



CREDIT CARD FRAUD DETECTION USING MACHINE LEARNING ALGORITHMS

Mrs. PINNAMARAJU.T.S.PRIYA ^{*1}, Ms. MADDI HINDUDEVI ^{*2}

^{*1} Assistant Professor, Department of Computer Science and Application,
Sanketika Vidhya Parishad Engineering College, P.M. Palem,
Visakhapatnam, Andhra Pradesh

^{*2} MCA Student, Department of Computer Science and Application,
Sanketika Vidhya Parishad Engineering College, P.M. Palem,
Visakhapatnam, Andhra Pradesh
maddihindudevi@gmail.com ^{#2}

ABSTRACT

Credit card fraud is a big issue in the financial services industry. Every year, billions of dollars are lost due to credit card theft. Due to confidentiality concerns, there is a scarcity of research studies analysing real-world credit card data. Machine learning techniques are employed to detect credit card fraud in this article. Standard models are initially employed. Then, hybrid methods based on Ada Boost and majority voting are used. A publicly available credit card data set is used to assess the model's performance. Following that, a real-world credit card data set obtained from a financial institution is analysed. In addition, noise is introduced into the data samples to test the robustness of the algorithms. The experimental results show that the majority voting method provides high accuracy rates in detecting credit card fraud.

KEY WORDS:

Credit card fraud, Ada Boost, Hybrid Models, Machine Learning Algorithm, Confidentiality, Financial Services.

1. INTRODUCTION

Fraud is defined as improper or criminal deception with the intent of gaining financial or personal gain. Two strategies can be employed to avoid fraud losses: fraud prevention and fraud detection. Fraud prevention is a proactive approach that prevents fraud from occurring in the first place. On the other hand, fraud detection is required when a fraudster attempts a fraudulent transaction. Credit card fraud is the unauthorized use of credit card information to make purchases. Credit card transactions can be done both physically and digitally. When conducting physical transactions, the credit card is used. This can happen over the phone or the internet in digital transactions. Cardholders commonly supply the card number, expiration date, and card verification number over the phone or in the mail.

With the advent of e-commerce over the last decade, the use of credit cards has skyrocketed. In 2011, there were over 320 million credit card transactions in Malaysia, which climbed to approximately 360 million in 2015. The number of fraud instances has steadily climbed in tandem with the rise in credit card usage. Despite numerous authorization techniques in place, credit card fraud cases have not been effectively thwarted. Fraudsters prefer the internet because it conceals their identity and location. The surge in credit card fraud has a significant influence on the financial sector. In 2015, global credit card theft totaled a whopping USD \$21.84 billion.

2. LITERATURE SURVEY

In this section we will mainly discuss about the background work that is carried out in order to prove the performance of our proposed Method. Literature survey is the most important step in software development process. For any software or application development, this step plays a very crucial role by determining the several factors like time, money, effort, lines of code and company strength. Once all these several factors are satisfied, then we need to determine which operating system and language used for developing the application. Once the programmers start building the application, they will first observe what are the pre-defined inventions that are done on same concept and then they will try to design the task in some innovated manner.

MOTIVATION

BORA MEHAR SRI SATYA TEJA et al (2022) have reported a study using supervised ML techniques. They have used random forest technique to find fraudulent transactions and Accuracy. For classification of the dataset decision trees are used they also defined that random forest would work efficiently by increasing the Accuracy after oversampling and SMOTE is used to increase classes in data set in balanced way to increase the Accuracy.

Zahra Faraji (2022) has examined Logistic Regression, Decision Tree, Random Forest, XGBoost, KNN and Ensemble, reported that Logistic Regression execution is superior to Decision Tree. XGBoost is quickest and have superior execution. KNN and Logistic Regression have better execution in detecting better and concluded that the model complexity doesn't ensure great performance.

B. N. V. Madhubabu et al (2021) have reported using supervised ML techniques. They have used random forest algorithm for tracking the fraudulent transactions and the Accuracy. For classification of the dataset decision trees are used. Using confusion matrix performance of RFA is evaluated. Although RFA produced good results with small data set, it is not accurate for imbalanced dataset.

Andhavarapu Bhanusri et al (2020) have proposed a study where various machine learning algorithms are compared with each other to evaluate the best classifier. They implemented various machine learning techniques on an imbalanced dataset such as Logistic Regression, naive bayes, random forest with ensemble classifiers utilizing boosting technique. This evaluation is done in view of the calculation's quantitative estimates like Accuracy, Precision, Recall, f1 score, support, confusion matrix. By contrasting the above three techniques, it was concluded that random forest classifier with boosting technique is superior to Logistic Regression and naive bayes techniques. In spite of the fact that Random Forest with Boosting technique performs best in this situation,

utilizing these three techniques we cannot decide the names of fraud and legitimate exchanges for the given dataset using machine learning.

S P Maniraj et al (2019) have reported a study using ML techniques. They have used Local outlier factor to measure the local length of the sample with respect to its neighbors and isolation forest algorithm for the feature selection and to split the data set between the least and greatest qualities for finding fraudulent transactions and we compare the techniques for the testing purpose to determine the Accuracy and Precision of those transactions. They used Jupyter notebook to program in python.

3. EXISTING METHODOLOGY

Three fraud detection approaches are described. To begin, a clustering model is employed to categorise legitimate and fraudulent transactions based on data parameter values. Second, Gaussian mixture models of past and present behaviour can be calculated to detect any deviations from the past. Finally, Bayesian networks are utilised to describe the statistics of a certain user as well as the statistics of other fraud scenarios.

LIMITATION OF EXISTING SYSTEM

1. The high amount of losses due to fraud and the awareness of the relation between loss and the available limit has to be reduced.
2. Testing credit card FDSs using real data set is a difficult task.
3. The fraud has to be deducted in real time and the number of false alert.

4. PROPOSED SYSTEM & ITS ADVANTAGES

For identifying credit card fraud, a total of twelve machine learning algorithms are deployed. Standard neural networks and deep learning models are among the algorithms used. Furthermore, the AdaBoost and majority voting methods are used to create hybrid models. The main contribution of this study is the evaluation of a range of machine learning models for fraud detection using a real-world credit card data set.

ADVANTAGES OF PROPOSED SYSTEM:

The following are the benefits of the proposed system. They are:

1. The results obtained by the several ML algorithms are accurate.
2. The performance of each and every individual ML algorithm differs with one another.
3. The proposed method greatly optimize the processing time.
4. We can finally predict the best algorithm after comparing several ML algorithms.

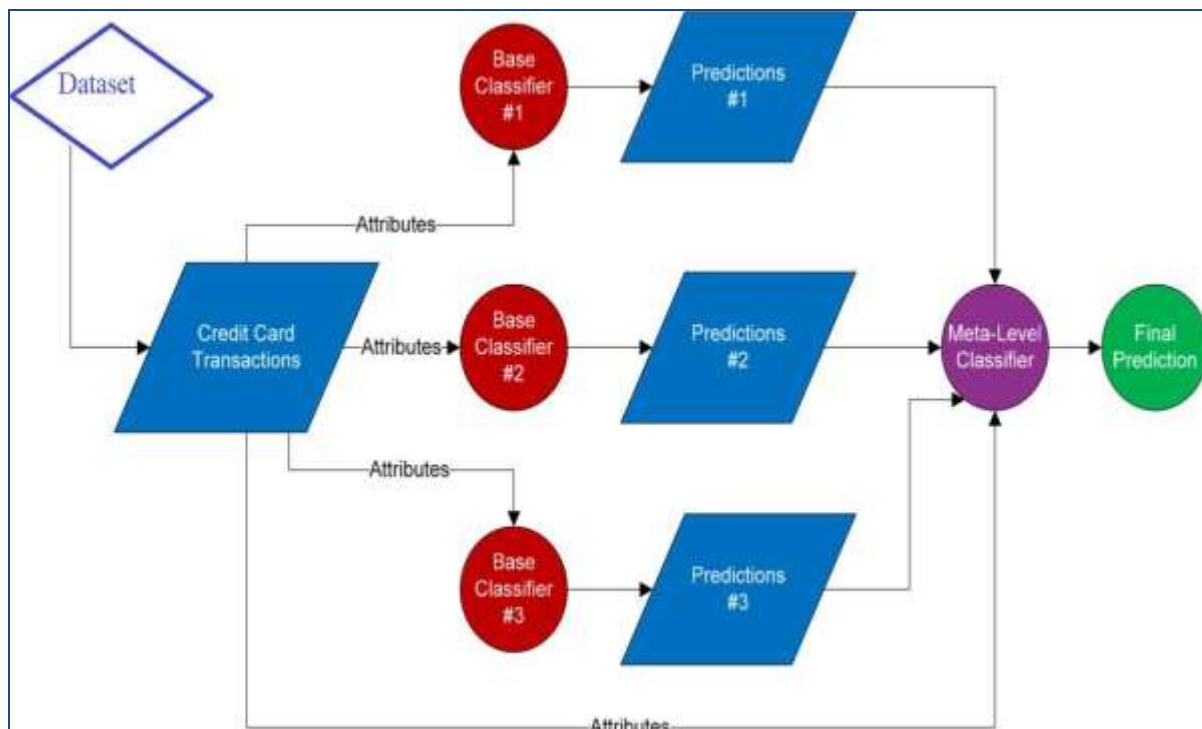


Figure 1. Denotes the Proposed Architecture

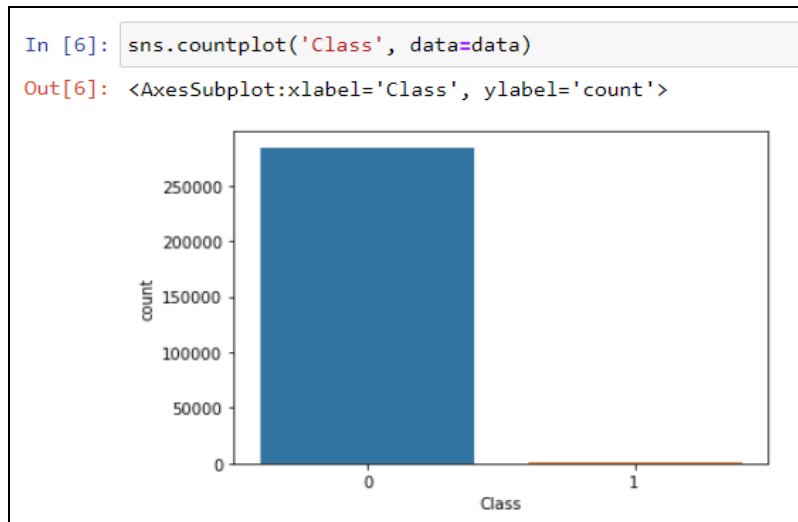
Figure 1 clearly represents the proposed architecture in which credit card fraud detection is predicted based on several classifiers. Here multiple classifiers are used to predict and based on the predictions we finally calculated the final prediction.

5. PROPOSED METHODOLOGY

In this proposed application we try to implement some ML algorithms to find out the credit card fraud transaction from public dataset.

A. ABOUT THE DATASET

We gathered the data from a public website, Kaggle. This dataset contains 284807 total records. These records are based on two days duration time data. Only a small amount of record is fraud. In the dataset 28 columns are already transformed by principal component analysis (PCA). The amount and time are the features that need to be scaled for acquiring better model. The dataset consists of not fraud class is 284315 and fraud class is 492. The count plot visualizes the imbalance data. The class 1 comprises of 0.172% of total dataset.



B. TRAINING AND TESTING

The highly imbalanced data that is typical in such applications, data from the two classes are sampled at different rates to obtain training data with reasonable proportion of fraud to non-fraud cases. As noted earlier, random under sampling of the majority class has been found to be generally better than other sampling approaches. We use random under sampling to obtain training dataset with varying proportions of fraud cases. Performance is observed on a separate Test dataset having 0.5% fraudulent transactions. As described in the data section, dataset has 492 observed fraudulent transactions. We Sampled legitimate transactions from dataset to create varying fraud rates in the modeling and test datasets. In other words, we kept the same number of fraudulent transactions in the modeling datasets, but varied the number of legitimate transactions from dataset to create varying fraud rates. Similarly, the actual fraud rates in the test dataset is 0.5%

ADAPTIVE BOOSTING

Adaptive Boosting is a machine learning meta-algorithm. It can be used in conjunction with many other types of learning algorithms to improve performance. The output of the other learning algorithms ('weak learners') is combined into a weighted sum that represents the final output of the boosted classifier. AdaBoost is adaptive in the sense that subsequent weak learners are tweaked in favor of those instances misclassified by previous classifiers. AdaBoost is sensitive to noisy data and outliers. In some problems it can be less susceptible to the overfitting problem than other learning algorithms. The individual learners can be weak, but as long as the performance of each one is slightly better than random guessing, the final model can be proven to converge to a strong learner. Every learning algorithm tends to suit some problem types better than others, and typically has many different parameters and configurations to adjust before it achieves optimal performance on a dataset, AdaBoost with decision trees as the weak learners is often referred to as the best out-of-the-box classifier. When used with decision tree learning, information gathered at each stage of the AdaBoost algorithm about the relative 'hardness' of each training sample is fed into the tree growing algorithm.

```

In [68]: from sklearn.ensemble import AdaBoostClassifier
adaboost = AdaBoostClassifier()
adaboost.fit(X_train,y_train)

Out[68]:
AdaBoostClassifier
AdaBoostClassifier()

In [69]: Y_pred_ad = adaboost.predict(X_test)

In [74]: ad_acc=sklearn.metrics.accuracy_score(y_test, Y_pred_ad)*100
ad_prec = sklearn.metrics.precision_score(y_test,Y_pred_ad)*100
ad_recall = sklearn.metrics.recall_score(y_test,Y_pred_ad)*100
ad_f1 = sklearn.metrics.f1_score(y_test,Y_pred_ad)*100
ad_Roc_auc = sklearn.metrics.roc_auc_score(y_test,Y_pred_ad)
print("ADABOOST:")
print("Accuracy:",ad_acc);
print("Precision:",ad_prec);
print("recall:",ad_recall);
print("F1-score:",ad_f1);
print("roc_auc:",ad_Roc_auc);
cm5 = confusion_matrix(y_test,Y_pred_ad)
print(cm5)
specificity = cm5[0,0]/(cm5[0,0]+cm5[0,1])
sensitivity = cm5[1,1]/(cm5[1,0]+cm5[1,1])
print("sensitivity:",sensitivity)
print("specificity:",specificity)

ADABOOST:
Accuracy: 96.41666666666666
Precision: 96.67832167832168
recall: 95.84055459272098
f1-score: 96.25761531766753
roc_auc: 0.9639539768159323
[[604 19]
 [ 24 553]]
sensitivity: 0.9584055459272097
specificity: 0.9695024077046549

```

NAÏVE BAYES

Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. There is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable. For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter. A naive Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the color, roundness, and diameter features.

LOGISTIC REGRESSION Logistic Regression is a supervised technique helps in classification. It provides output with the help of probabilistic values. The value will be 0 or 1. Logistic Regression follows the cost function called Sigmoid Function.

$$\text{Sigmoid Function } f(x) = \frac{1}{1 + e^{-x}}$$

SUPPORT VECTOR MACHINE

One of the supervised learning method is SVM. We use SVM classification technique called SVC. SVC is called Support Vector Classifier. For the given dataset Support Vector Machine technique creates its hyperplane. The main objective of SVM is creating hyperplane. For the given data, there will be many hyperplanes. Finding the optimized hyperplane is the task of SVM. Hyperplane is created with closest points of different classes. The points which are finalized are called support vectors. Every transaction undergoes through this hyperplane equation to classify the class it belongs to. This is about Support Vector Machine.

CATBOOST

In CatBoost, cat is category and boost is boosting. CatBoost does not required extensive training for providing better output. For credit card fraud detection we need classification so we use CatBoost Classifier. CatBoost develops symmetric trees. When we use the CatBoost Classifier without any parameters it creates 1000 decision trees from 0 to 999. These decision trees are used when we want to predict the output. Each input undergoes these decision trees to predict the class.

```
In [68]: from catboost import CatBoostClassifier
model = CatBoostClassifier()
catboost=model.fit(X_train, y_train)
```

```
Learning rate set to 0.015991
0:   learn: 0.6582874   total: 168ms   remaining: 2m 47s
1:   learn: 0.6288097   total: 176ms   remaining: 1m 27s
2:   learn: 0.5998417   total: 184ms   remaining: 1m 1s
3:   learn: 0.5740697   total: 192ms   remaining: 47.8s
4:   learn: 0.5500055   total: 198ms   remaining: 39.5s
5:   learn: 0.5253607   total: 205ms   remaining: 33.9s
6:   learn: 0.5046221   total: 211ms   remaining: 29.9s
7:   learn: 0.4855414   total: 217ms   remaining: 26.9s
8:   learn: 0.4655448   total: 225ms   remaining: 24.8s
9:   learn: 0.4463715   total: 235ms   remaining: 23.2s
10:  learn: 0.4304139   total: 242ms   remaining: 21.7s
11:  learn: 0.4154122   total: 248ms   remaining: 20.4s
12:  learn: 0.3991800   total: 256ms   remaining: 19.4s
13:  learn: 0.3845007   total: 263ms   remaining: 18.5s
14:  learn: 0.3692573   total: 271ms   remaining: 17.8s
15:  learn: 0.3543913   total: 279ms   remaining: 17.1s
16:  learn: 0.3410315   total: 286ms   remaining: 16.5s
17:  learn: 0.3289320   total: 293ms   remaining: 16s
```

```
In [69]: y_pred_cb=model.predict(X_test)
```

```
In [77]: cb_acc=sklearn.metrics.accuracy_score(y_test, y_pred_cb)*100
cb_prec = sklearn.metrics.precision_score(y_test,y_pred_cb)*100
cb_recall = sklearn.metrics.recall_score(y_test,y_pred_cb)*100
cb_f1 = sklearn.metrics.f1_score(y_test,y_pred_cb)*100
cb_Roc_auc = sklearn.metrics.roc_auc_score(y_test,y_pred_cb)
print("CATBOOSTING:")
print("Accuracy:",cb_acc);
print("Precision:",cb_prec);
print("recall:",cb_recall);
print("f1-score:",cb_f1);
print("roc_auc:",cb_Roc_auc);
cm4 = confusion_matrix(y_test,y_pred_cb)
print(cm4)
specificity = cm4[0,0]/(cm4[0,0]+cm4[0,1])
sensitivity = cm4[1,1]/(cm4[1,0]+cm4[1,1])
print("sensitivity:",sensitivity)
print("specificity:",specificity)
```

```
CATBOOSTING:
Accuracy: 97.58333333333333
Precision: 98.40989399293287
recall: 96.53379549393414
f1-score: 97.46281714785651
roc_auc: 0.9754458635049836
[[614  9]
 [ 20 557]]
sensitivity: 0.9653379549393414
specificity: 0.985553772070626
```

RESULTS

S. No	Algorithm	Accuracy	Precision	Recall	F1-Score	ROC_AUC
1	LOGISTIC REGRESSION	97.25	99.09	95.15	97.08	97.17
2	SUPPORT VECTOR MACHINE	96.42	98.9	93.59	96.17	96.31
3	ADABOOST	96.42	96.68	95.84	96.26	96.4
4	CATBOOSTING	97.58	98.41	96.53	97.46	97.54

6. EXPERIMENTAL REPORTS

In this proposed application, we try to use google collab as working platform and try to show the performance of our proposed application.

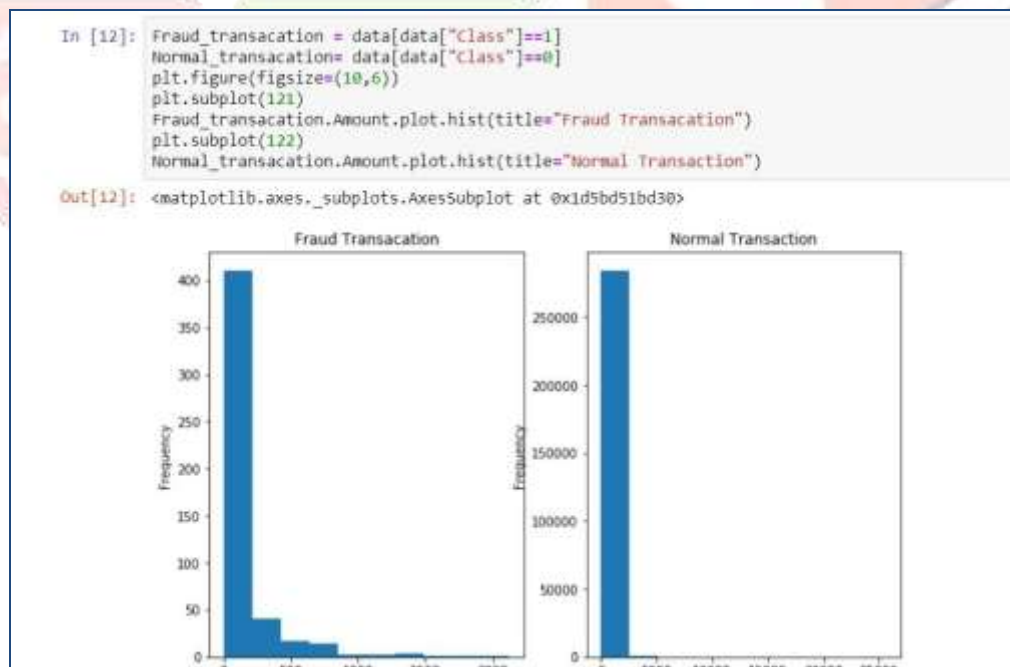
1) DATASET

Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	
1	0	-1.25981	-0.07278	2.530347	1.781153	-0.35832	0.461268	0.238999	0.098698	0.367387	0.267994	-0.35106	-0.51278	-0.99129	-0.111217	1.468177	-0.4704	0.287971	0.025791	0.481993	0.255421	-0.03835	0.177838
2	1	1.19187	0.304153	0.09448	0.448134	0.903818	-0.26235	-0.3788	0.085102	-0.25943	-0.39697	1.812177	1.962315	0.489395	-0.14477	0.625358	0.483917	-0.11448	-0.18136	-0.34578	-0.09388	-0.23578	-0.53887
3	1	-1.58835	-1.94026	1.773269	0.57676	-0.5812	1.808489	0.795463	0.247676	-1.55465	0.267943	0.624521	0.866884	0.717281	-0.16595	1.145585	-0.80028	1.189969	-0.12136	-1.26388	-0.52488	0.247988	0.779579
4	1	0.96677	-0.18521	1.782969	0.36328	-0.01831	1.247283	0.217493	0.377435	-1.88703	-0.85485	-0.22846	1.782318	0.507757	0.38792	-0.61421	-1.03945	0.68489	1.665775	-1.23262	-0.20889	-0.1883	0.055274
5	2	-1.25813	0.87737	1.548718	0.403034	-0.40719	0.099911	0.582941	-0.27953	0.817739	0.753874	-0.82284	0.538296	1.340852	-1.11967	0.175121	-0.45145	-0.23769	-0.03859	0.883487	0.408542	-0.09493	0.796276
6	3	0.42567	0.868251	1.411189	-0.18825	0.433887	-0.23913	0.474303	0.262394	-0.58667	-0.17442	1.842283	0.289894	-0.35889	-0.11871	0.557627	0.403176	0.85813	0.088651	-0.01819	0.884489	0.30285	-0.55982
7	4	1.129628	0.143084	0.061111	1.303653	0.181883	0.777708	-0.00516	0.081213	0.48486	-0.89815	-0.43681	0.15383	0.75116	0.167172	0.050144	-0.44319	0.001821	-0.61189	0.84558	-0.21961	-0.16772	-0.17071
8	7	-0.84437	1.417984	1.87438	-0.49221	0.848834	0.428138	1.128261	-0.80788	0.615375	1.249374	-0.62247	0.291474	1.757984	-1.23387	0.686133	-0.07613	-1.22123	-0.35822	0.54505	-0.15074	1.84565	-0.10245
9	7	-0.89409	0.888157	-0.11319	-0.27151	0.609599	1.721838	0.170145	0.853884	-0.88303	0.41944	-0.79512	-0.11845	-0.28625	0.814455	-0.33878	0.17008	-0.49077	0.118765	0.57038	0.052738	-0.07943	-0.26889
10	8	0.13816	1.118589	0.643657	0.222129	0.409365	0.24634	0.451583	0.668539	-0.79673	0.36885	1.017654	0.81828	1.088844	0.44052	0.150219	0.739453	0.54888	0.476677	0.451773	0.337721	1.24895	-0.63375
11	14	0.49944	-1.17634	0.81386	-0.37947	-1.97118	-0.62915	-1.42324	0.048265	-1.72041	0.263619	1.99844	-0.67244	-0.51395	-0.09500	0.23883	0.811967	0.253413	0.894344	-0.22137	-0.38723	-0.0889	0.153484
12	10	0.894938	0.618139	-0.8743	-0.89482	1.934584	0.113027	0.470355	0.538247	-0.50889	0.889765	-0.29512	0.32814	-0.89885	0.363852	0.528984	-0.12949	0.80398	0.189985	0.781664	0.125941	0.849624	0.188422
13	11	1.340899	-1.21264	0.38393	-1.2349	-1.48542	-0.75313	-0.8894	-0.22749	-2.09403	1.32370	0.27566	-0.34368	1.854517	-0.11763	0.725575	-0.81561	0.873036	-0.84776	-0.68319	-1.01076	-0.21385	-0.483259
14	11	1.088974	0.287721	0.828813	2.71252	-0.1784	0.137544	-0.09672	0.155882	-0.22108	0.46013	-0.77196	0.323887	-0.02188	1.17848	-0.83509	0.18993	1.124885	-0.9885	-0.88292	-0.3531	-0.05888	0.074422
15	12	-2.79285	-0.32777	1.86175	1.767473	-0.13885	0.802596	-0.42291	-1.96711	0.755713	1.151887	0.848955	0.782844	0.170849	-0.73498	0.486796	-0.38308	-0.25587	0.782829	0.221888	-0.81211	0.151863	0.122182
16	12	0.75242	0.145845	2.057313	-0.66864	-1.15818	-0.07785	-0.68858	0.083603	-0.48627	0.747311	-0.71888	-0.77841	1.047827	-0.2668	1.136919	1.862114	-0.27817	-0.42399	0.821525	0.263491	0.499625	0.16385
17	12	1.181215	-0.8401	1.267132	1.889091	-0.716	0.388369	-0.58606	0.18838	0.782333	-0.26798	-0.49331	0.936708	0.78838	-0.48865	0.284074	0.14863	0.80021	-0.54951	0.57588	-0.11381	-0.02465	0.196012
18	13	-0.43881	0.918996	0.504591	-0.72722	0.915479	-0.17787	0.707962	0.087982	-0.60527	-0.25788	0.324898	0.277292	0.252624	-0.2819	-0.18852	1.343174	-0.91871	0.68067	0.025436	-0.04702	-0.1548	-0.67284
19	14	-0.40236	-0.45635	1.180385	1.786289	1.849208	-1.76141	-1.53974	0.168842	1.13308	0.146519	0.193731	0.800217	-0.26657	-0.47913	-0.52881	0.472004	-0.72548	0.187081	-0.40807	-0.19648	-0.3618	0.88446
20	15	1.082936	-1.02935	0.454795	0.43881	-1.55543	0.78396	1.08066	-0.05183	-1.97868	0.638879	1.077542	0.63205	-0.40586	0.852011	-0.04288	0.11643	0.384241	0.554411	0.81443	-0.38791	0.17795	-0.17567
21	16	0.694885	-1.38181	1.025211	0.894139	-1.19121	1.309189	-0.67859	0.44529	-0.4462	0.588523	1.819131	1.288519	0.42386	-0.17205	-0.80798	-0.04438	0.505867	0.625847	-1.30841	-0.13831	-0.29518	-0.57196
22	17	0.963496	0.118481	-0.17148	2.390384	1.129968	1.698318	0.107712	0.521812	-1.18131	0.174396	1.69031	-0.400774	-0.93942	0.863739	0.710011	-0.60221	0.881484	-1.73716	-1.81761	-0.48911	0.141897	0.402487
23	18	1.168616	0.58213	-0.0673	1.262589	0.438804	0.688438	0.241412	0.138882	-0.08926	0.921575	0.144786	-0.53238	-2.10535	1.12887	0.003075	0.414425	-0.45448	-0.08887	-0.8086	-0.30717	0.018702	-0.86187
24	18	0.147491	0.277966	1.184711	-0.8928	-1.11418	-0.15012	-0.94936	-1.61784	1.544071	-0.81388	-0.8833	0.524813	-0.45338	0.881293	1.555204	-1.19689	0.781131	0.436621	2.17807	-0.19088	1.6818	0.008454
25	21	-1.94818	-0.8449	-0.40257	-1.02188	2.911898	1.957973	-0.08308	0.820549	0.849987	0.571743	-0.08128	-0.21575	0.044181	0.818898	1.187138	0.578841	-0.57067	0.244061	0.888803	-0.22672	-0.17953	-0.79833
26	21	-0.37429	-0.12149	1.321011	0.418808	0.265108	-0.89564	0.543885	-0.10463	0.478684	1.444951	-0.86957	-0.18852	-0.65521	-0.2798	-0.21167	-0.13113	0.020751	-0.48847	0.585751	-0.38668	-0.40364	-0.1274
27	21	1.173285	0.333488	0.281805	1.133461	-0.17258	0.91865	0.369025	-0.32736	-0.18465	0.84814	1.01443	0.87935	1.492385	0.101418	0.781478	0.04438	0.51164	-0.32586	0.39993	0.027878	0.867803	0.127812
28	21	1.320707	-0.17044	0.493235	0.73008	-0.83878	-0.81188	-0.2648	-0.23898	-1.07142	0.888839	-0.94151	-0.11192	0.362485	0.175455	0.781187	-1.03087	-0.21894	1.271705	-1.04862	-0.52189	-0.28818	-0.32336

2) PLOT HISTOGRAM FOR EACH PARAMETER



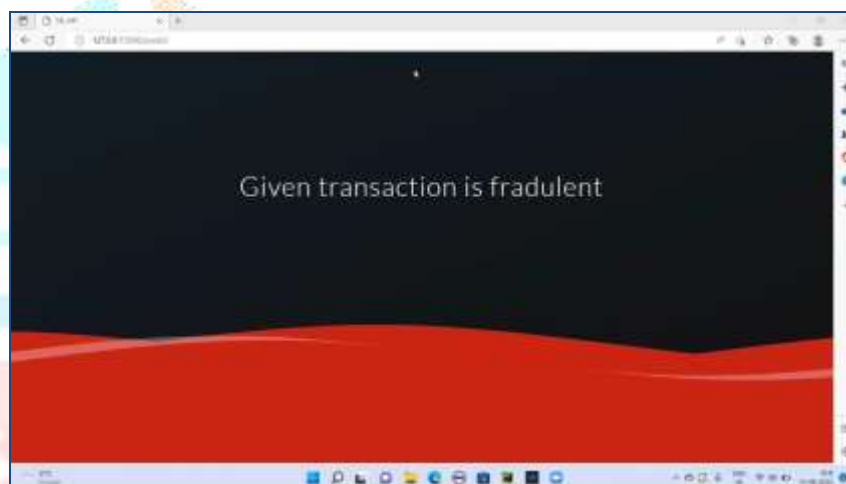
3) PLOT FOR NORMAL AND FRAUD TRANSACTIONS



4) MAJORITY VOTING



5) TEST THE INPUT



7. CONCLUSION

This project offered credit card fraud detection using machine learning techniques. A variety of standard models, including Nave Bayes and decision trees, were applied. Individual (standard) models and hybrid models using AdaBoost and majority voting combination methods were evaluated using a publicly available credit card data set. Because it considers true and false positive and negative predicted outcomes, the confusion metric has been adopted as a performance measure. For evaluation, an actual credit card data set from a financial institution was also utilised. The same individual and hybrid models were used. Using AdaBoost and majority voting methods, a perfect score of 1 was obtained.

8. REFERENCES

- [1] <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8292883>
- [2] J. T. Quah and M. Sriganesh, —Real-time credit card fraud detection using computational intelligence, Expert Systems with Applications, pp. 1721-1732, 2008.

[3] https://en.wikipedia.org/wiki/Machine_learning

[4] <https://en.wikipedia.org/wiki/AdaBoost>

[5] https://en.wikipedia.org/wiki/Naive_Bayes_classifier

[6] <https://www.kdnuggets.com/2016/04/datacamp-learning-python-data-analysis-data-science.html>

[7] <https://opensource.com/article/18/9/top-3-python-libraries-data-science>

[8] https://www.researchgate.net/publication/323213023_Credit_card_fraud_detection_using_AdaBoost_and_majority_voting

[9] N. Laleh and A. M. Azgomi, —A Taxonomy of Frauds and Fraud Detection Techniques, IICISTM, vol. 31, pp. 256-267, 2009.

[10] E. Aleskerov, B. Fieisleben and R. Bharat, —CARDWATCH: A Neural Network Based Database Mining System for Credit Card Fraud Detection, Department of Electrical Engineering and Computer Science, University of Siegen, pp. 220-226, 1997.

[11] N. F. R. Centre, —Bank card and cheque fraud, National Fraud Authority, UK.

9. ABOUT THE AUTHORS



Mrs. PINNAMARAJU.T.S.PRIYA working as Assistant professor in Master of computer application(MCA) in Sanketika Vidya Parishad Engineering College, Visakhapatnam Andhra Pradesh. with 6 years of experience in Masters of Computer Applications (MCA) , accredited by NAAC. With her area of interests in C, DBMS, Computer Organization, Software Engineering, IOT.



Ms. MADDI HINDUDEVI is currently pursuing her 2 years MCA in Department of Computer Science and Applications at Sanketika Vidhya Parishad Engineering College, P.M. Palem, Visakhapatnam, Andhra Pradesh. Her area of interest includes Python, Java, C, and C++.