

# Research on Data Prediction Based on Data Mining Combined Model

Dong Han and Chunhua Wang

*College of information engineering, Huanghuai University  
Corresponding E-mail: Flyequn@163.com*

## **Abstract**

*Given modelling method of combination forecast model based on the Data mining. Based on Data mining's combination forecast model's modelling method can reduce the serious influence that the variable substitution brings and has fully used useful information in the primary data. It obviously improved the accuracy of the prediction model. With data mining technology as the entry point and in combination with the analysis on data prediction characteristics, data prediction implementation access based on data mining combination model has been explored in the thesis. Research on variable substitution to non-linear regression forecast model precision's influence, and seek the modelling method that can improve the forecast precision. Based on the Data mining, the transform in space and the weighted processing combined method, make full use of information that the primary data provide.*

**Keywords:** *Data mining technology; data prediction; Non-linear model; Space transformation*

## **1. Introduction**

Data mining technology is generally to process, observe and analyze large-scale data by use of statistical analysis software, aiming to find the hidden rule or effective information existing in data. The application core of the data is large-scale data prediction which completely depends on data source. With the advent of big data era, among the top 10 technologies that have been internationally recognized with the most powerful influence and development potential in the future, data mining technologies ranks the third as the further expansion and deepening of statistics and data base technology [1]. The rapid development of computer technology and the wide application of data base, people have seen a great progress in their data collection and storage ability; especially with the generalization of the internet in recent years, various data and information have begun to explode. In face of the mass and complex data, people tend to feel helpless and confused, hard to effectively analysis and deal with them, with some managers even making decisions purely by intuition, instead of making analysis and judgment based on the historical data, which is undoubtedly a loss [2]. It is under such circumstances that data mining technology comes out in response to the demands of the times.

## **2. Overview of Data Mining Technology**

Data mining refers to the process of mining the potentially useful and regularly existing information and knowledge that hides in the large magnitude of practical data that people are not clear about in the first place [3]. While judged from the perspective of finance, data mining is a whole new financial information processing technology, mainly characterized by the analysis and exploration for the mass and complex data in financial data base to discover the hidden key rule that could help individual or institute with investment decision, and thus help people to make accurate judgment or decision.

### **3. The Characteristics of Data Prediction**

Big data requires a new process model to develop mass, high-growth-rate and diversified information assets with a stronger decision-making power, perception and procedure optimization ability [4]. Data prediction and analysis relies on data source, therefore, the characteristics of the data source also decides the characteristics of the big data prediction.

#### **3.1. Real Samples Rather than Sampling**

Cloud computation and data base can make it quite easy to acquire a sample data big enough and the entire data. Google can provide Google flue trend just because it has covered over 70% of American search market, no longer necessary to investigate the data by sampling but just to mine and analyze big data recording base. However the big data also have their flaws and the systematic deviation remains possible, for the real sample is not equal to the entire sample [5]. Therefore, there exists an issue on the threshold value of data scale. In case the number of data is less than the value, the questions can never be solves, while in case it reaches the value, there will come the solutions to the originally insurmountable questions [6]; even though the number exceeds the value, there won't be any more help to solve the questions.

#### **3.2. Efficient Rather than Accurate**

The traditional sampling requires a high accuracy in specific operation, because the slightest error may cause a grave consequence. Just imagine randomly selecting 1000 from the entire sample of the global 100mn people, in case of any errors in the operation on the 1000, there will produce a huge deviation among the 100mn[7]. While in case of full sample, the deviation will always be the same without the risk of magnifying. Google artificial intelligence expert Novarg ever wrote that the simple algorithm based on big data could be more effective than the complex algorithm based on small data. Data analysis is not simply for the sake of data analysis, instead it has many decision purposes and thus the timeliness also matters. Accurate calculation is conducted at the expense of time consumption, and in the era of small data, perusing accuracy is the forced method to avoid deviation expansion. In this big data era, rapidly acquiring a rough outline and development vein is far more important than the strict requirement for accuracy [8].

#### **3.3. Relevancy Rather than Causality**

Different from traditional logic reasoning, big data research requires a series of analysis & conclusion operations like statistic search, comparison, clustering and classification, and therefore, it inherits some characteristics from statistic science. Statistics pays attention to data relevancy or correlation. The so-called correlation means some rules existing between the values of two or more variables. "The analysis on correlation" aims to discover the correlation networking hidden in the data set which can generally be represented by support degree, reliability and degree of interest. The recommended algorithm of Amazon is quite well-known for telling what users might like by their consumption records which can be the historical records of the users or someone else. However, it cannot tell the reasons why they like. Just understanding the correlation is far from introducing the recommended algorithm to Amazon logistics and warehouse layout, or else some extra loss might be brought about. That is also the boundary line between relevancy and causality prediction.

### **4. Data Prediction Access Based on Data Mining Combination Model**

In health statistical research, it is necessary to discover a hidden rule from a lot of data, and it is best to present it in a mathematical model. Obviously the vast majority of these mathematical models are nonlinear. Because nonlinear regression models are more complex than linear regression models, it is not easy to calculate the regression parameters. On the premise of meeting the needs of the actual situation, sometimes non-linear models are approximated to regression models [9] to solve practical problems. By approximating a regression model by a nonlinear model, generally first substitute variables of the non-linear function and convert it into a linear model; afterwards, implement a linear regression, and then revert to the nonlinear model. Wherein, the calculation process of converting from a nonlinear model to a linear model, and then from a linear model to a nonlinear model, some interference information is added while the original information is lost, which will sometimes seriously affect the prediction accuracy of the nonlinear regression model obtained, whereas the combination forecasting model based on data mining methods can overcome this defect.

#### 4.1. Principles and Methods

##### 4.1.1. Form of the Nonlinear Mathematical Model

The nonlinear mathematical model can be expressed as follows:

$$y = f(x_1, x_2, \dots, x_m, a_1, a_2, \dots, a_l) + e \quad (1)$$

where  $e \sim N(0, \sigma^2)$ . The independent variable  $x = (x_1, x_2, \dots, x_m) \in R^m$  in Model (1) is a point in a m-dimensional space; parameter  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l) \in R^l$  is a point in an l-dimensional space; the dependent variable  $y \in R^1$  is a point in a one-dimensional space. Multivariate function  $f(x_1, x_2, \dots, x_m; \alpha_1, \alpha_2, \dots, \alpha_l)$  with a parameter  $\alpha$  is the nonlinear function for the independent variable  $x = (x_1, x_2, \dots, x_m) \in R^m$ . For a nonlinear regression analysis, the first problem to be solved is how to obtain the best estimate of the l-dimensional parameter  $\alpha$ .

##### 4.1.2. Method to Approximate Common Nonlinear Mathematical Models to Regression Model Parameters

In health statistics, the most widely used non-linear mathematical models include exponential parameter, power parameter, S-growth parameter, special power parameter and exponential parameter  $y = [g(x)]^\alpha \exp[\beta h(x)]$ . For the nonlinear mathematical model obtained from the experiment, assume its data set in the m+1 -dimensional space

$$X - Y \text{ is } \{((x_1, x_2, \dots, x_m)_i, y_i) | i = 1, 2, \dots, n\}.$$

As per the experimental data of the nonlinear mathematical model that has not been fitted in the m+1-dimensional space  $X - Y$ , the theoretical prediction data set corresponding to it is  $\{((x_1, x_2, \dots, x_m)_i, y_i) | i = 1, 2, \dots, n\}$ . It is difficult to directly find out the theoretical prediction value in the m+1-dimensional space  $X - Y$  due to the nonlinearity of the model, so commonly alternative methods are used to make variable substitution to the nonlinear function: convert into a linear model, carry out linear regression and then revert to a nonlinear model.

Variable substitution  $z = F(y)$  can be used to convert the data in the m+1-dimensional space  $X - Y$  into the data in the m+1-dimensional space  $X - Z$ . Afterwards, the image collection of the data set in the new m+1-dimensional space X-Z is

$\{((x_1, x_2, \dots, x_m)_i, z_i) \mid i = 1, 2, \dots, n\} = \{((x_1, x_2, \dots, x_m)_i, F(y_i)) \mid i = 1, 2, \dots, n\}$   
 . As thus, its theoretical prediction data set in the new  $m+1$ -dimensional space  $X-Z$  is  
 $\{((x_1, x_2, \dots, x_m)_i, z_i) \mid i = 1, 2, \dots, n\} = \{((x_1, x_2, \dots, x_m)_i, F(y_i)) \mid i = 1, 2, \dots, n\}$ .

When determining the corresponding nonlinear mathematical model as per the experimental data set, some textbooks and papers usually first get the residual sum of squares in the new  $m + 1$  -dimensional variable space  $X-Z$ :

$$s_1 = \sum_{i=1}^n (z_i - \hat{z}_i)^2 \quad (2)$$

Then the least square method is used to determine the best estimate  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l) \in R^l$  of 1-dimensional parameter  $\alpha$ . Finally, substitute the estimate into Equation (1) to obtain the nonlinear mathematical model.

The above method is used to determine the nonlinear mathematical model, which, however, has naturally hidden a serious unnoticeable defect, which is the residual sum of squares  $S_1$  in the new  $m+1$ -dimensional variable space  $X-Z$  and the minimal 1-dimensional parameter  $\alpha$  does not necessarily ensure that the residual sum of squares  $S_2$  in the original  $m+1$ -dimensional variable space  $X-Z$  is minimal, where

$$s_2 = \sum_{i=1}^n (y_i - \tilde{y}_i)^2 \quad (3)$$

It is this defect that has led to an error in the regression parameter of the nonlinear mathematical model, as obtained using the above method. In severe cases, it will even make the nonlinear model ineffective completely.

#### 4.1.3. Improvement of Approximating Common Nonlinear Mathematical Models to Regression Model Parameters

From the above analysis, it can be seen that in order to derive an ideal nonlinear mathematical model, the data mining method is employed. Further, the information provided by the original data is made full use of to ensure that the residual sum of squares  $S_2$  of the original variable space  $X - Y$  is minimal. Expand the function  $\hat{z}_i = F(\hat{y}_i) = F[y(x_i)]$  which contains the unknown the 1-dimensional parameter  $\alpha$  on the function  $y^i$  by Taylor series, it can be found that

$$\hat{z}_i = F(\hat{y}_i) = z_i + F'(y_i)(\hat{y}_i - y_i) + \frac{1}{2} F''(y_i) (\hat{y}_i - y_i)^2 + \dots$$

When  $\hat{y}_i \rightarrow y_i$ , exclude the infinitesimal  $(\hat{y}_i - y_i)$ , all higher-order infinite small can be deducted as:

$$(\hat{y}_i - y_i) \approx \frac{(\hat{z}_i - z_i)}{F'(y_i)} \quad (4)$$

Then it can be obtained that approximate expression for the residual sum of squares  $S_2$  in the  $m+1$ -dimensional space  $X - Y$ .

$$S_2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \approx \sum_{i=1}^n \frac{(z_i - \hat{z}_i)^2}{[F'(y_i)]^2} \quad (5)$$

The least square method is used to Equation (5) to find out the normal equations corresponding to the nonlinear mathematical model and find out the best estimate  $\alpha^* = (\alpha_1, \alpha_2, \dots, \alpha_l) \in R^l$  of the parameter  $\alpha$  in 1-dimensional space. Eventually, substitute the best estimate into Equation (1) to get the nonlinear mathematical model in the original variable dimensional space  $X - Y$ .

#### 4.1.4. Method of the Combination Forecasting Model Based on Data Mining

Use the improvement method of approximating the commonly used regression model parameter of the nonlinear model; apply the least squares method and obtain the normal equation model, so as to find out the approximation regression model for the original nonlinear function. The impact of the model on prediction accuracy is closely related to the original data, and the approximate expression of  $S_2$  may be considered to be deduced by increasing the weight of the large-value original data. Therefore, the forecasting model derived with  $S_2$  demonstrates high forecasting accuracy of the large-value original data, but low forecasting accuracy of the small-value original data.

In medical statistics, focusing on solving specific practical problems, when the data are large and the forecasting model coincides with the problem solving target, it can perfectly solve this problem. When the data are small, to solve this problem perfectly, the following equation can be applied to derive the forecasting model. Obviously,  $S_3$  is to increase the weight of the small-value original data,

$$S_3 = \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{y_i} \right)^2 \approx \sum_{i=1}^n \frac{(z_i - \hat{z}_i)^2}{[y_i F'(y_i)]^2}$$

While weakening the weight of the large-value original data. Similarly, a forecasting model with high forecasting accuracy can also be constructed for the original data whose value is between large and small.

#### 4.2. In Empirical Analysis

Example 1 finds out the estimated value for parameters  $\alpha$  and  $\beta$  in the given nonlinear mathematical model  $y = \beta x / (\alpha + x)$ , as shown in Table 1

**Table 1. The Experimental Data for the Amount of Drug Absorption Y and Plasma Concentration X Per Gram of Protein**

x	12.7	21.2	51.7	77.2	212.49.5	22.5	42.3	67.8	243.8	
y	0.103	0.466	0.767	1.573	2.462	0.083	0.399	0.899	1.735	2.360

The original approximation regression method, improved approximation regression method and the Gauss-Newton algorithm are used to find out the estimated values and residual sum of squares for parameters  $\alpha$  and  $\beta$  in the nonlinear mathematical model, as follows: for the original approximation regression method:  $\alpha^{\wedge} = -103.302$ ,  $\beta^{\wedge} = -0.865$ , and the residual sum of square is 33.6095; for the improved approximation regression method:  $\alpha^{\wedge} = 85.413$ ,  $\beta^{\wedge} = 3.386$ , and the residual sum of square is 0.6988; for the Gauss-Newton algorithm:  $\alpha = 144.660$ ,  $\beta = 4.050$ , and residual sum of square is 0.4393. In accordance with the relational expression between the amount of drug absorption y and plasma concentration x per gram of protein, if the plasma concentration x tends towards the positive infinite (that's not the case for real experiments), y will tend to approach parameter  $\beta$  in the model. Parameter  $\beta$  is the theoretical saturation value for y. Experimental data indicate that parameter  $\beta$  in the given nonlinear mathematical model should at least meet the condition  $\beta > 2.462$ . Calculated by the original approximation regression method, the estimated value for parameter  $\beta$  is -0.865, which is widely divergent from the real experiment.

By analyzing the residual sum of square or observing the standard residual plot and the fitting figure, it can be found that the nonlinear mathematical model determined by the original approximation regression method fails because of the poor signal to noise ratio of the experimental data. The improved approximation regression method is used to find out

the estimated value for parameter  $\beta$  is 3.386, which is basically consistent with the real experiment process. By analyzing the residual sum of square or observing the standard residual plot and the fitting figure, it can be found that the nonlinear mathematical model determined by the improved approximation regression algorithm is in good agreement with the actual situation. According to the requirements of the mathematical statistical theory, when the plasma concentration  $x \in [9.5, 234.8]$  is within experimental control, it is applicable to use the improved approximation regression method to calculate the estimated values for parameters  $\alpha$  and  $\beta$ . Most notably, when calculating the estimated values for parameters  $\alpha$  and  $\beta$  by employing the improved and the original approximation regression methods, the calculation method and time of both methods are basically the same.

Example 2: the relationship between times of parasitic disease treatment  $x$  and the positive review rates  $y$  is  $y = \exp(\alpha x + \beta)$ . As per the following eight sets of experimental data, find out the estimated values for parameters  $\alpha$  and  $\beta$  in the nonlinear mathematical model under given conditions: (1) forecasting model of positive review rates with less than four times of treatment; (2) forecasting model of positive review rates with more than four times of treatment.

**Table 2. The Original Experimental Data**

x	1	2	3	4	5	6	7	8
y	63.9	36.0	17.1	10.5	7.3	4.5	2.8	1.7

For the original experimental data in Table 2, the combination forecasting model is employed based on the data mining method. It has been calculated that the forecasting model of positive review rates with less than four times of treatment is`

$$y_1 = \exp(-0.583x + 4.763)$$

The forecasting model of positive review rates with more than four times of treatment is  $y_2 = \exp(-0.506x + 4.526)$

The original data and combination forecasting data are listed in Table 3, from which it can be seen that the combination forecasting result is ideal.

**Table 3. The Original Data and Combination Forecasting Data**

x	1	2	3	4	5	6	7	8
y	63.9	36.0	17.1	10.5	7.3	4.5	2.8	1.7
y1	63.6	35.5	19.8	11.1	6.2	3.5	1.9	1.1
y2	55.7	33.6	20.3	12.2	7.4	4.4	2.7	1.6

## 5. Discussions

The improved approximation regression method and the original approximation regression method share a similarity that: the data in the  $m+1$ -dimensional space  $X - Y$  are converted into data in the  $m+1$ -dimensional space  $X - Z$ , and variable substitution  $z = F(y)$  is introduced. The difference is that in regard to the original approximation regression method, use the residual sum of squares  $s_1 = \sum_{i=1}^n (z_i - \hat{z}_i)^2$  in the new  $m+1$ -dimensional variable space  $X - Z$ , and then find out the best estimate

$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l) \in R^l$  of I-dimensional parameter  $\alpha$  by using the least squares method. Substitute the best estimate into Equation (1) to obtain the nonlinear mathematical model. As for the improved method, use the residual sum of squares  $s_2 = \sum_{i=1}^n (y_i - \tilde{y}_i)^2$  in the original m+1-dimensional variable space X-Y, and then find out the best estimate  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l) \in R^l$  of I-dimensional parameter  $\alpha$  by using the least squares method. Substitute the best estimate into Equation (1) to obtain the nonlinear mathematical model. Precisely because of this difference, compared to the nonlinear mathematical model obtained by the original approximation regression method, the nonlinear mathematical model obtained by the improved method can more significantly improve the regression accuracy, and the normal equation derived by the improved approximation regression method is merely the weighted normal equation derived by the original approximation approach and retains the advantages of its ease of use.

The combination forecasting model based on data mining can dig out more information from the original data, which is conducive to solving practical problems in different situations. According to statistic investigation, data mining has a potential huge market value and will form a new industry in China in near future, with the increase of data volume and the wide application of computer.

## References

- [1] F. L. Krause and U. Kaufmann, "Meta-modelling for interoperability in data mining", CIRP Annals. vol. 145, (2007), pp. 191-196.
- [2] J. Trujillo and L. M. Sergio, "A UML based approach for modeling ETL processes in data warehouses", Conceptual Modeling 2003, Chicago, Notes in Computer Science, (2003), pp. 277-282.
- [3] J. Lines, L. M. Davis and J. Hills, "A shapelet transform for time series classification", Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining, Beijing, China, (2012), pp. 289-297.
- [4] A. Alzghoul, M. Lofstrand and B. Backe, "Data stream forecasting for system fault prediction", Computers and Industrial Engineering, vol. 62, no. 4, (2012), pp. 972-978.
- [5] S. Hashemi and Y. Yang, "Flexible decision tree for data stream classification in the presence of concept change, noise and missing values", Data Mining and Knowledge Discovery, vol. 19, no. 1, (2009), pp. 95-131.
- [6] S. Junier and E. Mostert, "A decision support system for the implementation of the Water Framework Directive in the Netherlands: Process, validity and useful information", Environmental Science & Policy, vol. 40, (2014), pp. 49-56.
- [7] E. Sunea, "Improving Decision Making Process in Universities: A Conceptual Model of Intelligent Decision Support System", Proceeding-Social and Behavioral Sciences, vol. 76, (2013), pp. 795-800.
- [8] Z. Y. Zhuang, C. L. Wilkin and A. Ceglowski, "A framework for an intelligent decision support system: A case in pathology test ordering", Decision Support Systems, vol. 55, no. 2, (2013), pp. 476-487.
- [9] Z. Shiqiang, "Approach on the Fitting Optimization Index of Curve Regression", Chinese Journal of Health Statistics, (2002), pp. 9-13.

## Authors



**Dong Han**, professor, Research direction: Computer analysis. Data mining



**Chunhua Wang**, assisted professor, Research direction: Computer analysis. Data mining.