

CSC 74020 Machine Learning Midterm project

Bayes Rule Classifier

xinyue.chang

October 2020

1 Abstract

This model is created for the midterm project for machine learning class at CUNY graduate center, fall semester, 2020. The purpose of this model is to explore the implementation of Bayes rules in machine learning. This model is mainly based on Bayes rule to train the dataset and consider the expected gain at the same time while building up the decision rule. The program language used for coding this model is python 3.0. The data in the dataset used in this model is followed by the cumulative probability distribution created randomly by the dataset generator built in the model. There are totally 5 dimensions, and 3 different classes in the dataset and two possible values for each dimension. Besides creating the decision rule, the model also performed a cross validation and optimization, and both helped to improve the performance of the model and produce a better result.

2 Introduction

The Bayes rules are widely used in calculating joint probabilities and conditional probabilities. And it is a very basic and important concept in statistical and data analyst machine learning. It is very common to use Bayes rule in both regressions and classification to calculate all the different probabilities and use for training the dataset to generate the information needed for the decision rule. With enough training data, using Bayes' Theorem can be used to build a very powerful model for classifiers and prediction for the unknown.

I am building this classifier for my midterm project and I think this whole process really helps me have a deeper understanding about the basic concept of Bayes rules and also applied the theory into building up the model to solving the real problem. The main goal of the classifier is to create a decision rule that can achieve the total economic expected gains as higher as possible. The input for this model are a dataset which created by the dataset generator; a 3*3 economic gain matrix, because there are a total 3 classes in the dataset; the prior probabilities for each class, and the data set $((x_1, \dots, x_Z), (c_1, \dots, c_Z))$

i)). And the output from this model are the Discrete Bayes Rule $fd(C)$, a 3 by 3 confusion matrix, as well as the expected gain. After the model was built, I also did data shuffling, cross valuations and optimization to explore the data and to discover what effects can improve the performance of the model to achieve a better result. By shuffling the data set and doing cross valuations, allows the model to decrease the noise created by the order of dataset, and exposed the learning process to different part of the dataset. By doing the optimization, it helps the model to modify and adjust the conditional probability by a small adjustment, so that increases the performance of the model. I will elaborate more details about my project, and the problem that I had during the process in the following paper.

3 Technical Process

3.1 The data set generator

The first step for this project is to create the data set. To do so, I created a one-dimensional sample data generator. This generator first will create a cumulative probability distribution based on how many possible values are in the sample, and then generate a certain size of sample data from the cumulative probability distribution that was created before.

For example, in my model, there are 2 possible values for X1 dimension, the generator first creates 2 uniform[0,1] pseudo random numbers, and by normalizing this set of numbers, it creates the Cumulative Probability Distribution for X1. I set up the sample size for X1 100,000, then the generator randomly creates a list of 100,000 number of uniform [0, 1] pseudo random number, and for each number in this list, if the its fall into the first part of the probability distribution, it would be given value 1, otherwise, it would be given value 2. I used this 1-D sample generator to create a total 5 lists of one-D data lists, each dimension has 2 possible values, and combined them together as my 5-dimension data set. For the data set class tag, I used the same generator to create a list of class tags which has 3 values with the size same as the data set I created in the last step. And by applying this class tag list to the 5-dimensional data, the data set generate process is completed.

The second step is to separate the training data and testing data, in my model I used the first 75% of the data I created as my training data set and the rest 25% as my testing data set.

3.2 The data training process

The calculation in this training process is based on the Bayes Theorem equation:

$$P(d|Ci) = P(Ci, d)/P(Ci), \quad (1)$$

To assign any C_k to d such that the $\sum_{j=1}^k e(C_j, C_K)PT(C_j, d)$ is maximized, it has to find out the largest joint probability $Pr(Ci, d)$ for each d , and in order to

find the joint probability $\Pr(C_i, d)$, we need to calculate the prior probabilities, $\Pr(C_i)$, and the conditional probabilities, $\Pr(d | C_i)$ from the training data set. From the data set that I created in the last step, it is easy to calculate the prior probabilities: First, separates the whole data set into 3 subsets based on the classes, then the prior probabilities, $\Pr(C_i)$, for each class are calculated by the length of each subsets divided by the total length of the training data set. The conditional probabilities are the probabilities of each tuple in the given class. In each class's subset, I counted the frequency of each d in each class, and divided by the total class length, by doing so, I have all the conditional probabilities for each d in each of the classes. I make my data set to a DataFrame and use the `groupby()` function in pandas to separate the class, and `collections` function from `counter` to organize the summarize all the d in each class.

3.3 Bayes decision rule

The Bayes decision rule is simply to find the assigned class C_k that maximize the expected economic gain, since I set the economic gain value identical for all the right assignments in my model, the model only needs to find out in which class that measurement d in the testing set has the largest joint probability, and then assigns the class tag to that measurement d . The model will separate the testing data set into 2 lists, one is the measurement data set list, the other one is the true class tag list. Then the model will go through each data point and find the same data point in the 3 different data subsets to find the conditional probabilities, and the multiple by each class's prior probabilities, and the true assignment economic gain value. After this calculation, the model will pick the class which has the largest product and assign that class to the measurement d . After the model go through all the testing measurements, it created a new list with all assigned classes for each d in the testing data set.

Since the 3 classes data subsets are all data frame, I use the `loc()` function in panda to locate the data point location and find the conditional probability for each d in each class.

3.4 Confusion matrix

The confusion matrix shows how the decision rule performed. the model combines the true class list and assigned class list that decision rule created in last step together, and create a new list of class tag combination for testing set, and then the model counts the frequency for each combination and divided it with totally length of the list to find out the probabilities for each combination to conduct the confusion matrix.

The sum of the diagonal of the confusion matrix is the identification accuracy rate, which indicates how accurate the decision rule is.

3.5 Economic gain

The expected economic gain is the value of the possible gain for all situations. and the confusion matrix shows all the probabilities for each combination, so we take the probabilities from the confusion matrix and multiply by the economic gain under that combination of true class and assign the class. The expected economic gain is 0.9082 at this point.

4 Model improvement

4.1 data shuffling

In the real world, sometimes the given data has some certain order, for this kind of data order can cause bias since the training data set and testing data set are not evenly mixed so the training data set and testing data set might both contain some special feature that cause the training process to be less efficient and decision rule less accurate. So shuffling data can help to evenly distribute the data set and increase the possibility of training the data set including as much information as possible. Even the data set in my model was generated randomly, after shuffling, and repeating the training process, it still shows an improvement of the model, and economic gain increased from original 0.9082 to 0.9126.

I used the `DataFram.sample(frac=1)` function to return a random sample of items from an axis of object to shuffle the entire data set.

4.2 Cross Validation

Cross validation allows us to use all the data for training at least once, the model has 5 folders, and chooses one each time as the testing set and rest of them as the training set. the results show when using the second and third folder as testing set and rest of the data as training set, the model performs the best, which increase the economic gain to 0.99.(figure 1)

```

[[0.463, 0.0, 0.0105], [0.1545, 0.0, 0.0025], [0.365, 0.0, 0.0045]]
in fold 0 economic gain is: 0.935
p(c1)= 0.489625
p(c2)= 0.161
p(c3)= 0.349375
[[0.495, 0.0, 0.0], [0.159, 0.0, 0.0], [0.346, 0.0, 0.0]]
in fold 1 economic gain is: 0.99
p(c1)= 0.4885
p(c2)= 0.159625
p(c3)= 0.351875
[[0.4995, 0.0, 0.0], [0.1645, 0.0, 0.0], [0.336, 0.0, 0.0]]
in fold 2 economic gain is: 0.999
p(c1)= 0.488625
p(c2)= 0.162375
p(c3)= 0.349
[[0.499, 0.0, 0.0], [0.1535, 0.0, 0.0], [0.3475, 0.0, 0.0]]
in fold 3 economic gain is: 0.998
p(c1)= 0.49175
p(c2)= 0.1585
p(c3)= 0.34975
[[0.4865, 0.0, 0.0], [0.169, 0.0, 0.0], [0.3445, 0.0, 0.0]]
in fold 4 economic gain is: 0.973

```

Figure 1: this is the result from cross validation

This model has a 5-fold cross validation, but it can be easily changed to the number of the fold that is wanted. The shuffling and cross validation are both improving the model's performance by adjusting how to use the data set itself.

4.3 Optimization

The other optimization I did for the model is modifying the class conditional probabilities. I used the first true classes and assigned classes list for this optimization. For each wrong assignment, I increase the conditional probability for the d given the true class, for example, the true class for d_i is 2, but it's assigned class is 3, the model will increase the $\Pr(d_i | C_2)$ by a small delta. after the adjustment for all wrong assignments, the model normalized the conditional probabilities for each class to 1 again, and then performed the decision rule again with the new conditional probabilities. By doing so, it increases the conditional probabilities in the right class for all wrong assigned measurement data, and the identification accuracy is increased at the same time. I also used `DataFrame.loc()` function to locate the conditional probabilities in each class. here shows how economic gain changes with the increase of delta (Figure 2).

```
In [66]: from matplotlib import pyplot as plt
plt.plot(d_value, eg_d)
plt.show()
```

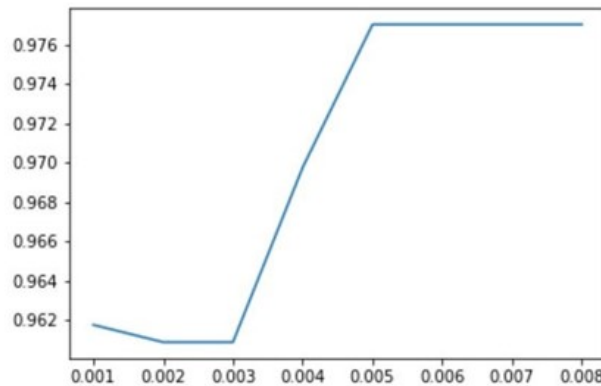


Figure 2: The relationship between economic gain and delta

5 Conclusion

Following the Bayes rule and learning from the data set, one can build a model for the data analysis and help to predicted unknown. there are several effects that affect the performance of the model, for example the data quality, data size, the distribution of the data, and the technique that used during the discover process. The Data shuffling, cross validation and modifying the conditional probabilities are all helpful for improving the model's performance, however this encasement of the improvement will slow down and stop when the performance reaches a certain point. I believe this limitation is caused by the information that carried in the data set itself. In the future I would love to explore more of how Bayes rule works in regression.

6 Reference:

http://haralick.org/ML/discrete_bayes.pdf
http://haralick.org/ML/midterm_project.pdf