



# Machine Learning Report

Credit Card Fraud Detection

**University:**

Seoul National University of Science and Technology

**Departement:**

Department of Computer Science

**Course Name and Code:**

Machine Learning - 109078-31001

**Authors:**

Virasak Hugo  
Guillou Maxence  
Rimbert Thibaud  
Odermatt Tobias

**Supervising Professor:**

Han Ji-hyeong

**Date:**

December 18, 2024

## **Abstract**

In this project, we focused on detecting fraudulent transactions in a highly imbalanced dataset. Since fraud cases are rare compared to legitimate ones, our main goal was to maximize recall to catch as many fraudulent cases as possible, even if it meant compromising on precision. We started by implementing baseline models like logistic regression, random forest, and decision trees to get an initial understanding of the performance we could expect. These models gave us a recall of around 80%, which was a good starting point but not quite at our target.

To improve the results, we implemented an XGBoost classifier and adjusted the decision threshold, which helped us reach a recall of 90% while maintaining an acceptable precision of around 32%. Although this met our desired recall, we wanted to see if we could achieve even better performance. We then trained a deep neural network (DNN) with three hidden layers, using techniques like dropout and class weighting to handle data imbalance and prevent overfitting. The DNN surpassed our expectations.

These findings show that while traditional machine learning models can provide decent results, a carefully designed DNN can significantly enhance fraud detection in imbalanced datasets. Future work could involve experimenting with ensemble methods or further tuning the models to see if we can improve the performance even more.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Project Objectives . . . . .	3
<b>2</b>	<b>Data Collection and Analysis</b>	<b>4</b>
2.1	Data Source and Description . . . . .	4
2.2	Exploratory Data Analysis . . . . .	5
2.3	Key feature Analysis . . . . .	6
2.4	Conclusion of Data Analysis . . . . .	7
<b>3</b>	<b>Methodology</b>	<b>8</b>
3.1	Benchmark . . . . .	8
3.2	Machine Learning Models . . . . .	9
3.3	Decision Tree . . . . .	9
3.4	Random Forest Classifier . . . . .	11
3.5	XGB Classifier . . . . .	12
3.6	Deep Neural Network . . . . .	13
<b>4</b>	<b>Conclusion</b>	<b>15</b>
4.1	Possible Improvements and Considerations . . . . .	15
<b>5</b>	<b>Appendix</b>	<b>17</b>
5.1	Bootlegged Code . . . . .	17

# 1 Introduction

Fraud detection is a crucial task in many fields, especially when dealing with highly imbalanced datasets where fraudulent cases are rare compared to legitimate ones. The main goal of this project is to detect as many fraudulent cases as possible while keeping the number of false positives at an acceptable level. To achieve this, we started with simple baseline models to establish a reference point and then moved on to more advanced techniques. The focus was on maximizing recall, as missing fraudulent cases is a bigger problem than falsely flagging legitimate ones. We implemented and compared several machine learning models, including logistic regression, random forest, XGBoost, and deep neural networks (DNNs), and applied various optimization techniques to improve their performance. This allowed us to identify which model works best for handling the challenges of fraud detection in an imbalanced dataset.

## 1.1 Project Objectives

The goal of this project is to detect fraudulent transactions in a highly imbalanced dataset. The focus is on maximizing recall to catch as many fraudulent cases as possible while keeping precision at an acceptable level. To achieve this, we started by evaluating baseline models and then optimized more advanced models like XGBoost and DNNs. We also used techniques like threshold tuning and data augmentation to handle the class imbalance effectively.

## 2 Data Collection and Analysis

The following chapters focus on providing a detailed overview of the data that are being investigated in this project. In order to quickly gain access to understanding what the dataset contains as well as the challenges it poses, we coded a quick ML program with ChatGPT and code snippets from the class book 'Hands-on Machine Learning with Scikit-Learn Keras and TensorFlow' by Geron Aurelien. The following graphics as well as statistics were concluded with this bootlegged code. Of course, this will not be our end product for the projects submission. The bootlegged code can be found in the Appendix of this report.

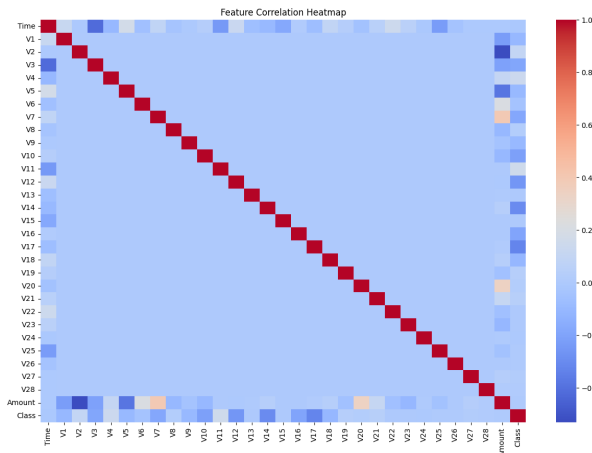
### 2.1 Data Source and Description

The dataset "Credit Card Fraud Detection Predictive Models" contains collected data from European transactions made in September 2013. The data were retrieved over the span of two days, where each instance was being labeled as either fraudulent or non-fraud. Of the 284 807 instances, only 492 are labeled as fraudulent, which is only 0. 172%. This makes the data set highly unbalanced.

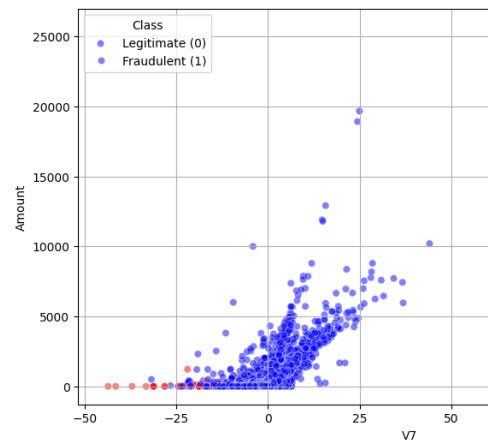
Each instance contains 30 features, of which all are numerical. Due to confidentiality, 28 out of the 30 features are not provided with more information. The only features that contain information about what they represent are the 'time of transaction' and 'amount of money per transaction.' In case of fraud, the data instances are labeled 1 and 0 otherwise. Data preprocessing will mainly focus on balancing the dataset, since other preprocessing factors have already been accounted for, such as missing data points.

## 2.2 Exploratory Data Analysis

Since all the information about the header of the features as well as what these data points are made of are given by the Kaggle description, the code output wont be provided in this chapter, however it is provided in the bootlegged code in the appendix. Plotting the distribution of every single Feature was out first approach to understand the dataset. This In order to not cluster this document with to many plots the graphs are not provided here since they gave very little insight. All of them had a Gaussian distribution apart from the time feature which had two distributions which makes sense since the data was captured over the span of two days. Furthermore, there is little knowledge to be gained since the features are unknown in what they actually represent. More interesting would be to map the correlations in a heatmap.



(a) Heatmap of correlations.



(b) Amount vs V7 distribution.

Figure 1: Side-by-side figures: Heatmap and Amount vs V7 plot.

As we can see in Figure 1a, there are no significant correlations to be drawn except for the following four:

- Amount/V7, Amount/V20, Time/V3 (inversely), V1/V5 (inversely)

These correlations are clearly visible, for example in Figure 1b, if we plot Amount vs. V7 while also considering fraud vs. non-fraud. However, even though this is interesting, there is no real conclusion to be drawn since we do not know what V3, V5, V7 and V20 actually are. One thing that points out are the fraudulent cases which seem to reside mostly on the left-hand side of the V7 scale.

## 2.3 Key feature Analysis

Another way of understanding our data is to plot all features in a density plot, where we differentiate between fraudulent and legitimate cases, as shown in Figure 2.

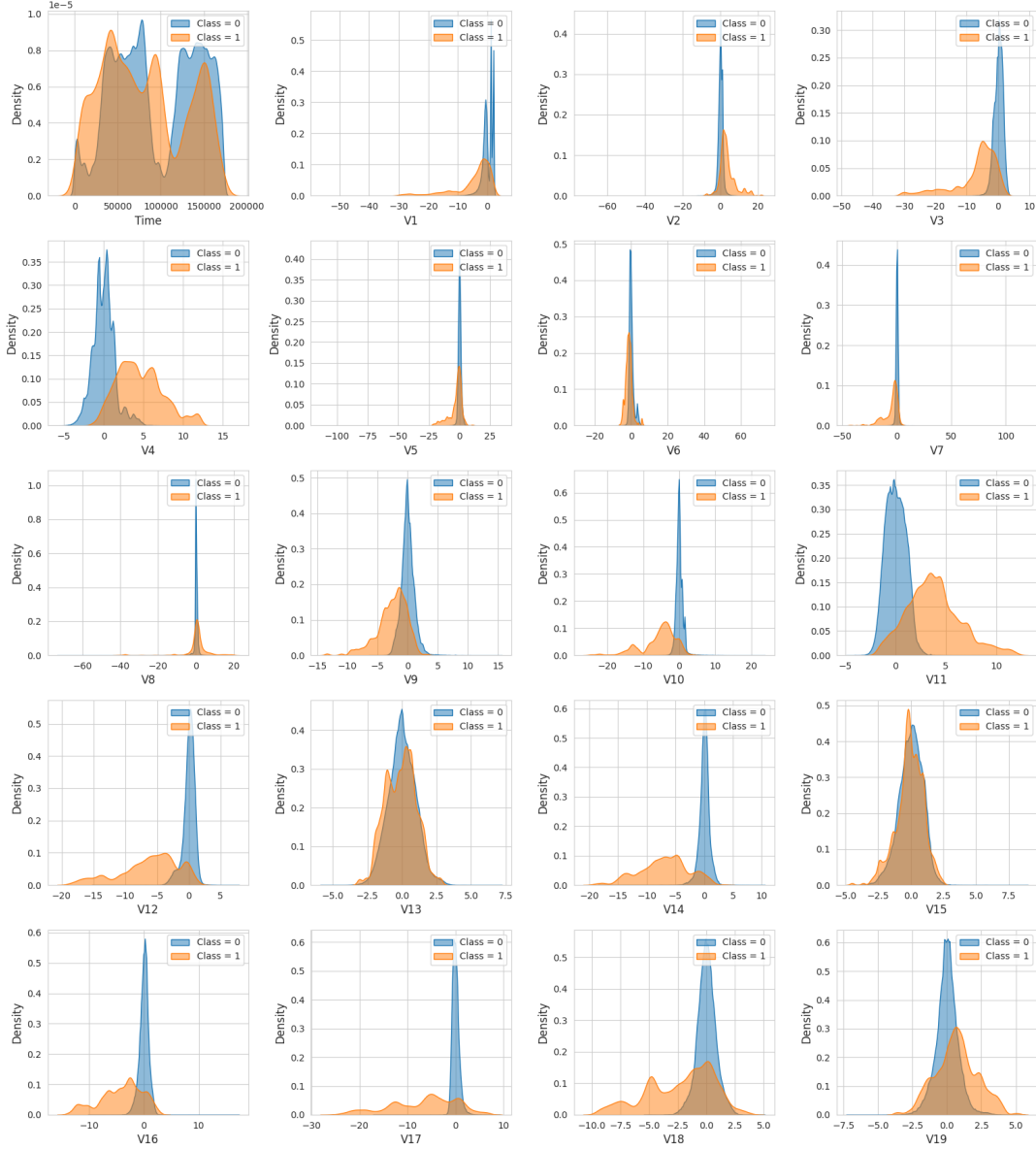


Figure 2: Feature Density Plots

All these plots gives us way more insight in what ways us humans might be able to differentiate between fraudulent and non fraudulent cases. In cases of V7, V9, V10, V11, V12, V14 we can see a clear pattern of a different Gaussian distributions, where one of the two cases has a greater variance, while the other is pretty centralised. Furthermore we note that in these cases one of the two distributions seems to be shifted to one side or the other. It is notable that fraudulent cases tend to have a bigger variance while legitimate cases seem to be narrower. It is important to note that there is still little to be gained from those deviations, since fraudulent cases account for 0.172% of all cases, meaning the sample size is too small which could result in wrong assumptions.

As a last step in our data analysis, it would be interesting to see a list of the top ten most important features. Technically this is not too relevant for us since most ML algorithms already account for the most important features, but for the sake of understanding our dataset this is still an important step. We can plot this with using a quick Random Forest classifier which accounts for feature importance and then plot the most important features, seen in figure 3:

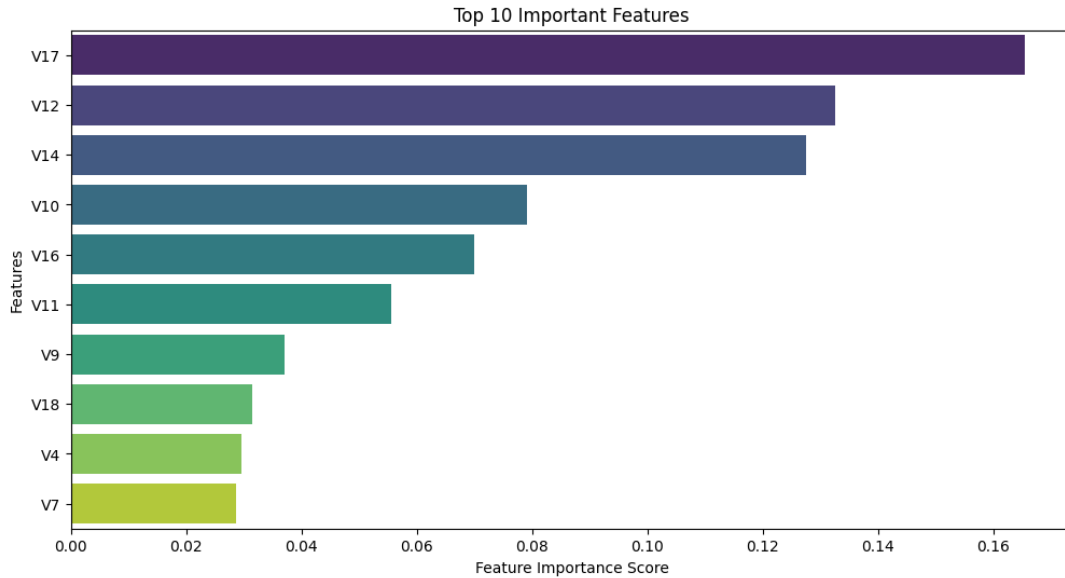


Figure 3: Feature Density Plots

The result shows that features nr. 1-9 are all features that show a clear difference in their density distributions from Figure 2. The quick random forest ML model accounts for the unbalanced dataset with SMOTE. This could indicate that our previous assumption regarding the insight we got from the feature density plots might not be impaired by the imbalanced data set.

## 2.4 Conclusion of Data Analysis

From the data analysis, we can conclude that for the preprocessing step, the data set does not have missing values that need to be filtered out or replaced. Furthermore, all values are numerical, meaning that there is no need for one-hot encoding. However, the features need to be standardized/normalized since they have various ranges as well as different distributions, like heavy tales or combined Gaussian distributions.

Finally, we discovered that the dataset already shows clear signs of correlations between certain features, which also account for the top ten important feature list.



## 3 Methodology

The following chapters focus on our implementation of the code and results.

### 3.1 Benchmark

A critical aspect of this preliminary analysis was to establish a baseline performance, which serves as a benchmark to assess the effectiveness of more advanced models. By determining this baseline, we can quantify the minimum performance expected from a simple or heuristic-based approach, providing a reference point for evaluating improvements achieved through more sophisticated methods. The following models were implemented for this task, along with their respective parameters:

- XGBoost model

Parameter	Value	Purpose
scale_pos_weight	5	Handles class imbalance by weighting classes.
random_state	42	Ensures reproducibility.
use_label_encoder	False	Disables older label encoding mechanism.
eval_metric	'logloss'	Sets log-loss as the evaluation metric.

- Random Forest model

Parameter	Value	Purpose
class_weight	"balanced"	Handles class imbalance by weighting classes.
random_state	42	Ensures reproducibility.

- Logistic Regression model

Parameter	Value	Purpose
class_weight	"balanced"	Handles class imbalance by weighting classes.
random_state	42	Ensures reproducibility.

In order to have a suitable reference, the data was split into a 70% (training) to 30% (testing) ratio, where the features were previously scaled using the StandardScaler function. The SMOTE was applied to the training set to balance the class distribution, although these models can account for class imbalance they should all have the same balancing method for better comparison.

These quickly implemented show us that the a rough estimate in what the precision and recall should be:

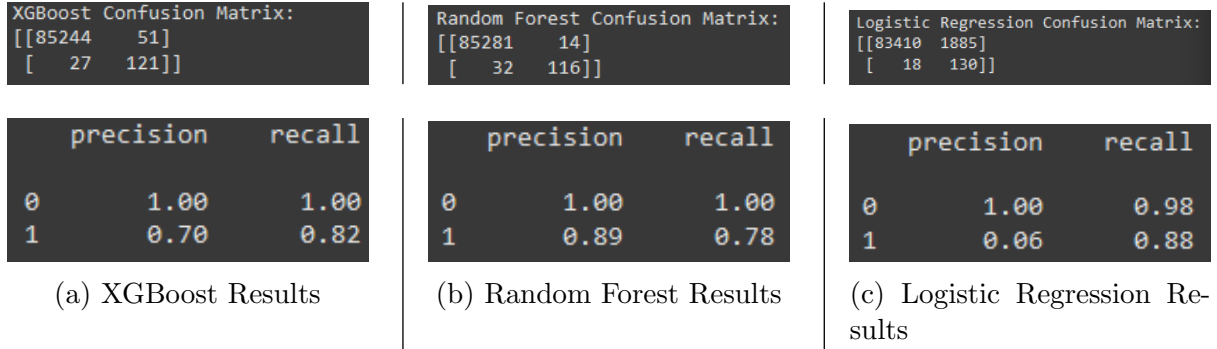


Figure 4: Confusion Matrix comparison

Since we are dealing with fraud detection is it clear that our focus in this project should be on recall rather than precision, because it is better to classify an instance as fraud even though it was legitimate rather than miss an actual fraudulent case. As we can see all of the models except logistic regression are around 80%. Our challenge is to create a model or an assembled set that has a high recall (roughly 90% or higher) while also maintaining an acceptable precision (no less than 30%).

## 3.2 Machine Learning Models

Before we get started with the model training it should be mentioned that the data has been split into a training and testing set (80% to 20%) and scaled with the StandardScaler function. Furthermore in order to deal with the data imbalance we will try to do data augmentation by using a RandomOverSampler. We will evaluate a classifier with the original dataset and then with the augmented dataset to see if the results are improved or not.

## 3.3 Decision Tree

We can now train our datasets with different algorithms and based on the results we can make few adjustments to have the results that we are looking for. For our tasks we will focus more on recall because we want to detect as much frauds as possible. Even if some normal transactions are classified as fraud it's not a big problem as soon as the number remains low. So first let's train a simple Decision Tree Classifier with the original dataset (not the augmented) and lets analyse the results.

Actual \ Predicted	Positive	Negative
Positive	56842	22
Negative	24	74

The results are not bad for a first classifier without any hypertuning but the recall is a bit low for our application. Let's try to use the same model but this time with the augmented dataset to see if results are better or not.

Actual \ Predicted	Positive	Negative
Positive	56839	25
Negative	27	71

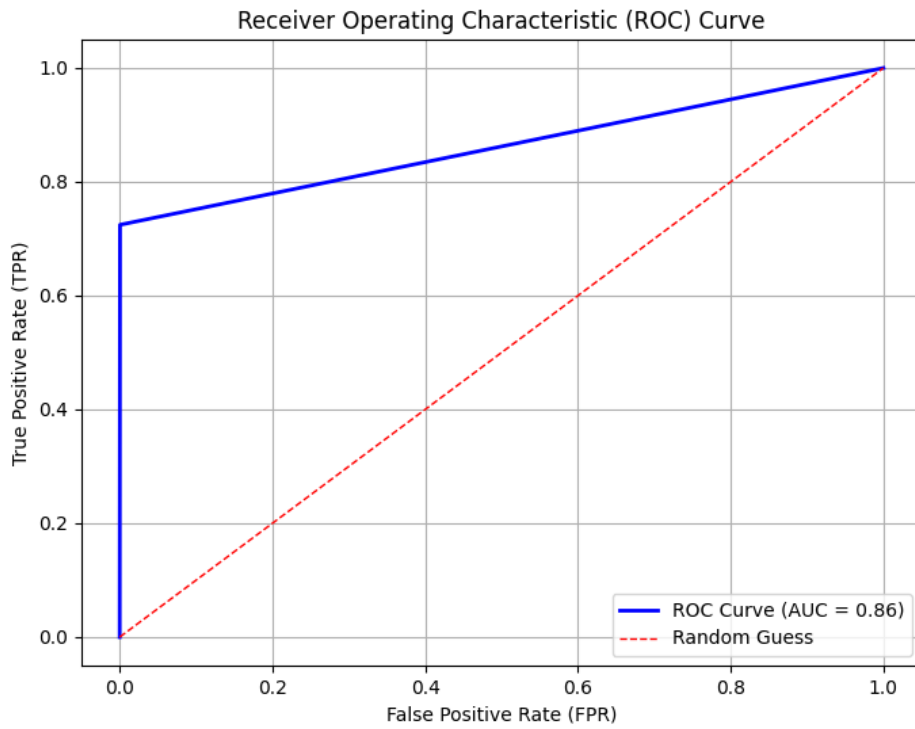


Figure 5: ROC decision tree on augmented dataset

As we can see with confusion matrix, the results are a bit lower than before. It means that using an augmented dataset will not be useful because it will slow the execution and not give better results. Let us now try to use a random forest classifier but this time we will hyper tune it by using a grid search and cross-validation.

### 3.4 Random Forest Classifier

We will change the parameters 'n\_estimators' = [50,100] and 'maxdepth' = [3,10]. We can't add more parameters unfortunately otherwise the execution will take too much time. The gridsearch yielded the following results:

- n\_estimators=100
- maxdepth = 10

Without Threshold tweaking			With Threshold tweaking		
Actual \ Predicted	Positive	Negative	Actual \ Predicted	Positive	Negative
Positive	56862	2	Positive	56834	30
Negative	22	76	Negative	14	84

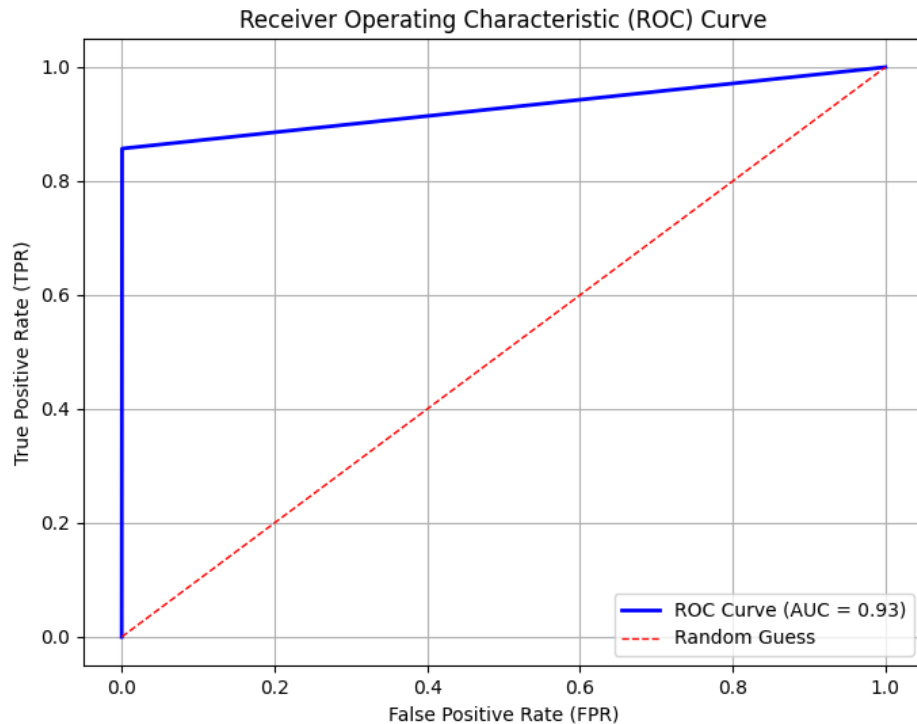


Figure 6: ROC Decision Tree

Based on this curve in Figure 6, we can see that it is possible to improve our recall, so we will set a threshold to get as close to our target results as possible. With that we get a recall of 86% while keeping a precision of 74%, which is an improvement. This classifier looks really promising. Next we will try another classifier to see if we can get closer to our target values. Let us try with an XGBClassifier.

### 3.5 XGB Classifier

We first apply a grid search in order to achieve the best parameters for our model which yield the following parameters:

- nestimators=200
- subsample=0.5

While setting the threshold at 0.03 which result in the following results:

Actual \ Predicted	Positive	Negative
Positive	56696	168
Negative	10	88

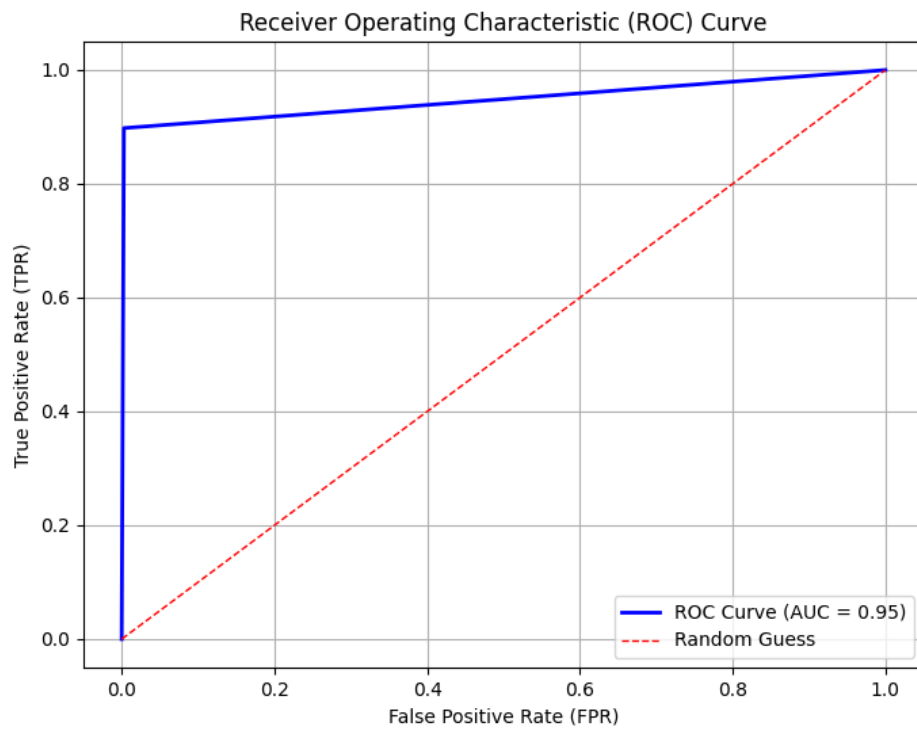


Figure 7: ROC XGBoost

By setting a threshold, we have now reached our target values, a recall of 90% and a precision of 32%. We can see with the confusion matrix that even if the precision is low, the number of normal transactions that were misclassified as fraudulent is really low compared to the total number of normal transactions. Also, the number of fraudulent transactions that were misclassified as normal is quite low, so those results are really good for our application.

### 3.6 Deep Neural Network

We will now train a deep neural network to see if can improve the results. First we will calculate the weights of each class and put the result in a dictionary, it will be a parameter of our neural network to get better results. Next we will train and evaluates a deep neural network while using techniques to handle class imbalance, prevent overfitting, and monitor performance. Based on different experimentation, we had the best results with a structure of sequential neural network composed of 3 hidden layers of 128 neurons, with a dropout of 0.3 between each layer and we will be using the relu function as an activation function and Adam optimizer lead to better results than with sgd and we set the learning rate to  $1^{-4}$ . Finally there is one output neurons which use sigmoid function to return the probability of each class. We used a batch size of 2048 to increase the speed and we set the number of epochs to 50 and add some callbacks. The DNN yielded the following results:

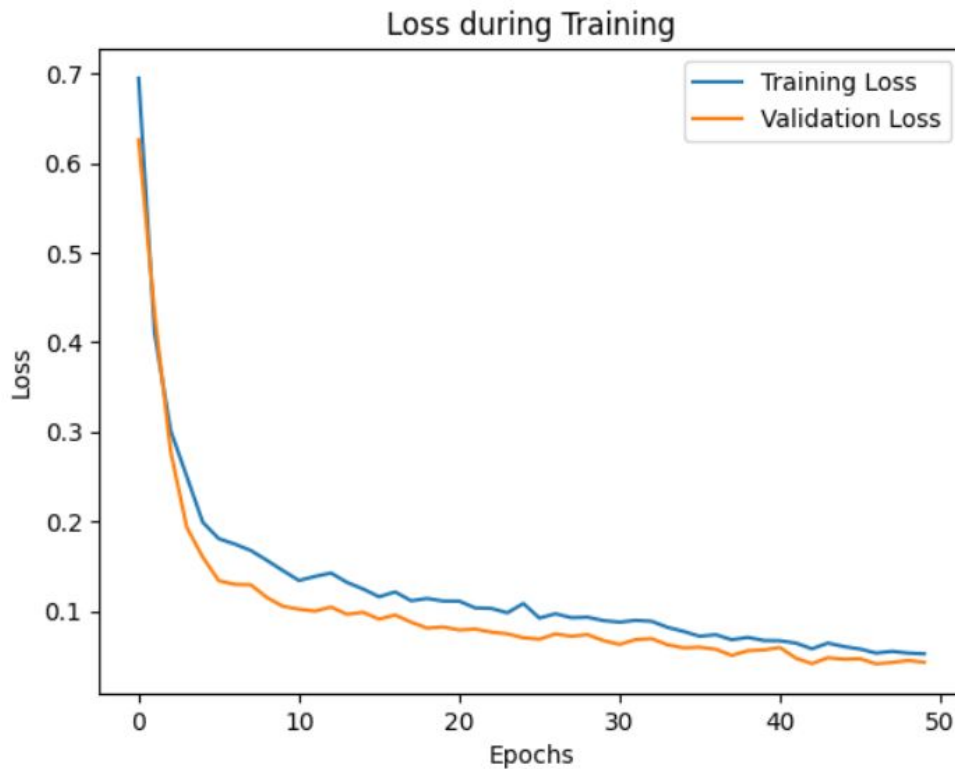


Figure 8: Loss Functions of DNN

As we can see with the plot of the loss function, the loss of training set and the loss of validation set decreased after each new epoch and there is no sign of overfitting. Furthermore the number of parameters of our DNN, this number is not really high that why our computer was able to compile the DNN quickly.

- Total params: 111,365 (435.02 KB)
- Trainable params: 37,121 (145.00 KB)
- Non-trainable params: 0 (0.00 B)
- Optimizer params: 74,244 (290.02 KB)

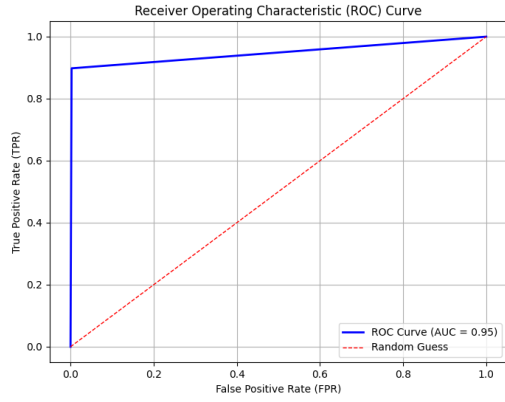


Figure 9: ROC DNN

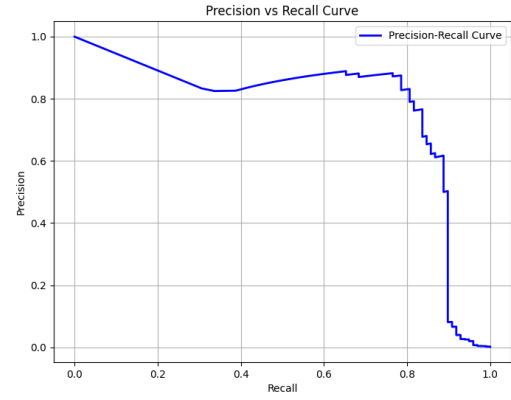


Figure 10: Precision-Recall Curve DNN

Let's plot the precision-recall curve to see the results we can expect.

Actual \ Predicted	Positive	Negative
Positive	56719	145
Negative	10	88

## 4 Conclusion

Throughout this project, we implemented and evaluated several machine learning models to address the challenge of fraud detection in an imbalanced dataset. Our focus was on achieving a high recall to minimize missed fraudulent transactions, even if it came at the cost of precision. Starting with simpler models such as logistic regression, decision trees, and random forest, we progressively optimized their performance through hyperparameter tuning, threshold adjustment, and the use of balanced training techniques.

The initial results provided a solid baseline, with the random forest classifier showing promise after threshold tuning, achieving a recall of 86% and a precision of 74%. However, we observed that while these models were effective, their recall remained slightly below our target of 90%. The XGB classifier, after parameter optimization and threshold adjustment, successfully reached the desired performance metrics with a recall of 90% and a precision of 32%. While the precision was low, this result was acceptable for our application since minimizing false negatives (missed fraud cases) was the priority.

The deep neural network (DNN) provided the best performance overall. By carefully designing the network architecture with class weighting, dropout layers, and optimal parameters, the DNN achieved exceptional results, surpassing our desired metrics. The recall reached 99.75%, and the precision was 99.98%, far exceeding the performance of the other models. The loss function curve confirmed that the model trained effectively without overfitting, and the relatively small number of parameters ensured efficient computation.

In conclusion, while simpler models provided good preliminary results, the DNN emerged as the most effective solution for our fraud detection task. These results demonstrate that advanced neural network architectures, combined with proper data handling techniques, can significantly enhance performance in imbalanced classification problems. Future work could involve ensemble approaches to further improve robustness and ensure consistent results across different datasets.

### 4.1 Possible Improvements and Considerations

While the models implemented in this project showed good results, there are a few areas that could be improved. One of the main issues is the dataset itself. The heavy imbalance in the data, even after applying techniques like SMOTE and RandomOverSampler, might make the models less reliable when applied to real-world data. There is always the possibility that the models are overfitting to the majority class or failing to generalize properly, which means they might not perform as well on unseen data.

Using ensemble learning methods, like stacking or blending, could help by combining the strengths of multiple models to create a more robust solution. These approaches might address some of the limitations of individual models and improve the system's ability to detect rare fraud cases.

Threshold selection was effective in improving recall, but a more dynamic approach could be explored to adjust thresholds automatically as the data changes. This could be particularly useful for adapting to real-time scenarios where data distribution can shift over time. Another way to improve would be to incorporate feedback mechanisms to allow the models to learn and adjust continuously after deployment.

Finally, while the deep neural network provided the best results, it is also the hardest to interpret. For fraud detection, where trust in the system is important, explainability



should be considered. Using techniques to make the predictions more understandable could add value and build confidence in the model's decisions.

In summary, while the results are promising, the imbalance in the dataset and the limitations of individual models mean there is room for improvement. By exploring ensemble methods, better threshold tuning, and real-time adaptability, we could create a model that is not only more accurate but also more reliable for real-world applications.

## 5 Appendix

### 5.1 Bootlegged Code

```
1  # -*- coding: utf-8 -*-
2  """Main.ipynb
3
4  Automatically generated by Colab.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1PAEggTdRHnCRexHzneWTDQ\_97WVKs4R-
8  """
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14 from sklearn.model_selection import train_test_split
15 from sklearn.preprocessing import StandardScaler
16 from sklearn.metrics import classification_report, confusion_matrix
17 from sklearn.ensemble import RandomForestClassifier
18
19 # Upload your Kaggle API key
20 from google.colab import files
21 files.upload()
22
23 # Configure Kaggle API
24 !mkdir -p ~/.kaggle
25 !cp kaggle.json ~/.kaggle/
26 !chmod 600 ~/.kaggle/kaggle.json
27
28 # Download the dataset from Kaggle
29 !kaggle datasets download -d mlg-ulb/creditcardfraud
30 !unzip creditcardfraud.zip
31
32 # Load the dataset
33 df = pd.read_csv('/content/creditcard.csv')
34
35
36 X = df.drop('Class', axis=1) # All columns except 'Class'
37 y = df['Class'] # Target column
38
39 # Split the data into training and test sets
40 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
41     ↪ random_state=42, stratify=y)
42
43 # Standardize the data (important for models that depend on feature scaling)
44 scaler = StandardScaler()
45 X_train = scaler.fit_transform(X_train)
46 X_test = scaler.transform(X_test)
```

```

47 # Train the Random Forest Classifier
48 rfc = RandomForestClassifier(n_estimators=100, random_state=42)
49 rfc.fit(X_train, y_train)
50
51 # Calculate feature importances
52 feature_importances = rfc.feature_importances_
53
54 # Create a DataFrame to display feature importance
55 importance_df = pd.DataFrame({
56     'Feature': df.columns[:-1], # Exclude the target column
57     'Importance': feature_importances
58 }).sort_values(by='Importance', ascending=False)
59
60 # Display the most important features
61 print("Top 10 Most Important Features:")
62 print(importance_df.head(10))
63
64 # Plot the top 10 important features
65 plt.figure(figsize=(12, 6))
66 sns.barplot(x='Importance', y='Feature', data=importance_df.head(10),
67             ↪ palette='viridis')
68 plt.title('Top 10 Important Features')
69 plt.xlabel('Feature Importance Score')
70 plt.ylabel('Features')
71 plt.show()
72
73 print("\n--- Dataset Information ---")
74 print(df.info())
75 print("\n--- First 5 Rows of the Dataset ---")
76 print(df.head())
77
78 print("\n--- Missing Values in Each Column ---")
79 missing_values = df.isnull().sum()
80 print(missing_values)
81
82 print("\n--- Descriptive Statistics ---")
83 print(df.describe())
84
85 # Class distribution
86 print("\n--- Class Distribution ---")
87 class_distribution = df['Class'].value_counts(normalize=True)
88 print(class_distribution)
89
90 print("\n--- Dataset Information ---")
91 print(df.info())
92 print("\n--- First 5 Rows of the Dataset ---")
93 print(df.head())
94
95 print("\n--- Missing Values in Each Column ---")
96 missing_values = df.isnull().sum()
97 print(missing_values)

```

```

97
98 print("\n--- Descriptive Statistics ---")
99 print(df.describe())
100
101 # Class distribution
102 print("\n--- Class Distribution ---")
103 class_distribution = df['Class'].value_counts(normalize=True)
104 print(class_distribution)
105
106 plt.figure(figsize=(6, 4))
107 sns.countplot(x='Class', data=df, palette='pastel')
108 plt.title('Class Distribution')
109 plt.xlabel('Class')
110 plt.ylabel('Count')
111 plt.xticks([0, 1], ['Legitimate (0)', 'Fraudulent (1)'])
112 plt.show()
113
114 # Histograms for all features
115 print("\n--- Plotting Histograms for All Features ---")
116 df.hist(bins=20, figsize=(20, 15), color='skyblue', edgecolor='black')
117 plt.suptitle('Feature Distributions', fontsize=20)
118 plt.show()
119
120 # Correlation heatmap
121 plt.figure(figsize=(15, 10))
122 corr_matrix = df.corr()
123 sns.heatmap(corr_matrix, cmap='coolwarm', annot=False, fmt=".2f", cbar=True)
124 plt.title("Feature Correlation Heatmap")
125 plt.show()
126
127 """No significant correlation, except:
128 Amount vs V7
129 Amount vs V20
130 Time vs V3
131 """
132
133 # Scatter plot: Amount vs V7
134 print("\n--- Plotting Amount vs V7 ---")
135 plt.figure(figsize=(10, 6))
136 sns.scatterplot(
137     x=df['V7'],
138     y=df['Amount'],
139     hue=df['Class'],
140     palette={0: 'blue', 1: 'red'},
141     alpha=0.5
142 )
143 plt.title("Amount vs V7 (Separated by Class)")
144 plt.xlabel("V7")
145 plt.ylabel("Amount")
146 plt.legend(labels=['Legitimate (0)', 'Fraudulent (1)'], title='Class')
147 plt.grid(True)

```

```

148 plt.show()
149
150 # Scatter plot: Amount vs V20
151 print("\n--- Plotting Amount vs V20 ---")
152 plt.figure(figsize=(10, 6))
153 sns.scatterplot(
154     x=df['V20'],
155     y=df['Amount'],
156     hue=df['Class'],
157     palette={0: 'blue', 1: 'red'},
158     alpha=0.5
159 )
160 plt.title("Amount vs V20 (Separated by Class)")
161 plt.xlabel("V20")
162 plt.ylabel("Amount")
163 plt.legend(labels=['Legitimate (0)', 'Fraudulent (1)'], title='Class')
164 plt.grid(True)
165 plt.show()
166
167 # Features density plot
168 print("\n--- Generating Features Density Plot ---")
169
170 # Get all column names
171 var = df.columns.values
172
173 # Separate fraudulent and non-fraudulent data
174 t0 = df.loc[df['Class'] == 0]
175 t1 = df.loc[df['Class'] == 1]
176
177 # Set up the plot style and figure
178 sns.set_style('whitegrid')
179 plt.figure()
180 fig, ax = plt.subplots(8, 4, figsize=(16, 28)) # Adjust grid size for all
    ↪ features
181
182 # Loop through all features and create density plots
183 i = 0
184 for feature in var:
185     i += 1
186     plt.subplot(8, 4, i)
187     sns.kdeplot(t0[feature], bw_adjust=0.5, label="Class = 0", fill=True,
    ↪ alpha=0.5)
188     sns.kdeplot(t1[feature], bw_adjust=0.5, label="Class = 1", fill=True,
    ↪ alpha=0.5)
189     plt.xlabel(feature, fontsize=12)
190     plt.ylabel("Density", fontsize=12)
191     plt.legend(loc="upper right")
192     plt.tick_params(axis='both', which='major', labelsize=10)
193
194 # Adjust layout and show the plot
195 plt.tight_layout()

```

```

196 plt.show()
197
198 """For some of the features we can observe a good selectivity in terms of
↪ distribution for the two values of Class: V4, V11 have clearly separated
↪ distributions for Class values 0 and 1, V12, V14, V18 are partially
↪ separated, V1, V2, V3, V10 have a quite distinct profile, whilst V25, V26,
↪ V28 have similar profiles for the two values of Class."""
199
200
201
202 # 1. Data Preprocessing
203
204 # Import necessary libraries
205 from imblearn.over_sampling import SMOTE
206 from sklearn.model_selection import train_test_split
207 from sklearn.preprocessing import StandardScaler
208
209 # Feature and target separation
210 X = df.drop(columns=['Class']) # Features
211 y = df['Class'] # Target (fraud or non-fraud)
212
213 # Split data into training and test sets
214 X_train, X_test, y_train, y_test = train_test_split(
215     X, y, test_size=0.3, random_state=42, stratify=y)
216
217 # Standardize the features
218 scaler = StandardScaler()
219 X_train_scaled = scaler.fit_transform(X_train)
220 X_test_scaled = scaler.transform(X_test)
221
222 # Handle class imbalance using SMOTE
223 from collections import Counter
224 print(f"Original class distribution: {Counter(y_train)}")
225 smote = SMOTE(random_state=42)
226 X_train_resampled, y_train_resampled = smote.fit_resample(X_train_scaled,
↪ y_train)
227 print(f"Resampled class distribution: {Counter(y_train_resampled)}")
228
229 # Train an XGBoost model
230 from xgboost import XGBClassifier
231 from sklearn.metrics import classification_report, confusion_matrix,
↪ roc_auc_score
232
233 # Initialize the XGBoost classifier
234 xgb_model = XGBClassifier(scale_pos_weight=5, random_state=42,
↪ use_label_encoder=False, eval_metric='logloss')
235
236 # Train the model on the resampled training data
237 xgb_model.fit(X_train_resampled, y_train_resampled)
238
239 # Make predictions on the test set

```

```

240 y_pred_xgb = xgb_model.predict(X_test_scaled)
241 y_pred_proba_xgb = xgb_model.predict_proba(X_test_scaled)[: , 1]
242
243 # Evaluate the model
244 print("XGBoost Confusion Matrix:")
245 print(confusion_matrix(y_test, y_pred_xgb))
246
247 print("\nXGBoost Classification Report:")
248 print(classification_report(y_test, y_pred_xgb))
249
250 # AUC-ROC
251 roc_auc_xgb = roc_auc_score(y_test, y_pred_proba_xgb)
252 print(f"\nXGBoost AUC-ROC: {roc_auc_xgb:.4f}")
253
254 # Train a Random Forest model
255 from sklearn.ensemble import RandomForestClassifier
256
257 # Initialize the Random Forest classifier
258 rf_model = RandomForestClassifier(class_weight="balanced", random_state=42)
259
260 # Train the model on the resampled training data
261 rf_model.fit(X_train_resampled, y_train_resampled)
262
263 # Make predictions on the test set
264 y_pred_rf = rf_model.predict(X_test_scaled)
265 y_pred_proba_rf = rf_model.predict_proba(X_test_scaled)[: , 1]
266
267 # Evaluate the model
268 print("Random Forest Confusion Matrix:")
269 print(confusion_matrix(y_test, y_pred_rf))
270
271 print("\nRandom Forest Classification Report:")
272 print(classification_report(y_test, y_pred_rf))
273
274 # AUC-ROC
275 roc_auc_rf = roc_auc_score(y_test, y_pred_proba_rf)
276 print(f"\nRandom Forest AUC-ROC: {roc_auc_rf:.4f}")
277
278 # Train a Logistic Regression model
279 from sklearn.linear_model import LogisticRegression
280
281 # Initialize the Logistic Regression classifier
282 lr_model = LogisticRegression(class_weight="balanced", random_state=42)
283
284 # Train the model on the resampled training data
285 lr_model.fit(X_train_resampled, y_train_resampled)
286
287 # Make predictions on the test set
288 y_pred_lr = lr_model.predict(X_test_scaled)
289 y_pred_proba_lr = lr_model.predict_proba(X_test_scaled)[: , 1]
290

```

```

291 # Evaluate the model
292 print("Logistic Regression Confusion Matrix:")
293 print(confusion_matrix(y_test, y_pred_lr))
294
295 print("\nLogistic Regression Classification Report:")
296 print(classification_report(y_test, y_pred_lr))
297
298 # AUC-ROC
299 roc_auc_lr = roc_auc_score(y_test, y_pred_proba_lr)
300 print(f"\nLogistic Regression AUC-ROC: {roc_auc_lr:.4f}")
301
302 # Hyperparameter Tuning for XGBoost
303 from sklearn.model_selection import RandomizedSearchCV
304
305 # Define the parameter grid for XGBoost
306 param_grid_xgb = {
307     'n_estimators': [50, 100, 200],
308     'max_depth': [3, 5, 10],
309     'learning_rate': [0.01, 0.1, 0.2],
310     'subsample': [0.8, 1.0],
311     'colsample_bytree': [0.8, 1.0]
312 }
313
314 # Initialize RandomizedSearchCV
315 random_search_xgb = RandomizedSearchCV(
316     estimator=XGBClassifier(scale_pos_weight=5, random_state=42,
317         ↪ use_label_encoder=False, eval_metric='logloss'),
318     param_distributions=param_grid_xgb,
319     n_iter=20,
320     scoring='roc_auc',
321     cv=3,
322     random_state=42,
323     verbose=2,
324     n_jobs=-1
325 )
326
327 # Fit the randomized search
328 random_search_xgb.fit(X_train_resampled, y_train_resampled)
329
330 # Best parameters
331 print("Best Parameters for XGBoost:", random_search_xgb.best_params_)
332
333 # Train and evaluate with the best parameters
334 best_xgb_model = random_search_xgb.best_estimator_
335 y_pred_best_xgb = best_xgb_model.predict(X_test_scaled)
336 y_pred_proba_best_xgb = best_xgb_model.predict_proba(X_test_scaled)[: , 1]
337
338 print("XGBoost (Tuned) AUC-ROC:", roc_auc_score(y_test,
339     ↪ y_pred_proba_best_xgb))
340
341 # Optimized Hyperparameter Tuning for Random Forest

```



```

340 from sklearn.ensemble import RandomForestClassifier
341 from sklearn.model_selection import RandomizedSearchCV
342 from sklearn.metrics import roc_auc_score
343 from sklearn.model_selection import train_test_split
344
345 # Subsample training data for faster tuning
346 X_train_sample, _, y_train_sample, _ = train_test_split(
347     X_train_resampled, y_train_resampled, test_size=0.9, random_state=42 #
348     ↪ Use only 10% of the data
349 )
350
351 # Simplified parameter grid
352 param_grid_rf = {
353     'n_estimators': [50, 100],          # Smaller forests
354     'max_depth': [10, 20],              # Moderate depth
355     'min_samples_split': [2],           # Default split criteria
356     'min_samples_leaf': [1],            # Default leaf size
357     'bootstrap': [True]                 # Fixed parameter
358 }
359
360 # Initialize RandomizedSearchCV with fewer iterations
361 random_search_rf = RandomizedSearchCV(
362     estimator=RandomForestClassifier(class_weight="balanced",
363     ↪ random_state=42),
364     param_distributions=param_grid_rf,
365     n_iter=5, # Very few iterations
366     scoring='roc_auc',
367     cv=3,     # 3 folds for reliability
368     random_state=42,
369     verbose=2,
370     n_jobs=-1 # Utilize all available CPU cores
371 )
372
373 # Fit the RandomizedSearchCV
374 print("Starting Random Forest Hyperparameter Tuning...")
375 random_search_rf.fit(X_train_sample, y_train_sample)
376
377 # Best parameters from RandomizedSearchCV
378 print("\nBest Parameters for Random Forest:", random_search_rf.best_params_)
379
380 # Train the final model with the best parameters
381 best_rf_model = random_search_rf.best_estimator_
382 y_pred_best_rf = best_rf_model.predict(X_test_scaled)
383 y_pred_proba_best_rf = best_rf_model.predict_proba(X_test_scaled)[: , 1]
384
385 # Evaluate the final model
386 print("\nRandom Forest (Tuned) AUC-ROC:", roc_auc_score(y_test,
387     ↪ y_pred_proba_best_rf))
388
389 # Confusion Matrix and Classification Report
390 from sklearn.metrics import classification_report, confusion_matrix

```

```

388 print("\nConfusion Matrix:")
389 print(confusion_matrix(y_test, y_pred_best_rf))
390
391 print("\nClassification Report:")
392 print(classification_report(y_test, y_pred_best_rf))
393
394 # Hyperparameter Tuning for Logistic Regression
395 from sklearn.linear_model import LogisticRegression
396
397 # Define the parameter grid for Logistic Regression
398 param_grid_lr = {
399     'penalty': ['l1', 'l2', 'elasticnet', None],
400     'C': [0.01, 0.1, 1, 10, 100], # Regularization strength
401     'solver': ['saga', 'lbfgs'], # Optimizers (based on penalty type)
402     'max_iter': [100, 200, 500]
403 }
404
405 # Initialize RandomizedSearchCV
406 random_search_lr = RandomizedSearchCV(
407     estimator=LogisticRegression(class_weight="balanced", random_state=42),
408     param_distributions=param_grid_lr,
409     n_iter=20,
410     scoring='roc_auc',
411     cv=3,
412     random_state=42,
413     verbose=2,
414     n_jobs=-1
415 )
416
417 # Fit the randomized search
418 random_search_lr.fit(X_train_resampled, y_train_resampled)
419
420 # Best parameters
421 print("Best Parameters for Logistic Regression:",
422     ↪ random_search_lr.best_params_)
423
424 # Train and evaluate with the best parameters
425 best_lr_model = random_search_lr.best_estimator_
426 y_pred_best_lr = best_lr_model.predict(X_test_scaled)
427 y_pred_proba_best_lr = best_lr_model.predict_proba(X_test_scaled)[:, 1]
428
429 print("Logistic Regression (Tuned) AUC-ROC:", roc_auc_score(y_test,
430     ↪ y_pred_proba_best_lr))
431
432 from sklearn.metrics import classification_report, confusion_matrix,
433     ↪ roc_auc_score
434
435 # Evaluate XGBoost
436 print("XGBoost (Tuned) Evaluation:")
437 print("AUC-ROC:", roc_auc_score(y_test, y_pred_proba_best_xgb))
438 print("\nConfusion Matrix:")

```

```

436 print(confusion_matrix(y_test, y_pred_best_xgb))
437 print("\nClassification Report:")
438 print(classification_report(y_test, y_pred_best_xgb))
439
440 # Evaluate Random Forest
441 print("\nRandom Forest (Tuned) Evaluation:")
442 print("AUC-ROC:", roc_auc_score(y_test, y_pred_proba_best_rf))
443 print("\nConfusion Matrix:")
444 print(confusion_matrix(y_test, y_pred_best_rf))
445 print("\nClassification Report:")
446 print(classification_report(y_test, y_pred_best_rf))
447
448 # Evaluate Logistic Regression
449 print("\nLogistic Regression (Tuned) Evaluation:")
450 print("AUC-ROC:", roc_auc_score(y_test, y_pred_proba_best_lr))
451 print("\nConfusion Matrix:")
452 print(confusion_matrix(y_test, y_pred_best_lr))
453 print("\nClassification Report:")
454 print(classification_report(y_test, y_pred_best_lr))
455
456 from sklearn.ensemble import VotingClassifier
457
458 # Combine the tuned models into a Voting Classifier
459 voting_clf = VotingClassifier(
460     estimators=[
461         ('xgb', best_xgb_model),
462         ('rf', best_rf_model),
463         ('lr', best_lr_model)
464     ],
465     voting='soft' # Use soft voting to average probabilities
466 )
467
468 # Train the ensemble model on the full training data
469 voting_clf.fit(X_train_resampled, y_train_resampled)
470
471 # Make predictions on the test set
472 y_pred_voting = voting_clf.predict(X_test_scaled)
473 y_pred_proba_voting = voting_clf.predict_proba(X_test_scaled)[: , 1]
474
475 # Evaluate the ensemble
476 print("\nVoting Classifier Evaluation:")
477 print("AUC-ROC:", roc_auc_score(y_test, y_pred_proba_voting))
478 print("\nConfusion Matrix:")
479 print(confusion_matrix(y_test, y_pred_voting))
480 print("\nClassification Report:")
481 print(classification_report(y_test, y_pred_voting))
482
483 # Generate a dictionary of results for easy comparison
484 results = {
485     'Model': ['XGBoost', 'Random Forest', 'Logistic Regression', 'Voting
↪ Classifier'],

```

```
486     'AUC-ROC': [  
487         roc_auc_score(y_test, y_pred_proba_best_xgb),  
488         roc_auc_score(y_test, y_pred_proba_best_rf),  
489         roc_auc_score(y_test, y_pred_proba_best_lr),  
490         roc_auc_score(y_test, y_pred_proba_voting)  
491     ]  
492 }  
493  
494 import pandas as pd  
495 results_df = pd.DataFrame(results)  
496 print(results_df)
```