

CAR INSURANCE PLAN

Machine Learning I Project



Lianjie Shan

1. Introduction and description

In our research, I finished the regression section using MLP method. The final predictive model was built through my work.

For a brief introduction of background of MLP, in the part of model selection, I hope to train a neural network to achieve our goal because of the complexity of data and the relatively vague correlation between variables. MLP has a high degree of parallel processing, a high degree of nonlinear global function, good fault tolerance, associative memory function, very strong adaptive, self-learning function, so we finally decided to use MLP multilayer perceptron.

2. Individual work and Results

2.1. Regression (MLP) and Results

After separating potential claimants with SVM, we can predict accurately whether insurance customers will pay claims in the future, but we're not satisfied. We hope to get more concrete results. Therefore, our goal is to find a model that is able to predict the specific amount of customer claims.

2.1.1. Data preparation

In this part, we use the new dataset ‘CLM1value.csv’ that we split from original data. It contains 2,413 observations with non-null values and 5 customer characteristics - BLUEBOOK, TRAVTIME, INCOME, MVR_PTS and AGE as variables, CLM_AMT with specific integers as our target.

	BLUEBOOK	TRAVTIME	INCOME	MVR_PTS	AGE	CLM_AMT
6	17430.0	46.0	125301.0	0.0	34.0	2946.0
8	18930.0	21.0	50815.0	2.0	40.0	6477.0
10	16970.0	44.0	107961.0	10.0	37.0	4021.0
11	11200.0	34.0	62978.0	0.0	34.0	2501.0
13	18300.0	15.0	77100.0	0.0	53.0	6077.0

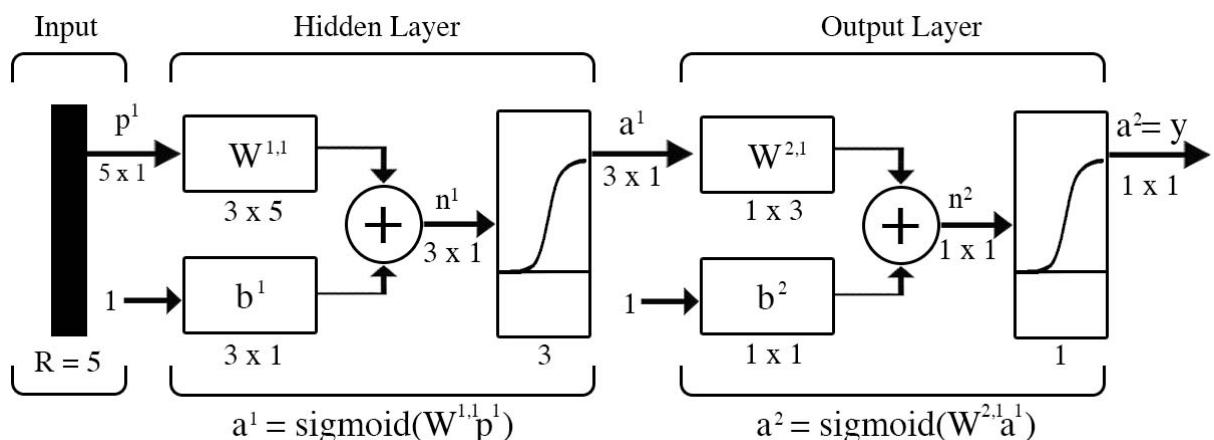
In practice, we stacked by encountered overflow problem. Obviously, a wide range of data makes MLP troubles to operate, so we standardized the data.

	BLUEBOOK	TRAVTIME	INCOME	MVR PTS	AGE	CLM AMT
6	0.074086	0.467391	0.000000	0.230769	0.483333	0.291478
8	0.118703	0.260870	0.137686	0.000000	0.200000	0.688553
10	0.161179	0.119565	0.118515	0.000000	0.733333	0.526172
11	0.140435	0.228261	0.097696	0.000000	0.266667	0.449420
13	0.176655	0.163043	0.053070	1.000000	0.400000	0.612204

Similar to SVM, remove missing values. We respectively divided data into training data (70%) and test data (30%) after standardization and removed the missing values. Because there were not too many observations in the original data and some missing values appear in the standardization process, only 403 values of training data can be obtained after these operations, while only 174 values of test data can be obtained.

2.1.2 Network Architecture

For this research, we used the multilayer perceptron (MLP) method to train the dataset and to build the predictive model. Our MLP model included 2 layers which can also call the hidden layer and output layer. To be more specific, the log sigmoid function was used in both of two layers. At the same time, we selected 5 features as our input data. The plot of network architecture is as follows.



$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

As we can see from the plot above, when we have training dataset of inputs (feature), we can obtain the outputs after two layers of training.

2.1.3 Process of Building Model and Results

As we have done in the data preparation section, we split data into the training dataset and test dataset. And training dataset was used to produce a predictive model after MLP training. For the initialization of MLP training, we set the number of neurons is equal to 5 and the times of epochs are equal to 1000. After initialization, we put our inputs (training dataset) in the MLP model to obtain the weight matrix to get the predictive model. And the weight vector we obtained from the MLP is as following (neurons = 5, epochs = 1000)

```

Stage 1) Random starting synaptic weights:
    Layer 1 (3 neurons, each with 5 inputs):
[[ -0.16595599  0.44064899 -0.99977125 -0.39533485 -0.70648822]
 [-0.81532281 -0.62747958 -0.30887855 -0.20646505  0.07763347]
 [-0.16161097  0.370439  -0.5910955   0.75623487 -0.94522481]
 [ 0.34093502 -0.1653904   0.11737966 -0.71922612 -0.60379702]
 [ 0.60148914  0.93652315 -0.37315164  0.38464523  0.7527783 ]]

    Layer 2 (1 neuron, with 4 inputs):
[[ 0.78921333]
 [-0.82991158]
 [-0.92189043]
 [-0.66033916]
 [ 0.75628501]]

Stage 2) New synaptic weights after training:
    Layer 1 (3 neurons, each with 5 inputs):
[[ -3.14944429 -0.10220113 -1.53256736 -1.85949674 -3.23175782]
 [-5.46013835 -2.02573787 -0.84631291 -2.92134375 -4.29557013]
 [-2.81680502  0.55810412  1.90918765 -0.39694533 -3.16541547]
 [-2.57080477 -4.80057779 -1.66770337 -4.49542986 -3.80295988]
 [-5.62322547 -7.56311993 -1.47108119 -6.35276698 -5.70353273]]

    Layer 2 (1 neuron, with 4 inputs):
[[ -2.02821766]
 [ 4.03192123]
 [-2.65357125]
 [ 1.79833417]
 [-0.83030681]]

```

According to the weight vector, we can get the predive model with the assistance of the formula:

$$n = f(W.P + b)$$

In this equation, n is neural network output; f is transformed function; W is weight vector; P is neural network input; b is biased in the network; “.” means dot product. Therefore, based on the formula above, we can obtain the predictive function:

$$N = \text{sigmoid}(W.P)$$

However, according to the output value of our model, we cannot intuitively see the specific value of the claim that the customers may make in the future because of the standardized data.

To solve this problem, we apply the inverse function in Python to convert the data into its original form. In this way, we can see the specific value of the claim. The final results of predictions are as followings.

Inverse Prediction Values from Standardized to Real

[0.3673214]	[3671.99163825]
[0.42892033]	[4282.74509204]
[0.38475559]	[3844.85163628]
[0.37078258]	[3706.30930536]
[0.32756746]	[3277.8313754]
[0.38619124]	[3859.08617053]
[0.37190939]	[3717.48156009]
[0.44247311]	[4417.12087744]
[0.39065171]	[3903.31168905]
[0.37662705]	→ [3764.25721683]
[0.37076745]	[3706.15926393]
[0.37946452]	[3792.39073447]
[0.3826891]	[3824.36247572]
[0.36139664]	[3613.24770296]
[0.39366315]	[3933.17009324]
[0.37221891]	[3720.55045839]
[0.32236985]	[3226.29706708]
[0.31275087]	[3130.92492173]
[0.43665823]	[4359.46630453]
.	.
.	.
.	.

2.1.4 Evaluation and Improvement of MLP Model

Introduction of the Mean Square Error

The mean square error(MSE), as a popular method of calculating errors, is often used as a standard to evaluate a model. The following is the MSE equation.

$$MSE = \frac{1}{N} \sum_{i=1}^N (e_i)^2 = \frac{1}{N} \sum_{i=1}^N (t_i - a_i)^2$$

In this function, N is equal to the number of the inputs; e is the errors; t is the target; a is the network output; i is equal to ith input of the dataset.

To define the network performance, MLP is supervised learning which indicates the targets and errors are included in this model. Thus we can use MSE to evaluate the network quality.

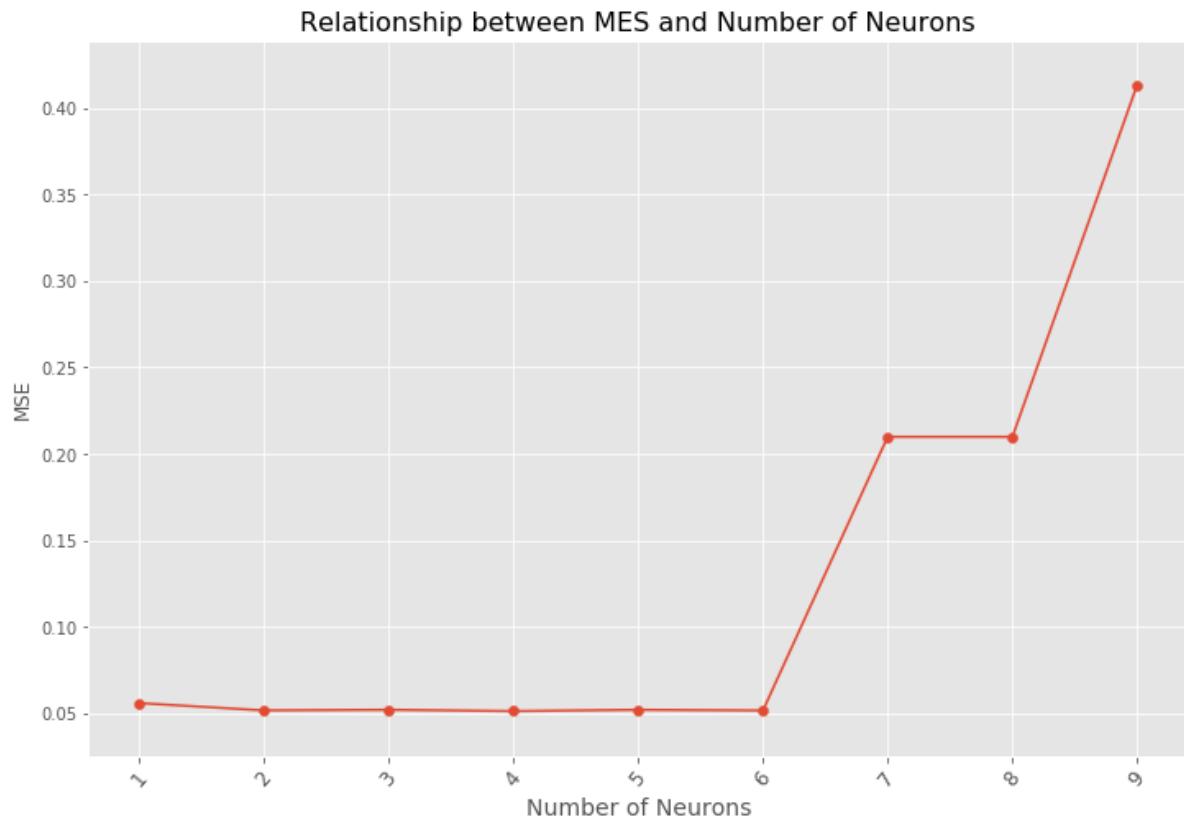
Evaluating the Model with MSE

In order to evaluate the predictive model, we use the test dataset which we have split from the original dataset to calculate the value of MSE. And the value of MSE is equal to 0.3 indicating that our predictive model is acceptable. Thus, our model can be used in the real-life that to assist the insurance company to predict the value of the claim the customers may produce in the future.

Improving the Accurate of the Model

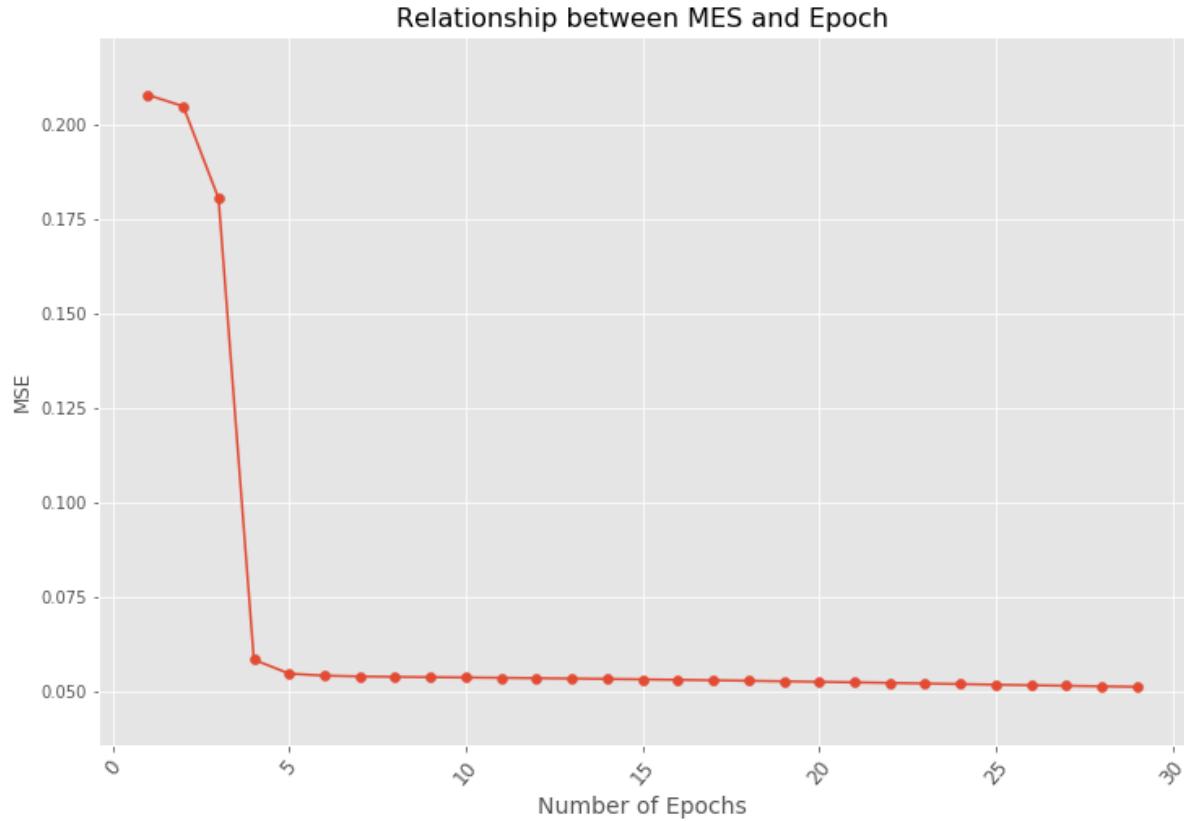
After building the predictive model, we want to adjust the number of neurons and epochs to improve the accuracy of the model. And we can adjust the MSE value to the minimum through the control variable method to enhance the quality of the model.

For the neurons section, we fix the value of epochs is equal to 100. And the relationship between the MSE and neurons is as following.



From this plot, when the value of neurons is greater than 6, an obvious increase in the value of MSE indicates the problem of overfitting. Thus, we selected the number of neurons to minimize the value of MSE. And the ideal number of neurons is 3.

For the epochs section, we fix the value of neurons is equal to 3. And the relationship between the MSE and epochs is as following.



From this plot, as the epoch was greater than 4, the MSE value became minimal and stable. Thus, the optimal value of epoch is equal to 4.

When we established the best values of epoch and neurons, we can re-calculate the MSE values of the model. After calculation, we get that the value of MSE is about 0.03 which is a significant decline compared with before. Thus, the accuracy of the model has been improved in this way.

3. Summary and Improvement

3.1. Summary

For this study, we used SVM method to make a classification which would able to divide customers by whether they will claim after buying car

insurance. At the same time, we applied MLP method to build predictive model in order to predict specific claim amount. Consequently, both MLP and SVM provided satisfying results indicating the feasibility of this research method.

3.2. Improvement of the Model

Due to the huge and complex data volume of users in real life, the accuracy of our prediction model will be reduced to some extent. In this case, we would use K-means algorithm to classify our original data and use the MLP method for each subclass to make regression.

3.2.1 Introduction of K-means Algorithm

Clustering is a kind of unsupervised learning, that is, grouping similar data into the same cluster. The more similar the objects in the cluster are, the better the clustering effect is. In this case, the K-means algorithm is a method of dividing data into K classes without any supervisory signals. And the general steps of the algorithm are as follows:

Step1: Choose K center points at random

Step2: Assign each data point to its nearest central point;

Step3: The average distance from the point in each class to the center of the class is recalculated

Step4: Assign each data to its nearest central point;

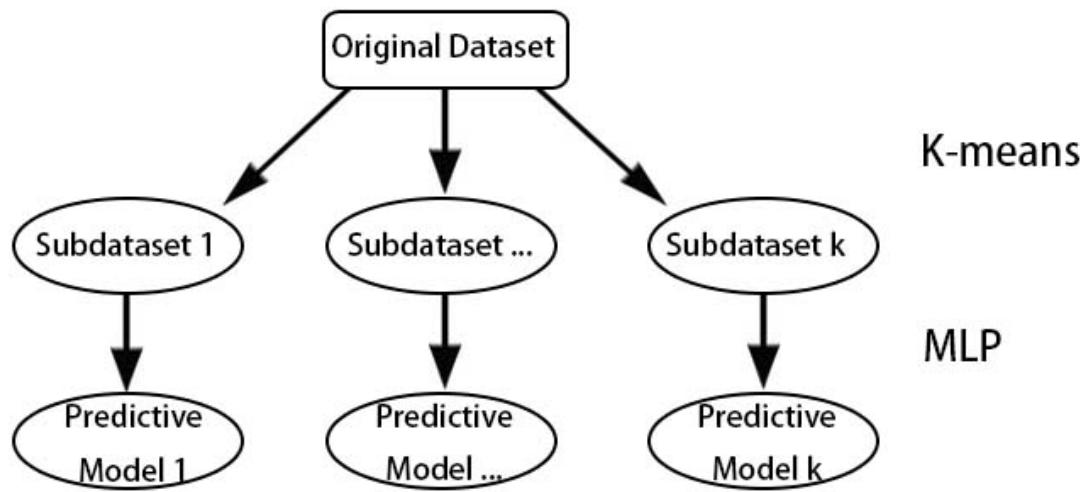
Step5: Repeat steps 3 and 4 until all observations are no longer assigned or the maximum number of iterations is reached.

3.2.2 Application of K-means Algorithm

In our research, we want to use the k-means algorithm to find an optimal classification group number, that is to say, the classification group number that can make the value of MSE become the smallest. In this case, we are breaking down the original data into k classes, and within each of those classes we will re-using MLP to build a predictive model. That brings our total number of predictive models to K.

When we use the established K models to make predictions, we first find the cluster to which the input (customer) belongs, and then use the model to predict the value of claim.

With the assistance of this method, we can minimize the value of MSE and improve the accuracy of the model.



In our research, we did not use K-means algorithm to improve the accuracy of our model because the amount of data was not large, and the final model evaluation was acceptable. However, we believe our approach is feasible. In the future research, when we face with a larger and more complex dataset, we will use K-means algorithm to verify the feasibility of this method.

4. Code Calculation

Copy: 77 Modified:15 Own:53, Result = 47.69%

5. References

Sato, Kaz. “*Using Machine Learning for Insurance Pricing Optimization | Google Cloud Blog.*” *Google*, Google Cloud Platform, 19 Mar. 2017, cloud.google.com/blog/big-data/2017/03/using-machine-learning-for-insurance-pricing-optimization.

Malhotra, Ravi, and Swati Sharma. *MACHINE LEARNING IN INSURANCE - Accenture.com*. Accenture, 2018, www.accenture.com/t20180822T093440Z_w_en/_acnmedia/PDF-84/Accenture-Machine-Leaning-Insurance.pdf.

Raschka, Sebastian, and Vahid Mirjalili. *Python Machine Learning: Machine Learning and Deep Learning with Python, Scikit-Learn, and TensorFlow*. Pack Publishing, 2018.

6. Appendix

Columns Explanation

#	Columns	Description	Data type
1	ID	Customers' ID Number	int
2	KIDSDRIV	The number of kids the customers need to drive for.	int
3	BIRTH	Date of birth.	Factor
4	AGE	Age	int
5	HOMEKIDS	The number of kids the customers have.	int
6	YOJ	Working age.	int
7	INCOME	Income.	int
8	PARENT1	Whether the customers' parents pay for their insurance.	Factor
9	HOME_VAL	Property value.	Factor
10	MSTATUS	Marriage status.	Factor
11	GENDER	Gender.	Factor
12	EDUCATION	Education level.	Factor
13	OCCUPATION	Occupation.	Factor
14	TRAVTIME	Travel times.	int
15	CAR_USE	Whether customers' cars are used for private or commercial.	Factor
16	BLUEBOOK	Blue Book is a guidebook that compiles and quotes prices for new and used automobiles and other vehicles of all makes, models and types.	Factor
17	TIF	The insured vehicle.	int
18	CAR_TYPE	Car Type.	Factor
19	RED_CAR	Red cars are more expensive to insure.	Factor
20	OLDCALIM	Cumulative claim amount.	Factor
21	CLM_FREQ	Claim frequency.	int
22	REVOKED	Revoked time.	Factor

23	MVR PTS	Motor vehicle record points. The Violation/Accident Guidelines and Points Columns on the right are used to assign points to each accident or violation over a three-year period.	int
24	CLM AMT	Claim amount.	Factor
25	CAR AGE	Car age.	int
26	CLAIM FLAG	Claim flag.	int
27	URBANICITY	Car insurance varies according to the area of activity.	Factor