validation criterion, our forecasting model is an
 379 optimal with k=2.

380 In this model the variables ranking are selected by error. Variables are dropped after rank 600. The
 381 neural-type method used as a core of algorithm in our model. The summary of the results of our model
 382 depict that model complexity (it informs about the number of coefficients in the model and the number of
 383 layers) is 2 of 6. It means that the model has two layers and six coefficients or weight of polynomial.
 384 Maximum number of layer selected in our model are 33 with initial layer^{vii} width 1. The Criterion value of this
 385 model is 0.060354 which informs about the value of validation criterion configured in the Solver module^{viii}.
 386 Top-ranked model has the smallest criterion value. Our model's low criterion value indicates that the model
 387 is suitable for this data.

388 The formula of suggested forecasting model under GMDH –neural network is given by

$$389 \quad Y_t = 6677.04 + 1.036 Y_{t-15} + 0.005 Y_{t-23}$$

390 Accuracy of model shows different accuracy metrics for the model selected in the model browser.
 391 Model contains accuracy measures calculated for observations used to create the model. Error measure is
 392 used to choose a metric for calculation of the mean and the root mean errors. Available metrics are the
 393 absolute (MAE and RMSE), which outputs mean error values “as is” and the target percentage (MAPE),
 394 where for each model value we calculate percentage deviation from the actual value and then the
 395 percentage deviations are averaged [20]. The model statistics of GMDH - neural network are presented in

396 Table 5.

397 Table 5: RMSE, MAPE, MAE values of fitted GMDH neural network models

Sl. No	Variable	Model	RMSE	MAPE	MAE	R ²
1	Production	GMDH	1833.72	5.275	1473.56	0.99

398
 399 Calculation of magnitude of predicted variable involves only the observations that are used for
 400 training and testing. The forecasting values are presented in Table 6. In our study, GMDH neural networks
 401 model forecasting oilseeds production will be 28176 thousand tonnes, 22145 thousand tonnes, 32864

402 thousand tonnes, 32008 thousand tonnes and 35751 thousand tonnes in 2016-17, 2017-18, 2018-19,
 403 2019-20, 2020-21, respectively.

Model	Variable	Value	Years				
			2016-17	2017-18	2018-19	2019-20	2020-21
GMDH	Production (000 tonnes)						
		Forecast	28176	22145	32864	32008	35751
		Lower	24508	18477	29196	28340	32083
		Upper	31844	25813	36532	35676	39419

404 Table 6: Forecast values with GMDH neural network model

405

406 **5. COMPARISON BETWEEN ARIMA and GMDH-Neural Network Model:**

407 Now the question that arises is which model is better and appropriate for forecasting the oilseeds
 408 production. To find the solution, we compare the model statistics of ARIMA and GMDH-neural network in
 409 terms of RMSE, MAE and MAPE. Model with lower values of RMSE, MAE and RMPE as compare to the
 410 other model, is better. The model statistics of GMDH-neural network and ARIMA both are presented in
 411 Table 7. The table indicates that GMDH-neural network is better model than ARIMA in all respect.

Variable	Model	RMSE	MAE	MAPE	R ²
Production	ARIMA (0,1,1)	2811.85	2008.66	11.71	0.88
	GMDH	1833.72	1473.89	5.275	0.99

412 Table 7: RMSE, MAPE, MAE statistics of fitted ARIMA models and GMDH

413

414 To verify our results, we considered similar research works such as Srinivasan, 2008, [22] and Xu
 415 et.al. 2012, [24]. Srinivasan (2008) used a GMDH-type neural network and traditional time series models to
 416 forecast predicted energy demand. It was shown that a GMDH-type neural network was superior in
 417 forecasting energy demand compared to traditional time series models with respect to MAPE. In another
 418 study, Xu et.al. (2012) applied a GMDH algorithm and ARIMA models to forecast the daily power load.
 419 According to their results, GMDH-based results were superior to the results of ARIMA models in terms of
 420 MAPE for forecasting performance.

421 *Since the above analysis lends support to the choice of GMDH-neural network over ARIMA type*
 422 *modeling we would propose the values obtained from GMDH-neural network as the forecast outcome.*

423 **6. FINAL FORECASTING:**

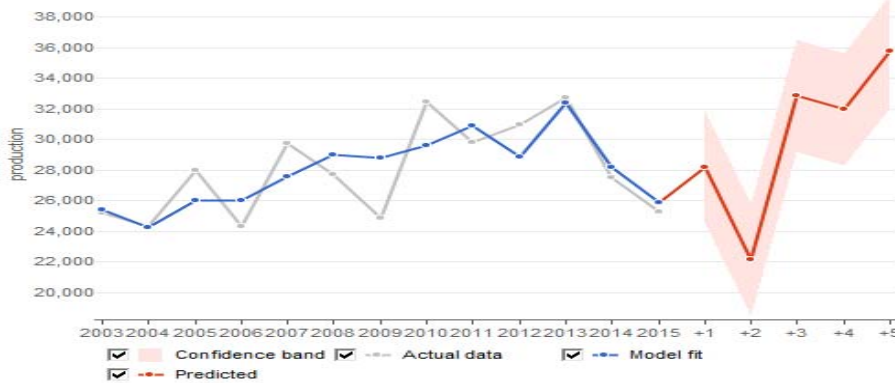
424 The outcome of GMDH model are presented precisely in Table8,

Table 8: Forecast values with GMDH- neural network model

Variable	Model	Predicted				
Production (000 tonnes)	GMDH	2016-17	2017-18	2018-19	2019-20	2020-21
		28176	22145	32864	32008	35751

426

427 The graphical presentation of forecasted value of oilseeds production under GMDH- neural network
 428 is depicted in Figure 7. In the diagram below, time is measured along the horizontal axis and the vertical
 429 axis measures level of production. Actual value is presented by black line and fitted value is shown in blue.
 430 The red line indicates the forecast value of oilseeds production whereas confidence band has been
 431 represented by shaded area.



432

433 Fig.7:Actual value, fitted value and forecast value and confidence band in GMDH model

434

435 From both from Table 8 and Fig 7, it is clear that the expected oilseeds production will increase in India in
 436 near future which will reduce the gap between demand and supply of oilseeds. Alternatively, it can be said
 437 that this rise in supply will be helpful in meet in the growing domestic demand for edible oil due to increase
 438 in population. As a result, the dependence on imported edible oil will reduce substantially, preventing the
 439 huge expenditure of already scarce foreign exchange.

440 7. CONCLUSION

441 ARIMA models are not always adequate for the time series that contains non-linear structures. In this
 442 context, a nonlinear GMDH can be an effective way to improve forecasting performance. Based on the
 443 results obtained in our study, one can infer that application of GMDH techniques in modeling and
 444 forecasting of time series can increase the forecasting accuracy. More specifically, the GMDH-neural

445 network model performed better for forecasting oilseed production of India as compared to ARIMA models.
446 The results of forecasting in GMDH-neural network methods reveals that India's oilseeds production will be
447 28176 thousand tonnes in 2016-17. It will decline to 22145 thousand tonnes in 2017-18 and thereafter it will
448 increase to 32864 thousand tonnes in 2018-19, 32008 thousand tonnes in 2019-20 and 35751 thousand
449 tonnes in 2020-21. This production of oilseeds may not be adequate to make our country self-sufficient. This
450 is because the demand for oilseeds grows faster along with rising population. Still the gap between demand
451 and supply of oilseeds will reduce, resulting in reduced dependence on imported of edible oil and drain of
452 foreign exchange from India will be under control.

453 **8. REFERENCES:**

- 454 1. Adhikari Ratnadip, Agarwal RK (2013) *An Introductory Study on Time Series Modeling and Forecasting*,
455 Lambert Academic Publishing, <https://arxiv.org/abs/1302.6613>.
- 456 2. Agropedia (2011) Oilseeds Scenario of India, submitted by YadavKiran
457 <http://agropedia.iitk.ac.in/content/oilseeds-scenario-india>
- 458 3. Ahmadi MH, Ahmadi MA, Mehrpooya M, Rosen M A (2015) Using GMDH Neural Networks to Model the
459 Power and Torque of a Stirling Engine. *Sustainability*, 7: 2243-2255.
- 460 4. Amanifard N, Nariman-Zadeh N, Borji M, Khalkhali A, Habibdoust A (2008a) Modelling and Pareto
461 optimization of heat transfer and flow coefficient in microchannels using GMDH type neural networks.
462 *Energy Conversion and Management*. 49 (2):311-325.
- 463 5. Amanifard N, Nariman-Zadeh, N, Farahani, M H, Khalkhali, A (2008b) Modelling of multiple short-length-
464 scale stall cells in an axial compressor using evolved GMDH neural networks, *Energy Conversion and*
465 *Management*. 49(10): 2588–2594.
- 466 6. Arivarasi R and Madhavi Ganesan (2015) Time Series Analysis of Vegetable Production and Production
467 and Forecating using ARIMA Model, *Asian Journal of Science and Technology*, 6(10):1844-
468 1848, <http://www.journalajst.com/sites/default/files/2456.pdf>
- 469 7. Borkar Prema & Bodade VM (2017) Application of ARIMA Model for Forecasting Pulses Productivity in
470 India, *Journal of Agricultural Engineering and Food Technology*, 4(1).
- 471 8. Chang FJ and Hwang YY (1999) A self-organization algorithm for real-time flood forecast, *Hydrological*
472 *processes*, 13: 123-138.
- 473 9. Datta LK (2009) Study on Analysis of Financial Data Using Multivariate and Time Series Technique,
474 Ph.D thesis, Department of Statistics, Saurashtra University, Rajkot, Gujarat.
- 475 10. Iba H, De Garish. H & Sato T(1995) A numerical approach to genetic programming for system
476 identification, *Evolutionary Computation*, 3 (4): 417-452.
- 477 11. Ivakhnenko AG (1971) Polynomial theory of complex system, *IEEE Trans. Syst., Man Cybern.* SMCI-1,
478 1: 364-378.
- 479 12. Kumar Manoj & Anand Madhu (2014) An Application of Time Series ARIMA Forecasting Model for
480 Predicting Sugarcane Production In India, *Faculty of Economic Sciences*, 9 (1): 81-94.
- 481 13. Latysh Elena, Koshulko Oleksiy (2012) Testing k-value in k-fold Cross Validation of Forecasting Models
482 for Time Series Analysis of G-spreads of Top-quality RUB Bonds, *The 5th International Workshop on*
483 *Inductive Modelling IWIM'2012*.
- 484 14. Nasiru MO and Olanrewaju SO.(2015) Forecasting Airline Fatalities in the World Using a Univariate
485 Time Series Model, *International Journal of Statistics and Applications*, 5(5): 223-230,
486 <http://article.sapub.org/10.5923.j.statistics.20150505.06.html>

- 487 15. Nariman-Zadeh N, Darvizeh A, Ahmad-Zadeh R (2003) Hybrid Genetic Design of GMDH-type Neural
488 Networks Using Singular Value Decomposition for Modelling and Prediction of the Explosive Cutting
489 Process, *Journal of Engineering Manufacture*, Proceedings of the IMechE Part` B, 217: 779 – 790.
- 490 16. Osman Dag and CeylanYozgatligil (2016) GMDH approach can help energy experts to design Stirling
491 heat engines with high levels of performance, reliability and robustness and with a low degree of
492 uncertainty. *The R Journal*, 8(1): August, 2016.
- 493 17. Padhan Purna Chandra (2012) Application of ARIMA Model for Forecasting Agricultural Productivity in
494 India India. *Journal of Agriculture and Social Science*, 8(2): 50–56.
- 495 18. Samsudin R, Saad P, Shabri A(2011) A hybrid GMDH and least squares support vector machine in time
496 series forecasting, *Neural Network World*, 3(11): 251-268
- 497 19. Shabri A and Samsudin R (2014) A Hybrid GMDH and Box-Jenkins Models in Time Series
498 Forecasting, *Applied Mathematical Sciences*, 8(62): 3051 -3062.
- 499 20. Simulation results, GMDH shell documents,
500 http://d.gmdhshell.com/bf3/doku.php?id=processing_results
- 501 21. Solver, GMDH shell documents. <http://d.gmdhshell.com/docs/solver>
- 502 22. Srinivasan D (2008) Energy demand prediction using GMDH networks. *Neurocomputing*, 72(1): 625–
503 629,
- 504 23. Rathod Santosha, Singh KN, Patil SG, Naik Ravindrakumar H, Ray Mrinmoy, and MeenaVikram Singh
505 (2018) Modeling and forecasting of oilseed production of India through artificial intelligence techniques,
506 *Indian Journal of Agricultural Sciences*, 88(1): 22-27.
- 507 24. Xu H, Dong Y, Wu J, and Zhao W (2012) Application of GMDH to short-term load forecasting. *In*
508 *Advances in Intelligent Systems*, pages. 27–32. Springer-Verlag.
- 509 25. Wang X, Li L, Lockington D, Pullar D and Jeng DS(2005) Organizing Polynomial Neural Network for
510 Modelling Complex Hydrological Processes, Research Report, R861, Department of Civil Engineering,
511 The University of Sydney.
- 512 26. WikiVisually - Group method of data handling(2018)
513 https://wikivisually.com/wiki/Group_method_of_data_handling

514

515 List of Abbreviations

516

517	AR:	Autoregressive
518	MA:	Moving Average
519	ARMA:	Autoregressive Moving Average
520	ARIMA:	Autoregressive Integrated Moving Average
521	ACF:	Autocorrelations Functions
522	PACF:	Partial Autocorrelations Functions
523	GMDH:	Group Method Data Handling
524	ANN:	Artificial Neural Network
525	RMSE:	Root Mean Square Error
526	MAE:	Mean Absolute Error
527	MAPE:	Mean Absolute Percentage Error
528	AIC:	Akaike Information Criteria
529	BIC:	Bayesian Information Criteria
530	Q Statistics:	Box-Pierce
531	LB:	Ljung-Box
532	TMO:	Technology Mission Oilseeds
533	ISOPOM:	Integrated Scheme on Oilseeds, Pulses, Oil Palm and Maize
534	PD :	Partial Descriptions

535

536

537

End note

ⁱ**Extrapolation** techniques make forecasts using only the past data.

ⁱⁱThe Ljung-Box Q statistic to test whether a series of observations over time are random and independent. If observations are not independent, one observation can be correlated with a different observation k time units later, a relationship called autocorrelation. Autocorrelation can decrease the accuracy of a time-based predictive model, such as time series plot, and lead to misinterpretation of the data.

ⁱⁱⁱ**ANN**: The basic objective of ANNs was to construct a model for mimicking the intelligence of human brain into machine. Similar to the work of a human brain, ANNs try to recognize regularities and patterns in the input data, learn from experience and then provide generalized results based on their known previous knowledge. Although the development of ANNs was mainly biologically motivated, but afterwards they have been applied in many different areas, especially for forecasting and classification purposes [1].

^{iv}**GMDH-type neural network** algorithms are modeling techniques which learn the relations among the variables. In the perspective of time series, the algorithm learns the relationship among the lags. After learning the relations, it automatically selects the way to follow in algorithm.

^v**Core algorithms** perform generation and selection of model structures. Then model coefficients are fitted using the least squares method.

^{vi}**Variables ranking** turns on preliminary ranking and reduction of variables. Ranking of variables according to their individual ability to predict testing data.

^{vii}**Initial layer** width means how many neurons are added to the set of inputs at each new layer.

^{viii}Solver [21] module produces predictive models for target variables.