# DETECTION OF FRAUDULENT FINANCIAL STATEMENTS USING THE HYBRID DATA MINING APPROACH

## A PROJECT REPORT

**Submitted by**

**ARAVINDS B R**
**14ITR007**

**KISHORE N**
**14ITR041**

**MANISHA B**
**14ITR046**

*in partial fulfilment of the requirements*

*for the award of the degree*

*of*

## BACHELOR OF TECHNOLOGY

## IN

## INFORMATION TECHNOLOGY

**DEPARTMENT OF INFORMATION TECHNOLOGY**

**SCHOOL OF COMMUNICATION AND COMPUTER SCIENCES**



Estd : 1984

## KONGU ENGINEERING COLLEGE

**(Autonomous)**

**PERUNDURAI  ERODE – 638 060**

**APRIL 2018**

**DEPARTMENT OF INFORMATION TECHNOLOGY**

**KONGU ENGINEERING COLLEGE**

**(Autonomous)**

**PERUNDURAI ERODE – 638060**

**APRIL 2018**

**BONAFIDE CERTIFICATE**

This is to certify that the Project Report entitled **DETECTION OF FRAUDULENT FINANCIAL STATEMENTS USING THE HYBRID DATA MINING APPROACH** is bonafide record of project work done by **ARAVINDS B R** (14ITR007), **KISHORE N** (14ITR041) and **MANISHA B** (14ITR046) in partial fulfillment of the requirements for the award of the Degree of Bachelor of Technology in Information Technology of Anna University, Chennai during the year 2017 - 2018.

**SUPERVISOR**                                      **HEAD OF THE DEPARTMENT**

                                                              **(Signature with seal)**

**Date:**

Submitted for the end semester viva voce examination held on _____

**INTERNAL EXAMINER**                                      **EXTERNAL EXAMINER**

# DEPARTMENT OF INFORMATION TECHNOLOGY

# KONGU ENGINEERING COLLEGE

**(Autonomous)**

**PERUNDURAI ERODE – 638060**

**APRIL 2018**

# DECLARATION

We affirm that the Project Report entitled **DETECTION OF FRAUDULENT FINANCIAL STATEMENTS USING THE HYBRID DATA MINING APPROACH** being submitted in partial fulfillment of the requirements for the award of Bachelor of Technology is the original work carried out by us. It has not formed the part of any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

<div align="right">

**ARAVINDS B R**
**14ITR007**

</div>

Date:

<div align="right">

**KISHORE N**
**14ITR041**

</div>

<div align="right">

**MANISHA B**
**14ITR046**

</div>

I certify that the declaration made by the above candidates is true to the best of my knowledge.

Date:                                       Name and Signature of the Supervisor with seal

# ABSTRACT

Financial statements are basic documents represent financial status of a company. Fraudulent financial statements are intentional and illegal acts that result in misleading financial statements or misleading financial disclosure. It is important to detect the fraudulent financial statement prior to its occurrence and prevent company's financial loss. The purpose of the project is to construct fraudulent financial statement detection model which identifies whether the financial statement of an organization is either fraud or not fraud. The percentage of data considered as training data, test data and validation data in Random Forest,Artificial Neural Network(ANN), Gradient Boosting Machine(GBM) and Generalized Linear Model(GLM)  will differ. The detection of fraudulent financial statement of an organization will reduce the loss incurred due to the fraud made with the help of financial statement.

The credit card dataset is given as the input for all the implemented methods. Initially, variable selection process takes place where along with 28 variables amount and class are also used. Random forest or random decision forest method is used which is an ensemble learning method for classification, regression and other tasks. In the next stage, in order to construct fraudulent financial statement detection models the above implemented Artificial Neural Network, Gradient Boosting Machine,Generalized Linear Model are used for the selection of more optimized result within these models.

.

# ACKNOWLEDGEMENT

First and foremost, we acknowledge the abundant grace and presence of Almighty throughout different phases of the project and its successful completion.

We wish to express our extreme gratefulness to our beloved Correspondent **Thiru.A.VENKATACHALAM B.Sc.,** and all the trust members of Kongu Vellalar Institute of Technology Trust for providing all the necessary facilities to complete the project successfully.

We express our deep sense of gratitude to our beloved Principal **Prof.S.KUPPUSWAMI B.E., M.Sc(Engg)., Dr.Ing(France).,**for providing us an opportunity to complete the project.

We express our gratitude and heartfelt sincere thanks to **Dr.S.VARADHAGANAPATHY M.S., M.E., Ph.D.,** Head, Department of Information Technology, for his valuable suggestions.

We are thankful to our project coordinators **Dr.C.NALINI M.E., Ph.D.,** and **Mrs.S.ANITHA M.E.,** for their valuable guidance and support to complete our project successfully.

We are highly indebted to **Dr.P.SURESH, M.E., Ph.D.,** Department of Information Technology for his valuable supervision and advice for the fruitful completion of the project.

We are very much thankful to our Department Faculty members for their valuable support and suggestions.

**TABLE OF CONTENTS**

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| | | |
|---|---|---|
| CHAID | - | Chi-square Automatic Interaction Detector |
| RF | - | Random Forest |
| BBN | - | Bayesian Belief Network |
| SVM | - | Support Vector Machine |
| GBM | - | Gradient Boosting Machine |
| GLM | - | Generalized Linear Model |
| ANN | - | Artificial Neural Network |
| CART | - | Classification And Regression Trees |

# CHAPTER 1

# INTRODUCTION

## 1.1 DATA MINING

Data mining is a process used by companies to turn raw data into useful information. By using software to look for patterns in large batches of data, businesses can learn more about their customers and develop more effective marketing strategies as well as increase sales and decrease costs. Data mining depends on effective data collection and warehousing.

### 1.1.1 The Scope of Data mining

Data mining derives its name from the similarities between searching for valuable business information in a large database for example, finding linked products in gigabytes of store scanner data and mining a mountain for a vein of valuable ore. Both processes require either shifting through an immense amount of material, or intelligently probing it to find exactly where the value resides. Given databases of sufficient size and quality, data mining technology can generate new business opportunities by providing these capabilities:

- **Automated prediction of trends and behaviours :** Data mining automates the process of finding predictive information in large databases. Questions that traditionally required extensive hands-on analysis can now be answered directly from the data quickly.

- **Automated discovery of previously unknown patterns :** Data mining tools sweep through databases and identify previously hidden patterns in one step. An example of pattern discovery is the analysis of retail sales data to identify seemingly unrelated products that are often purchased together.

### 1.1.2 Data mining process

The data mining process breaks down into some steps. First, organizations collect data and load it into their data warehouses. Next, they store and manage the data. Business analysts, management teams and information technology professionals access the data and determine how they want to organize it. Then, application software sorts the data based on the user's results, and finally, the end user presents the data in the form of graph

### 1.1.3 Classification

Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data. For example, a classification model could be used to identify loan applicants as low, medium, or high credit risks. A classification task begins with a data set in which the class assignments are known. For example, a classification model that predicts credit risk could be developed based on observed data for many loan applicants over a period of time. In addition to the historical credit rating, the data might track employment history, home ownership or rental, years of residence, number and type of investments, and so on. Credit rating would be the target, the other attributes would be the predictors, and the data for each customer would constitute a case. Classifications are discrete and do not imply order. Continuous, floating-point values would indicate a numerical, rather than a categorical, target. A predictive model with a numerical target uses a regression algorithm, not a classification algorithm.

### 1.1.4 Prediction

Predictive analytics encompasses a variety of statistical techniques from predictive modeling, machine learning, and data mining that analyze current and historical facts to make predictions about future or otherwise unknown events. In business, predictive models exploit patterns found in historical and transactional data to identify risks and opportunities. Models capture relationships among many factors to allow assessment of risk or potential associated with a particular set of conditions, guiding decision making for candidate transactions. The

defining functional effect of these technical approaches is that predictive analytics provides a predictive score (probability) for each individual (customer, employee, healthcare patient, vehicle, component, machine, or other organizational unit) in order to determine, inform, or influence organizational processes that pertain across large numbers of individuals, such as in marketing, credit risk assessment, fraud detection, manufacturing, healthcare, and government operations including law enforcement. The Predictor resource mines data collected in the Data Warehouse. The Data Warehouse database can store all site data, including user profile data, click-history data, and transaction data. Although the Data Warehouse database is quite large, the capabilities of the Predictor resource can efficiently mine the data.

## 1.2 OBJECTIVE

The objective of this project is to detect the fraudulent financial statements of an organization using an effective and accurate fraudulent financial statement detection model. Financial statements are company's basic documents that reflect its financial status. The financial statement is the main basis for decision making on the part of a vast number of investors, creditors and other persons in need of accounting information, as well as a concrete expression of business performance, financial status and the social responsibility of listed companies. The process strives to bring out the methodology and implementation of these classification techniques and stress upon the results and conclusion induced on the basis of accuracy and time complexity.

# CHAPTER 2

# SYSTEM ANALYSIS

## 2.1 LITERATURE REVIEW

In research studies Lenard and Alam, (2009) proposed that the analysis of ratios is chosen as one of the methods to determine fraud. Financial difficulties may be motivation for managers to engage in fraudulent activities. According to Ravisankar et al. (2011), the higher levels of debt may increase the probability of the fraudulent financial statements too.

Wei Zhou (2011) proposed the Detecting evolutionary financial statement fraud in which data mining is used widely in financial statements especially when new schemes of financial statement fraud adapt to the detection techniques. Therefore they concluded by specifying the future research is needed to design the active discovery module that is both effective and efficient. Furthermore, although they have suggested the response surface method to extract the domain knowledge and to adaptively learn the changes from fraudsters, there may exist alternatives that have equal or better performance than the RSM.

Wang et al., (2012) proposed the hybrid system consists of machine learning and rule-based classifiers. For the machine learning classifier, we investigate a variety of lexical, syntactic and knowledge-based features, and show how much these features contribute to the performance of the classifier through experiments. Na et al. (2012) proposed a rule-based system for polarity classification. Clause-level sentiment classification algorithm is developed and applied to fraudulent reviews on a discussion forum.

## 2.2 SUMMARY

The literature survey helps to understand more about the trending data mining model for fraudulent financial statement detection rather than the conventional models such as CHAID and CART. The literary references indicate that the usage of the financial ratios for determining the fraudulent statements of financial reports is a convenient and straightforward means. However, the problem of interpreting the result interpretation arises, i.e. what value of financial ratio indicates that the reports are fraudulent.

# CHAPTER 3

# PROBLEM DEFINITION

## 3.1 PROBLEM STATEMENT

The regression analysis method was initially used in organizations to detect the fraudulent financial statements. But the regression analysis method causes a number of disadvantages in terms of judgment and its error rate is relatively high. Regression analysis is a form of predictive modeling technique which investigates the relationship between a dependent and independent variable. The judgment accuracy rates of using data mining techniques to detect fraudulent financial statements of any organization vary, and the construction of the model is neither complete nor perfect. There are four elements of fraud that are found in the case of financial statements. They are,

1) serious erroneous expressions of the nature of transactions

2) knowingly violating rules

3) the victim accepting a misstatement as fact

4) damage due to financial losses caused by the above three situations.

Figure 3.1 shows the working of the model where the original data of the year before the year of fraud is given as the input and then the variables are classified as financial and non financial variables and then a two-stage statistical treatment is followed. Data mining techniques are then used to create the detection model for comparison of model accuracy. Studies applying DT techniques to detect fraudulent financial statements. The judgment accuracy rates of using data mining techniques to detect fraudulent financial statements vary and the construction of the model is neither complete nor perfect. As stated above, most studies only use 1–2 data mining techniques, without offering model comparison and most use one-stage statistical treatment to establish the detection model.

**Figure 3.1 Research design and procedure**

## 3.2 PROPOSED METHOD

The proposed method to detect the fraudulent financial statements using the credit card dataset with various models such as Random Forest, Artificial Neural Network(ANN), Gradient Boosting Machine(GBM), Generalized Linear method(GLM). The purpose of the project is to construct a valid and rigorous fraudulent financial statement detection model. Analysis of financial details is one of those simple methods to identify frauds. The aim of the research is to distinguish financial details, the values of which could indicate the fraud in financial statements. Moreover, the logistic regression model of fraud detection in financial statements has been developed. In the first stage, for the variable selection process appropriate deviations are used with certain constrains. Random forest or random decision forest method is used which is an ensemble learning method for classification, regression and other tasks. In the second stage, in order to construct fraudulent financial statement detection models ANN, GBM, GLM are used for the selection of more optimized result within these models. Data mining tools allow enterprises to predict future trends.

# CHAPTER 4

## SYSTEM REQUIREMENTS

### 4.1 SOFTWARE REQUIREMENTS

Operating system : Windows 8.1

Language : Python 3.6.4

### 4.2 HARDWARE REQUIREMENTS

Processor : Dual Core

Storage : 100GB

RAM : 4GB

### 4.3 PACKAGE DESCRIPTION

pandas : Provides data structures to work with structured, time series data.

seaborn : Library for making statistical graphics

H2O : Used for data modeling and general computing

matplotlib : Package publishes quality 2D graphics

numpy : General-purpose array-processing package

# CHAPTER 5

## SYSTEM DESIGN

### 5.1 OVERVIEW

In the existing system CHAID and CART were used. Since they are conventional algorithm the recent methods are used. They are Random Forest(RF), Artificial Neural Network(ANN), Gradient Boosting Machine(GBM) and Generalized Linear Model(GLM). The percentage of data to be used as training data, validation data and test data will differ for each and every method.



**Figure 5.1 Flow Diagram**

**5.2 MODULE DESCRIPTION**

**5.2.1 Random Forest**

Random forests is an ensemble learning method for classification and regression, that operate by constructing a multitude of decision tree at training time and outputting the class that is the mode of the class or mean prediction of the individual trees .Random forest uses full decision trees. It has low bias and high variance where variance can be defined as the error from the small fluctuations in the training data set and bias can be defined as the error from the spontaneous assumptions. The prediction describes a method of building a forest of uncorrelated trees using a CART like procedure, combined with randomized node optimization and bagging.

Important points are,

    i.  Using out-of-bag error as an estimate of the generalization error.

    ii.  Measuring variable importance through permutation.

The algorithm working is explained as follows,

1. The input for the Random Forest Algorithm is the credit card dataset which consists of 30 variables along with the class and amount variable.
2. The data will be split into training and test data. Then, the Random Forest model is constructed and then the testing of data follows.
3. The confusion matrix is built from which the accuracy can be calculated.

**5.2.2 Gradient Boosting Machine**

Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. The trees will be built with the help of previously built trees.

The algorithm working is explained as follows,

1. The credit card dataset is given as the input.
2. The data will be split into three parts such as train, test and valid data.

3. In training data, the model is built. In test data, the output will be predicted with the help of training data. The validation data will be used to validate whether the given output is valid or not.

### 5.2.3 Generalized Linear Model

In statistics, the generalized linear model (GLM) is a flexible generalization of ordinary linear regression that allows for response variables that have error distribution models other than a normal distribution.

The algorithm working is explained as follows,

1. The input for the Generalized Linear Model is the credit card dataset which consists of 30 variables which excludes the class and amount variable.
2. The data will be split into three parts such as train, test and valid data.
3. The default number of trees is set to 50 initially and the learning rate is increased from the default value.
4. The default number of trees is further decreased and the learning rate is increased from the previous value.
5. In training data, the model is built. In test data, the output will be predicted with the help of training data. The validation data will be used to validate whether the given output is valid or not.

### 5.2.5 Artificial Neural Network

Neural computing refers to a pattern recognition methodology for machine learning. The resulting model from neural computing is often called an Artificial Neural Network (ANN) or a neural network. Figure 5.1 represents that the ANN is composed of multiple nodes such as $X_1, X_2 \ldots X_i$, The result of these operations is passed to other neurons.

The output at each node is called its activation or node value. Each link is associated with weight $W_{1j}, W_{2j} \ldots W_{ij}$. ANNs are capable of learning, which takes place by altering weight values. Neural networks have been used in many business applications for pattern recognition, forecasting, prediction, and classification. Neural network computing is a key component of any data mining tool kit.

The algorithm working is explained as follows,

1. The credit card dataset is given as the input for Artificial Neural Network

2. The data will be split into three parts such as train, test and valid data.

3. Tuning is performed to reduce the error rate.

4. In training data, the model is built. In test data, the output will be predicted with the help of training data. The validation data will be used to validate whether the given output is valid or not.

5. The number of hidden layers are increased to predict the improved accuracy.



**Figure 5.2 Processing in an Artificial Neuron**

# CHAPTER 6

# SYSTEM IMPLEMENTATION

## 6.1 DATASET

The credit card details dataset is taken from the web resource www.kaggle.com. Figure 6.1 shows the sample details format from the www.kaggle.com. The credit card dataset consists of 30 variables in total. Since the credit card dataset deals with confidential information, the names of the attributes used in the credit card dataset is not revealed for the security purposes.

| Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | V10 | V11 | V12 | V13 | V14 | V15 | V16 | V17 | V18 | V19 | V20 | V21 | V22 | V23 | V24 | V25 | V26 | V27 | V28 | Amount | Class |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Figure 6.1 Credit Card Dataset**

Each method used in this model takes input from the same credit card dataset. Instead of the names, 28 variables out of the total 30 are named from V1 to V28 and the last two variables are named as amount and class. Along with the credit card dataset, spiral and grid are two other datasets which are used for Artificial Neural Network(ANN) method. The value used in class variable is either 0 or 1. If the financial statement is fraud then the predicted value should be set to 1 and if the financial statement is not fraud then the predicted value is set to 0.

## 6.2 RESULT ANALYSIS

The credit card dataset is given as input for all the four methods separately and the accuracy of each method will be compared with the accuracy of other methods of the fraudulent financial statement detection model. When the attribute named class has the predicted value as 1 then the financial statement is considered as a fraudulent financial statement and if the value is predicted to be 0 then financial statement is fraud free.

Random Forest algorithm splits the dataset into two parts for processing i.e., 70% as training data and remaining 30% as test data. Since the Random Forest method deals with fully grown decision trees, the accuracy of this method is lower than other methods. Random Forest Algorithm has low bias value and higher variance value.

Artificial Neural Network method splits the dataset into the following manner i.e., 60% as training data, 20% as validation data and remaining 20% as test data. The implementation of this method gives us the maximum accuracy of 87%. Gradient Boosting Machine Algorithm is considered as a better alternative algorithm for random forest because Gradient Boosting Machine constructs a new tree with the help of previously constructed tree and also this algorithm has high bias and low variance but the Random Forest Algorithm has low bias and high variance. The values predicted in the Gradient Boosting Machine Algorithm are mean square error , root mean square error, mean absolute error and root mean squared logarithmic error.

Generalized Linear Model is mainly used in statistical analysis and the data is split into three such as 60% as training data, 20% as validation data and remaining 20% as test data. The multinomial model will be constructed initially with the default parameters with the help of Limited-memory Broyden–Fletcher–Goldfarb–Shanno Algorithm. The default lambda

value is 0.001. In order to update the Generalized Linear Model Estimator the lambda value is lowered to 0.0001 which leads to improvement in the mean square error value, mean absolute error value and then the confusion matrix is built. Two classes are used in this model which implies fraud and not fraud respectively. A binomial classifier is built and features are added to the binomial model to improve its predictive capacity and then the final multinomial classifier with a lambda value of 0.001 is constructed. The hit ratio gradually increases since the first multinomial model.



**Figure 6.2 Training Data Result**

**Figure 6.3 Test Data Result**

Figure 6.2 and Figure 6.3 represents the histogram representation of training data and test data respectively. With the help of the true positive rate and false positive rate calculated in the confusion matrix the output graph which consists of receiver operating characteristic curve is produced.

# CHAPTER 7

# CONCLUSION AND FUTURE ENHANCEMENT

## 7.1 CONCLUSION

A company's financial statement is the key basis for all investor judgments, and is the last line of defense for investor interests. If management attempts to withhold information, even when independent CPAs, investment banking and securities analysts are involved, investors can experience significant losses. Based on the empirical results of this study, the accuracy of the DT CHAID, combined with CART, in detecting fraudulent financial statements, is relatively high. It can therefore be used as a tool to help auditors in the detection of fraudulent financial statements. The research findings can provide a reference for investors, shareholders, company managers, credit rating institutions, auditors, CPAs (certified public accountants), securities analysts, financial regulatory authorities, and relevant academic institutions.

## 7.2 FUTURE ENHANCEMENT

In our future scope we planned to detect the fraudulent financial statements of any organization using artificial intelligence concept which prevents the financial loss incurred by the organization and improves the accuracy. The project can be extended in such a way that methods with higher accuracy will be used to detect the fraudulent financial statements effectively which helps to detect the fraudulent behavior prior to its occurrence.

# APPENDIX I - CODING

**RANDOM FOREST:**

```python
import pandas as pd

import matplotlib.pyplot as plt

import matplotlib.gridspec as gridspec

import seaborn as sns

import os

import h2o

import numpy as np

data = reviews = pd.read_csv('dataset.csv')

data['Class'] = data['Class'].astype('category')

cols = list(data.columns)

data.groupby('Class')['Class'].count()

h2o.init()

creditcard = h2o.upload_file(path ='dataset.csv')

creditcard['Class'] = creditcard['Class'].asfactor()

df = creditcard.as_data_frame()

cols = list(df.columns)

df.isnull().sum()

df.query('Amount == 0').groupby('Class').count()

df.query('Amount != 0').to_csv('creditcard_amount_positive.csv')

creditcard = h2o.upload_file(path ='creditcard_amount_positive.csv')

creditcard['Class'] = creditcard['Class'].asfactor()

df = creditcard.as_data_frame()

from scipy import stats

variables = cols[:-1]
```

```
gs = gridspec.GridSpec(num_cols, 1)

for i, cn in enumerate(df[cols[:-1]]):

    ax = plt.subplot(gs[i])

    sns.distplot(df[cn][df.Class == 1], bins=50)

    sns.distplot(df[cn][df.Class == 0], bins=50)

    ax.set_xlabel('')

    ax.set_title('histogram of feature: ' + str(cn))

plt.show()

from scipy import stats

variables = cols[:-1]

keep = []

p_value_alpha = 0.05 #defult p-value for statistical significance

for variable in variables:

    fraud_v = df[variable][df.Class == 1]

    not_fraud_v = df[variable][df.Class == 0].sample(len(fraud_v))

    p_value = stats.ttest_ind(not_fraud_v, fraud_v).pvalue

    if p_value >= p_value_alpha:

        print("Distributions are equal. Discard {} variable".format(variable))

    else:

        print("Distributions are diferent. Keep {} variable".format(variable))

        keep.append(variable)

from h2o.estimators.random_forest import H2ORandomForestEstimator

train, valid = creditcard.split_frame(ratios=[0.7])

response_var = 'Class'

features = [col for col in cols if col != response_var]

naive_rf_model = H2ORandomForestEstimator()
```

```python
performance_train = naive_rf_model.model_performance(train=True)
import itertools
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_curve, auc
def plot_confusion_matrix(cm, classes,
                          normalize=False,
                          title='Confusion matrix',
                          cmap=plt.cm.Blues):
    if normalize:
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
        print("Normalized confusion matrix")
    else:
        print('Confusion matrix, without normalization')

      print(cm)
    plt.imshow(cm, interpolation='nearest', cmap=cmap)
    plt.title(title)
    plt.colorbar()
    tick_marks = np.arange(len(classes))
    plt.xticks(tick_marks, classes, rotation=45)
    plt.yticks(tick_marks, classes)
    fmt = '.2f' if normalize else 'd'
    thresh = cm.max() / 2.
    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
        plt.text(j, i, format(cm[i, j], fmt),
                 horizontalalignment="center",
                 color="white" if cm[i, j] > thresh else "black")
    plt.tight_layout()
```

```
    plt.ylabel('True label')

    plt.xlabel('Predicted label')

preds = naive_rf_model.predict(train)

cm = confusion_matrix(train.as_data_frame()['Class'], preds.as_data_frame()['predict'])

plot_confusion_matrix(cm, ['Non-Fraud', 'Fraud'], False)

fpr, tpr, threshold = roc_curve(train.as_data_frame()['Class'],
preds.as_data_frame()['predict'])

plt.plot(fpr, tpr, color='darkorange',

      lw=2, label='ROC curve')

plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')

plt.xlim([-0.04, 1.0])

plt.ylim([-0.04, 1.05])

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate')

plt.title('Curva ROC')

plt.legend(loc="lower right")

plt.show()

predictions = naive_rf_model.predict(valid)

cm = confusion_matrix(valid.as_data_frame()['Class'], predictions.as_data_frame()['predict'])

plot_confusion_matrix(cm, ['Non-Fraud', 'Fraud'], False)

fpr, tpr, threshold = roc_curve(valid.as_data_frame()['Class'],
predictions.as_data_frame()['predict'])

plt.plot(fpr, tpr, color='darkorange',

      lw=2, label='ROC curve')

plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')

plt.xlim([-0.04, 1.0])

plt.ylim([-0.04, 1.05])

plt.xlabel('False Positive Rate')
```

```
plt.ylabel('True Positive Rate')

plt.title('Curva ROC')

plt.legend(loc="lower right")

plt.show()
```

**DEEP LEARNING:**

```
import h2o

h2o.init(max_mem_size = 2)            #uses all cores by default

h2o.remove_all()                      #clean slate, in case cluster was already running

from h2o.estimators.deeplearning import H2OAutoEncoderEstimator,
H2ODeepLearningEstimator

import matplotlib.pyplot as plt

import numpy as np

import pandas as pd

from h2o.estimators.deeplearning import H2OAutoEncoderEstimator,
H2ODeepLearningEstimator

from h2o.estimators.gbm import H2OGradientBoostingEstimator

import os

spiral = h2o.import_file(path = os.path.realpath("input/spiral.csv"))

grid  = h2o.import_file(path = os.path.realpath("input/grid.csv"))

spiral_df = spiral.as_data_frame(use_pandas=True)

grid_df = grid.as_data_frame(use_pandas=True)

grid_x, grid_y = grid_df.x.reshape(201,201), grid_df.y.reshape(201,201)

spiral_r = spiral_df[spiral_df.color == "Red"]

spiral_k = spiral_df[spiral_df.color == "Black"]

spiral_xr, spiral_yr = spiral_r[spiral_r.columns[0]], spiral_r[spiral_r.columns[1]]

spiral_xk, spiral_yk = spiral_k[spiral_k.columns[0]], spiral_k[spiral_k.columns[1]]
```

```
markersize_ = 7**2

plt.figure(figsize = (5,5))

plt.scatter(spiral_xr, spiral_yr, c = 'r', s=markersize_)

plt.scatter(spiral_xk, spiral_yk, c = 'k', s=markersize_)

plt.axis([-1.5, 1.5, -1.5, 1.5])

plt.title("Spiral");

X = spiral.col_names[0:2]

y = spiral.col_names[2]

dl_model = H2ODeepLearningEstimator(epochs=1000)

dl_model.train(X, y, spiral)

gbm_model = H2OGradientBoostingEstimator()

gbm_model.train(X, y, spiral)

drf_model = H2ORandomForestEstimator()

drf_model.train(X, y, spiral)

glm_model = H2OGeneralizedLinearEstimator(family="binomial")

glm_model.fit(spiral[X], spiral[y])                    #model.fit() example

models = [dl_model, gbm_model, drf_model, glm_model]

m_names = ["Deep Learning", "Gradient Boosted Method", "Distributed Random Forest",
"Generalized Linear Model"]

def plot_spirals(models, model_names):

    fig, ax = plt.subplots(2,2, figsize=(12,12))

    for k, subplot in enumerate(ax.flatten()):

        subplot.scatter(spiral_xr, spiral_yr, c = 'r', s=markersize_)

        subplot.scatter(spiral_xk, spiral_yk, c = 'k', s=markersize_)

        subplot.axis([-1.5, 1.5, -1.5, 1.5])

        subplot.set_title(model_names[k])

        subplot.set_ylabel('y')

        pred_z = models[k].predict(grid).as_data_frame(True)
```

```
    subplot.contour(grid_x, grid_y, (pred_z['predict'] ==
'Black').astype(np.int).reshape(201,201), colors='b')
```

```
dl_1 = H2ODeepLearningEstimator(epochs=1)
```

```
dl_1.train(X, y, spiral)
```

```
dl_250 = H2ODeepLearningEstimator(checkpoint=dl_1, epochs=250)
```

```
dl_250.train(X, y, spiral)
```

```
dl_500 = H2ODeepLearningEstimator(checkpoint=dl_250, epochs=500)
```

```
dl_500.train(X, y, spiral)
```

```
dl_750 = H2ODeepLearningEstimator(checkpoint=dl_500, epochs=750)
```

```
dl_750.train(X, y, spiral)
```

```
models_dl = [dl_1, dl_250, dl_500, dl_750]
```

```
m_names_dl = ["DL " + str(int(model.get_params()['epochs']['actual_value'])) + " Epochs"
for model in models_dl]
```

```
plot_spirals(models_dl, m_names_dl)
```

```
dl_1 = H2ODeepLearningEstimator(hidden=[1000], epochs=500)
```

```
dl_1.train(X, y, spiral)
```

```
dl_2 = H2ODeepLearningEstimator(hidden=[200,200], epochs=500)
```

```
dl_2.train(X, y, spiral)
```

```
dl_3 = H2ODeepLearningEstimator(hidden=[42,42,42], epochs=500)
```

```
dl_3.train(X, y, spiral)
```

```
dl_4 = H2ODeepLearningEstimator(hidden=[11,13,17,19], epochs = 1000)
```

```
dl_4.train(X, y, spiral)
```

```
# It is clear that different configurations can achieve similar performance, and that tuning will
be required for optimal performance.
```

```
models_network = [dl_1, dl_2, dl_3, dl_4]
```

```
m_names_network = ["1000", "200 x 200", "42 x 42 x 42", "11 x 13 x 17 x 19"]
```

```
plot_spirals(models_network, m_names_network)
```

```
models_act = []
```

```
m_names_act = []
```

```python
for i,method in enumerate(["Tanh","Maxout","Rectifier","RectifierWithDropout"]):

    models_act.append(H2ODeepLearningEstimator(activation=method, hidden=[100,100],
epochs=1000))

    models_act[i].train(X, y, spiral)

    m_names_act.append("DL "+ method + " Activation")

plot_spirals(models_act, m_names_act)

creditcard_df = h2o.import_file(path = os.path.realpath("datast.csv"))

train, valid, test = creditcard_df.split_frame([0.6, 0.2], seed=1234)

creditcard_X = creditcard_df.col_names[:-1]     #last column is cover_type,

creditcard_y = creditcard_df.col_names[-1]

creditcard_model_v1 = H2ODeepLearningEstimator(model_id="creditcard_v1", epochs=1,
variable_importances=True)

creditcard_model_v1.train(creditcard_X, creditcard_y, training_frame = train,
validation_frame = valid)

print creditcard_model_v1

var_df = pd.DataFrame(creditcard_model_v1.varimp(),

        columns=["Variable", "Relative Importance", "Scaled Importance", "Percentage"])

print var_df.shape

var_df.head(10)

creditcard_model_v2 = H2ODeepLearningEstimator(

    model_id="creditcard_v2",

    hidden=[32,32,32],              ## small network, runs faster

    epochs=1000000,                  ## hopefully converges earlier...

    score_validation_samples=10000,     ## sample the validation dataset (faster)

    stopping_rounds=2,

    stopping_metric="MSE", ## alternatives: "MSE","logloss","r2"

    stopping_tolerance=0.01)

creditcard_model_v2.train(creditcard_X, creditcard_y, training_frame=train,
validation_frame=valid)
```

```
cov_v2_df = creditcard_model_v2.score_history()

cov_v2_df

plt.plot(cov_v2_df['training_classification_error'], label="training_classification_error")

plt.plot(cov_v2_df['validation_classification_error'], label="validation_classification_error")

plt.title("Covertype Deep Learner (Early Stop)")

plt.legend();

# ### Adaptive Learning Rate

creditcard_model_tuned = H2ODeepLearningEstimator(

  model_id="creditcard_tuned",

  overwrite_with_best_model=False,

  hidden=[128,128,128],          # more hidden layers -> more complex interactions

  epochs=10,                # to keep it short enough

  score_validation_samples=10000,  # downsample validation set for faster scoring

  score_duty_cycle=0.025,        # don't score more than 2.5% of the wall time

  adaptive_rate=False,          # manually tuned learning rate

  rate=0.01,

  rate_annealing=0.000002,

  momentum_start=0.2,          # manually tuned momentum

  momentum_stable=0.4,

  momentum_ramp=10000000,

  l1=0.00001,                # add some L1/L2 regularization

  l2=0.00001,

  max_w2=10                  # helps stability for Rectifier

)

creditcard_model_tuned.train(creditcard_X, creditcard_y, training_frame=train,
validation_frame=valid)

cov_tuned_df = creditcard_model_tuned.scoring_history()

cov_tuned_df
```

```python
plt.plot(cov_tuned_df['training_classification_error'], label="training_classification_error")

plt.plot(cov_tuned_df['validation_classification_error'],
label="validation_classification_error")

plt.title("Covertype Deep Learner (Tuned)")

plt.legend();

pred = creditcard_model_tuned.predict(test[0:-1]).as_data_frame(use_pandas=True)

test_actual = test.as_data_frame(use_pandas=True)['Cover_Type']

(test_actual == pred['predict']).mean()
```

## GRADIENT BOOSTING METHOD:

```python
import h2o

import os

h2o.init(max_mem_size = "2G")          #specify max number of bytes. uses all cores by default.

h2o.remove_all() #clean slate, in case cluster was already running

from h2o.estimators.gbm import H2OGradientBoostingEstimator

creditcard_df = h2o.import_file(os.path.realpath("dataset.csv"))

# 60% for training

# 20% for validation (hyper parameter tuning)

# 20% for final testing

#split the data as described above

train, valid, test = creditcard_df.split_frame([0.6, 0.2], seed=1234)

#Prepare predictors and response columns

creditcard_X = creditcard_df.col_names[:-1]     #last column is Class, our desired response variable

creditcard_y = creditcard_df.col_names[-1]

gbm_v1 = H2OGradientBoostingEstimator(

    model_id="gbm_creditcard_v1",
```

```
gbm_v1.train(creditcard_X, creditcard_y, training_frame=train, validation_frame=valid)

gbm_v1.score_history()

gbm_v1.hit_ratio_table(valid, train = FALSE, valid = FALSE, xval = FALSE)

# ### GBM Round 2

gbm_v2 = H2OGradientBoostingEstimator(

    ntrees=20,

    learn_rate=0.2,

    max_depth=10,

    stopping_tolerance=0.01, #10-fold increase in threshold as defined in rf_v1

    stopping_rounds=2,

    score_each_iteration=True,

    model_id="gbm_creditcard_v2",

    seed=2000000

gbm_v2.train(creditcard_X, creditcard_y, training_frame=train, validation_frame=valid)

gbm_v2.hit_ratio_table(valid=True)

# ### GBM: Third Time

gbm_v3 = H2OGradientBoostingEstimator(

    ntrees=30,

    learn_rate=0.3,

    max_depth=10,

    sample_rate=0.7,

    col_sample_rate=0.7,

    stopping_rounds=2,

    stopping_tolerance=0.01, #10-fold increase in threshold as defined in rf_v1

    score_each_iteration=True,

    model_id="gbm_creditcard_v3",

    seed=2000000
```

```
gbm_v3.train(creditcard_X, creditcard_y, training_frame=train, validation_frame=valid)

gbm_v3.hit_ratio_table(valid=True)
```

## GENERALIZED LINEAR METHOD

```
import h2o

h2o.init(max_mem_size = "2G")           #specify max number of bytes. uses all cores by
default.

h2o.remove_all()                        #clean slate, in case cluster was already running

from h2o.estimators.glm import H2OGeneralizedLinearEstimator

import pandas as pd

import numpy as np

creditcard_df = h2o.import_file(os.path.realpath("dataset.csv"))

train, valid, test = creditcard_df.split_frame([0.7, 0.15], seed=1234)

creditcard_X = creditcard_df.col_names[:-1]     #last column is Cover_Type, our desired
response variable

creditcard_y = creditcard_df.col_names[-1]

glm_multi_v1 = H2OGeneralizedLinearEstimator(

        model_id='glm_v1',          #allows us to easily locate this model in Flow

        family='multinomial',

        solver='L_BFGS')

glm_multi_v1.train(creditcard_X, creditcard_y, training_frame=train,
validation_frame=valid)

glm_multi_v1

glm_multi_v1.hit_ratio_table(valid=True)

glm_multi_v2 = H2OGeneralizedLinearEstimator(

        model_id='glm_v2',

        family='multinomial',

        solver='L_BFGS',
```

```
          Lambda=0.0001          #default value 0.001
)

glm_multi_v2.train(creditcard_X, creditcard_y, training_frame=train,
validation_frame=valid)

glm_multi_v2

glm_multi_v2.hit_ratio_table(valid=True)

glm_multi_v2.confusion_matrix(valid)

c1 = creditcard_df[creditcard_df['Cover_Type'] == 'class_1']

c2 = creditcard_df[creditcard_df['Cover_Type'] == 'class_2']

df_b = c1.rbind(c2)

train_b, valid_b, test_b = df_b.split_frame([0.7, 0.15], seed=1234)

glm_binom_v1 = H2OGeneralizedLinearEstimator(

          model_id="glm_v3",

          solver="L_BFGS",

          family="binomial")

glm_binom_v1.train(creditcard_X, creditcard_y, training_frame=train_b,
validation_frame=valid_b)

glm_binom_v1.accuracy()

def cut_column(train_df, train, valid, test, col):

   only_col= train_df[col]                    #Isolate the column in question from the training
frame

   counts, breaks = np.histogram(only_col, bins=20)   #Generate counts and breaks for our
histogram

   min_val = min(only_col)-1                 #Establish min and max values

   max_val = max(only_col)+1

   new_b = [min_val]                         #Redefine breaks such that each bucket has
enough support

   for i in xrange(19):

      if counts[i] > 1000 and counts[i+1] > 1000:
```

```
        new_b.append(breaks[i+1])

    new_b.append(max_val)

    names = [col + '_' + str(x) for x in xrange(len(new_b)-1)]  #Generate names for buckets,
these will be categorical names

    train[col+"_cut"] = train[col].cut(breaks=new_b, labels=names)

    valid[col+"_cut"] = valid[col].cut(breaks=new_b, labels=names)

    test[col+"_cut"] = test[col].cut(breaks=new_b, labels=names)

def add_features(train, valid, test):

    train_df = train.as_data_frame(True)

    cut_column(train_df, train, valid, test, "Elevation")

    cut_column(train_df, train, valid, test, "Hillshade_Noon")

    cut_column(train_df, train, valid, test, "Hillshade_9am")

    cut_column(train_df, train, valid, test, "Hillshade_3pm")

    cut_column(train_dfs, train, valid, test, "Horizontal_Distance_To_Hydrology")

    cut_column(train_df, train, valid, test, "Slope")

    cut_column(train_df, train, valid, test, "Horizontal_Distance_To_Roadways")

    cut_column(train_df, train, valid, test, "Aspect")

    interaction_cols1 = ["Elevation_cut",

                "Wilderness_Area",

                "Soil_Type",

                "Hillshade_Noon_cut",

                "Hillshade_9am_cut",

                "Hillshade_3pm_cut",

                "Horizontal_Distance_To_Hydrology_cut",

                "Slope_cut",

                "Horizontal_Distance_To_Roadways_cut",

                "Aspect_cut"]

    train_cols = train.interaction(factors=interaction_cols1,    #Generate pairwise columns
```

```
                    pairwise=True,

                    max_factors=1000,

                    min_occurrence=100,

                    destination_frame="itrain")

valid_cols = valid.interaction(factors=interaction_cols1,

                    pairwise=True,

                    max_factors=1000,

                    min_occurrence=100,

                    destination_frame="ivalid")

test_cols = test.interaction(factors=interaction_cols1,

                    pairwise=True,

                    max_factors=1000,

                    min_occurrence=100,

                    destination_frame="itest")

train = train.cbind(train_cols)                #Append pairwise columns to H2OFrames

valid = valid.cbind(valid_cols)

test = test.cbind(test_cols)

#Add a three-way interaction for Hillshade

interaction_cols2 = ["Hillshade_Noon_cut","Hillshade_9am_cut","Hillshade_3pm_cut"]

train_cols = train.interaction(factors=interaction_cols2,    #Generate pairwise columns

                    pairwise=False,

                    max_factors=1000,

                    min_occurrence=100,

                    destination_frame="itrain")

valid_cols = valid.interaction(factors=interaction_cols2,

                    max_factors=1000,
```

```
                    min_occurrence=100,

                    destination_frame="ivalid")

   test_cols = test.interaction(factors=interaction_cols2,

   train = train.cbind(train_cols)                    #Append pairwise columns to H2OFrames

   valid = valid.cbind(valid_cols)

   test = test.cbind(test_cols)

   return train, valid, test

train_bf, valid_bf, test_bf = add_features(train_b, valid_b, test_b)

glm_binom_feat_1 = H2OGeneralizedLinearEstimator(family='binomial', solver='L_BFGS',
model_id='glm_v4')

glm_binom_feat_1.train(creditcard_X, creditcard_y, training_frame=train_bf,
validation_frame=valid_bf)

glm_binom_feat_1.accuracy(valid=True)

glm_binom_feat_2 = H2OGeneralizedLinearEstimator(family='binomial', solver='L_BFGS',
model_id='glm_v5', Lambda=0.001)

glm_binom_feat_2.train(creditcard_X, creditcard_y, training_frame=train_bf,
validation_frame=valid_bf)

glm_binom_feat_3 = H2OGeneralizedLinearEstimator(family='binomial',
model_id='glm_v6', lambda_search=True)

glm_binom_feat_3.train(creditcard_X, creditcard_y, training_frame=train_bf,
validation_frame=valid_bf)

glm_binom_feat_3.accuracy(valid=True)

train_f, valid_f, test_f = add_features(train, valid, test)

glm_multi_v3 = H2OGeneralizedLinearEstimator(

            solver='L_BFGS',

            Lambda=0.0001)

glm_multi_v3.train(creditcard_X, creditcard_y, training_frame=train_f,
validation_frame=valid_f)

glm_multi_v3.hit_ratio_table(valid=True)
```
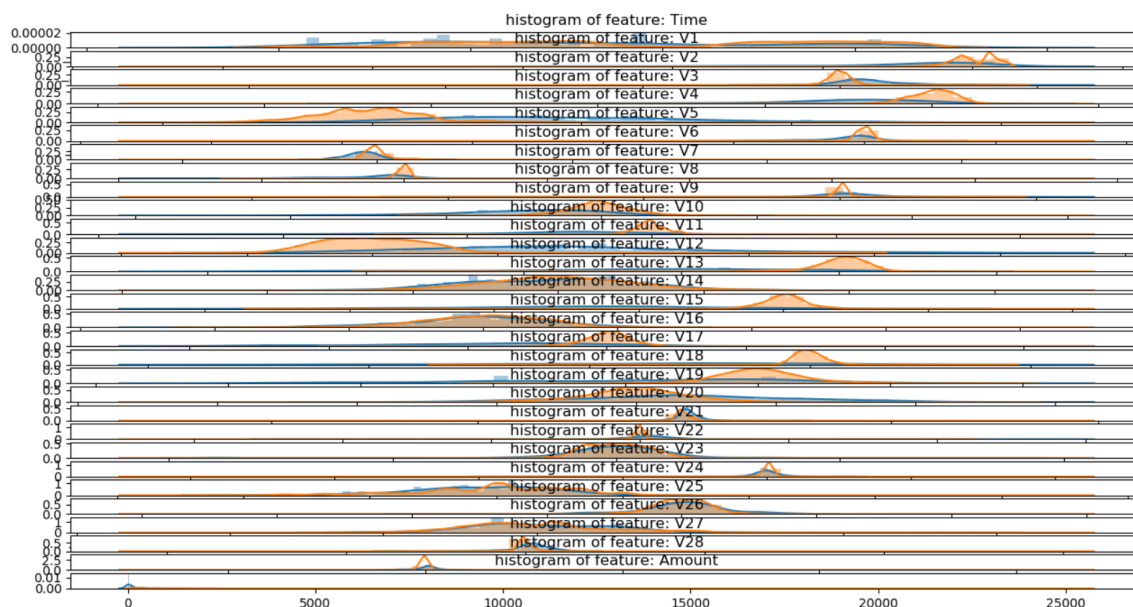
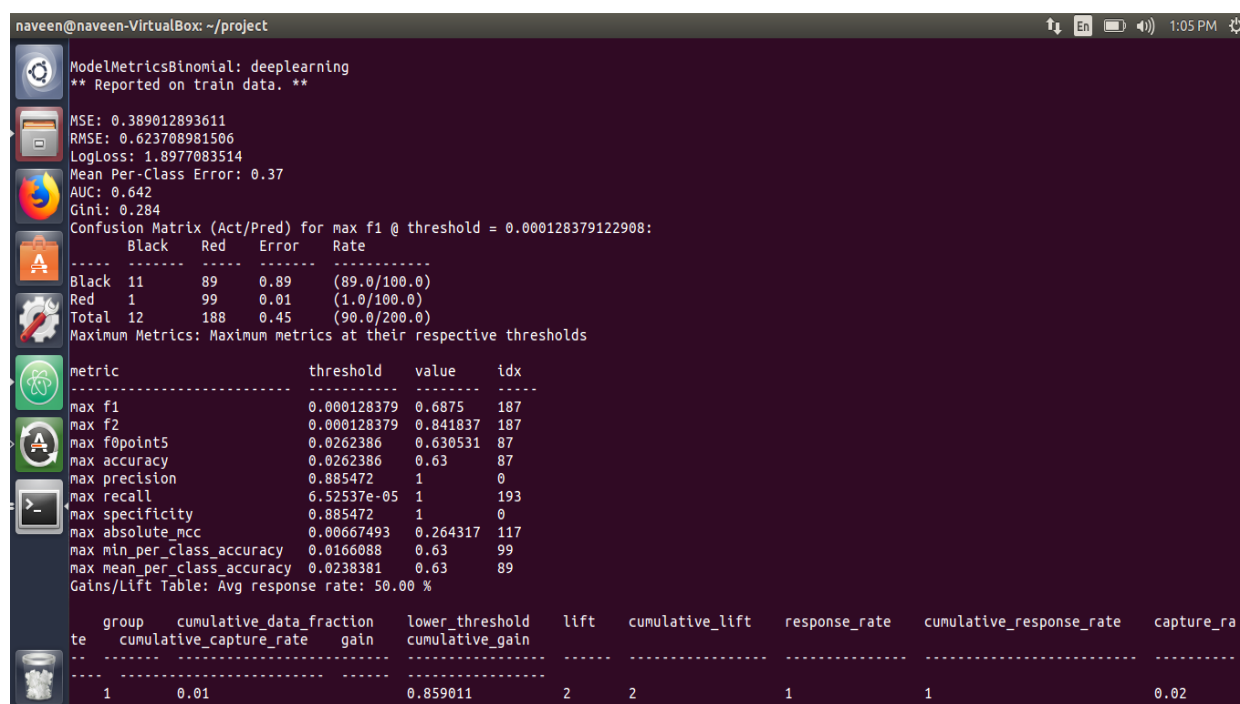# APPENDIX 2 - SCREENSHOTS



**Figure A.2.1 Variable Selection**



**Figure A.2.2 Accuracy of ANN**

**Figure A.2.3 GBM Result**

**REFERENCES**

1. Chen S (2016), 'Detection of fraudulent financial statements using the hybrid data mining approach', SpringerPlus, Vol.5, No.1, pp. 89.

2. Kamarudin K A et al., (2012), 'Aggressive financial reporting and corporate fraud', Procedia-Social and Behavioral Sciences,Vol.65, pp.638-643.

3. Ravisankar P et al., (2011), 'Detection of financial statement fraud and feature selection using data mining techniques', Decision Support Systems, Vol.50,No.2, pp. 491-500.

4. Yeh C C et al., (2010), 'A hybrid approach of DEA, rough set and support vector machines for business failure prediction', Expert Systems with Applications, Vol.37,No.2, pp.1535-1541.

5. Zhou W et al., (2011), 'Detecting evolutionary financial statement fraud', Decision Support Systems, Vol.50,No.3, pp.570-575.