

# 7410 Mini Case Study

## Real-life Fraud Detection Scenario

Name: Chen, Xingyi(陈星熠)

UID: xychenn

University Number: 3036198102

Email: xychenn@connect.hku.hk

December 2, 2023

### 1 Introduction

ENRON dataset[2] is a widely used dataset in the field of machine learning and natural language processing. It is named after the energy company Enron Corporation, which became infamous due to a massive accounting scandal that led to its bankruptcy in 2001. The dataset consists of around 500,000 emails sent by Enron employees, which were made public during the subsequent investigations.[4]

This mini-case study utilizes the ENRON dataset for fraud detection. Due to its large size and diverse content, this study uses a simplified version of the ENRON dataset.

In this study, I define designing a fraud data analytics plan that can identify committing people according to the inside records as the scope of the study. I first performed Exploratory Data Analysis (EDA) on the dataset and then used the pre-processed data for model training. I trained two models, Random Forest[1] and Neural Network[3], and compared the performance of both. I then tested the Random Forest model, which has better performance among the two, using a fraud scenario dataset developed by myself, and analyzed and summarized the above results. You can get all the relevant source code and datasets from the [link](#).

### 2 Exploratory Data Analysis

#### 2.1 Summary Description of the Dataset

The dataset[2] contends 22 variables with 20 in numeric type and 2 are character type. So I excluded these 2 char-

acter variables, namely poi and email address, and then attempted to find out appropriate candidate variables to do the further analysis. I calculated the deviations of those numeric variables, and 12 of them are distributed in relatively large ranges as shown in Figure 1.

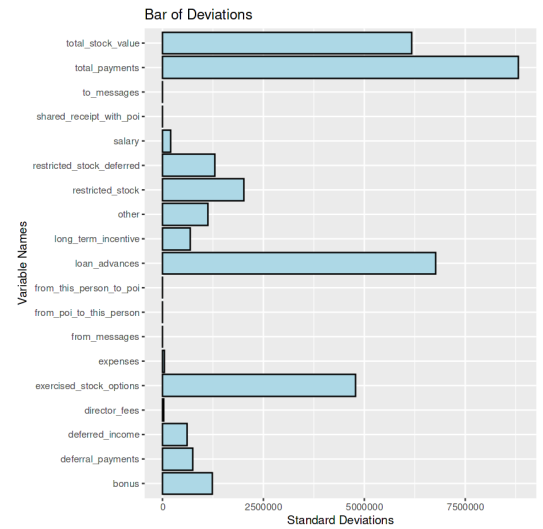


Figure 1: Bar of Deviations

#### 2.2 Univariate Analysis

In this section, I would like to take 3 variables, namely salary, exercised stock options, and bonus, into consideration. I will give my reasons in the following paragraphs.

Firstly, I assume that salary relates to people's rank in their company since a higher position means a higher salary in most cases.

Based on the assumption above, I may further infer

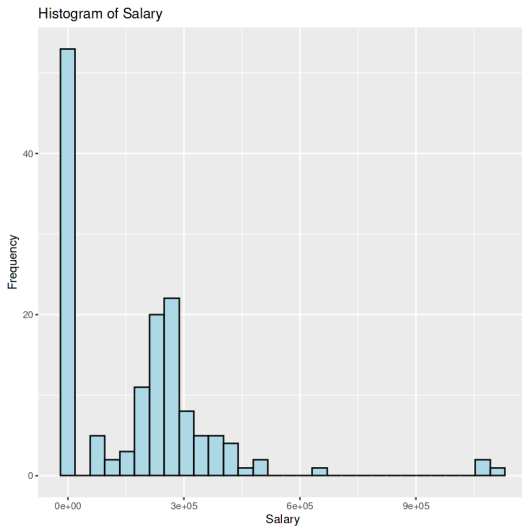


Figure 2: Histogram of Salary

that it should present a relatively falling trend of exercised stock options as the salary data goes up. Leaders of the company should exercise their stock options as little as possible since it is wise for them to keep their stock options as long as they can if they feel confident about the company's development under most circumstances.

Lastly, bonus data should not be too high compared with the value of its corresponding salary data from my perspective.

Additionally, the standard deviations of these 3 variables are relatively higher than those of others according to the previous section, which is also suspicious. Hence, I choose them to do univariate analysis.

**Salary.** As shown in Figure 3, there is some outliers' value over 9e05, and the amount of these outliers is really small in Figure 22 So, I infer that these outliers of salary data may come from leaders of the company. I printed their names, and it was a very list of ENRON leaders' names, which validated my previous assumption.

**Exercised Stock Options.** As shown in Figure 4, there is some outliers' value over 2.5e06. However, it was difficult to identify the dividing line between outliers and maxima from Figure 4, so I drew a histogram with a smaller scale (Figure 5) to find the exact range of outliers, which is above 4e6.

**Bonus.** As shown in Figure 6, there is some outliers' value over 2e06. I consider these outliers as unusual bonus receivers.

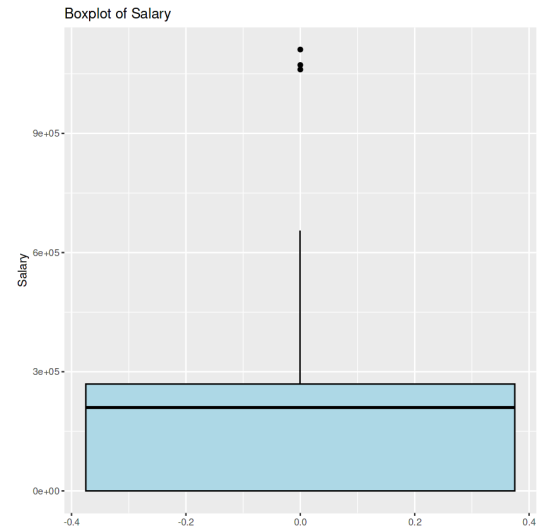


Figure 3: Boxplot of Salary

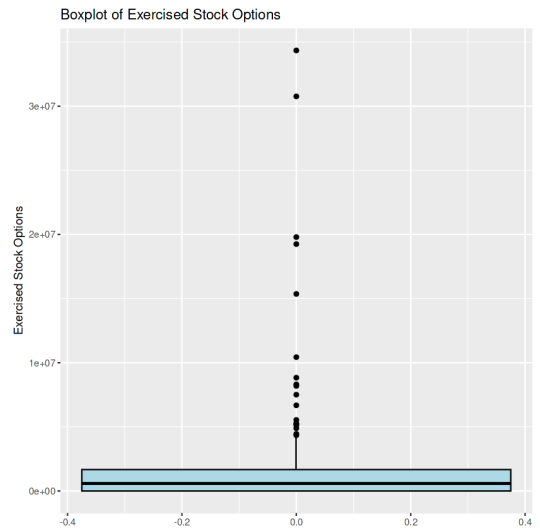


Figure 4: Boxplot of Exercised Stock Options

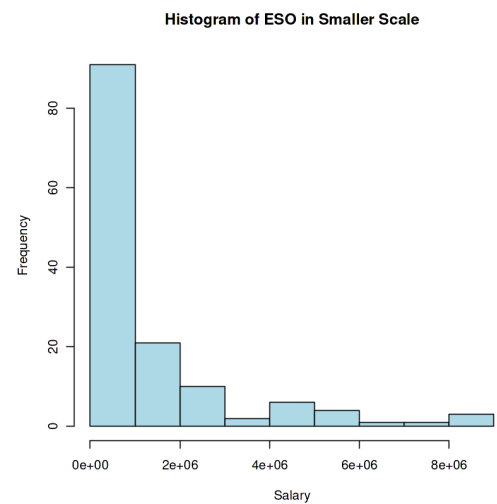


Figure 5: Histogram of ESO in Smaller Scale

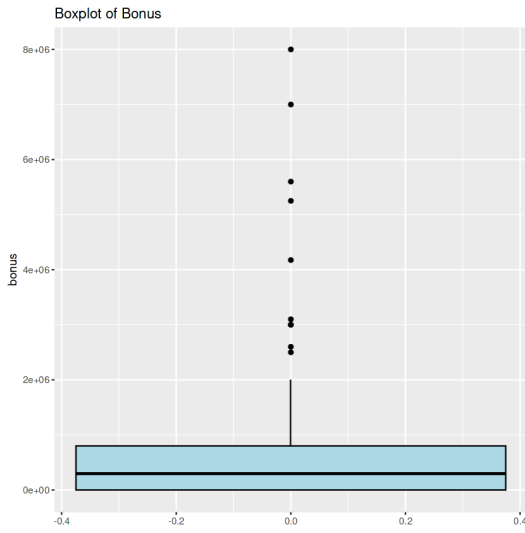


Figure 6: Boxplot of Bonus

## 2.3 Bivariate Analysis

In this section, I plotted scatter plots for the above three variables two by two to further analyze the relationship between the two variables and to identify bivariate outliers based on my following assumptions.

**Association Between Salary and ESO.** I have assumed that it should present a relatively falling trend of exercised stock options as the salary data goes up. Leaders of the company should exercise their stock options as little as possible since it is wise for them to keep their stock options as long as they can if they feel confident about the company's development under most circumstances.

However, the scatterplot(Figure 7) shows a number of data points that are not only salary outliers but also ESO(Exercised Stock Options) outliers. It is likely an indication that company leaders are not confident in the company's growth position and may even have engaged in illegal stock transactions.

**Association Between Salary and Bonus.** Based on the inference that company leaders may sell a large number of their stock, which breaks federal laws, I further assumed that they had accomplices in their employees. Figure 8 shows that there are some people who received bonuses over 10 times their salaries. Hence, I consider these outliers as employees receiving unusual bonuses instead of POI in case some of them are outstanding employees.

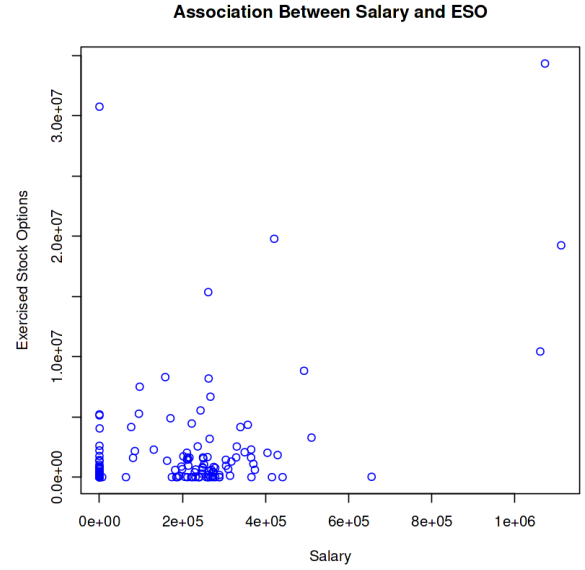


Figure 7: Association Between Salary and ESO

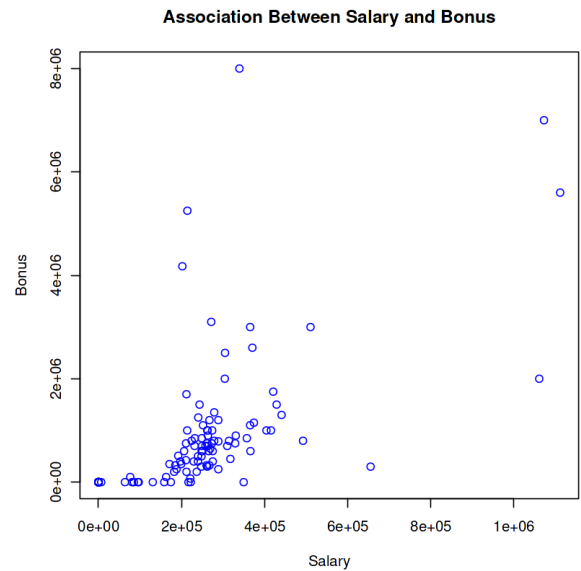


Figure 8: Association Between Salary and Bonus

## 2.4 Multivariate Analysis

In this section, I plotted correlations among selected variables(Figure 10) where Name and Email Address were removed from. In addition, I transformed poi data from character to numeric type so that we can see associations between poi data and others.

Figure 10 shows that Salary, Bonus, Total Stock Value, and ESO are significantly more strongly correlated with poi than the other variables. After perusing the dataset meticulously, I found that Total Stock Value is equal to the sum of ESO and Restricted Stock, thus the correlation between Total Stock Value and and ESO is strong. So essentially the strong correlation between Total Stock Value and poi is the strong correlation between ESO and poi. This further validates the validity of the assumptions I made above.

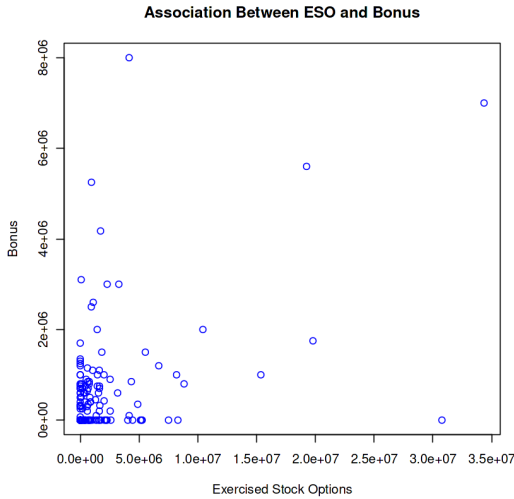


Figure 9: Association Between ESO and Bonus

## 2.5 Outlier Analysis

In this section, I aggregate univariate outliers, bivariate outliers, and multivariate outliers as a subset of the ENRON dataset and compare their poi values side-by-side.

Figure 11 shows that the more data points that are outliers in more variables at the same time, the greater the probability that their corresponding poi values are true. The joint Salary and ESO variables are used to identify poi are quite effective in the ENRON dataset.

It is worth noting that there is a person in the outlier set named LAVORATO JOHN J whose ESO and Bonus

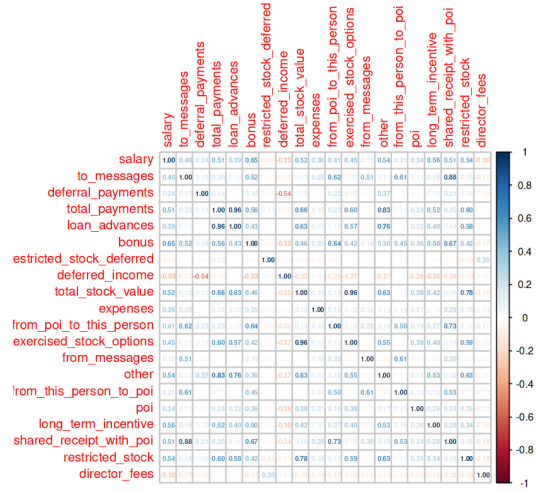


Figure 10: Correlations Among Variables

data are outliers but are not poi. I guess that he may be an excellent employee. Data points like this may have an impact on model training, and a decision should be made to either upsample such data or just remove such data based on the parameter capacity of the model.

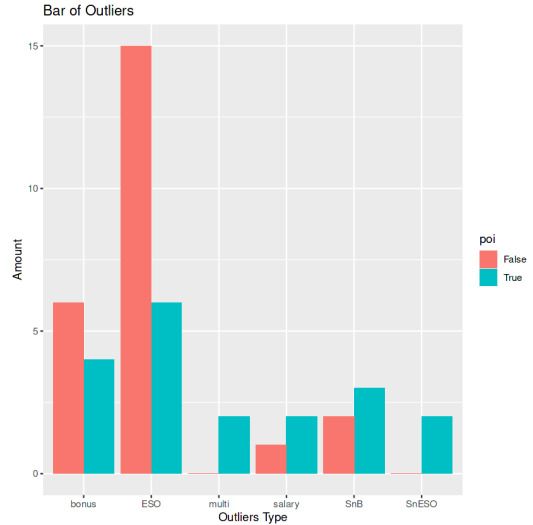


Figure 11: Bar of Outliers

## 3 Model Building and Validation

### 3.1 Comparison of Performances of the Two Models

In this section, two supervised learning algorithms, Random Forest[1] and Neural Network[3], are selected to conduct fraud detection modeling on ENRON dataset. Prior

to model training, I applied SMOTE (Synthetic Minority Over-Sampling Technique) to enhance the imbalanced dataset, shown as Figure 12 and Figure 13.

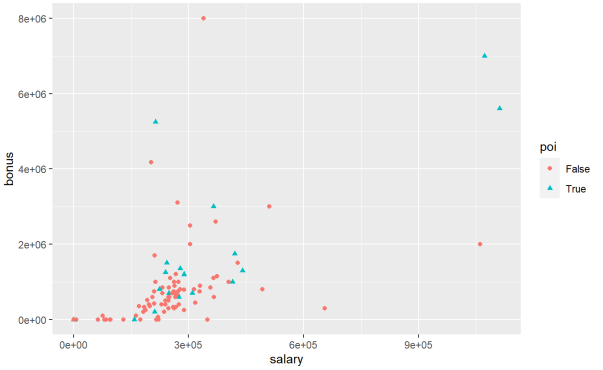


Figure 12: RAW ENRON

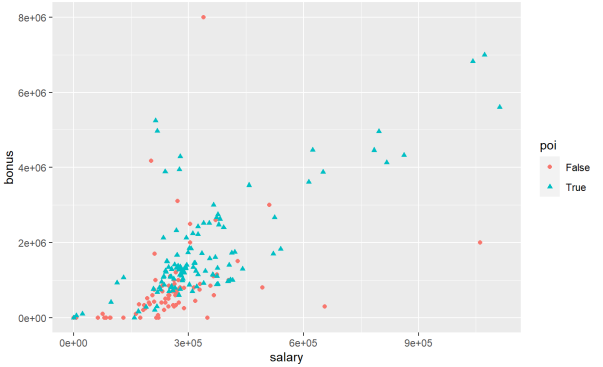


Figure 13: Enhance Imbalanced Dataset

I first trained a random forest model, and the final outcome shows an outstanding index of accuracy, which is up to 0.9742. The distribution of this RF model’s variable importance is shown in Figure 14 . I then trained a neural network model with 800 epochs, however, within taking more training time, the model didn’t get higher performance. The accuracy of the neural network is just 0.9142.

There is a clear performance gap between random forest and neural network, as shown in Figure 16 and Figure 17. Green dots indicate data points that are consistent with the ground truth, and the red ones indicate the wrong prediction. We can see that random forest model did a better job in this task.

### 3.2 Result Analysis

According to the experiment results above, I then further discuss the shortage and advantage of deploying these

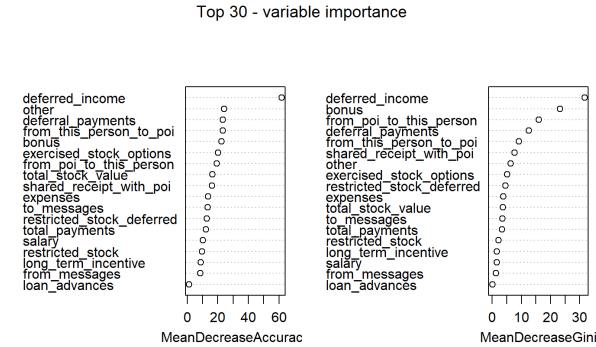


Figure 14: RF Variable Importance

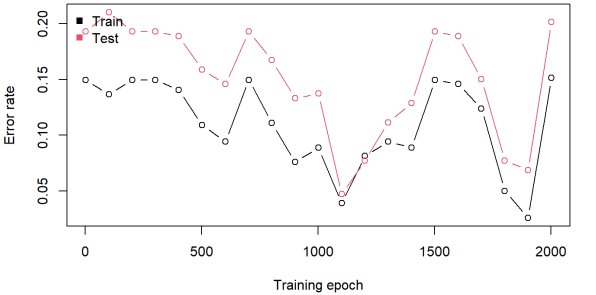


Figure 15: Process of NNet Training

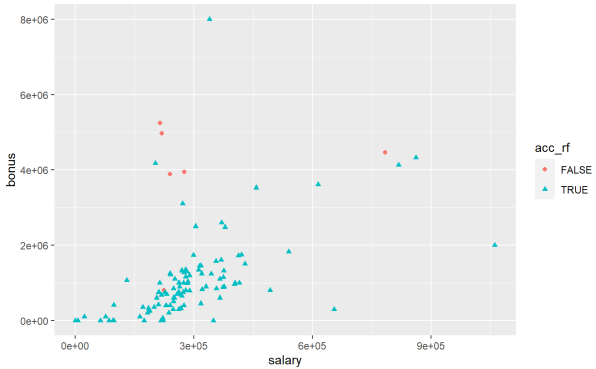


Figure 16: Accuracy of RandomForest

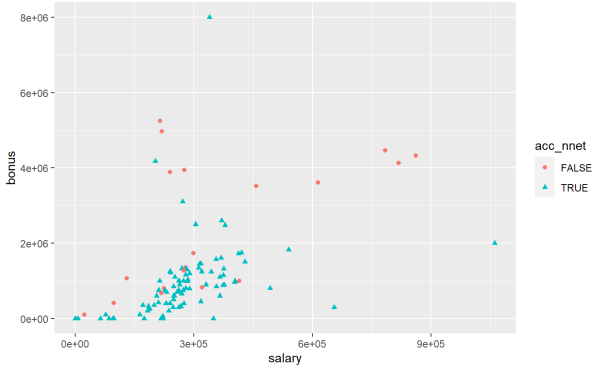


Figure 17: Accuracy of Neural Network

two models in real-life situations. In this experiment, the training time of the two models is short and the models' sizes are small due to the small amount of model parameters. So it's possible to run a random forest model or neural network model online, and they could be trained online for several certain trading cases. However, fewer model parameters means weaker learning ability. The random forest model can only handle some simple tasks while it's risky to apply random forest models to some complex fraud-detecting tasks. Neural network models although perform weaker than random forest models in this experiment, can improve their learning ability by extending more neural network layers. In general, it's more efficient to deploy random forest models in simple fraud detection, and deploying deeper neural network models on the complex fraud detection task could be more reliable.



Figure 18: Groundtruth



Figure 19: Prediction of RandomForest

## 4 Fraud Scenario Identification

In this section, I first give the permutation of fraud scenarios in the ENRON case. I then choose one of them, "In-

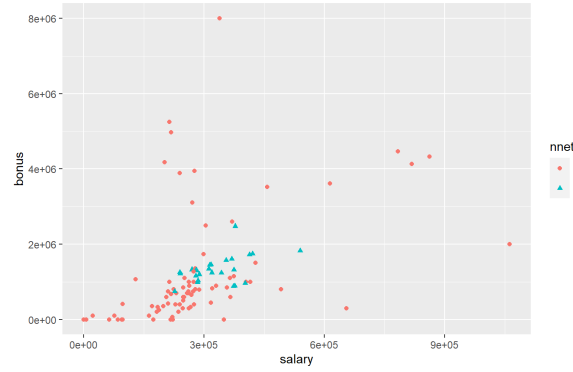


Figure 20: Prediction of Neural Network

sider Trading", as a reference to develop a new fraud scenario as well as a corresponding dataset. I use the dataset as a test dataset to evaluate the performance of the random forest model trained before and analyze results.

### 4.1 Permutation of Fraud Scenarios

As shown in Figure 21, there are 10 different specific fraud schemes in the ENRON case. Due to the volume of data used to analyze most of them being huge, I refer to the one of "Insider Trading" which is mainly identified by numeric data instead of text to simplify the test dataset.

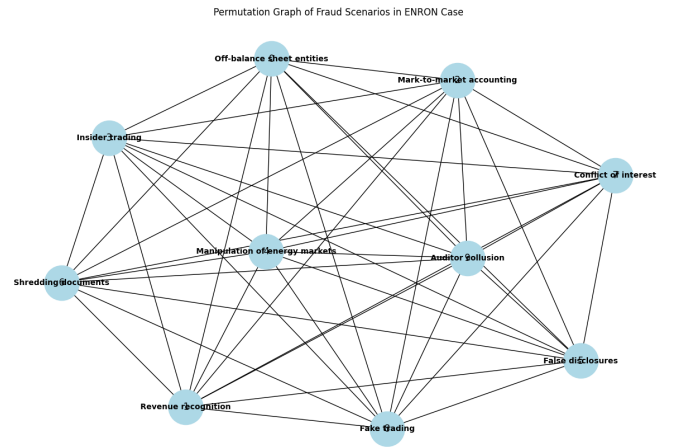


Figure 21: Permutation of All ENRON Fraud Scenarios

### 4.2 Test Financial Fraud Scenario Definition

**Scenario Definition.** Insider trading refers to the buying or selling of securities, such as stocks or options, based on non-public information about a company. In the Enron case, several executives engaged in illegal insider trading,

including CEO Jeffrey Skilling and CFO Andrew Fastow. They sold their Enron shares based on non-public information about the company’s financial problems, avoiding significant losses.

Insider trading is illegal under securities laws as it undermines the fairness and integrity of the financial markets. Skilling and Fastow violated these laws by trading Enron stock based on non-public information, gaining an unfair advantage over other investors.They were eventually charged and convicted for their involvement in the Enron scandal. Skilling was convicted of multiple counts of securities fraud, insider trading, and conspiracy, while Fastow pleaded guilty to charges of conspiracy and securities fraud. They both received significant prison sentences and financial penalties for their actions.[4]

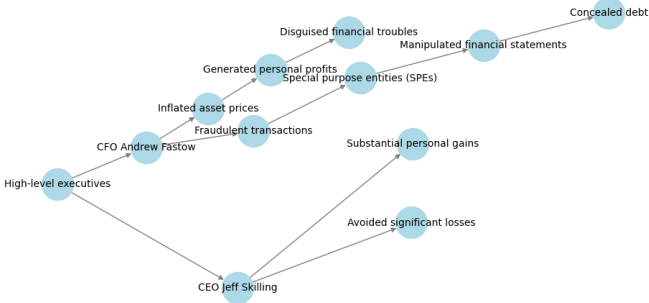


Figure 22: Permutation Graph of Inside Trading Scenario

**Dataset Creation.** I develop a new test dataset according to the ENRON dataset, e.g. randomly sample the people of interest(POI) data, without repeated samples, for random sampling times in a range of 2 to 4 and calculate the mean of those to produce a new record. I have produced 510 records in total, with 60 POI records and 450 non-POI records.

### 4.3 Test Results of RF Model

The prediction results of the trained random forest model and comparisons with groundtruth are shown in this subsection, further analysis of these results is in the next subsection. The accuracy of prediction is 0.9022.

### 4.4 Result Analysis

The variable importance of trained random forest model, shown in Figure 14, shows that the model focus a lot on the

Table 1: Prediction vs Groundtruth

		RF Prediction	
		False	True
GT	False	406	44
	True	11	49

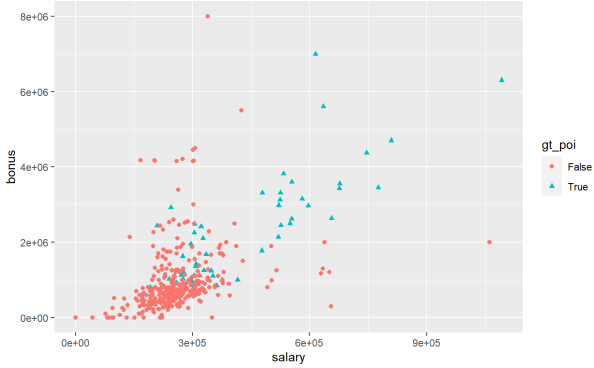


Figure 23: Groundtruth of Testset

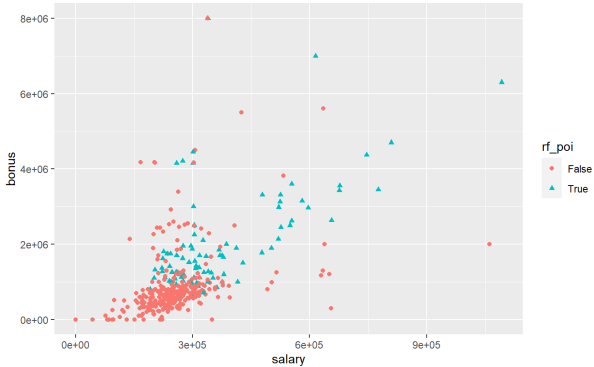


Figure 24: Prediction Results of RF Model on Testset

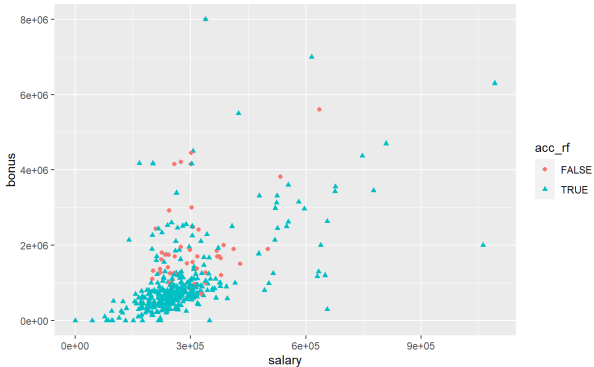


Figure 25: Prediction Accuracy of RF Model on Testset



value of someone’s bonus, exercised stock options, deferred income and etc. which are exactly the main features of a inside trading fraud scenario. Hence the model has a pretty well performance on the test dataset.

However, Table 1 shows that there is an imbalance in accuracy between POI detection and non-POI detection. It demonstrates that when the model faces non-POI, it can identify the data correctly with a probability of 90.22%, while its accuracy of identifying POI is only 81.67%. From my perspective, I speculate that maybe the importance of deferred income sits in the wrong position. Although it may somehow concern the possibility of inside trading fraud, it’s not the most vital factor to detect inside trading. So the model will make the wrong decision when the POI data comes to all inside-trading-related variable value is high however with a low value of deferred income.

## 5 Summary and Recommendation

### 5.1 Data-based Analysis

The help of a data-driven fraud detection model will significantly improve the work efficiency when we are doing data analysis of fraud detection, because a data-driven model is quite sensitive to the changes of vital data, sometimes even more sensitive than an experienced account, that can save us a lot of energy to staring at rows in the data to find the red flags. However, it can not replace human work yet, since financial operations are low-risk-tolerated. For my recommendation, we should use it as a tool for the present.

### 5.2 Non-data-based Analysis

Out of data, we are also required to focus on non-data analytic elements for more robust financial fraud detection. Instead of solely relying on data analytics, consider incorporating behavioral analysis techniques. This involves monitoring and analyzing the behavior of individuals or entities to identify patterns or anomalies that may indicate fraudulent activities. For example, monitoring transaction patterns or changes in spending habits can provide valuable insights.

## 5.3 Recommendations

To prevent ENRON inside trading financial fraud in the future, I would like to give some suggestions according to the experiment results in this mini-case study. Establishing robust internal controls and procedures to prevent unauthorized access to sensitive information may help a lot. In the ENRON case, the abuse of non-public information on company performance is a vital factor leading to the tragedy. This includes implementing segregation of duties, dual control systems, and regular monitoring of financial transactions. In addition, companies should perform regular internal audits to review trading activities, financial statements, and compliance with regulatory requirements. These audits can help identify any potential red flags or areas of concern that may indicate insider trading.

In general, financial fraud prevention is an ongoing effort that requires constant vigilance and commitment to maintaining a culture of integrity and compliance.

## References

- [1] BREIMAN, L. Random forests. *Machine Learning* 45, 1 (Oct 2001), 5–32.
- [2] KLIMT, B., AND YANG, Y. The enron corpus: A new dataset for email classification research. In *Machine Learning: ECML 2004* (Berlin, Heidelberg, 2004), J.-F. Boulicaut, F. Esposito, F. Giannotti, and D. Pedreschi, Eds., Springer Berlin Heidelberg, pp. 217–226.
- [3] MCCULLOCH, W. S., AND PITTS, W. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics* 5, 4 (Dec 1943), 115–133.
- [4] OPENAI. Chatgpt. <https://openai.com/research/chatgpt>, 2021.



6 Appendix

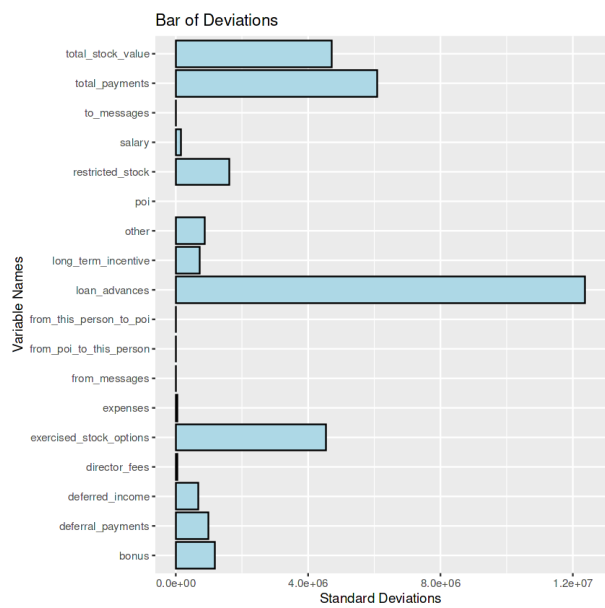


Figure 26: Deviation Distribution of Test Dataset

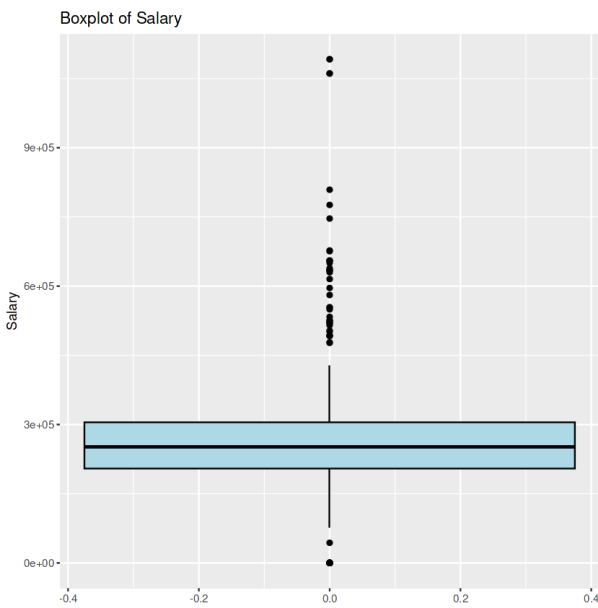


Figure 28: Boxplot of Salary on Testset

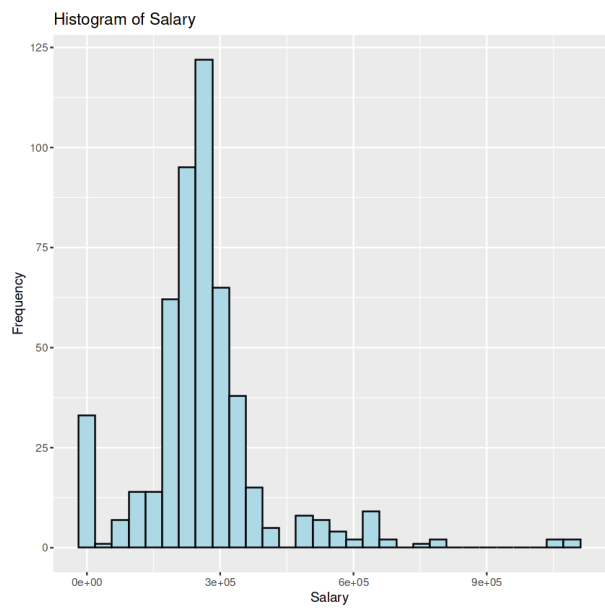


Figure 27: Histogram of Salary on Testset

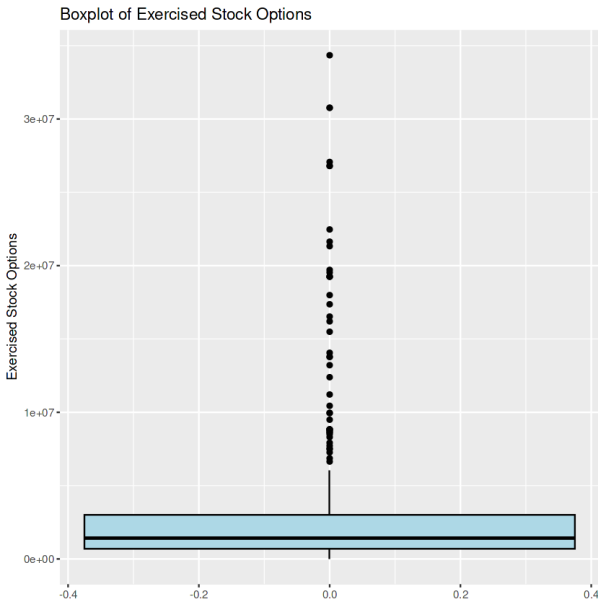


Figure 29: Boxplot of Exercised Stock Options on Testset

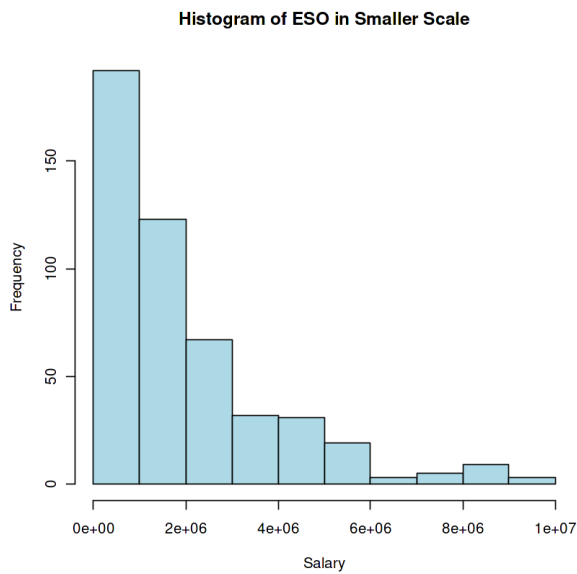


Figure 30: Histogram of ESO in Smaller Scale on Testset

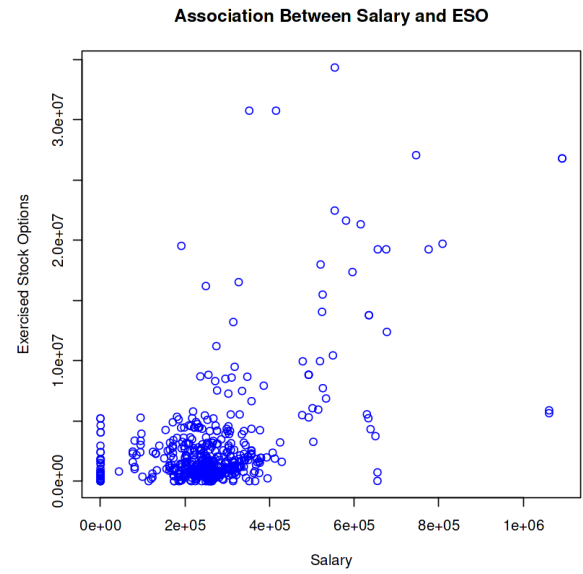


Figure 32: Association Between Salary and ESO on Testset

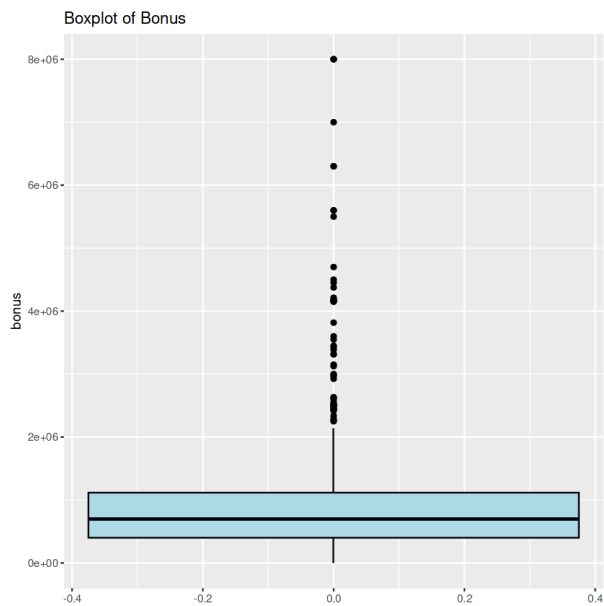


Figure 31: Boxplot of Bonus on Testset

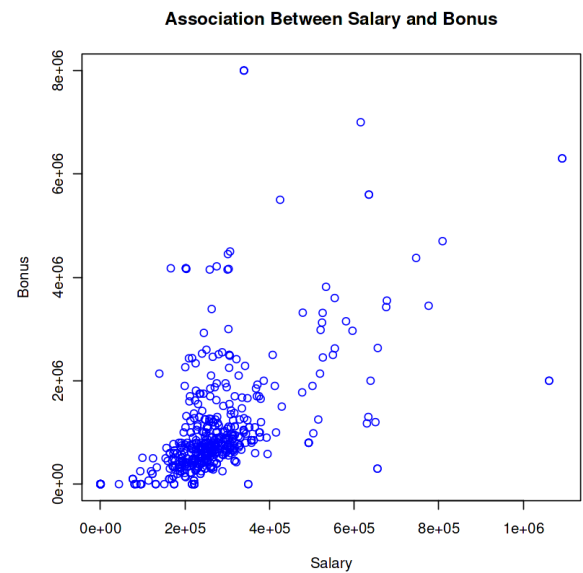


Figure 33: Association Between Salary and Bonus on Testset

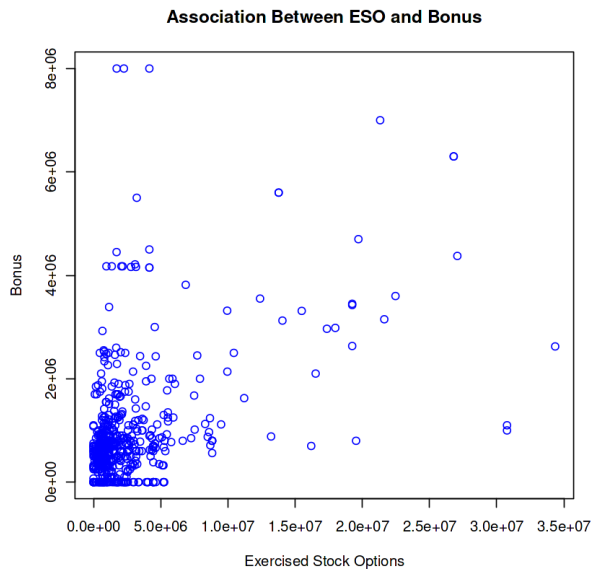


Figure 34: Association Between ESO and Bonus on Testset

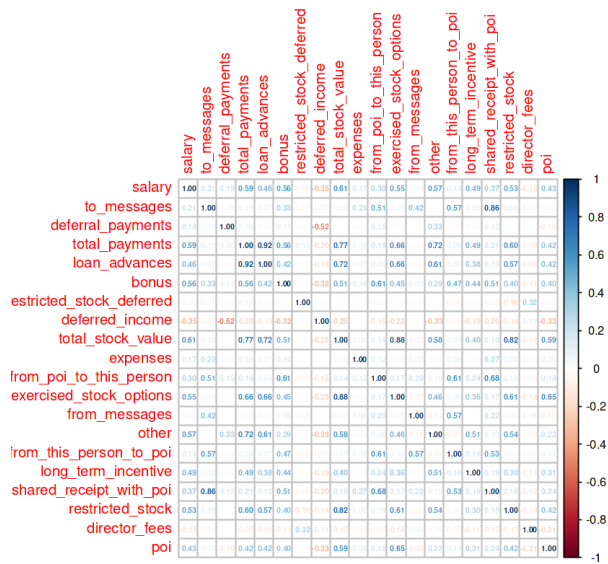


Figure 35: Correlations Among Variables on Testset

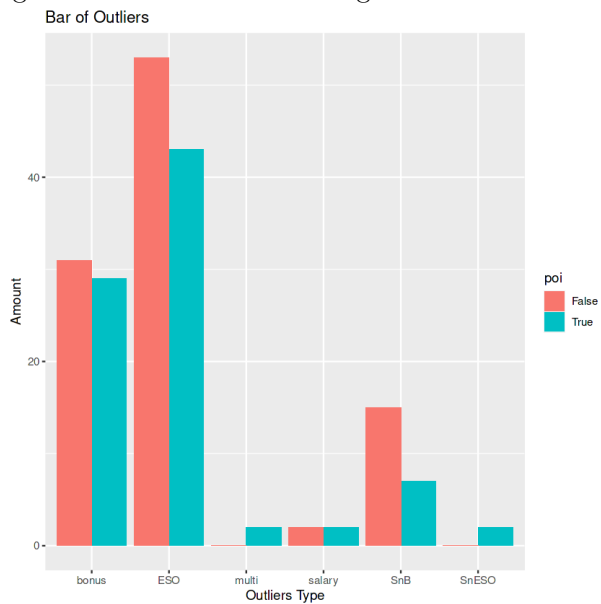


Figure 36: Bar of Outliers on Testset