

Research on Information Forecasting Based on Different Data Mining Techniques

Yiran Wang and Guang Zheng

College of Network Engineering, Zhoukou Normal University, Henan, Zhoukou
466001, China

Corresponding-Email: Wangyr@126.com

Abstract

This paper has been explored information data prediction implementation access based on data mining combination model. With data mining technology as the entry point and in combination with the analysis on information data prediction characteristics. Research on variable substitution to non-linear regression forecast model precision's influence, and seek the modeling 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. Given modeling method of combination forecast model based on the Data mining. Based on Data mining's combination forecast model's modeling 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.

Keywords: Data mining technology; data prediction; Non-linear model; Space trans-Formation

1. Introduction

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 [Krause F. L. *et al.* 2007]. 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 [Alzghoul A. *et al.* 2012]. Information data prediction and analysis relies on data source, therefore, the characteristics of the data source also decides the characteristics of the big data prediction. 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 [Trujillo J. *et al.* 2003]. It is under such circumstances that data mining technology comes out in response to the demands of the times. 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 [Lines J. *et al.* 2012]. 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.

This paper chooses Heze food enterprise to collect information. 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 [Hashemi S. *et al.* 2009]. 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 [Sandra Junier. *et al.* 2014]; even though the number exceeds the value, there won't be any more help to solve the questions.

2. Data Clustering, Classification and Feature Selection Algorithm

For real time IOT data stream, fast density based clustering is required which includes density based data stream clustering. This clustering grouped as density-grid based method and density based microclustering method. The main advantage of density-grid approach is its fast processing time that is independent of the number of data points and dependent only on number of cells. On other hand, density based microclustering method keep summary of clusters in microclusters and form a final clusters from them. proposed real experiment with HDC (hybrid density-based clustering for data stream)-stream algorithm. HDC-Stream only searches in potential list and if it cannot find the suitable microcluster, the data point is mapped to the grid, which keeps the outlier buffer. In future, author will focus on distributed HDC-stream density-based clustering to improve performance in IOT.

1. On the basis of feature values, decision tree classifies instances. For a given set S cases, C4.5 first grows an initial A. tree using divide-and-conquer algorithm as follow: tree is leaf labelled if all cases belongs to same class S or S is small. B. Otherwise, select test based on single attribute. Some limitation this algorithm pertains are: Empty branches, insignificant branches and over fitting. Most decision tree algorithm cannot perform well with problem that require diagonal partitioning.
2. The most common role of data mining is to find frequent itemsets from transaction datasets and derive association rules. Once itemsets are obtained, it's upfront to generate association rules. To achieve this Apriori algorithm is helpful. This algorithm is assumes that items within transaction or itemsets are sorted in lexicographic order. The Apriori algorithm generally perform in 2steps join and prune step. It then calculates frequency only for those candidates generated by scanning the database.

As growing pressure for classification of data in urgency situation: data classification for breach response, for e-discovery, for business unity as moving towards cloud.

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 [Elena Sunea, *et al.* 2013]. 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.

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.

Information 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 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 information forecasting model based on data mining methods can overcome this defect.

Form of the Nonlinear Mathematical Model

The nonlinear mathematical model can be expressed as follows:

$$y = f(x_1, x_2, \dots, x_m, \alpha_1, \alpha_2, \dots, \alpha_l) + e \quad (1)$$

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 α 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 α .

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

$$\{ ((x_1, x_2, \dots, x_m)_i, z_i) | i = 1, 2, \dots, n \} = \{ ((x_1, x_2, \dots, x_m)_i, F(y_i)) | i = 1, 2, \dots, n \}.$$

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

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 $X-Z$:

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

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 α . 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 α 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

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

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.

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 α 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:

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

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 α in l -dimensional space. Eventually, substitute the best estimate into Equation (1) to get the nonlinear mathematical model in the original variable dimensional space $X - Y$.

Method of the Combination Information 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 deducted by increasing the weight of the large-value original data. Therefore, the information forecasting model derived with S_2 demonstrates high information forecasting accuracy of the large-value original data, but low information forecasting accuracy of the small-value original data.

In medical statistics, focusing on solving specific practical problems, when the data are large and the information 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 information forecasting model. Obviously, S_3 is to increase the weight of the small-value original data,

$$S_3 = \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{y_i}\right)^2 \approx \sum_{i=1}^n \frac{(z_i - \hat{z}_i)^2}{[y_i F'(y_i)]^2} \quad (6)$$

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

In Empirical Analysis

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

Table 1. The Estimated Value 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	0.1030	4.660	7.671	5.732	4.620	0.830	3.990	8.991	7.352

Note: X stands for plasma concentration. Y stands for the amount of drug absorption

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 α and β 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 β in the model. Parameter β is the theoretical saturation value for y. Experimental data indicate that parameter β 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 β 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 β 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 α and β . Most notably, when calculating the estimated values for parameters α and β 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 α and β in the nonlinear mathematical model under given conditions: (1) information forecasting model of positive review rates with less than four times of treatment; (2) information 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	63.936	0.17	1.110	57.34	52.81	7		

Note: X stands for plasma concentration. Y stands for the amount of drug absorption

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

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

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

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

Table 3. The Original Data and Combination Information Forecasting Data

x	1	2	3	4	5	6	7	8
y	63.936	0.17	110.57	34.52	81.7			
y1	63.635	5.19	8.11	16.23	51.91	1.1		
y2	55.733	6.20	3.12	27.44	42.71	6		

Note: X stands for plasma concentration. Y stands for the amount of drug absorption
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 α 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 α 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.

Conclusion

The combination information 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.

Conflict of Interests

The authors declared that they have no conflicts of interest to this work.

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 information 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.

Authors



Wang Yiran, is lecture in the network engineering, China. His research and teaching interests focus on computer network, computer safety.