# COMPARING DECISION TREE AND RANDOM FOREST DATA MINING METHODS, THE IMPACT ON THE PREDICTION RESULTS OF BANK CUSTOMER CREDIT CLASSIFICATION

*E.I. Gubin, PhD, Associate Professor*
*Yubo Jin, student gr. 8PM0I*
*Tomsk Polytechnic University*
E-mail: yuybo2@tpu.ru

## Introduction

Data mining is the process of extracting information and knowledge that people do not know in advance but that has potential usefulness from a large amount of incomplete, noisy, fuzzy, and random actual data. Generally speaking, the results of data mining do not require completely accurate knowledge, but rather a general trend. As far as specific applications are concerned, data mining is the process of using various analytical tools to discover relationships between models and data in massive data sets, and these models and relationships can be used to make predictions. One of the primary risks faced by banks is credit risk, of which loan risk is the main element. This paper aims to screen credit customers according to the personal information on the finance dataset, finding the right person (the good customer).

## Research methods

Using classification techniques in data mining, customers can be divided into different categories. The target variable of the customer's personal data we obtained is "GB", which means 1 is good, 0 is bad. It can be seen that this is a classification problem. We use decision trees and random forests in our classification technology to build credit risk models. The steps we need to do: 1) Data reading and preprocessing. 2) Model building. 3) Parameter optimization. 4) Model prediction and feature importance. 5) Visualize decision trees and random forests. 6) Compare the classification accuracy of decision trees and random forests. By comparing the two data mining classification methods, we analyze the reasons for the differences and select the most suitable and accurate model.

## Results

First, use Python to read the raw data in CSV format, convert it into dataframe format, and view the dataset description. After that, clean the data, fill in the missing value with the mean of the same category, delete the noise data with excessive error, check that there are no duplicate values, and perform digitalization of the data in string format.

Here we select the prepared data sheet, which contains the credit information of 3000 people. We define the rest of the data, except "GB" as x (feature) and the categorical variable "GB" as y (target). Usually, 75% of the data set is used as the training set and the remaining 25% as the test set. Next, use the sklearn libraries to build decision tree and random forest models, respectively, and use the grid search method for parameter optimization.

In Table 1, showing the importance of different features, the larger the value, or the closer it is to 1, the more important it is. On the contrary, the closer it is to 0, the less important it is. Here we will arrange them in reverse order.

Table 1. The importance of features

| Feature | Importance |
|---|---|
| Age | 0.227958 |
| TMJOB1 (Time at Job) | 0.101802 |
| INCOME | 0.077459 |
| cards_ (Credit Cards int format) | 0.072145 |
| PERS_H (Time at Job) | 0.069979 |
| CASH (Requested cash) | 0.065912 |
| EC_CARD (EC_card holders) | 0.049334 |
| INC (Salary) | 0.044683 |
| prof_ (Profession) | 0.040487 |
| NMBLOAN (Num Mybank Loans) | 0.039799 |

Samples: Total number of samples, Values: bad values/ good values.

We need to know from decision trees when to stop: One way of doing this is to set a minimum number of training inputs to be used on each leaf. Another way is to set the maximum depth of your model. Maximum depth refers to the length of the longest path from a root to a leaf. Here we choose one of the decision trees in the random forest for display.
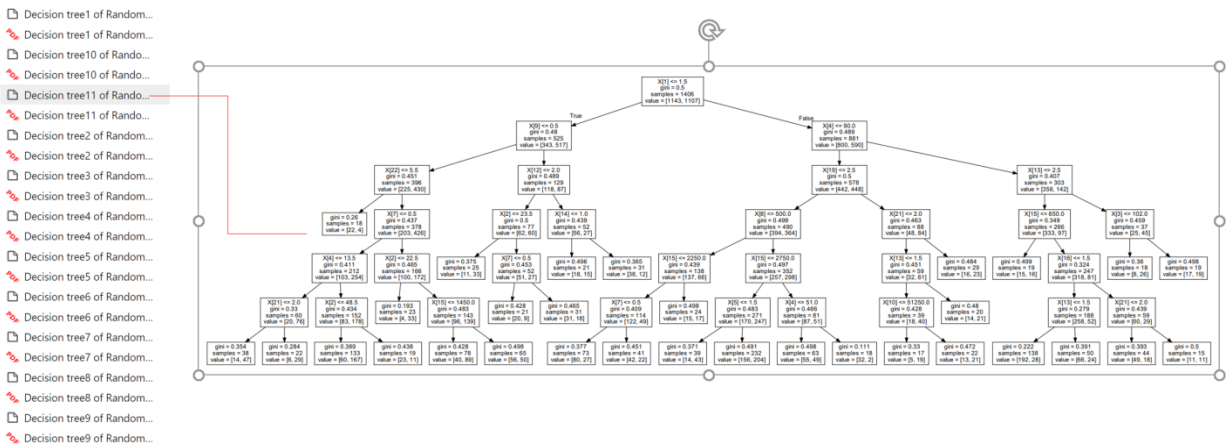


Fig. 1. Visualize one of the decision trees in a random forest

Compare the performance metrics scores of the two models in Table 2:

Table 2. Get classification performance metrics

| Performance metrics | Decision tree Score | Random forests Score |
|---|---|---|
| Accuracy | 0.57333 | 0.67066 |
| Precision | 0.63788 | 0.68238 |
| Recall | 0.54653 | 0.69796 |
| f1-score | 0.58868 | 0.69008 |

It can be seen that the classification of random forests is better than that of decision trees because random forest is based on the combination of multiple weak classifiers of decision trees into a strong classifier.

**Conclusion**

A single decision tree can achieve complete accuracy or a root mean square error on the training set, which is impossible for random forests, but the performance of random forests on the test set is significantly better than that of decision trees. Avert overfitting through random forests and achieve better generalization performance. The accuracy of random forests is 67.066%, significantly better than the accuracy of decision trees, at 57.333%. So, we finally chose the random forest model.

**References**

1. Onesmus Mbaabu. (09.04.2020) Introduction to Random Forest in Machine Learning [Electronic resource]. – URL: https://www.section.io/engineering-education/introduction-to-random-forest-in-machine-learning
2. The application of big data in the financial industry [Electronic resource]. – URL: http://c.biancheng.net/view/3736.html
3. Prashant Gupta. (18.05.2017) Decision Trees in Machine Learning [Electronic resource]. – URL: https://towardsdatascience.com/decision-trees-in-machine-learning-641b9c4e8052
4. Huang Shan, Gubin E. Data cleaning for data analysis. Томск, 2018г. - С. 387-389.