

Credit Card Fraud Detection Model

Machine Learning Models Comparison

Logistic Regression

Long Short Term Memory(LSTM)

Random Forest

Group-6

Scam Guards



Team Members:

Ayon Kumar Das, *Computer Science*, adas010@citymail.cuny.edu

Myesha Mahazabeen, *Computer Science*, mmahaza000@citymail.cuny.edu

Najia Jahan, *Computer Science*, njahan010@citymail.cuny.edu

Professor Yunhua Zhao

CSC 44800

Final Project

12/19/2023

Table of Contents

Abstract.....	3
Introduction.....	3
Dataset Overview.....	4
Machine Learning Models.....	5
Logistic Regression.....	5
Long Short Term Memory(LSTM).....	7
Random Forest.....	11
Comparison of the Models.....	13
Accuracy Scores Comparison.....	13
Comparison of the Graphs & Matrices.....	14
ROC Curve.....	14
Precision-Recall Curve.....	16
Correlation Matrix.....	18
Confusion Matrix.....	20
Challenges.....	23
Conclusion.....	24
References.....	25

Abstract

This project introduces a practical approach to detect credit card fraud in real-time electronic transactions. We employ Logistic Regression, Long Short Term Memory (LSTM), and Random Forest models, emphasizing precision and recall. The study involves a detailed comparison of these models, assessing their performance using key metrics such as Accuracy, Precision, Recall, F1-Score, and Matthews Correlation Coefficient. By leveraging these techniques, our aim is to provide insights for improving fraud detection in the ever-evolving landscape of electronic transactions

Introduction

In the age of digital transactions, the rise of credit card fraud urges us to find strong ways to stop it. Banks and financial companies face a big challenge in protecting people from tricky fraud tactics. This study goes a step further by checking not just one, but three computer methods—logistic regression, long short term memory(LSTM), and random forest—to see how good they are at telling apart real and fake credit card transactions. Our main goal is to share smart ideas that can help make even better systems to catch fraud. For this project, we're using a dataset with a substantial volume of transactions, ensuring our study covers a wide range of situations in the real world of electronic transactions. This extensive dataset allows us to comprehensively assess the performance of the models and draw meaningful insights to contribute to the ongoing efforts in enhancing credit card fraud detection.

Dataset Overview

The dataset utilized in this project provides a snapshot of credit card transactions conducted by European users in September 2013, spanning a concise two-day period. A total of 284,807 transactions are included, with only 492 identified as fraudulent—comprising a mere 0.172% of the dataset. To address confidentiality concerns, the dataset primarily includes features resulting from Principal Component Analysis (PCA) transformation, with 'Time' (representing seconds since the first transaction) and 'Amount' (indicating the transaction sum) being the exceptions. The crucial 'Class' feature serves as the label, denoting whether a transaction is fraudulent (Class 1) or non-fraudulent (Class 0) [1].

Below are two small snippets of the dataset used.

# Time	# V1	# V2	# V3	# V4	# V5	# V6	# V7	# V8	# V9	# V10	# V11	#
0	-1.359887133673 8	-0.872781173389 8497	2.5363467379691 4	1.3781552242744 3	-0.3383287679942 518	0.4623877777622 92	0.2395985548612 57	0.8986979812618 587	0.3637869696112 13	0.8987941719789 316	-0.551599533268 813	-3 3
0	1.1918571113148 6	0.2661587128596 3	0.1664881133532 1	0.4481548784689 11	0.8688176492822 243	-0.882368888815 5687	-0.878882983332 3113	0.8851816549148 184	-0.255425128189 186	-0.166974414884 614	1.6127266618547 9	1 7
1	-1.358354861598 23	-1.348163874736 89	1.7732893426311 9	0.3797795938343 28	-0.583198133318 193	1.8884993887926 3	0.7914689564584 22	0.2476757865889 91	-1.514654322685 83	0.2876428652166 96	0.6245814594248 95	6 3
1	-0.966271711572 887	-0.185226888882 898	1.792993395787 2	-0.863291275836 453	-0.818388879683 8823	1.2472831675248 6	0.2376889397717 8	0.3774358746522 62	-1.387824862781 97	-0.854951922471 3749	-0.226487263835 481	6 6
2	-1.158233093495 23	0.8777367548484 51	1.548717846511 21	0.4838339339551 653	-0.487193377311 653	0.8959214624684 256	0.5929487453855 45	-0.278532677192 282	0.8177393882352 94	0.7538744319763 54	-0.822842877946 363	6 6
2	-0.425965884412 454	0.9685238448829 85	1.1411893423221 9	-0.168252879768 382	0.4289868887721 9	-0.829727551663 9742	0.4762889487288 27	0.2683143338748 74	-0.568671375712 51	-0.371487196834 471	1.3412619880195 7	6 3
4	1.2296576345879 3	0.1418835878493 26	0.8453787735899 449	1.2826127367359 4	0.1918889885976 45	0.2727881228990 98	-0.885159882882 58983	0.8812129398830 894	0.4649599947838 86	-0.899254321128 9237	-1.416987243149 28	- 6
7	-0.644269442348 146	1.4179635454738 5	1.8743883763556 015	-0.492199818495 015	0.9489348947641 57	0.4281184628338 89	1.1286313583835 3	-3.887864238735 89	0.6153747386678 27	1.2493761781517 6	-0.619467796121 913	6 6
7	-0.894286882282 82	0.2861571962765 44	-0.113192212729 871	-0.271526138888 684	2.6695986595986 684	3.7218188611275 1	0.3781451276769 16	0.8518844432889 85	-0.392847586798 684	-0.418438432848 439	-0.785116586646 538	- 6
9	-0.338261752425 75	1.1195933764156 6	1.8443665515731 6	-0.222187276738 296	0.4993688864972 7	-0.246761188619 91	0.6515832864899 72	0.8695385865186 387	-0.736727316364 189	-0.366845639286 541	1.8176144678326 2	6 2

creditcard.csv (150.83 MB)

Detail Compact Column

31 of 31 columns

#	# V19	# V20	# V21	# V22	# V23	# V24	# V25	# V26	# V27	# V28	# Amount	# Class
1	0.483992968255733	0.251412898239785	-0.018386777944153	0.277837575558899	-0.118473918188767	0.8649288749146731	0.128539358273528	-0.189114843888824	0.133558376748387	-0.02185383534538215	149.62	0
1	-0.145783841325259	-0.0698831352238283	-0.225775248893138	-0.638671952771851	0.101288821253234	-0.339846475529127	0.167178484418143	0.125894532368176	-0.08898389914322813	0.0147241691924927	2.69	0
1	-2.26185789538414	0.524979725224484	0.247998153469754	0.771679481917229	0.989412262347719	-0.689288956498685	-0.327641833735251	-0.139896571514147	-0.955327948384261	-0.0597518485929284	378.66	0
1	-1.2326219788892	-0.288837781168366	-0.188388452835545	0.88527359678253453	-0.198328518742841	-1.17557533186321	0.647376834682838	-0.22192884458487	0.0627228487293833	0.0614576285886333	123.5	0
1	0.883486924968175	0.488542368392758	-0.08943869713232919	0.79827849458971	-0.137458879619863	0.141266983824769	-0.206889587619756	0.582292224181569	0.219422229513348	0.215153147499286	69.99	0
1	-0.8331937877876282	0.8849676728682849	-0.288253514656728	-0.559824796253248	-0.8263976679795373	-0.371426583174346	-0.232793816737834	0.185914779897957	-0.253844224739337	-0.8818882569229443	3.67	0
1	-0.8455758446637976	-0.21963255278686	-0.167716265815783	-0.278789726172363	-0.154183786889385	-0.788855415884671	0.75813693588659	-0.257236845917139	0.8345874297438413	0.88516776898624916	4.99	0
1	0.324584731321494	-0.156741852488285	1.94346533978412	-1.01545478979971	0.857583529867291	-0.64978980559993	-0.415266566234811	-0.0516342969262494	-1.20692188894258	-1.08533918832377	48.8	0
1	0.57832816746536	0.8527356691149697	-0.8734251881859225	-0.268891632235551	-0.284232669947878	1.011591881878582	0.3732846881462294	-0.384157387782294	0.0117473564581996	0.14248432992147	93.2	0
1	0.451772964394125	0.283711454727929	-0.246913936918808	-0.533752642486113	-0.12879488888185	-0.385849925313426	-0.0697338468416923	0.8941988339514961	0.246219384619926	0.8838756493473326	3.68	0

Machine Learning Models

Logistic Regression

Logistic Regression is a powerful statistical method commonly employed for binary classification tasks, such as predicting whether an email is spam or not, or, as in this project, determining whether a credit card transaction is fraudulent. Unlike linear regression, which predicts continuous outcomes, logistic regression models the probability that a given instance belongs to a particular category.

Mathematical Equation [2]:

The logistic regression model employs the logistic function (also known as the sigmoid function) to transform the linear combination of input features into a probability score between 0 and 1.

The logistic function is represented by the equation:

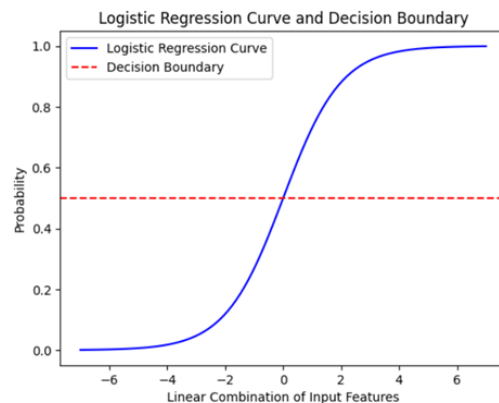
$$P(Y=1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

Here,

- $P(Y=1)$ is the probability of the instance belonging to Class 1 (fraudulent in our case).
- e is the base of the natural logarithm.
- β_0 is the intercept.
- $\beta_1, \beta_2, \dots, \beta_n$ are the coefficients associated with the input features X_1, X_2, \dots, X_n .

Graphical Representation:

The sigmoid function generates an S-shaped curve, mapping any real-valued number to the range $[0, 1]$. This curve is essential for logistic regression, as it allows the model to output probabilities. The logistic regression decision boundary is set where the sigmoid function equals 0.5. Instances falling on one side of this boundary are predicted as Class 1, while those on the other side are predicted as Class 0.



In the graph, the x-axis represents the linear combination of input features, and the y-axis represents the predicted probability. As the linear combination increases, the probability

approaches 1, and as it decreases, the probability approaches 0. The decision boundary (shown as a vertical line) determines the threshold for classification.

Model Training

Trained on a dataset with fraudulent and non-fraudulent transactions, Logistic Regression learns patterns using a sigmoid function. Coefficients and the function's output help set a decision boundary, making it effective for binary outcomes like credit card fraud detection. The model's interpretability and suitability for this scenario contribute to its valuable role in the project.

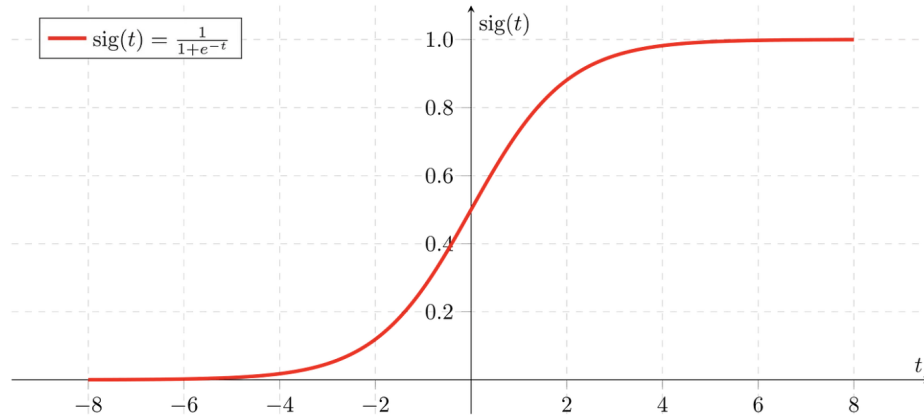
LSTM(Long Short Term Memory)

LSTM is a recurrent neural network (RNN) design primarily used for tasks involving sequence prediction problems. The single hidden state of a traditional RNN can make it challenging for the network to learn long-term dependencies. Long short term memory networks (LSTMs) solve this issue by adding a memory cell, which is a container that can store information for a long time.

They can learn long-term dependencies in sequential data, which makes them a good fit for tasks like price prediction, speech recognition, language translation, and time series forecasting.

The three gates, the input, forget, and output gates, control the memory cell. This memory cell's input, output, and addition of data is determined by these gates. The data added to the memory cell is managed by the input gate. What information is taken out of the memory cell is managed by the forget gate. Information that is output from the memory cell is also controlled by the output gate. As a result, long-term dependencies can be learned by LSTM networks, which are able to selectively store or discard information as it passes across the network[5].

Sigmoid Activation Function with Graph



LSTM gates are the sigmoid activation functions; that is, they output a value between 0 and 1, which is often one of the two. Because we want a gate to only produce positive values and to be able to clearly indicate whether we should maintain a certain feature or discard it, we use the sigmoid function for gates. "0" indicates that nothing can pass through the gates. A gate that reads "1" indicates that everything can pass through it [6].

LSTM Model Architecture [6]

The chain structure of the LSTM architecture is made up of neural networks and various memory building blocks known as cells. The gates perform memory manipulations, while the cells store information. Three gates are present in LSTM:

Forget Gate

The forget gate is used to remove information that is no longer useful in the cell state. It takes two inputs, x_t (input at that specific time) and h_{t-1} (previous cell output), multiplies them with weight matrices, and adds bias. The resulting product is then passed through an activation function that produces a binary output: if the output is 0, the piece of information is forgotten for that specific cell state, if the output is 1, it is kept for future use. The equation for forget gate is:

$$f_t = s(w_f[h_{t-1}, x_t] + b_f)$$

Input Gate

The input gate modifies the cell state by adding pertinent information. Using inputs h_{t-1} and x_t , the sigmoid function is first used to regulate the information before filtering the values to be remembered in a manner akin to a forget gate. Next, a vector containing all possible values from h_{t-1} and x_t is created using the tanh function, which produces an output ranging from -1 to +1. Finally, the useful information is obtained by multiplying the vector values by the regulated values. The equation for the input gate is:

$$i_t = s(w_i[h_{t-1}, x_t] + b_i)$$

Then the previous state is multiplied by f_t , ignoring the data that was previously decided to ignore. Then $i_t * C_t$ is included. This is a representation of the updated candidate values, with each state value's amount of update considered.

$$c_t = \tanh(w_c[h_{t-1}, x_t] + b_c)$$

$$c_t = f_t * c_{t-1} + i_t * c_t$$

$$h_t = o_t * \tanh(c_t)$$

Here, c_t denotes cell state(memory) and c_t refers to candidate for cell state.

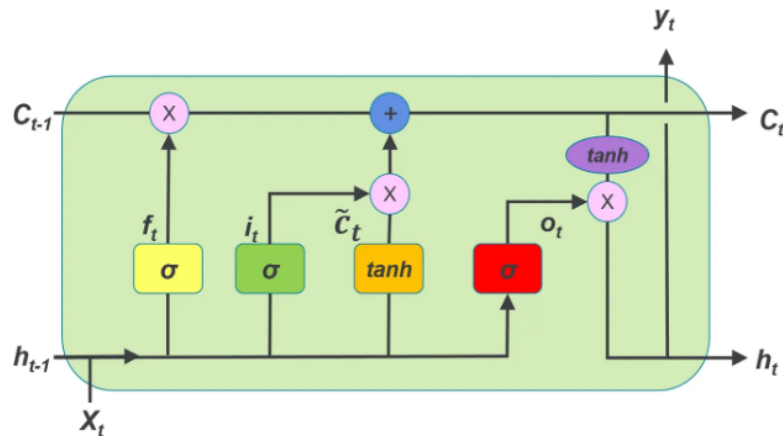
Output Gate

The output gate performs the function of extracting valuable information to be displayed as output from the current cell state. Using the tanh function on the cell, a vector is first created. Next, the data is filtered by the values that need to be remembered using inputs h_{t-1} and x_t , and the information is regulated using the sigmoid function. The vector's values and the regulated values are finally multiplied before being sent as an output and input to the following cell. The equation for output gate is:

$$o_t = s(w_o [h_{t-1}, x_t] + b_o)$$

Here, h_{t-1} denotes the output of the previous lstm block at t-1 time. x_t refers to input at current timestamp, b_x refers to biases for all the gates and w_x refers to the weight for all the gates. And s is the sigmoid function.

Below is the figure of a single LSTM block:



To control the information flow between 0 and 1, each gate has sigmoid activation functions that help it decide what to keep and what to discard. In addition, the hyperbolic tangent (tanh) activation function is used to modify the cell state. This function guarantees that values fall between -1 and 1, enabling more stable updates to the cell state. With this architecture, information can be updated or selectively retained over long sequences, allowing LSTMs to capture long-range dependencies in sequential data efficiently.

An LSTM Neural Network is constructed using TensorFlow's Keras API. The model consists of an LSTM layer with 100 units, followed by a Dropout layer with a dropout rate of 0.5 to prevent overfitting. The output layer is a dense layer with a sigmoid activation function, suitable for binary classification problems. The model is compiled using binary cross-entropy as the loss function and the Adam optimizer with a learning rate of 0.001. The accuracy metric is used for model evaluation

Training the Model

The LSTM model is trained on the preprocessed data using a training set, and the performance is evaluated on a validation set. The training process involves 10 epochs with a batch size of 64.

Random Forest [3]

Overview: In our credit card fraud detection project, we have employed the Random Forest model as one of the key machine learning algorithms. Random Forest is an ensemble learning method known for its effectiveness in classification tasks. Here's a concise overview of its role and characteristics in our project.

Ensemble Learning: Random Forest operates by creating multiple decision trees during the training phase. Each decision tree is constructed using a different subset of the dataset and a random subset of features. The final prediction is then determined by aggregating the individual predictions of these trees.

Decision Trees in Random Forest: The fundamental building block of Random Forest is the decision tree. These trees are simple models that make decisions based on a series of rules. In our project, decision trees are trained to discern patterns in credit card transactions, distinguishing between legitimate and fraudulent activities.

Bagging Technique: Random Forest employs a technique called bagging (Bootstrap Aggregating). This involves training each decision tree on a random subset of the dataset, allowing for diversity among the trees. This diversity helps improve the model's robustness and generalization to unseen data.

Feature Importance: One of the notable features of Random Forest is its ability to provide a measure of feature importance. This indicates the relevance of each feature in making predictions. For credit card fraud detection, understanding which features contribute most to the model's decision-making process is crucial.

Model Training: During the training phase, the Random Forest model is fed with our credit card transaction dataset. This dataset includes features derived from Principal Component Analysis (PCA) transformation, along with 'Time' and 'Amount' attributes. The model is trained to recognize patterns and intricacies associated with fraudulent transactions.

Comparison of the Models

Comparing the Accuracy

When comparing our models, we evaluated their performance using these key metrics: Accuracy, Precision, Recall, F1 Score, and Matthews Correlation Coefficient (MCC). Accuracy gives an overall measure, Precision gauges the accuracy of positive predictions, Recall assesses sensitivity to detecting true positives, F1 Score balances Precision and Recall, and MCC offers a comprehensive correlation measure. Together, these metrics guide our assessment, ensuring a thorough evaluation of the models in the context of credit card fraud detection.

Scores	Logistic Regression	Long Short Term Memory	Random Forest
Accuracy	0.99892	0.99956	0.99959
Precision	0.81356	0.96202	0.97468
Recall	0.48979	0.77551	0.78571
F1 Score	0.61166	0.85875	0.87005
Matthew Correlation Coefficient	0.63078	0.86354	0.87492

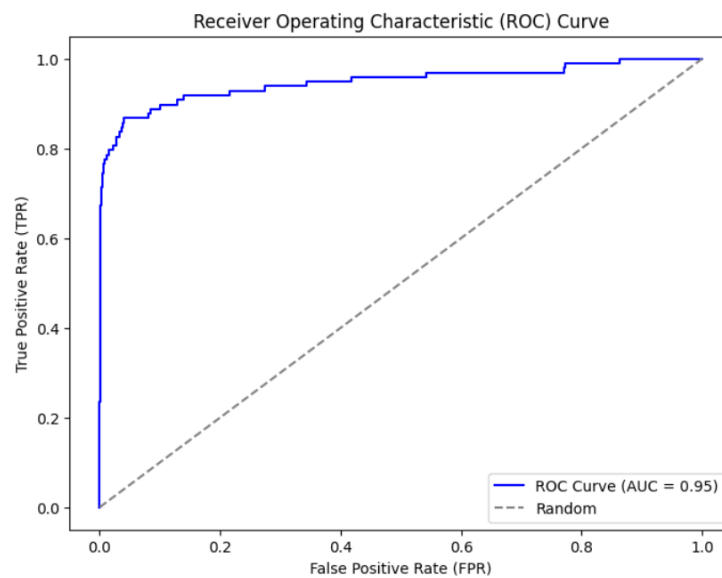
In comparing the accuracy of the three models based on the output values on the table above, it is evident that Random Forest outperforms both Logistic Regression and Long Short Term Memory (LSTM) across various metrics. Random Forest exhibits the highest accuracy, precision, recall, F1 Score, and Matthews Correlation Coefficient among the models, with values of 0.99959, 0.97468, 0.78571, 0.87005, and 0.87492, respectively. This indicates a robust performance in correctly classifying both fraudulent and non-fraudulent transactions. While LSTM follows

closely with competitive scores, Logistic Regression lags behind, particularly in recall, where it demonstrates the lowest sensitivity to fraudulent transactions. The superiority of Random Forest in accuracy metrics suggests its effectiveness in credit card fraud detection, making it a promising choice for real-time electronic transaction security.

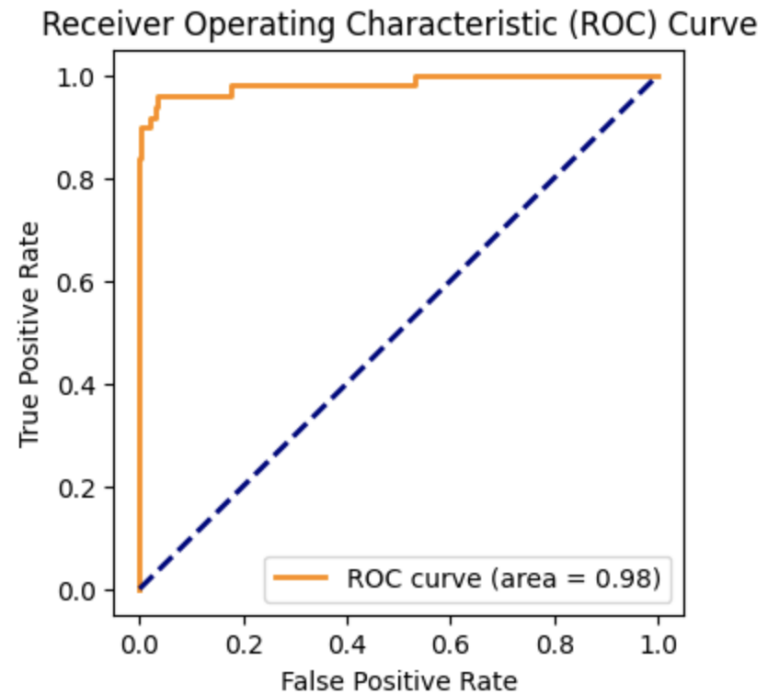
Comparing the Matrices & Graphs

Receiver Operating Characteristics (ROC Curve)

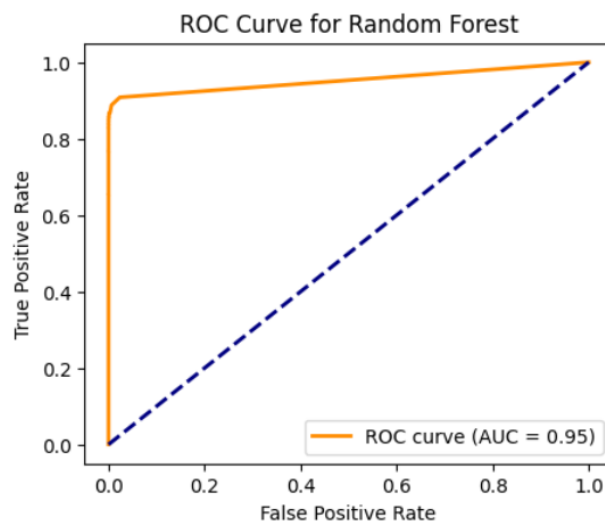
Logistic Regression



Long Short Term Memory



Random Forest



The Receiver Operating Characteristic (ROC) curve assesses how well a binary classification model distinguishes between positive and negative instances across different threshold levels.

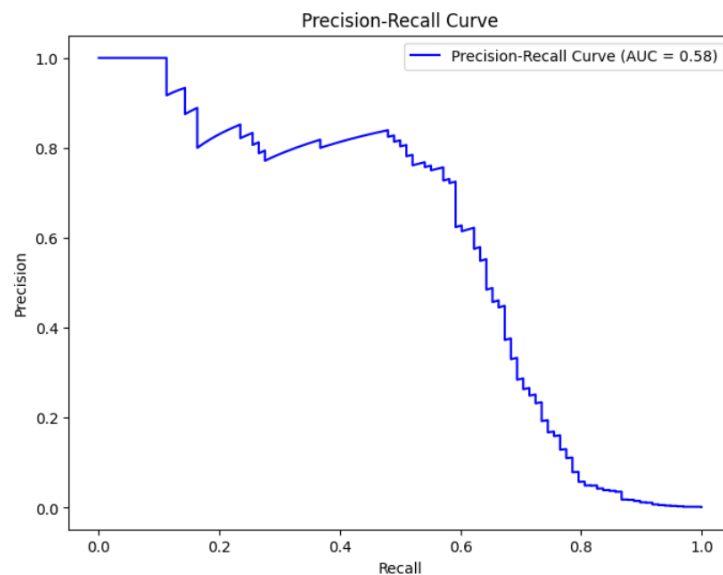
The Area Under the Curve (AUC) summarizes this performance. From our output:

- **Logistic Regression:** AUC = 0.95 - Good ability to distinguish with a balanced true positive rate and false positive rate.
- **LSTM:** AUC = 0.98 - Excellent discrimination, superior to Logistic Regression, indicating strong performance.
- **Random Forest:** AUC = 0.95 - Robust discrimination, similar to Logistic Regression.

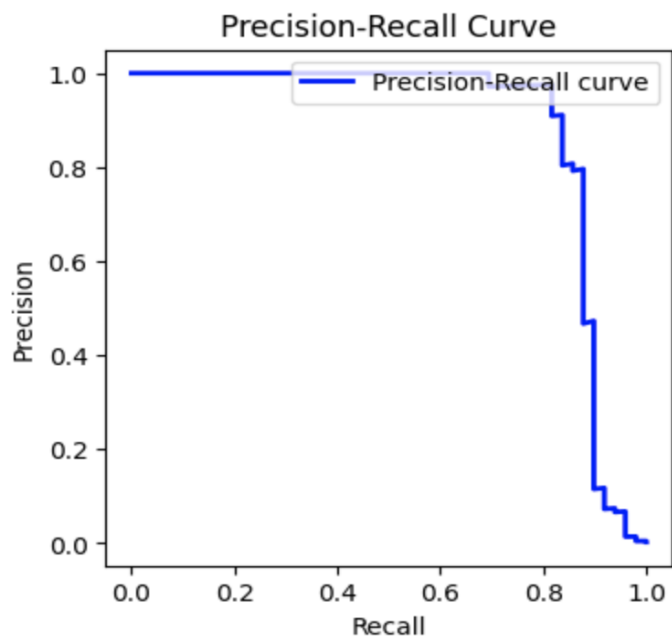
All models exhibit reasonably strong performance, with LSTM standing out for superior discrimination.

Precision-Recall Curve:

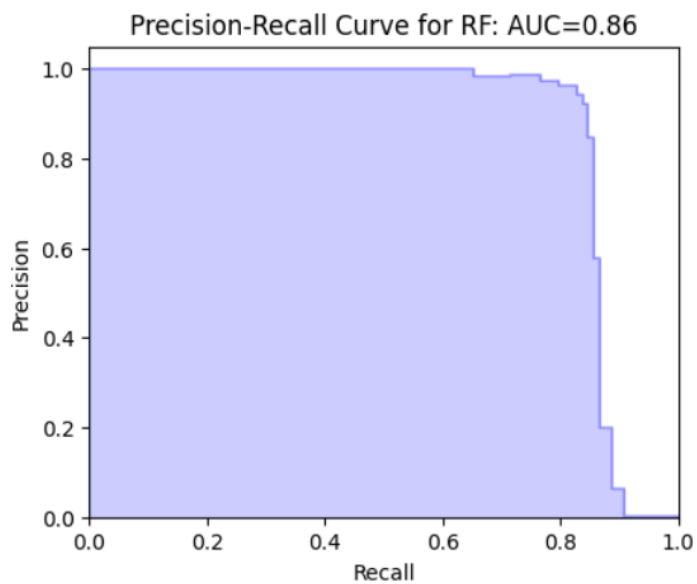
Logistic Regression



Long Short Term Memory



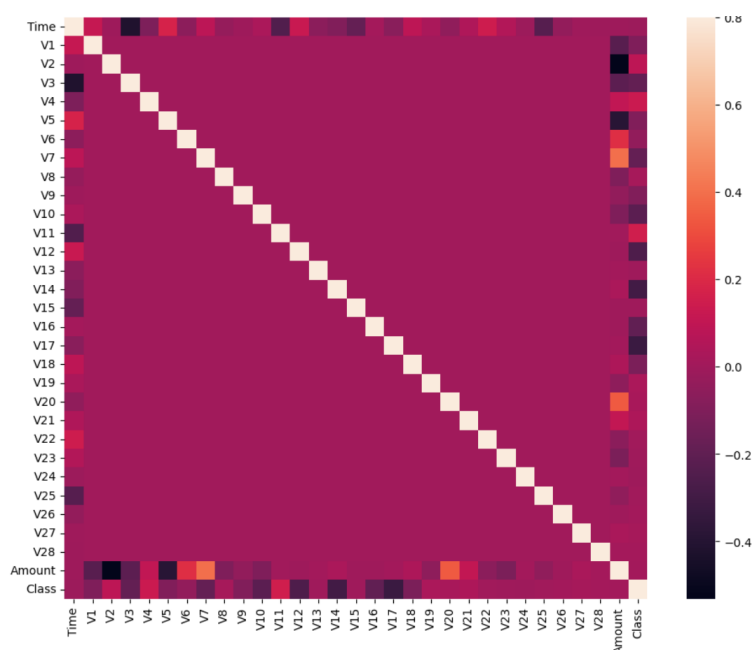
Random Forest



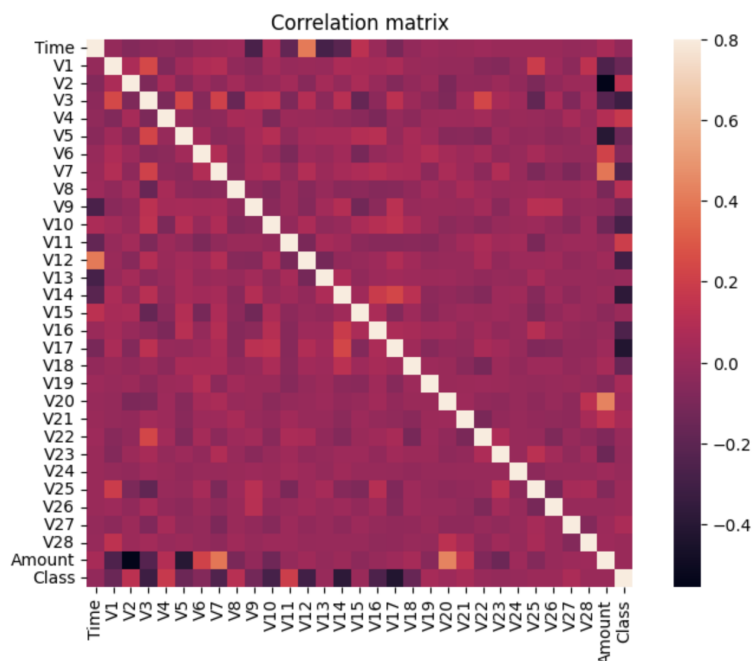
The Precision-Recall (PR) curve evaluates a binary classification model's precision and recall across various threshold levels, providing insights into its ability to correctly identify positive instances while minimizing false positives. The Area Under the Curve (AUC) summarizes the PR curve's performance. From the outputs above, we can see how LSTM and Random Forest show strong and robust precision-recall characteristics while Logistic Regression lags behind.

Correlation Matrix (Heatmap):

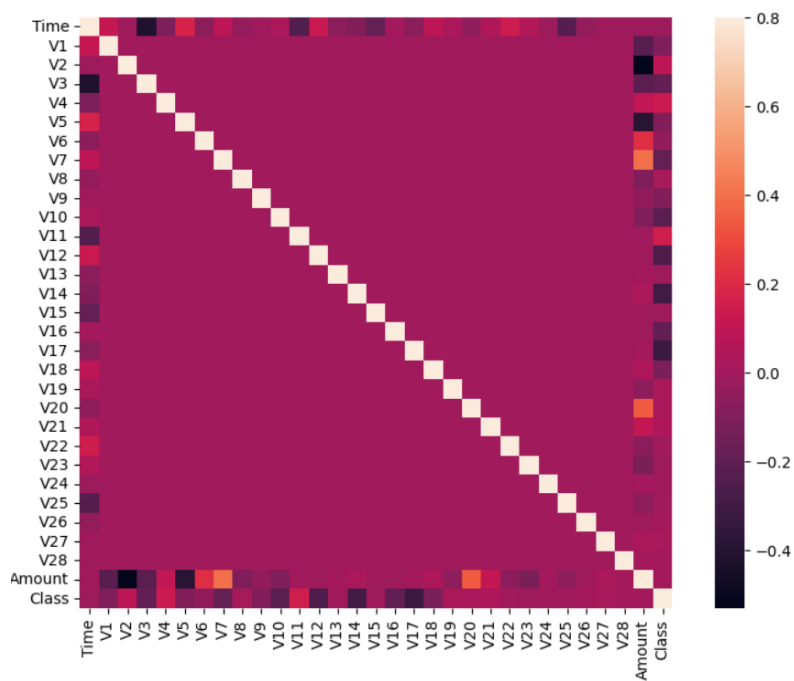
Logistic Regression



Long Short Term Memory



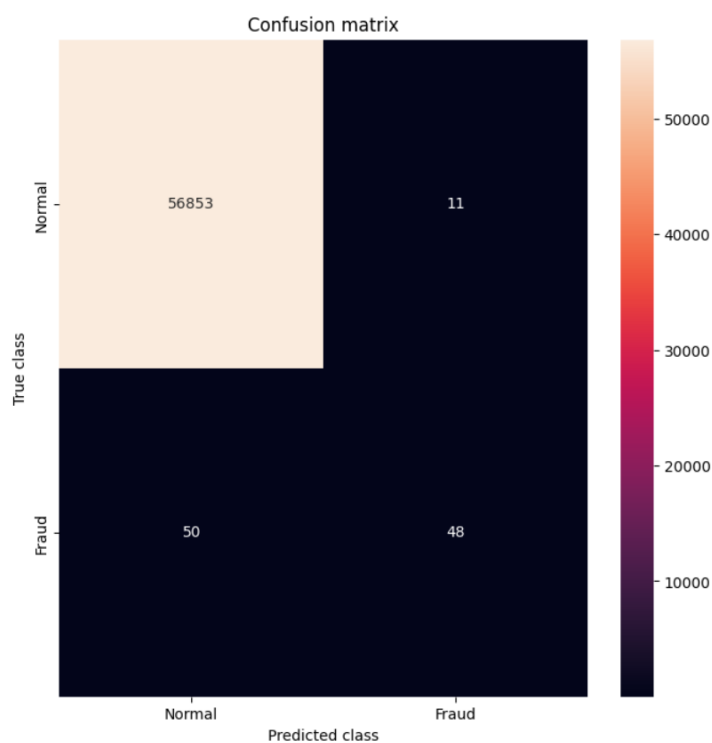
Random Forest



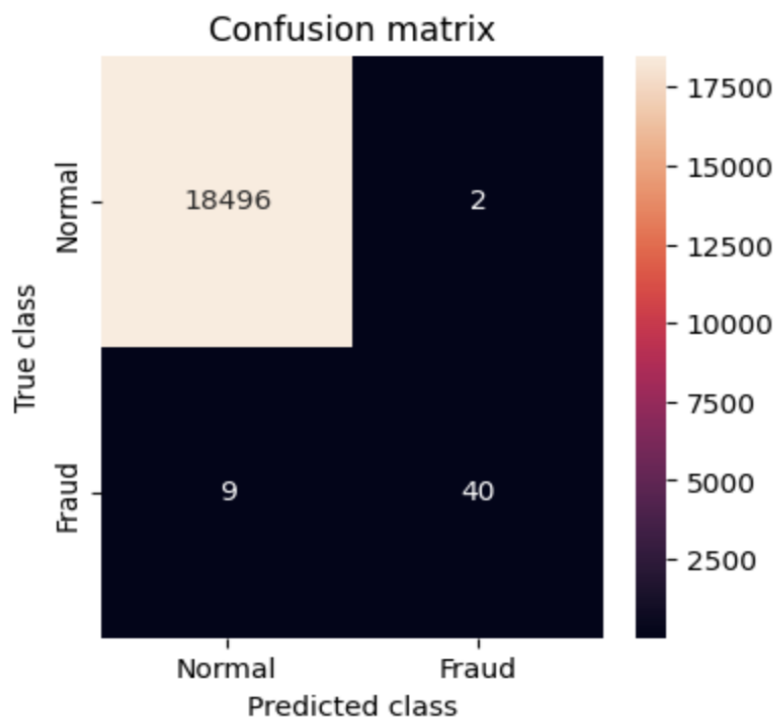
The correlation matrix is like a map that shows how different features in our fraud detection data relate to each other. Each entry in the matrix tells us the strength and direction of the relationship between two features. If the number is close to 1, it means they move together in the same direction, and if it's close to -1, it means they move in opposite directions. The correlation of 0 suggests no linear relationship. By examining the correlation matrices of the three models, we observe that the Logistic Regression and Random Forest models exhibit minimal correlation with other features. However, in the LSTM model, there is a more noticeable correlation between features.

Confusion Matrix:

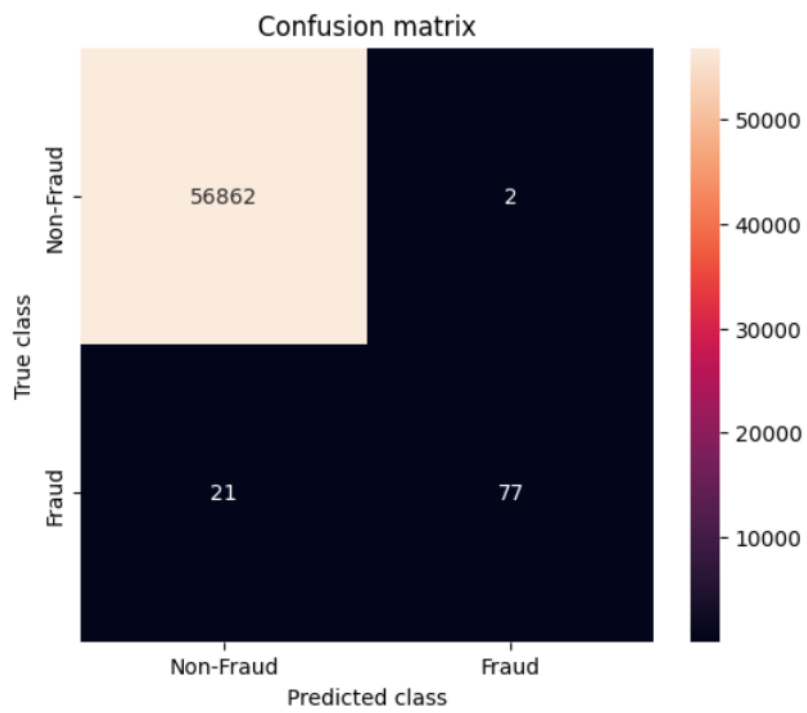
Logistic Regression



Long Short Term Memory



Random Forest



Model Names	True Negative (Non Fraud)	False Positive	False Negative	True Positive (Fraud)
Logistic Regression	56853	11	50	48
Long Short Term Memory	18496	2	9	40
Random Forest	56869	2	21	77

From the values above from the confusion matrix we can compare the efficiency of our models like below:

Logistic Regression: Logistic Regression demonstrates a good balance between correctly identifying non-fraudulent transactions (True Negatives) and correctly identifying fraudulent transactions (True Positives). However, as seen above, it has the highest number of false positives and false negatives compared to the other two models, indicating potential room for improvement.

Long Short Term Memory (LSTM): LSTM exhibits strong performance in correctly classifying non-fraudulent transactions (high True Negatives) and effectively identifying fraudulent transactions (high True Positives). The low values for false positives and false negatives suggest a robust ability to discriminate between the classes.

Random Forest: Random Forest, similar to Logistic Regression, maintains a balance between True Negatives and True Positives. It has a lower number of false positives and lower number of false negatives compared to Logistic Regression and a higher number of true positives. This suggests that Random Forest may be more sensitive to fraudulent transactions but generates fewer false alarms compared to Logistic Regression.

In summary, each model has its strengths and weaknesses. Logistic Regression and Random Forest demonstrate a balanced performance, while LSTM excels in accurately identifying both non-fraudulent and fraudulent transactions with fewer false predictions. The choice of the most suitable model depends on the specific requirements and trade-offs desired for the credit card fraud detection system.

Challenges

Building this project involves overcoming several challenges:

- **Imbalanced Data:** The dataset is heavily imbalanced, with only 0.2% of transactions labeled as fraudulent and 99.8% as non-fraudulent. This imbalance poses a significant challenge for the machine learning models leading to a bias towards predicting the majority class. Traditional metrics like accuracy may not be reliable, so we use precision, recall, and F1 score for better evaluation [4]. F1 score, considering both precision and recall, offers a nuanced view, crucial for accurately predicting the minority class alongside the majority.
- **Interpretability vs. Complexity:** Striking a balance between the interpretability of the model and its complexity proved to be a critical challenge. Achieving both a clear understanding of the model's decisions and high accuracy required thoughtful consideration and exploration.
- **Threshold Setting:** Determining the optimal threshold for classifying transactions played a vital role in influencing precision, recall, and overall model performance. Finding the right balance was essential to ensure the model's effectiveness in identifying fraud without generating excessive false positives.

- **Ethical and Regulatory Compliance:** Navigating ethical and regulatory considerations when handling sensitive financial data introduced complexity to the model development process. The dataset utilized in this project underwent an anonymized Principal Component Analysis (PCA) transformation for the features data which ensured privacy laws.
- **Adaptability to Fraud Tactics:** Constructing a model capable of adapting to evolving fraud tactics presented an ongoing challenge. Continuous monitoring and updates will be necessary to keep the model effective in the face of emerging fraud patterns and techniques.

Conclusion

In summary, our project aimed to detect credit card fraud using three models: Logistic Regression, LSTM, and Random Forest. After thorough comparison, Random Forest emerged as the most effective, surpassing the others in accuracy. We tackled challenges like imbalanced data and interpretability, emphasizing precision and recall. While each model had its strengths, Random Forest, with its ensemble approach, proved highly efficient in distinguishing between legitimate and fraudulent transactions. Ongoing monitoring is crucial for adapting to evolving fraud tactics. This study contributes valuable insights for developing strong, adaptive models to enhance electronic transaction security.

References

1. Machine Learning Group - ULB, "Credit Card Fraud Detection," Kaggle, March 23, 2018, <https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud/data>.
2. "Sigmoid Function." Wikipedia, November 27, 2023.
https://en.wikipedia.org/wiki/Sigmoid_function.
3. Melcher, Kathrin, and Rosaria Silipo. "Fraud Detection Using Random Forest, Neural Autoencoder, and Isolation Forest Techniques." *InfoQ*, InfoQ, 14 Aug. 2019, www.infoq.com/articles/fraud-detection-random-forest/.
4. Mazumder, Saikat. "5 Techniques to Handle Imbalanced Data for a Classification Problem." *Analytics Vidhya*, 27 Sept. 2023, www.analyticsvidhya.com/blog/2021/06/5-techniques-to-handle-imbalanced-data-for-a-classification-problem/.
5. <https://www.geeksforgeeks.org/deep-learning-introduction-to-long-short-term-memory/>
6. <https://medium.com/@divyanshu132/lstm-and-its-equations-5ee9246d04af>