

FLIP01 PROJECT

JINGBAO LUO

ABSTRACT. This article utilizes the Isolation Forest algorithm to achieve credit card fraud detection with a recognition accuracy of 0.97 by continuously fine-tuning the algorithm.

CONTENTS

1. Introduction	2
2. Method	2
3. Data Exploration	3
4. Experiment	5
5. Conclusions	6

Date: 2023-08-13.

Key words and phrases. Machine Learning, Credit Card Data, IsolationForest.

1. INTRODUCTION

Credit card fraud refers to intentional acts of using forged or invalidated credit cards, impersonating others to deceive for financial gain, or engaging in malicious overdraft behavior using one's own credit card. There are three main forms of credit card fraud: card lost or stolen fraud, application fraud, and counterfeit credit card fraud. Among these fraud cases, over 60% involve counterfeit credit cards, characterized by organized criminal activities that include stealing card information, manufacturing fake cards, selling counterfeit cards, and then using these fake cards to commit crimes for huge profits. Credit card fraud detection is a crucial method for banks to minimize losses.

This article utilizes a dataset sourced from Kaggle to implement credit card fraud detection. The dataset contains transaction information of European cardholders during September 2013, conducted through credit cards. The dataset covers transactions that occurred within a span of two days, with a total of 284,807 transactions. Among these transactions, there were 492 cases of fraud, making the dataset highly imbalanced, where the positive class (fraudulent transactions) accounts for only 0.172

The original dataset has undergone anonymization and Principal Component Analysis (PCA) processing. The anonymous variables V1, V2, ..., V28 represent the principal components obtained through PCA, and the only variables not subjected to PCA are Time and Amount. Time represents the time elapsed between each transaction and the first transaction in the dataset, measured in seconds, while Amount represents the transaction amount. The "Class" variable is a categorical variable, taking the value 1 when fraud occurs and 0 otherwise.

This article utilizes a dataset sourced from Kaggle to implement credit card fraud detection. The dataset contains transaction information of European cardholders during September 2013, conducted through credit cards. The dataset covers transactions that occurred within a span of two days, with a total of 284,807 transactions. Among these transactions, there were 492 cases of fraud, making the dataset highly imbalanced, where the positive class (fraudulent transactions) accounts for only 0.172

The original dataset has undergone anonymization and Principal Component Analysis (PCA) processing. The anonymous variables V1, V2, ..., V28 represent the principal components obtained through PCA, and the only variables not subjected to PCA are Time and Amount. Time represents the time elapsed between each transaction and the first transaction in the dataset, measured in seconds, while Amount represents the transaction amount. The "Class" variable is a categorical variable, taking the value 1 when fraud occurs and 0 otherwise.

The main implementation in this article is utilizing the Isolation Forest algorithm to achieve credit card fraud detection. By continuously fine-tuning the algorithm through parameter adjustments, the article aims to improve the algorithm's performance.

2. METHOD

The Isolation Forest algorithm, proposed in 2008 by Liu Fei, Zhou Zhihua, and others, does not rely on distance or density metrics to describe the differences between samples and other samples. Instead, it directly characterizes the so-called



[illegible]

The logic of the Isolation Forest algorithm is intuitive. It uses binary trees to split the data, and both sample selection and feature selection are performed using randomization. If a certain sample is an outlier, it may require very few iterations to be isolated.

Isolation Forest algorithm code is [algorithm 1](#):

Input: X -Input data, t -number of trees, φ - subsampling size

Output: a set of t *iTRees*

- ```

1: Initialize Forest
2: set height limit $l = \text{ceiling}(\log_2 \varphi)$
3: for $i = 1$ to t do
4: $X' \leftarrow \cup(X, \varphi)$
5: $\text{Forest} \leftarrow \cup i\text{Trees}(X', 0, l)$
6: end for
7: return Forest

```

In this session, basic tools from pandas were utilized to perform data analysis and exploration. By analyzing the data, a more intuitive understanding was gained. table 1 provides a basic overview of the data for each attribute, including measures such as the mean, maximum, minimum, and other indicators.

From the perspective of fraud occurrence, the data was visualized as shown in fig. 1. It can be observed that the instances of fraud are significantly fewer than non-fraudulent cases, indicating the presence of data imbalance.

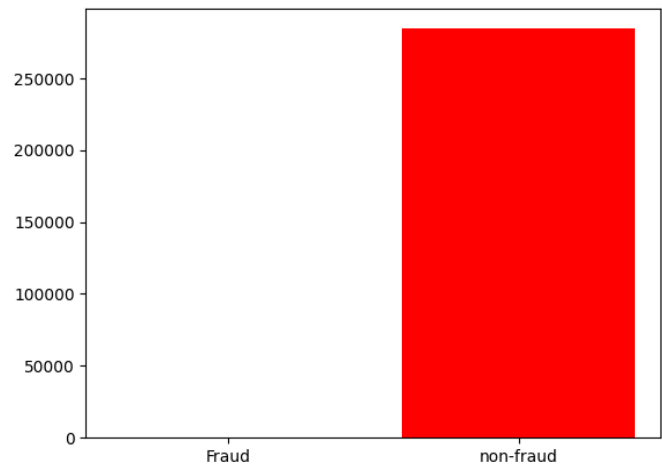


FIGURE 1. fraud or no fraud

fig. 2 is a heatmap of correlations, representing the correlation of each attribute with fraud. From the figure, it can be seen that the individual attributes have a relatively low impact on the class, and some even exhibit negative correlations.

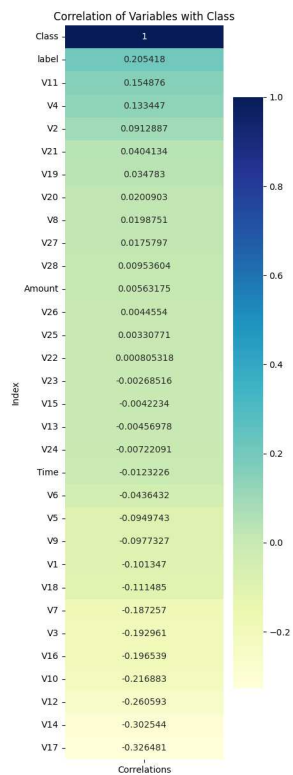


FIGURE 2. correlation heat plot

TABLE 2. parameter adjustment

| Parameter                | Accuracy | Recall | Parameter                                | Accuracy | Recall |
|--------------------------|----------|--------|------------------------------------------|----------|--------|
| none-parameter-adjust    | 0.963    | 0.821  | none-parameter-adjust(feature_delete)    | 0.953    | 0.663  |
| n_estimators-adjust(220) | 0.981    | 0.722  | n_estimators-adjust(feature_delete)(140) | 0.980    | 0.520  |
| n_features-adjust(11)    | 0.981    | 0.750  | n_features-adjust(feature_delete)(11)    | 0.981    | 0.522  |
| max_samples-adjust       | 0.981    | 0.829  | nmax_samples-adjust(feature_delete)      | 0.980    | 0.634  |

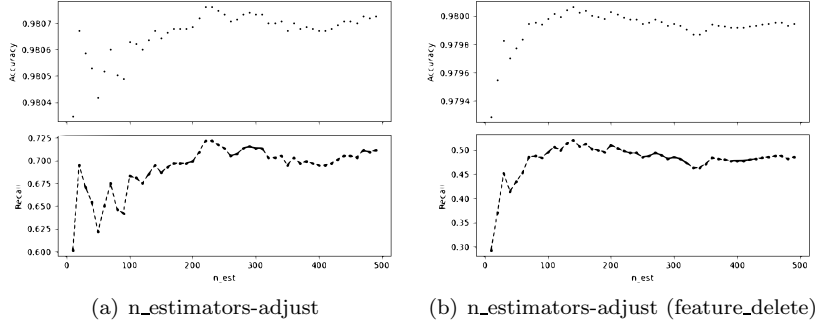


FIGURE 3. n\_estimators-adjust

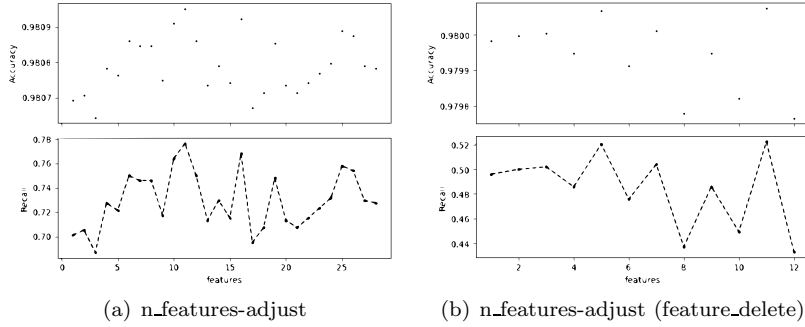


FIGURE 4. n\_efeatures-adjust

#### 4. EXPERIMENT

In the experimental design, we achieved the optimal performance of the model by continuously adjusting parameters. The evaluation of the training set's accuracy and recall mainly involved removing attributes and continuously changing the parameters of the isolation forest.

table 2 displays the accuracy and recall rates under various parameter adjustments. From the table, it can be observed that the recall rate deteriorates after removing some features. This implies that although these features may exhibit negative correlations with the outcome, they should not be indiscriminately deleted. By iteratively trying different parameter combinations, we ultimately identified the optimal parameters.

fig. 3 and fig. 4 and fig. 5 are visualizations of specific parameter adjustments.

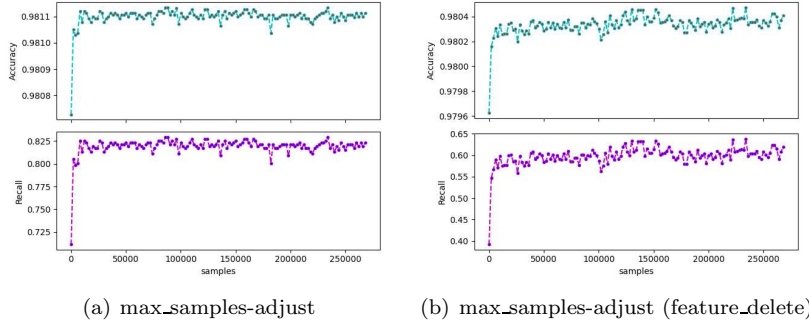


FIGURE 5. max\_samples-adjust

## 5. CONCLUSIONS

The detection of credit card fraud data was achieved through Isolation Forest, and by continuously adjusting the parameters of the algorithm, a high level of accuracy and recall was ultimately reached, at 0.98 and 0.8, respectively.

However, a notable drawback of this method is the lack of a more in-depth analysis of attributes. This could mean that in the data preprocessing stage or feature selection process, there was insufficient exploration, understanding, and utilization of attribute information relevant to fraud. While optimizing the accuracy and recall of the model is undoubtedly important, equal attention should also be given to the importance and correlation of attributes, as well as their role in fraud detection.

(A. 1) SCHOOL OF ECONOMICS AND MANAGEMENT,, NANJING UNIVERSITY OF SCIENCE AND TECHNOLOGY, NANJING 210094, CHINA

*Email address, A. 1:* jluo@tulip.academy