

## Research Article

# Claim Amount Forecasting and Pricing of Automobile Insurance Based on the BP Neural Network

Wenguang Yu <sup>1</sup>, Guofeng Guan,<sup>2</sup> Jingchao Li,<sup>3</sup> Qi Wang,<sup>2</sup> Xiaohan Xie,<sup>1</sup> Yu Zhang,<sup>1</sup> Yujuan Huang <sup>4</sup>, Xinliang Yu,<sup>1</sup> and Chaoran Cui<sup>5</sup>

<sup>1</sup>School of Insurance, Shandong University of Finance and Economics, Jinan 250014, China

<sup>2</sup>School of Mathematic and Quantitative Economics, Shandong University of Finance and Economics, Jinan 250014, China

<sup>3</sup>College of Mathematics and Statistics, Shenzhen University, Shenzhen 518060, China

<sup>4</sup>Office of Academic Research, Shandong Jiaotong University, Jinan 250357, China

<sup>5</sup>School of Computer Science & Technology, Shandong University of Finance and Economics, Jinan 250014, China

Correspondence should be addressed to Wenguang Yu; yuwg@sdufe.edu.cn

Received 31 October 2020; Revised 23 December 2020; Accepted 7 January 2021; Published 20 January 2021

Academic Editor: Benjamin Miranda Tabak

Copyright © 2021 Wenguang Yu et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The BP neural network model is a hot issue in recent academic research, and it has been successfully applied to many other fields, but few researchers apply the BP neural network model to the field of automobile insurance. The main method that has been used in the prediction of the total claim amount in automobile insurance is the generalized linear model, where the BP neural network model could provide a different approach to estimate the total claim loss. This paper uses a genetic algorithm to optimize the structure of the BP neural network at first, and the calculation speed is significantly improved. At the same time, by considering the overfitting problem, an early stop method is introduced to avoid the overfitting problem. In the model, a three-layer BP neural network model, which includes the input layer, hidden layer, and output layer, is trained. With consideration of various factors, a total claim amount prediction model is established, and the trained BP neural network model is used to predict the total claim amount of automobile insurance based on the data of the training set. The results show that the accuracy of the prediction by using the BP neural network model to both the data of Shandong Province and to the data of six cities is over 95%. Then, the predicted total claim amount is used to calculate premiums for five cities in Shandong Province according to credibility theory. The results show that the average premium of the five cities is slightly higher than the actual claim amount of the city. The combination of BP neural network and credibility theory can perform accurate claim amount estimation and pricing for automobile insurance, which can effectively improve the current situation of the automobile insurance business and promote the development of insurance industry.

## 1. Introduction

Insurance industry in China has made great progress along with the continuous development of Chinese economy, and insurance plays an increasingly important role in people's daily life. Therefore, a fair and comprehensive pricing system is essential to the development of insurance industry, which can effectively avoid the adverse selection problem, can maintain the insurance industry in a healthy competition, and can promote the development of insurance industry. According to the China Statistical Yearbook in 2015, the

premium income of automobile insurance was 619.9 billion yuan, accounting for 73.59% of the premium income of property insurance; in 2016, the premium income of automobile insurance was 683.42 billion yuan, accounting for 73.76% of the premium income of property insurance; in 2017, the premium income of automobile insurance was 752.11 billion yuan, accounting for 71.35% of the premium income of property insurance. It can be seen that the premium income of automobile insurance is steadily increasing, and the proportion of the premium income to the property insurance premium income is maintained at more than 70%.

Hence, the profitability of automobile insurance plays a decisive role in the operating efficiency of property insurance companies. However, the operating condition of the automobile insurance business in China is generally poor. According to the data of China Insurance Regulatory Commission in 2018, only seven property and automobile insurance companies made profit, where 48 out of 55 unlisted property and casualty insurance companies engaged in the automobile insurance business suffered losses to varying degrees. The total loss of automobile insurance is about 8.68 billion yuan, and the loss ratio is too high in most years.

The main problem of excessive high claims in automobile insurance is inadequate premium ratemaking. Inadequate premium ratemaking does not simply mean that the premium rates are too high or too low, but means that the premiums among different risks are not differentiated or the distinctions are inappropriate, which leads to a large number of adverse selection and causes poor business quality so that the premium income does not match with the risk the company takes.

The study on premium ratemaking of automobile insurance has been attracting many scholars' attention. Bailey and Simon first proposed the idea of classification pricing, which classified insurance policies according to a certain characteristic of the risk and priced each type of insurance policies separately [1]. Denneberg first proposed the Poisson-gamma model to study the frequency of nonhomogeneous insurance policy claims and obtained good fitting results in empirical research on auto insurance [2].

The generalized linear model (GLM) is a widely accepted model for premium ratemaking of automobile insurance in recent decades. In the last century, Nelder and Wedderburn first proposed the GLM, which has been widely used as once proposed [3]. Samson and Thomas used the GLM to perform pricing for the premium rate based on the data of a third-liability insurance from an insurance company in UK, and they found that no claim discount, automobile type, region, and age class of the automobile owner have significant impact on both the claim amount and claim frequency [4]. Smyth introduced the maximum likelihood estimation of the DGLM, considered the situation when the population obeys the normal and inverse Gaussian distribution, and analyzed the selection of the initial value of the iteration [5]. Stroiński and Currie evaluated the risk through the GLM based on the data of a third-liability insurance from an insurance company in UK and proved that the GLM plays an important role in premium ratemaking of automobile insurance [6]. Meng briefly analyzed the shortcomings of traditional nonlife insurance product rate determination methods, such as the single analysis method and minimum bias procedure. It also proved that the GLM can be applied to determine premiums for automobile products by using a group of automobile insurance data [7]. Draper showed that the GLM fitted better than the traditional model based on automobile insurance data of an insurance company in France by using SAS software [8]. Zhao and Chen applied the dual-generalized linear model to price the automobile insurance premium rate. Based on an empirical study of a group of automobile insurance loss data, they found that the

dual-generalized linear model is more reasonable to determine the premium rate compared with the GLM [9].

With the development of research on GLMs, scholars have found limitations of the GLM. Initially, the GLM only builds the regression relationship between the expected value of the response variable and explanatory variables and assumed that the dispersion parameter is constant. Although this assumption simplifies the model, it does not hold for some cases. Smyth and Jørgensen applied the DGLM to the premium ratemaking of auto insurance and directly predicted the premium rate of auto insurance. However, the regional factors were excluded in the empirical study, and the obtained rate structure did not reflect the regional differences [10]. Antonio and Beirlant combined the generalized linear-mixed model and Bayesian method to determine premium rates and found that this combination performs well [11]. According to the nonlife insurance loss data, Frees et al. firstly analyzed the claim frequency, claim type, and claim intensity under the framework of the hierarchical model, then applied the Bayesian method to determine the joint probability distribution among variables, and finally predicted the total claim loss in the future. Finally, they used the simulation method to predict the premium under the policy limit [12]. Wang et al. considered that the fat tail of automobile insurance loss data has significant impact on premium ratemaking; hence, they introduced the density function to describe the fat tail distribution, and based on it, they constructed the GAMLSS model under the type-two generalized beta distribution, which improved the limitations of assumptions and parameter modeling in the traditional GLM and increased the accuracy of the prediction for automobile insurance loss [13]. In the past, it was assumed that the total claim distribution was compound Poisson-gamma distributed. The GLM was used for the claim frequency and claim intensity, respectively, and then, the expected total claim was set to be equal to the expected value of claim frequency times the expected value of claim intensity. Zhang and Xie assumed that the total claim amount followed the Tweeite distribution, then directly established a GLM for the total claim amount, and obtained the average value of the total claim amount for each risk. Through the empirical analysis of the data, the above two methods are compared, and results showed that the data fitting degree based on Tweeite distribution was better [14].

With the development of science and technology, more and more attention has been paid to the driving behavior in premium ratemaking of automobile insurance. Ayuso et al. demonstrated how automobile insurance can be improved by incorporating mileage and driver behavior data. The key idea is that telemetry should facilitate the inclusion within insurance pricing of those factors that traffic authorities identify as being associated with risky drivers [15]. Huang and Meng studied the use of a wide range of driving behavior variables to predict the risk probability and claim frequency of an insured vehicle. The advantage of the model is that it can improve the interpretability and predictive accuracy of the model at the same time, thus providing a new solution for the classification pricing of UBI products [16].

Artificial neural network is a new model. It is a kind of computational model which imitates the structure of biological network after human beings have fully studied the structure of animals.

1. BP neural network model is applied to automobile insurance.
2. Genetic algorithm is used to optimize the structure of BP neural network.
3. An early stop method is introduced to avoid the overfitting problem.

The neural network is composed of multiple neurons, and each neuron is connected to each other to form a network. The network transmits and processes information and imitates the human brain structure. It has strong adaptability and can process linear and nonlinear data. The BP neural network algorithm based on backpropagation is one of the most mature and widely used neural network algorithms.

Many scholars have introduced neural network algorithms into the insurance industry. Brockett et al. applied Kohonen's self-organizing competitive networks to identify fraud problems in personal injury insurance [17]. Liu et al. compared the multiclass AdaBoost tree with the generalized linear model, two-layer BP (backpropagation) neural network, and support vector machine (SVM) to predict the effect of claim intensity, and they found that the AdaBoost method has the best prediction accuracy and relatively small variance [18]. Mzhavia applied neural networks to the risk classification of car drivers and found a set of neural networks with the best classification effect. The number of neurons in the input layer, hidden layer, and output layer of the network was 11, 12, and 2. The inspirit function was a hyperbolic tangent function [19]. Under the assumption of Poisson distribution, Wüthrich used the speed-acceleration data recorded by the Internet of vehicles to extract the driving behavior factor through the Bottleneck neural network learning algorithm and established a generalized additive model to predict the frequency of claims [20]. Zhang and Wang applied SOM to the claim prediction of automobile insurance, which provided a new way of premium ratemaking of automobile insurance [21]. In addition, the application of the neural network in other fields can be seen in Lin et al's research [22–25].

In summary, it can be found that there are abundant research studies on automobile insurance pricing and premium ratemaking, and scholars have been seeking new methods to price for automobile insurance more accurately. At the early stage, scholars mainly studied generalized linear models and continuously improved the generalized linear models so that the generalized linear models could be better applied to the premium ratemaking. However, the generalized linear models still have some shortcomings. The BP neural network has strong fault tolerance and high accuracy when fitting data. Aiming at the characteristics of the BP neural network, this paper tries to apply the BP neural network to price for automobile insurance rates and verifies the model with real data from insurance companies. The accuracy of the model is expected to provide new ideas for the premium ratemaking in the insurance industry.

The outline of the paper is organized as follows. In Section 2, we construct the BP neural network model. In Section 3, empirical analysis is carried out on the model. In Section 4, according to the characteristics of the data, we extend the model and verify that the model is also applicable to all regions of the country. In Section 5, we study the application of the model in automobile insurance ratemaking. Finally, conclusions and policy recommendations are given in Section 6.

## 2. The BP Neural Network Model and Optimization

*2.1. The BP Neural Network Model.* The BP neural network is currently the most widely used neural network. The learning rule is to adopt the algorithm of backpropagation, using the steepest descent method to adjust the coefficient in reverse according to the error between the actual output value and expected output value, until the coefficient is optimized to make the error within the acceptable range. The BP neural network can learn fixed patterns, use some data to determine the corresponding parameters, and then make predictions based on these parameters.

A three-layer BP neural network model is used in this paper, which is composed of an input layer, an output layer, and a hidden layer. Its structure is shown in Figure 1.

We assume that there are  $n$  neurons in the input layer, five neurons in the hidden layer, and two neurons in the output layer.  $X_k = (x_{1k}, x_{2k}, \dots, x_{nk})$  which are input values,  $k = 1, 2, \dots, m$ .  $\alpha_{ij}$  are the weights connecting the input layer and the hidden layer and  $\beta_{jl}$  are the weights connecting the hidden layer and the output layer value, where  $i = 1, 2, \dots, n$ ,  $j = 1, 2, \dots, 5$ , and  $l = 1, 2$ . Neurons in the same layer are not connected to each other, and each neuron between the input layer and hidden layer and hidden layer and output layer are connected.

The specific algorithm of the BP neural network model based on Figure 1 is as follows: assuming that the stimulus function uses the sigmoid function, sequentially input  $m$  sample data  $X_1, X_2, \dots, X_m$  and then randomly select the  $k$ th input sample  $X_k = (x_{1k}, x_{2k}, \dots, x_{nk})$ ; the hidden layer input vector is  $Y_k = (y_{1k}, y_{2k}, \dots, y_{5k})$ , the hidden layer output vector is  $Z_k = (z_{1k}, z_{2k}, \dots, z_{5k})$ , the input vector of the output layer is  $\tilde{Y}_k = (\tilde{y}_{1k}, \tilde{y}_{2k})$ ,  $\tilde{Z}_k = (\tilde{z}_{1k}, \tilde{z}_{2k})$  is the output vector of the output layer, the expected output vector is  $R_k = (r_{1k}, r_{2k})$ , the threshold of each neuron in the hidden layer is denoted as  $a_j$ , the threshold of each neuron in the output layer is denoted as  $b_l$ , the inspirit function is  $f(\cdot)$ , the learning parameter is  $\mu$ , and  $E = (1/2) \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk})^2$  is the error function.

The input and output of each neuron in the hidden layer and output layer can be calculated as follows:

$$y_{jk} = \sum_{i=1}^n \alpha_{ij} x_{ik} - a_j,$$

$$z_{jk} = f(y_{jk}),$$

$$\tilde{y}_{lk} = \sum_{j=1}^5 \beta_{jl} z_{jk} - b_l,$$

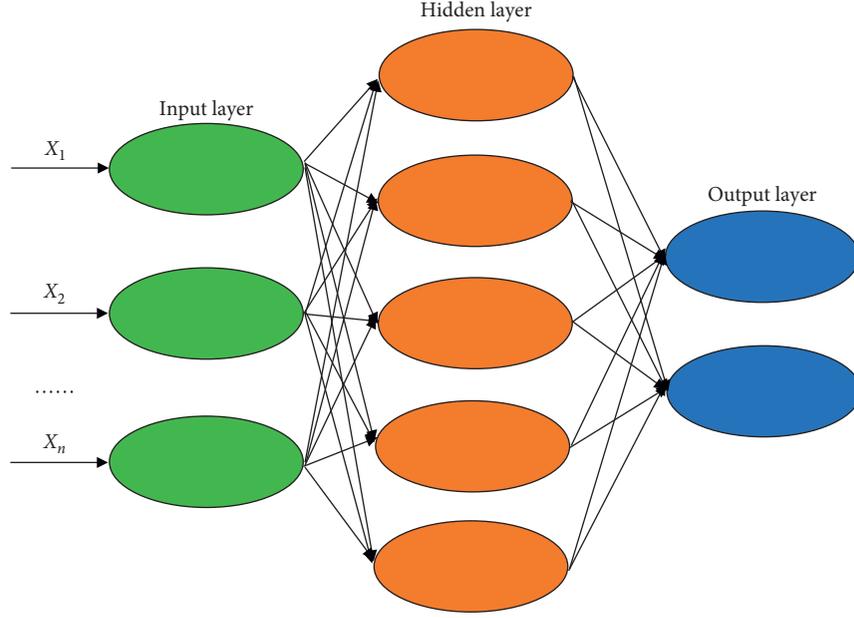


FIGURE 1: The three-layer BP neural network model.

$$\tilde{z}_{lk} = f(\tilde{y}_{lk}). \quad (1)$$

By using the expected output and actual output of the network, the partial derivative of the error function for each neuron in the output layer is as follows:

$$\begin{aligned} \frac{\partial E}{\partial \tilde{y}_{lk}} &= \frac{\partial [(1/2) \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk})^2]}{\partial \tilde{y}_{lk}} \\ &= \frac{\partial [(1/2) \sum_{l=1}^2 (r_{lk} - f(\tilde{y}_{lk}))^2]}{\partial \tilde{y}_{lk}} = - \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk}) f'(\tilde{y}_{lk}) \triangleq -\delta_{lk}. \end{aligned} \quad (2)$$

By using the connection weight from the hidden layer to the output layer, the output layer and output of the hidden layer to calculate the partial derivative of the error function for each neuron of the hidden layer is given as follows:

$$\begin{aligned} \frac{\partial E}{\partial y_{jk}} &= \frac{\partial [(1/2) \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk})^2]}{\partial y_{jk}}, \\ &= \frac{\partial [(1/2) \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk})^2]}{\partial z_{jk}} \cdot \frac{\partial z_{jk}}{\partial y_{jk}}, \\ &= \frac{\partial [(1/2) \sum_{l=1}^2 (r_{lk} - f(\sum_{j=1}^5 \beta_{jl} z_{jk} - b_l))^2]}{\partial z_{jk}} \cdot f'(y_{jk}), \\ &= - \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk}) f'(\tilde{y}_{lk}) \beta_{jl} \cdot f'(y_{jk}), \\ &= - \left( \sum_{l=1}^2 \delta_{lk} \beta_{jl} \right) f'(y_{jk}), \quad \triangleq -\rho_{jk}. \end{aligned} \quad (3)$$

Through the above two formulas, we can get the change of weight value  $\beta_{jl}$  in each adjustment as follows:

$$\Delta \beta_{jl} = -\mu \frac{\partial E}{\partial \beta_{jl}} = -\mu \frac{\partial E}{\partial \tilde{y}_{lk}} \cdot \frac{\partial \tilde{y}_{lk}}{\partial \beta_{jl}} = \mu \delta_{lk} z_{jk}. \quad (4)$$

After  $N$  adjustments, the  $(N+1)$ th value is as follows:

$$\beta_{jl}^{N+1} = \beta_{jl}^N + \Delta \beta_{jl}. \quad (5)$$

Similarly, we can get the change of weight value  $\alpha_{ij}$  in each adjustment and the  $(N+1)$ th value after  $N$  adjustments as follows:

$$\Delta \alpha_{ij} = -\mu \frac{\partial E}{\partial \alpha_{ij}} = \mu x_{ik} \rho_{jk}, \quad (6)$$

$$\alpha_{ij}^{N+1} = \alpha_{ij}^N + \Delta \alpha_{ij},$$

and the global error can be calculated as follows:

$$E = \frac{1}{2m} \sum_{k=1}^m \sum_{l=1}^2 (r_{lk} - \tilde{z}_{lk})^2. \quad (7)$$

Finally, compare the size of the global error with the setting error. If the global error exceeds the setting error, keep adjusting the weights until the setting error is met.

The calculation process of the BP neural network model can be presented graphically, as shown in Figure 2.

**2.2. Optimization and Control Problem of Overfitting Issue of the BP Neural Network.** The BP neural network performs local search according to the gradient descent method, and the weight adjustment in the network is realized by the local adjustment. However, the BP neural network has the following problems. Firstly, adjusting the weight locally makes

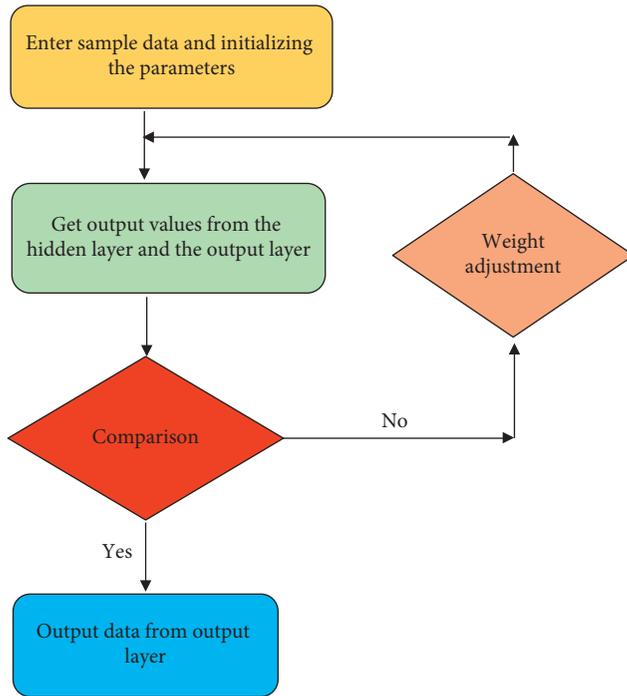


FIGURE 2: Back propagation flow chart of the BP neural network.

the weights fall into the local optimum rather than the global optimum. Secondly, when using the gradient descent method for optimization, there will be a flat area when the output neuron result is close to 0 or 1, where the error changes slightly with the weights in the flat zone, which causes the training process to be extremely slow and even be judged by the network that the global optimal solution has been found, and the network training will stop.

Considering the above shortcomings of the BP neural network, this paper chooses a genetic algorithm to improve the structure of the BP neural network. The genetic algorithm has the following characteristics:

- (1) The genetic algorithm starts searching from multiple starting points instead of starting from one single point so that the search range is larger when searching for the best, which is more conducive to search for the global optimal solution.
- (2) The genetic algorithm searches for the target solution in an adaptive manner, in which a series of operations such as selection, crossover, and mutation is operated based on a certain probability; hence, there are great flexibility and nondirectionality in the search process.
- (3) The genetic algorithm does not rely on search space knowledge and other auxiliary information. It only calculates the degree of individual pros and cons based on the fitness function and performs subsequent operations on this basis.
- (4) The genetic algorithm takes the coding of the required target variable as the operation target.

Based on the above characteristics, this method can use the parallel property of the global search to find the optimal

connection weight set of the neural network, and the use of gradientless optimization and random operators helps to evolve the initial weights of the artificial neural network so that the probability of the BP algorithm falling into a local minimum is minimized.

Therefore, this paper is based on the genetic algorithm to search the optimal value for the initial weight and threshold in the neural network to optimize the BP neural network. The process is set as follows. At the first stage, the genetic algorithm is used to train the neural network. The genetic algorithm is used to evolve the optimal initial weight and threshold set of neural network training. This is achieved by the genetic algorithm simultaneously searching in all possible directions in the search space and narrowing it down to the area where the best possible weight and deviation can be found. In the second stage, the neural network is trained using the BP algorithm. The training starts by initializing the BP algorithm and using the genetic algorithm to assist in training the initial weights and thresholds of the evolution. The neural network with optimal weights and thresholds is initialized by using the BP algorithm, and the global optimal solution started by the genetic algorithm is searched continuously by adjusting the weights and thresholds of the neural network.

Although the genetic algorithm is used to optimize the structure of the BP neural network, the network may still have an overfitting problem. The overfitting problem means that the accuracy of the network on the training set is very high, but the accuracy on the test set is relatively low, which affects the generalization of the model.

In order to avoid overfitting, the early stop method is applied in this paper. This method divides the sample into a confirmed subset and uses this subset to test the network error during the training process. In the initial stage of training, the network error will be reduced; but, when the network begins to overfit, the network error will rise; when the network error rises in a certain number of iterations, the network stops training. At this time, the weight and bias value of the network can be obtained when the network error is the smallest.

In addition, this paper also uses 26,100 pieces of data to train each parameter by increasing the amount of data as much as possible. The feature diversity of the data helps to make full use of all training parameters to result estimation and simulation more accurate.

Finally, this paper also reduces the number of BP neural network layers and the input layer and hidden layer neurons to prevent overfitting when building the model.

Figure 3 shows the error reduction graph in network training. The errors of the training set, prediction set, and confirmation subset have been decreasing until they become stable, indicating that the network has no overfitting problems. At the same time, the prediction accuracy of each prediction set below is also very high which verifies the conclusion.

### 3. Empirical Analysis of the BP Neural Network Model

**3.1. Data Source Analysis.** The data in this paper comes from real claim data of an insurance company in Shandong Province. The data contains eight columns, which are owner's age, owner's gender, number of seats, vehicle age,

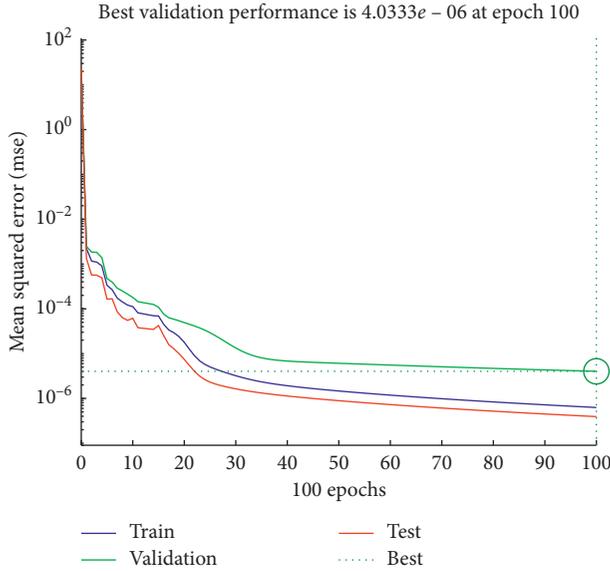


FIGURE 3: The error reduction graph under the early stop method.

purchase price, vehicle weight, NCD coefficient, and aggregate claim amounts, totaling 30,000 data. After removing outliers, there are 29,000 valid data remaining. The data covers 17 prefecture-level cities in Shandong Province, but the amount of data is not uniformly distributed in each prefecture-level city. Among them, the data of five prefecture-level cities such as Weifang, Binzhou, Jining, Yantai, and Weihai exceeds two thousand data, of which there are 6500 data in Yantai city, and the data of the remaining prefecture-level cities are less than one thousand. Some data are shown in Table 1.

**3.2. Overall Modeling and Result Analysis.** This paper takes the owner's age, owner's gender, number of seats, vehicle age, purchase price, vehicle weight, and NCD coefficient as input variables and the aggregate claim amount as the output variable. In addition, the data required by the BP neural network model is divided into a training set and a test set. The training set is used to train the network, adjust the parameters, and control the error within our goals, where the test set uses the trained network to predict the aggregate claim amount. In order to make the network training more accurate and to take the recognition degree of the output result graph into account, we select 90% of the data as a training set of the neural network, and the remaining 2,900 data are used as a test set to test prediction performance of the neural network.

Since there is a big difference in the magnitude of the data used in this paper, it will have a negative impact on the training of the network, reducing the learning ability of the network, and may even fail to reach the training goal. Therefore, all the data will be normalized at first, where all the data are mapped in the range of 0-1 to facilitate network training.

In this paper, the number of nodes in the input layer is set to 7, corresponding to 7 input variables. The number of nodes in the output layer is 1, which corresponds to the

aggregate claim amount. The number of nodes in the hidden layer has a significant impact on the performance of the network, and this paper adopts an empirical algorithm to determine the number of nodes in the hidden layer, that is,  $\sqrt{m+n+a}$ , where  $m$  and  $n$  are the number of nodes in the input layer and output layer, respectively, and  $a$  is a random number, and after many tests and adjustments, it is taken 1 in this paper. The inspirit function between the input layer and hidden layer adopts a tan *sig* function, which is a hyperbolic tangent inspirit function. The activation function between the hidden layer and output layer adopts a purelin inspirit function, which is a linear inspirit function. The training function uses a trainlm function. The network parameters in this paper are set as follows: the learning rate is 0.1, and the target accuracy is set to 0.00001. We use prediction accuracy to measure the results of network prediction, which is defined as follows:

$$\text{prediction accuracy} = 1 - \frac{\text{prediction value} - \text{real value}}{\text{real value}}. \quad (8)$$

We first let the network randomly select the training set and use the selected training set to train the network. The analysis of the training network result is shown in Figure 3.

Figure 4 is a state diagram of using the training set to complete the training for the network, and it indicates that the error between the output value of the training set and the actual output value is 0.000000628 after the iteration; hence, the network can accurately fit this type of data. This also shows that the BP neural network model can be used to fit and predict the total claim amount of automobile insurance, and the trained network is a network with superior performance.

Figure 5 shows the goodness of the fit, and it indicates that global goodness of fit  $R$  reaches 0.99952, and the coefficient is very close to 1, which strongly indicates that the input neuron variable has a strong explanatory effect on the output neuron variable. The BP neural network model is very suitable for the fitting of automobile insurance claim data, and it can be used to fit and predict the aggregate claim amount.

Figure 6 is an effect diagram of using the trained network to fit the training set data. Since the individuals with no claims in the data set account for the majority, the bottom of the graph is denser, and the actual value and fitted value are stacked together, indicating that the fitting effect for zero claim amount is very good. For some very small claims, although the plots are relatively dense, it can be seen that they can be fitted well too. For individuals with relatively large claims, the network can also be accurately fitted. Hence, the network fits the training set well, and it fits most individuals accurately.

Given that the BP neural network model can accurately fit the training set, the trained network model is used to predict the data of the test set. Figure 7 shows the comparison between predicted values and actual values of the test set. The test data set also accounts for the majority of individuals with no claim, so the phenomenon of

TABLE 1: Claim data.

Owner's age	Owner's gender	Number of seats	Vehicle age	Purchase price (yuan)	Vehicle weight (kg)	NCD coefficient	Aggregate claim amounts
48	1	5	4	62,900.0	1039	-0.4	0.0
31	1	5	3	61,900.0	1265	-0.3	0.0
50	1	8	5	30,000.0	975	-1.3	0.0
47	1	5	9	114,800.0	1255	-1.3	0.0
33	1	5	13	263,800.0	1580	-0.15	0.0
31	2	7	5	38,800.0	1205	-0.3	0.0
40	1	5	6	37,900.0	1210	-1.3	0.0
26	2	5	3	44,900.0	1020	-1.2	0.0
33	1	7	8	278,800.0	1855	-1	0.0
37	1	5	9	37,900.0	895	-1	0.0
49	1	5	3	119,800.0	1501	0.7	0.0
40	2	8	2	52,300.0	1305	0.15	0.0
49	1	5	17	55,500.0	1050	-0.4	1,200.0

Note. In owner's gender, 1 indicates male and 2 indicates female.

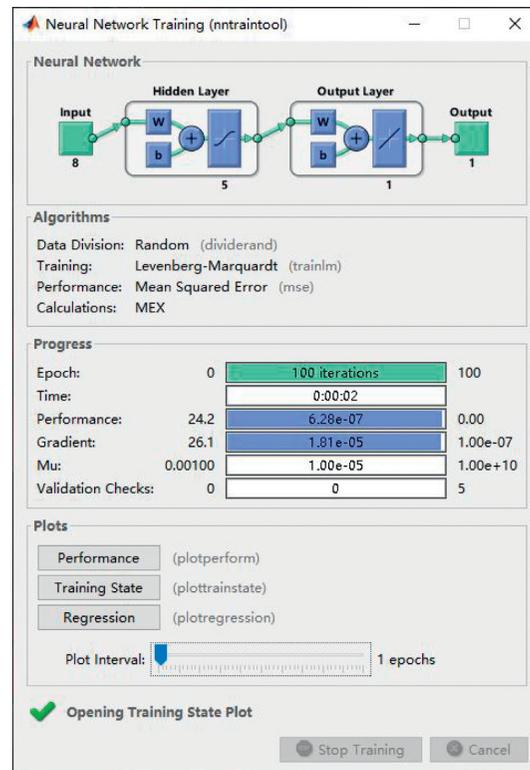


FIGURE 4: The overall modeling network state diagram.

accumulation appears at the bottom of the graph, which also shows that the network can accurately predict individuals with zero claims. In addition, the prediction of individuals whose total claim amount is not zero is also very accurate. The output shows that the prediction accuracy of the network on prediction data set is 99.68%, which proves the accuracy of the overall prediction results of the network, indicating that the BP neural network is very suitable for predicting the aggregate claim amount of automobile insurance.

#### 4. Generalization of the Model Based on Data

As the automobile claim data is confidential for each insurance company, it is difficult to obtain comprehensive data. The data used in this paper only includes the claim data of various cities in Shandong Province, so the data has certain geographical limitations. Since each place has its own unique characteristics in terms of topography, weather, and economy, the data of each prefecture-level city has its own uniqueness. It can be considered that the data characteristics

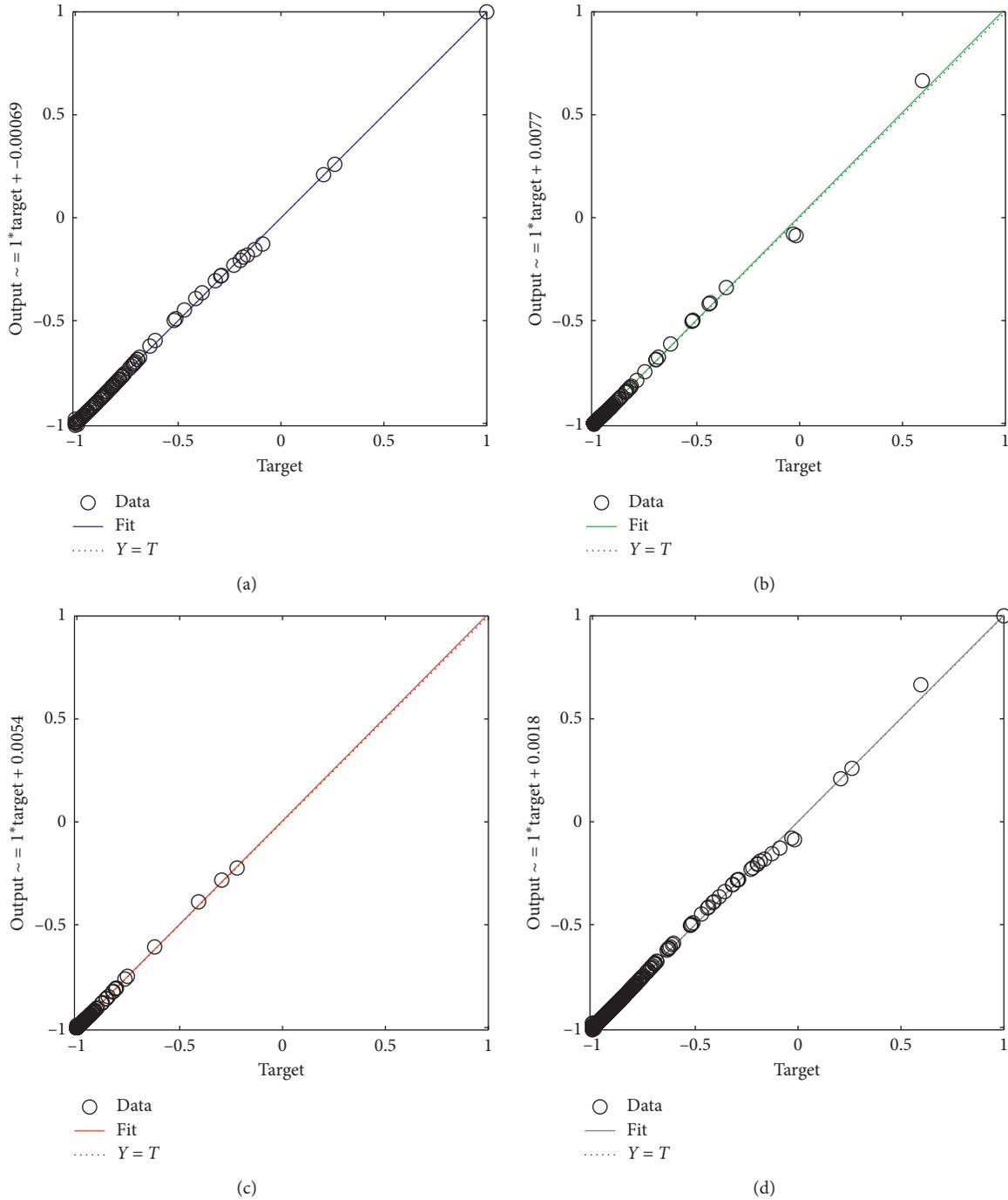


FIGURE 5: Overall goodness of the fit diagram: (a) training:  $R = 0.99972$ ; (b) validation:  $R = 0.99897$ ; (c) test:  $R = 0.99967$ ; (d) all:  $R = 0.99952$ .

of each prefecture-level city are different; hence, the accuracy and adaptability of the model will be tested and verified based on the data at the prefecture-level city level. If it can be verified that the BP neural network model can accurately fit and predict the data of each prefecture-level city, then the model is applicable to data national wide. Because some prefecture-level cities contain relatively few data and are not representative, this paper only uses prefecture-level cities with more than two thousand data to fit and to predict, such

as Binzhou City, Jining City, Weihai City, Weifang City, Jinan City, and Laiwu City.

When fitting and predicting the data of each prefecture-level city, this paper randomly divides the data into a training set and a prediction set by considering that a certain amount of data is required to train the network. The training set accounts for 90%, and the test set accounts for 10%. The target accuracy is set to 0.00001. The following describes the fitting and forecasting for the six prefecture-level cities in detail.

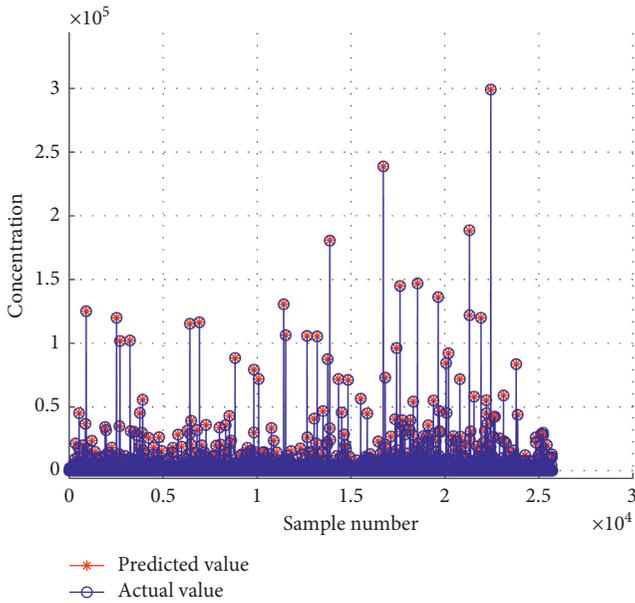


FIGURE 6: The overall modeling training set data fitting effect diagram.

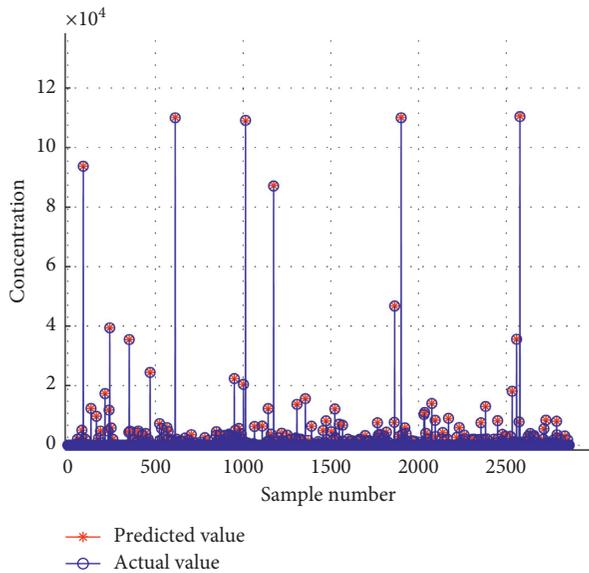


FIGURE 7: The prediction effect diagram of overall modeling test set data.

**4.1. Binzhou City.** Binzhou City has only 2300 pieces of data. By fitting and forecasting the data of Binzhou City, the following training results and prediction results are obtained, as shown in Figures 8 and 9.

Figure 8 shows the effect of fitting the data of the training set from Jinzhou City. Individuals with a claim amount of zero are accurately fitted and stacked at the bottom of the graph, and both small claims and large claims are also accurately fitted. From the above figure, it can be seen that the network accurately fits the Binzhou training set data as a whole.

Figure 9 gives the comparison between predicted values and actual values. It shows that the predicted values of the aggregate claim amount of zero is exactly the same as the true values, and small aggregate claim amounts can almost be predicted correctly. There is only a small error in the predicted value of the high claim data. From the output of the network, the prediction accuracy of the network on the prediction data set of Jinan City is 99.06%, which indicates that the network has made accurate predictions on the prediction data set of Binzhou City.

**4.2. Jining City.** There are only 2130 pieces of data in Jining City. By fitting and forecasting the data of Jining City, the following training result figure and prediction result figure are obtained, as shown in Figures 10 and 11.

Since this data set contains a large number of individuals with zero claims, there is a stacking phenomenon at the bottom of Figure 10, which also shows that the individuals with zero claim amount are accurately fitted. Figure 10 shows that the network accurately fits both small claims and large claims, meaning that the network fits well in general.

Figure 11 gives the comparison between the actual values and predicted values obtained by predicting the input data of the prediction set. It can be seen from Figure 11 that the predicted values of the small aggregate claim amounts are completely consistent with the actual values. The prediction accuracy of the network for the prediction set data is 99.49%, which shows that the network can also be applied to fit Jining automobile insurance claim data and accurately predict the accumulated claim amount.

**4.3. Weihai City.** There are 2020 pieces of data in Weihai City. By fitting and forecasting the data of Weihai City, the following training result figure and prediction result figure are obtained, as shown in Figures 12 and 13.

Figure 12 is a comparison diagram of the fitted values and actual values of the training set after the network is trained through the Weihai city training data set. It can be seen that there are many claims in the training set of this city, and the network accurately fits all the data with claims.

Figure 13 indicates that the network accurately predicts the data in the prediction data set with claims, and the predicted values are consistent with the actual values. The predicted and actual values of the data without claims are stacked at the bottom of the image, and the prediction accuracy of the prediction set of Weihai City by the network is 99.76%; hence, the prediction outcome of the network is very good.

**4.4. Weifang City.** There are 3000 pieces of data in Weifang City. By fitting and forecasting the data of Weifang City, the following training result figure and prediction result figure are obtained, as shown in Figures 14 and 15.

The accumulation phenomenon appears at the bottom of Figure 14, and the trained network accurately fits individuals

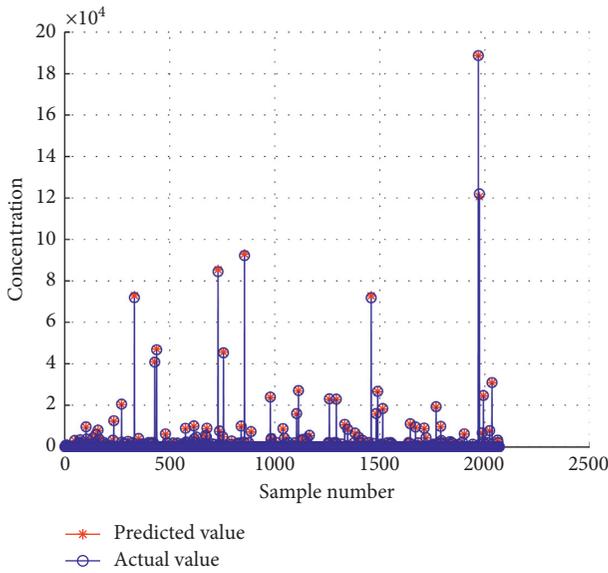


FIGURE 8: The fitting effect diagram of Binzhou training set data.

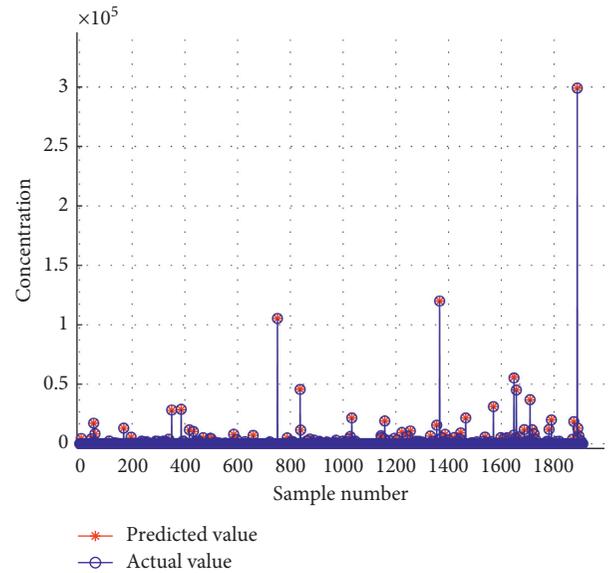


FIGURE 10: The fitting effect diagram of the Jining training set data.

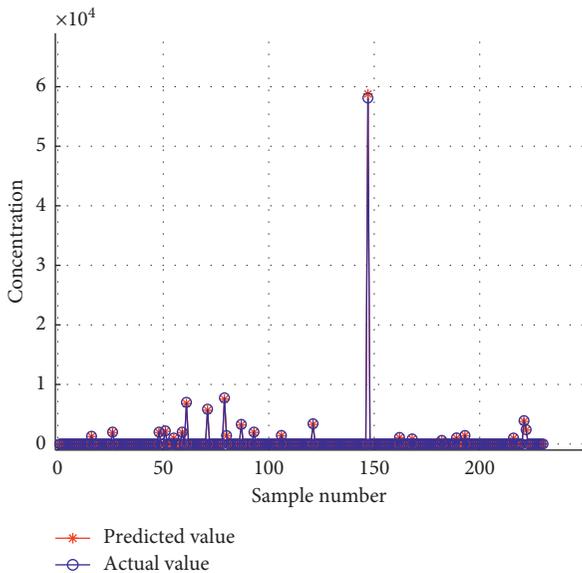


FIGURE 9: The Binzhou test set data prediction effect diagram.

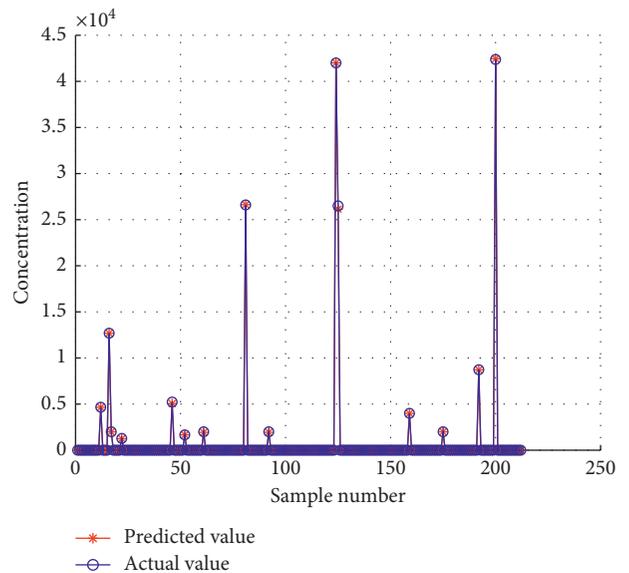


FIGURE 11: The Jining test set data prediction effect diagram.

with zero claims. Most of the data for claims in the training set of this city are small claim amounts and larger claim amounts. The fitting values of the network to claim data are consistent with the actual values, indicating that the network has a good fitting effect.

It can be seen from Figure 15 that claim amounts in Weihai’s prediction set vary. The predicted value of the network for each piece of claim data is consistent with the actual value. Again, the predicted values and actual values of the data for which no claim has occurred are stacked at the bottom of the image. The overall prediction accuracy of the network for the prediction set data is 98.83%, and the prediction performance is good.

4.5. *Jinan City*. There are 4370 pieces of data in Jinan City. By fitting and forecasting the data of Jinan City, the following training result figure and prediction result figure are obtained, as shown in Figures 16 and 17.

Figure 16 shows the comparison between fitted values and actual values of the trained network on the input data of the training set. Since there are more data in this prefecture-level city, the proportion of data without claims is larger, so the lower part of the graph appears to have a heavy accumulation phenomenon, and at the same time, fitted values of the claim data are consistent with actual values, and the network fits well. Besides, there are many claims in this data set, and the network still has a good fit for all claim data.

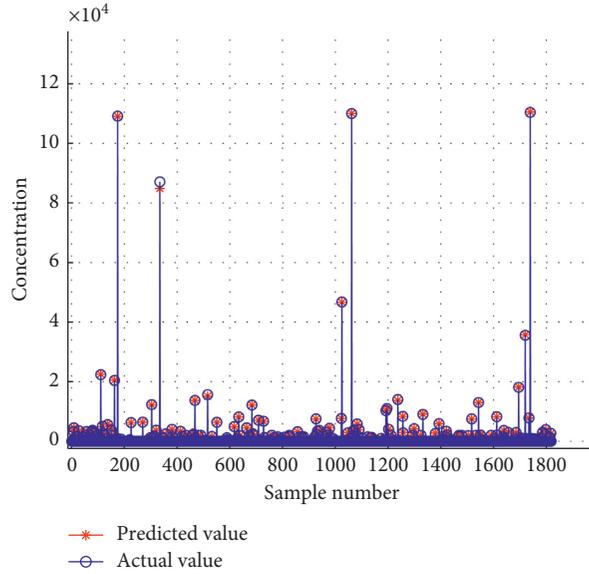


FIGURE 12: The fitting effect diagram of Weiwei training set data.

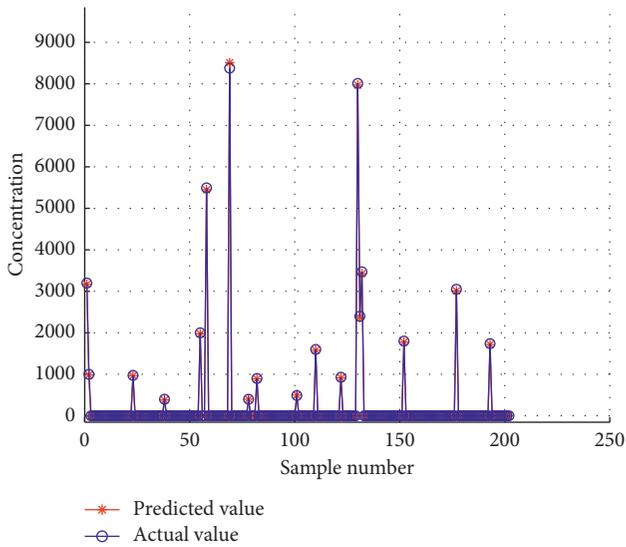


FIGURE 13: The Weiwei test set data prediction effect diagram.

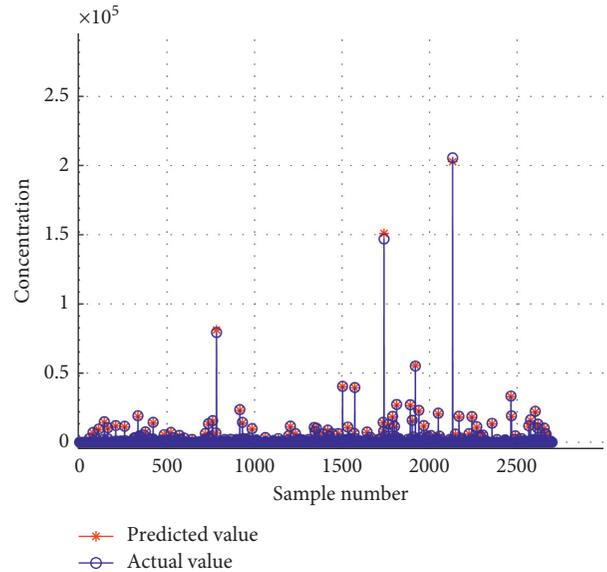


FIGURE 14: The fitting effect diagram of Weifang training set data.

Figure 17 shows the comparison between predicted values and actual values of the input data of the test set by the network. This prefecture-level city has more centralized claim data, and the network has made accurate predictions for most of the claim data. The prediction for zero claims is accumulated at the bottom of the graph. The overall prediction accuracy of the network for the prediction set data is 99.96%, and the prediction effect is very good.

4.6. *Laiwu City*. There are 2600 pieces of data in Laiwu City. By fitting and forecasting the data of Laiwu City, the following training result figure and prediction result figure are obtained, as shown in Figures 18 and 19.

Figure 18 shows the fitting effect of the training set in Laiwu City. The fitted value of the network for each piece of claim data can accurately correspond to the actual value. The predicted and actual values of the zero claim data are stacked at the bottom of the graph, and the network has a good fitting effect.

Figure 19 is the prediction output of the trained BP neural network on the Laiwu City prediction data set. It shows that the network can accurately predict the claim data in the prediction data set. The overall prediction accuracy of the network for this prediction set data is 99.98%, and the overall prediction performance of the network is good.

In summary, the BP neural network model not only can accurately fit and predict the claim data of the entire

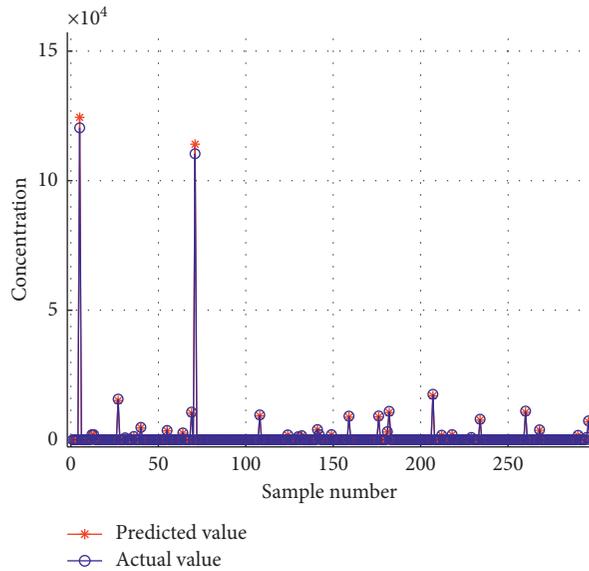


FIGURE 15: The Weifang test set data prediction effect diagram.

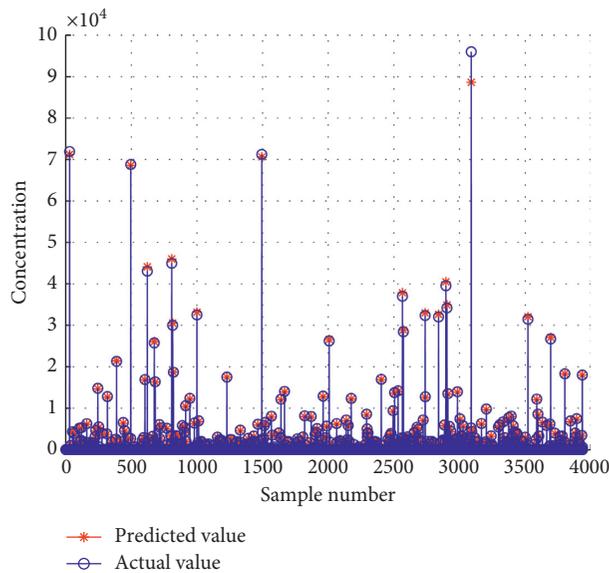


FIGURE 16: The fitting effect diagram of Jinan training set data.

Shandong Province but also can accurately fit and predict the data of six prefecture-level cities. Besides, the trained BP network has a prediction accuracy of over 95% for each prefecture-level city, indicating that it is reasonable and accurate to use the BP neural network model to fit and predict the aggregate claim amount of automobile insurance. Moreover, the BP neural network fits and predicts the data of six prefecture-level cities with different data characteristics well, indicating that the BP neural network can adapt to data with different characteristics in different regions. As the BP

neural network has a strong tolerance, the network can be used to fit and predict data nationwide.

### 5. Premium Ratemaking

5.1. Model Introduction. Credibility theory is the study of how to reasonably use prior information and individual claim experience to estimate, predict, and formulate posterior insurance premiums. The posterior premium estimate is calculated as follows:

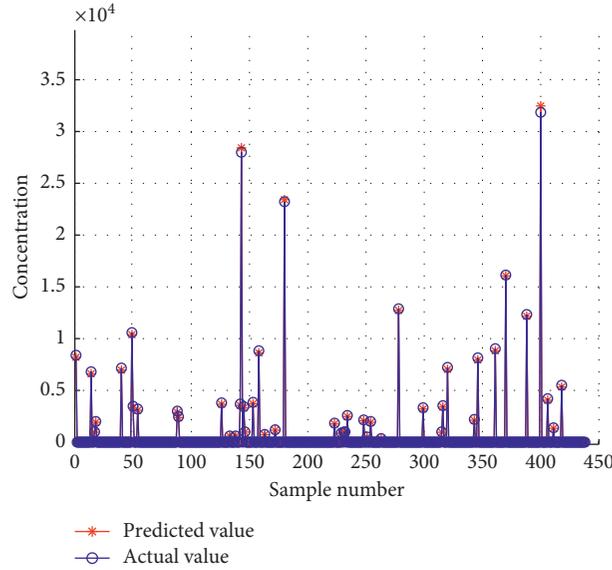


FIGURE 17: The Jinan test set data prediction effect diagram.

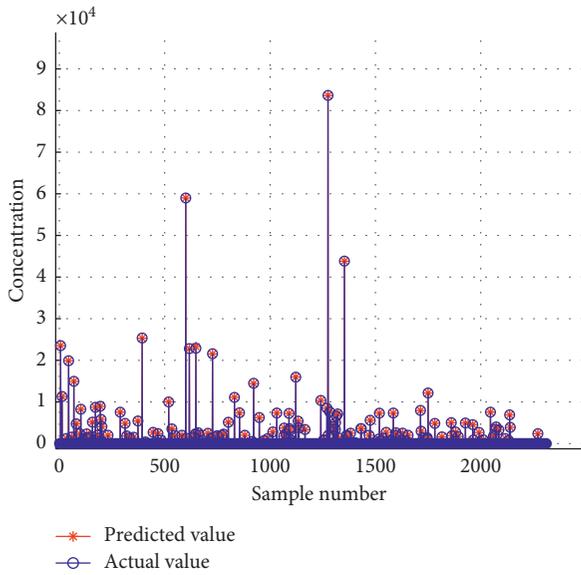


FIGURE 18: The fitting effect diagram of Laiwu training set data.

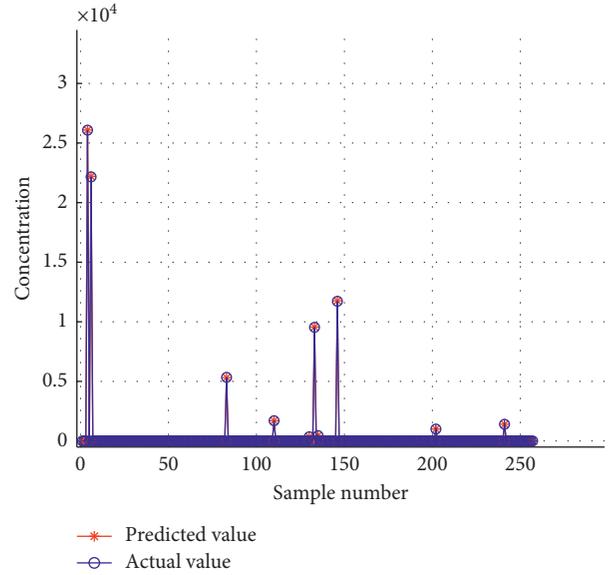


FIGURE 19: The Laiwu test set data prediction effect diagram.

$$\begin{aligned} \text{Estimation for posterior premium} &= Z * \text{Experience value} \\ &+ (1 - Z) * \text{Prior value,} \end{aligned} \tag{9}$$

where  $Z$  ( $0 < Z < 1$ ) is the reliability factor. The posterior premium estimate is called the reliability estimate. Only by choosing the reliability factor correctly, the adjusted insurance premium can then be closed to its actual risk level. In addition to insurance premiums, the credibility theory can also be used to estimate the number of claims, total claim amounts, loss ratio, and relative number of levels. There are two types of credibility models: the classical reliability model and most accurate reliability model.

The classical reliability model attempts to limit the influence of random fluctuations in the observed data on the estimated value, which is also called the limited-fluctuation reliability theory. In the classical reliability model, it is necessary to determine when the individual risk reaches a certain scale, and the reliability factor = 1, that is, the empirical data is fully credible. This scale becomes the “full credibility standard,” and the reliability factor less than 1 is called partial credibility. Suppose the data volume of the individual risk is  $n$ ; then,  $n_f$  is the full credibility standard. If  $n > n_f$ , then the reliability factor  $Z = 1$ . If  $n < n_f$ , then  $Z = \sqrt{n/n_f}$ .

The most accurate reliability model is the so-called least squares reliability model. This model determines the reliability factor by minimizing the sum of squared errors

between the estimated value and actual value and emphasizes the accuracy of the estimated results, mainly including the Buhlmann reliability model and Buhlmann–Straub reliability model. The Buhlmann reliability model assumes that the scale of the individual risk remains the same. If  $n$  represents the empirical period, that is, the number of years of observation of empirical data, the reliability factor  $Z$  is given as  $Z = n/(n + k)$ , and  $k$  is called the Buhlmann parameter, which is the ratio of the mean of the process variance (EPV) to the variance of the hypothetical mean (vhm). In the Buhlmann–Straub reliability model, the scale of the individual risk can be changed. The Buhlmann reliability model is a simplified form of it. In the most accurate reliability model, empirical data and prior data both have significant influence on the reliability factor. For equally important influences, the reliability factor can only be found when empirical data and prior data are determined.

### 5.2. The Buhlmann Model in Nonparametric Estimation.

Since nonparametric estimation does not require the use of overall information (the overall distribution and some parameter characteristics of the population), the distribution type of the population does not need to be assumed, and we can directly perform statistical testing on the distribution of the population, so, in this paper, we use the Buhlmann model from nonparametric estimation to calculate pure premiums.

Since the BP neural network can accurately predict the claim information for each of the insured and in the previous empirical analysis, we have proved that the BP neural network can accurately predict the total claim amount of the insured. Therefore, we use the obtained claim prediction value as empirical data to determine individual premiums. The following is based on the actual claim data and predicted data of the six prefecture-level cities to calculate the premium.

First of all, we use the average actual claim amounts  $\mu$  of the six prefecture-level cities as prior data and the claims of different individuals predicted by the BP neural network  $X_{ij}$  as empirical data, and by applying the formula  $P_i = Z_i \times \bar{X}_i + (1 - Z_i)\mu$ , the average net premiums of the corresponding six prefecture-level cities  $P_i$  are obtained, where  $i$  represents the  $i^{\text{th}}$  prefecture-level city,  $Z_i$  is the reliability factor of the  $i^{\text{th}}$  prefecture-level city,  $\bar{X}_i$  is the average predicted claim amount of the  $i^{\text{th}}$  prefecture-level city,  $X_{ij}$  is the predicted claim amount of the  $j^{\text{th}}$  insured from the  $i^{\text{th}}$  prefecture-level city,  $Y_{ij}$  indicates the actual claim amount of the  $j^{\text{th}}$  insured from the  $i^{\text{th}}$  prefecture-level city, and  $n_i$  indicates the number of people in the training set of the  $i^{\text{th}}$  prefecture-level city.  $i = 1, 2, \dots, 6$ ;  $j = 1, 2, \dots, n_i$ .

The average value of predicted claims for the  $i^{\text{th}}$  prefecture-level city in Shandong Province is calculated as  $\bar{X}_i = (9/n_i) \sum_{j=1}^{(n_i/9)} X_{ij}$ .

The average value of predicted claims of six prefecture-level cities in Shandong Province is  $\bar{X} = (1/6) \sum_{i=1}^6 \bar{X}_i$ .

The actual mean value of claims in the training set of the  $i^{\text{th}}$  prefecture-level city in Shandong Province is  $\mu_i = E(\bar{Y}_i) = E((1/n_i) \sum_{j=1}^{n_i} Y_{ij})$ .

The average value of actual claims in the training set of six prefecture-level cities in Shandong Province is  $\mu = E(\bar{Y}) = (1/6) \sum_{i=1}^6 E(\bar{Y}_i) = (1/6) \sum_{i=1}^6 \mu_i$ .

The variance of the predicted claim amount for the  $i^{\text{th}}$  prefecture-level city is  $v_i = \text{Var}(X_{ij})$ .

The variance of the forecast claim amount in Shandong Province is as follows:

$$v = \frac{1}{6} \sum_{i=1}^6 v_i = \frac{1}{6} \sum_{i=1}^6 \text{Var}(X_{ij}), \quad (10)$$

and we have

$$\text{Var}[\bar{X}_i] = E[\text{Var}(\bar{X}_i|\Theta_i)] + \text{Var}[E(\bar{X}_i|\Theta_i)] = \frac{v}{n} + a, \quad (11)$$

where the estimated value for the structural parameter is  $\mu_i$ , size of risk  $X_i$  is measured by  $\Theta_i$ , and  $v$  and  $a$  are calculated as follows:

- (1)  $\hat{\mu}_i = \bar{Y}_i = (1/n_i) \sum_{j=1}^{n_i} Y_{ij}$  is the estimated mean value of actual claims in the  $i^{\text{th}}$  prefecture-level city
- (2)  $\hat{v}_i = v_i = \text{Var}(X_{ij})$  is the variance of the predicted claim amount for the  $i^{\text{th}}$  prefecture-level city
- (3)  $\hat{v} = (1/6) \sum_{i=1}^6 \hat{v}_i$  is the variance of the forecast claim amount of six prefecture-level cities in Shandong Province

$$\hat{a} = \frac{1}{5} \sum_{i=1}^6 (\bar{X}_i - \bar{X})^2 - \frac{\hat{v}}{n}. \quad (12)$$

The pure premium of each prefecture-level city is calculated as follows:

$$P_i = \hat{Z}_i \bar{X}_i + (1 - \hat{Z}_i)\mu, \quad i = 1, 2, \dots, 6, \quad (13)$$

where

$$\hat{Z}_i = \frac{n_i}{n_i + \hat{k}}, \quad (14)$$

$$\hat{k} = \frac{\hat{v}}{\hat{a}}$$

**5.3. Empirical Analysis of Premium Ratemaking.** Through the Buhlmann model, we find the parameter values for each prefecture-level city and calculate the pure premium and corresponding aggregated pure premium for each prefecture-level city according to the corresponding parameters.

From Table 2, it can be seen that the training of the BP neural network requires a certain amount of data. Generally, the larger the amount of data, the better the performance of the trained network and the higher the accuracy of the predicted data. The  $Z$  value increases with the increase of the amount of data, which means that, as the amount of data increases, the predicted value obtained through the BP neural network weighs more in the determination of the

TABLE 2: The comparison between the predicted pure premium and actual total claim.

	Jinan City	Binzhou City	Jining City	Laiwu City	Weihai City	Weifang City
Z value	0.4627	0.2382	0.2245	0.2589	0.2486	0.2897
Size of prediction set	438	230	213	260	202	300
Average pure premium	521.14	630.85	629.44	602.09	599.8	574.78
Average claim amount in prediction set	434.19	633.87	527.16	531.09	480.37	468.83
Profit percentage	0.2	-0.004	0.19	0.13	0.24	0.22

TABLE 3: Pure premiums for some individuals.

Predicted claim amount $X_{ij}$	Premium calculated $P_{ij}$
0	449.462
0	449.462
6162	1917.2504
0	449.462
0	449.462
3533	1291.0226
0	449.462
0	449.462
1653	843.2066
.....	.....

premium, which makes the premium ratemaking more reasonable for individuals with different risks.

Besides, the average pure premium of each prefecture-level city calculated by the above model can ensure that insurance companies do not lose money in the automobile insurance business. The average net premiums calculated by Jinan and Binzhou are almost the same as the actual average claims, while the average net premiums calculated by Jining, Laiwu, Weihai, and Weifang are slightly higher than the actual average claims. This shows that our premium ratemaking model is reasonable and can be used as a new idea for automobile insurance companies to determine premium rates.

Based on the parameters obtained from the above model of prefecture-level cities, we can find the pure premiums that each individual should pay in the six prefecture-level city prediction sets. The calculation formula is as follows:

$$P_{ij} = \hat{Z}_i \times X_{ij} + (1 - \hat{Z}_i) \bar{X}_i, \quad (15)$$

where  $P_{ij}$  is the pure premium for the  $j^{\text{th}}$  individual in the  $i^{\text{th}}$  prefecture-level city,  $\hat{Z}_i$  is the reliability factor of the  $i^{\text{th}}$  prefecture-level city, and  $X_{ij}$  is the predicted claim amount of the  $j^{\text{th}}$  individual in the  $i^{\text{th}}$  prefecture-level city. Laiwu City has been taken as an example to find the pure premiums for individuals. Since the insured usually takes into account no claims preferential treatment rules in real life, this paper also takes this situation into account, and then, when the individual's predicted claim amount is less than 200, it is recorded as 0. The calculated results are shown in Table 3.

Table 3 shows the premium calculated for selected individuals with different risks. From Table 3, we can see that, for individuals with lower predicted claims, the calculated pure premiums are relatively low, while for individuals with higher predicted claims, the premiums are relatively high. This model can effectively identify risks and helps to charge the insured with premiums that are compatible with their

risks, making premium ratemaking more reasonable and fair.

## 6. Conclusion

A brand new method and idea to price for automobile insurance is attempted in this paper. Different from the traditional model, we did not separately model the number of claims and the individual claim amount, but directly modeled the aggregate claim amount. This idea is simple and straightforward, and the results are clear. Based on the BP neural network model in Matlab, this paper fits and predicts 29,000 valid data in Shandong Province, which shows that the BP neural network can predict the aggregate amount of automobile insurance claims very accurately. In addition, in order to break the regional limitations of the data, we fitted and predicted the data of each prefecture-level city separately with the prefecture-level city as a unit. It shows that the network is very inclusive of data and can break geographical restrictions; hence, it can be applied to nationwide data.

Given that the fitted and predicted aggregate claim amounts are accurate, this paper uses the credibility theory to calculate the average pure premiums for six prefecture-level cities. By using the average aggregate claim amount of the training set data of the entire Shandong Province, the average pure premium of each prefecture-level city is adjusted so that the overall automobile insurance compensation situation of Shandong Province can be considered, and the individual compensation situation of each prefecture-level city can be highlighted. The result shows that the pure premium calculated by this method can improve the profitability of the automobile insurance business.

Taking into account the different risks and aggregate claim amount of each of the insured, this paper aims to find appropriate net premiums for each of the insured corresponding to their risks. Using the reliability factor of individual's location-level city and the average claim amount of the prefecture-level city to calculate each individual's pure premium and trying to personalize the premium rate, through this calculation method, individual risks can be identified to a certain extent, and the premium rates can be differentiated to make the determination of car insurance rates more reasonable and fair.

The BP neural network is introduced and applied to the field of automobile insurance, which has the important application value for pricing of automobile insurance rates. The combination of the BP neural network and Buhlmann model can calculate the pure premium that matches the liability of the insurance company, which can effectively

improve the current situation of the loss of the insurance company's automobile insurance business, thereby boosting the enthusiasm and creativity of the insurance company to develop more and better products to promote the soundness of the country's automobile insurance industry.

To avoid the shortcomings of the BP neural network such as slow convergence speed and easy to fall into a local minimum, the genetic algorithm is selected to optimize the BP neural network in this paper. The genetic algorithm has a strong adaptive and optimal ability, which can effectively improve the convergence speed of the BP neural network and prevent the network from falling into a local minimum. At the same time, by considering the overfitting problem of the BP neural network, this paper chooses to use the early stop method to make the network error continue to decrease and stabilize, effectively avoiding the overfitting problem.

### Data Availability

The data used to support the findings of this work are available from the corresponding author upon request.

### Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

### Acknowledgments

This work was supported by the National Natural Science Foundation of China (no. 11301303), National Social Science Foundation of China (no. 15BJY007), Taishan Scholars Program of Shandong Province (no. tsqn20161041), Humanities and Social Sciences Project of the Ministry Education of China (no. 19YJA910002), Natural Science Foundation of Shandong Province (no. ZR2018MG002), Fostering Project of Dominant Discipline and Talent Team of Shandong Province Higher Education Institutions (no. 1716009), Shandong Provincial Social Science Project Planning Research Project (no. 19CQXJ08), Risk Management and Insurance Research Team of Shandong University of Finance and Economics, Excellent Talents Project of Shandong University of Finance and Economics, Collaborative Innovation Center Project of the Transformation of New and Old Kinetic Energy and Government Financial Allocation, Shenzhen Peacock Program (no. 000417), Shandong Jiaotong University "Climbing" Research Innovation Team Program, and 1251 Talent Cultivation Project of Shandong Jiaotong University.

### References

- [1] R. A. Bailey and L. J. Simon, "Two studies in automobile insurance ratemaking," *ASTIN Bulletin*, vol. 1, no. 4, pp. 192–217, 1960.
- [2] D. Denneberg, "Premium calculation: why standard deviation should be replaced by absolute deviation," *ASTIN Bulletin*, vol. 20, no. 2, pp. 181–190, 1990.
- [3] J. A. Nelder and R. W. M. Wedderburn, "Generalized linear models," *Journal of the Royal Statistical Society. Series A (General)*, vol. 135, no. 3, pp. 370–384, 1972.
- [4] D. Samson and H. Thomas, "Linear models as aids in insurance decision making: the estimation of automobile Insurance Claims," *Journal of Business Research*, vol. 15, no. 3, pp. 247–256, 1987.
- [5] G. K. Smyth, "Generalized linear models with varying dispersion," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 51, no. 1, pp. 47–60, 1989.
- [6] K. J. Stroinski and I. D. Currie, "Selection of variables for automobile insurance rating," *Insurance Mathematics and Economics*, vol. 8, no. 1, pp. 35–46, 1989.
- [7] S. W. Meng, "An application of generalized linear model to automotor insurance pricing," *Application of Statistics and Management*, vol. 26, no. 1, pp. 24–29, 2007.
- [8] N. R. Draper, "Generalized linear models for insurance data by Piet de Jong, Gillian Z. Heller," *International Statistical Review*, vol. 76, no. 2, p. 315, 2008.
- [9] M. Q. Zhao and Y. P. Chen, "The auto insurance ratemaking based on double generalized liner models and the comparison with generalized linear models," *Insurance Studies*, vol. 10, pp. 32–41, 2016.
- [10] G. K. Smyth and B. Jørgensen, "Fitting tweedie's compound Poisson model to insurance claims data: dispersion modelling," *ASTIN Bulletin*, vol. 32, no. 1, pp. 143–157, 2002.
- [11] K. Antonio and J. Beirlant, "Actuarial statistics with generalized linear mixed models," *Insurance: Mathematics and Economics*, vol. 40, no. 1, pp. 58–76, 2007.
- [12] E. W. Frees, P. Shi, and E. A. Valdez, "Actuarial applications of a hierarchical insurance claims model," *ASTIN Bulletin*, vol. 39, no. 1, pp. 165–197, 2009.
- [13] X. H. Wang, S. W. Meng, and Y. S. Wang, "Automobile insurance pricing models based on heavy-tailed loss distribution and its application," *Insurance Studies*, no. 4, pp. 67–78, 2017.
- [14] L. Z. Zhang and H. Y. Xie, "The application of Tweedie distribution in auto insurance ratemaking," *Insurance Studies*, no. 1, pp. 80–90, 2017.
- [15] M. Ayuso, M. Guillén, and J. P. Nielsen, "Improving automobile insurance ratemaking using telematics: incorporating mileage and driver behaviour data," *Transportation*, vol. 46, no. 3, pp. 735–752, 2019.
- [16] Y. F. Huang and S. W. Meng, "Automobile insurance classification ratemaking based on telematics driving data," *Decision Support Systems*, vol. 127, Article ID 113156, 2019.
- [17] P. L. Brockett, X. Xia, and R. A. Derrig, "Using kohonen's self-organizing feature map to uncover automobile bodily injury claims fraud," *The Journal of Risk and Insurance*, vol. 65, no. 2, pp. 245–274, 1998.
- [18] Y. Liu, B. J. Wang, and S. G. Lv, "Using multi-class adaboost tree for prediction frequency of auto insurance," *Journal of Applied Finance and Banking*, vol. 4, no. 5, pp. 45–53, 2014.
- [19] T. Mzhavia, *Vehicle Insurance Claim Data Study and Forecasting Model Using Artificial Neural Networks*, Tallinn University of Technology, Tallinn, Estonia, 2016.
- [20] M. V. Wüthrich, "Covariate selection from telematics car driving data," *European Actuarial Journal*, vol. 7, no. 1, pp. 89–108, 2017.
- [21] L. Z. Zhang and D. Wang, "A research on the rate making of automobile insurance with big data—modeling of automobile insurance claim severity based on SOM neural network," *Insurance Studies*, no. 9, pp. 56–65, 2018.

- [22] Y. C. Lin, J. Li, M.-S. Chen, Y.-X. Liu, and Y.-J. Liang, "A deep belief network to predict the hot deformation behavior of a ni-based superalloy," *Neural Computing and Applications*, vol. 29, no. 11, pp. 1015–1023, 2016.
- [23] Y. C. Lin, Y. J. Liang, M. S. Chen, and X. M. Chen, "A comparative study on phenomenon and deep belief network models for hot deformation behavior of an Al-Zn-Mg-Cu Alloy," *Applied Physics A*, vol. 123, p. 68, 2016.
- [24] Y. C. Lin, J. Huang, H.-B. Li, and D.-D. Chen, "Phase transformation and constitutive models of a hot compressed TC18 titanium alloy in the  $\alpha + \beta$  regime," *Vacuum*, vol. 157, pp. 83–91, 2018.
- [25] D.-D. Chen, Y. C. Lin, and F. Wu, "A design framework for optimizing forming processing parameters based on matrix cellular automaton and neural network-based model predictive control methods," *Applied Mathematical Modelling*, vol. 76, pp. 918–937, 2019.