



## RV Educational Institutions®

### RV Institute of Technology and Management

(Affiliated to VTU, Belagavi)

JP Nagar 8th Phase, Bengaluru – 560076

#### Department of Information Science and Engineering



Course Name: Data Mining and Data Warehousing

Course Code: 21CS643

Assignment 2

**Fraud Detection System for Financial Transactions**

VI Semester

2021 Scheme

Submitted By : Vijay Kumar Gowda K K(1RF21IS065)



## Abstract

This project aims to develop an effective fraud detection system for financial transactions using the "Credit Card Fraud Detection" dataset from Kaggle, which is characterized by a severe class imbalance with only 0.172% of transactions being fraudulent. To address this imbalance, we applied the NearMiss algorithm to create a balanced dataset by undersampling the majority class. A comprehensive exploratory data analysis (EDA) was conducted to understand the dataset's structure and characteristics. Various classifiers, including Logistic Regression, Decision Trees, Random Forests, Support Vector Machines (SVM), and k-Nearest Neighbors (k-NN), were trained and evaluated on the balanced dataset. Among these, the Random Forest classifier demonstrated superior performance, exhibiting high precision and recall, making it a robust choice for detecting fraudulent transactions.

Additionally, a simple neural network was developed and trained on the balanced dataset, showing competitive performance compared to traditional classifiers. The comparison revealed that while the neural network was promising, the Random Forest classifier slightly outperformed it in terms of precision and F1-score. The project addressed common challenges associated with imbalanced datasets, such as potential information loss during undersampling and computational costs. Various mitigation strategies, including oversampling with SMOTE and the use of alternative performance metrics, were explored. The findings underscore the importance of balancing techniques and robust classifiers in fraud detection systems, offering a solid foundation for future enhancements through advanced models and real-time data integration.



## Table of Content

	Page No
1.Introduction	1
2.Definition and Working	3
3.Dataset Overview	6
4.Implementation and Analysis	9
4.1 Gather Sense of Our Data	
4.1.1 Data Overview	
4.1.2 Data Normalization	
4.1.3 Class Distributions	
4.1.4 Data Distribution	
4.2 Scaling and Distributing	
4.2.1 Scaling Code	
4.2.2 Splitting Data	
4.3 Random Under-Sampling	
4.4 Distributing and Correlating Data	
4.5 Exploratory Data Analysis (EDA)	
4.5.1 Matrix Code	
4.5.2 Imbalanced Correlation Matrix	
4.5.3 Imbalanced Subsample Matrix	
4.5.4 Class Negative Correlation	
4.5.5 Class Positive Correlation	
4.5.6 Distributions	
4.5.7 V-Features	
4.6 Dimensionality Reduction and Clustering	
4.6.1 Dimensionality Reduction	
4.6.2 Dimensionality Clustering	
4.6.3 Dimensionality Clustering Mapping	
4.7 Model Training and Evaluation	
4.7.1 Learning Curves	
4.7.2 Logistic Regression Curve	
4.7.3 SMOTE Technique	
5.Conclusion	20
6.Appendix	22
7.References	24



## List Of Figures

### Figures

Figures	Description
Fig 4.1.1	Data Overview
Fig 4.1.2	Data Normalization
Fig 4.1.3	Class Distributions
Fig 4.1.4	Data Distribution
Fig 4.2	Scaling Code
Fig 4.3	Splitting data
Fig 4.4.1	Random Under-Sampling
Fig 4.4.2	Distributing and Correlating Data
Fig 4.5.1	Matrix Code
Fig 4.5.2	Imbalanced Correlation matrix
Fig 4.5.3	Imbalanced Subsample Matrix
Fig 4.5.4	Class Negative Correlation
Fig 4.5.5	Class Positive Correlation
Fig 4.6.1	Distributions
Fig 4.6.2	V-Features
Fig 4.7.1	Dimensionality Reduction
Fig 4.7.2	Dimensionality Clustering
Fig 4.7.3	Dimensionality Clustering Mapping
Fig 4.8	Learning Curves
Fig 4.9	Logistic Regression Curve
Fig 4.10	SMOTE Technique



## CHAPTER 1

# INTRODUCTION

In this project, we aim to develop a robust system for detecting fraudulent financial transactions using various predictive models. Financial fraud detection is a critical task in the banking and finance sector, as it helps prevent substantial monetary losses and protects customers from unauthorized transactions. The primary challenge in this domain is the highly imbalanced nature of the data, where fraudulent transactions are significantly fewer than non-fraudulent ones.

### Data Description

The dataset we are using for this project has been sourced from Kaggle and is specifically tailored for fraud detection. Due to privacy concerns, the features in the dataset have been anonymized and scaled. This means that while we do not know the exact nature of the features (e.g., whether they represent account balance, transaction type, etc.), we do know that they have been preprocessed to ensure uniformity and to protect sensitive information.

### Goals

The primary goals of this project are as follows:

1. Data Understanding: Gain a comprehensive understanding of the dataset, including the distribution of fraudulent and non-fraudulent transactions, and the characteristics of each feature.
2. Data Balancing: Address the imbalance in the dataset by creating a balanced sub-dataframe with an equal ratio of "Fraud" and "Non-Fraud" transactions. We will employ techniques like the NearMiss algorithm to achieve this.
3. Model Selection: Evaluate a range of classifiers to determine which one provides the highest accuracy in detecting fraudulent transactions. This will include traditional machine learning models as well as more complex algorithms.
4. Neural Network Implementation: Develop and train a neural network model to compare its performance against the best-performing classifier. Neural networks have shown great promise in various classification tasks, and we aim to leverage their capabilities for fraud detection.
5. Imbalanced Dataset Challenges: Identify and mitigate common issues associated with imbalanced datasets, such as misleading accuracy scores and the risk of overfitting. We will employ strategies like resampling, cross-validation, and appropriate performance metrics to ensure our models are robust and reliable.



### Objectives:

- Understand the Data: Analyze the distribution of transactions in the dataset.
- Create Balanced Data: Generate a sub-dataframe with a balanced ratio of "Fraud" and "Non-Fraud" transactions using the NearMiss algorithm.
- Evaluate Classifiers: Determine and compare the accuracy of various classifiers.
- Neural Network Comparison: Develop a neural network and compare its accuracy against the best-performing classifier.
- Address Imbalanced Dataset Challenges: Understand and mitigate common issues in handling imbalanced datasets.

The primary objective of this project is to thoroughly analyze and address the challenges posed by an imbalanced credit card fraud detection dataset, leveraging the steps outlined in a Kaggle notebook. First, we aim to understand the data by loading the dataset and performing initial inspections to grasp its structure, size, and feature types. Exploratory Data Analysis (EDA) will follow, focusing on the distribution of the target variable (fraud vs. non-fraud) and other features using visual tools like histograms and scatter plots to identify patterns or anomalies.

To address the significant class imbalance, we will implement the NearMiss algorithm, generating a sub-dataset with a balanced ratio of fraud and non-fraud transactions. This balanced dataset will undergo further EDA to ensure the balancing process preserves the underlying data patterns. Next, we will evaluate various classifiers, including Logistic Regression, Decision Trees, Random Forests, Support Vector Machines (SVM), and k-Nearest Neighbors (k-NN). These models will be trained and tested on the balanced dataset, with their performance assessed using accuracy, precision, recall, F1-score, and ROC-AUC metrics. This evaluation will identify the best-performing classifier.

Additionally, we will develop a simple neural network and train it on the balanced dataset, comparing its performance against the top traditional classifier using the same evaluation metrics. This comparison will highlight the strengths and weaknesses of neural networks versus traditional classifiers in handling this specific problem.

Understanding and mitigating common issues in handling imbalanced datasets is another critical objective. We will discuss challenges like bias towards the majority class, poor generalization, and misleading accuracy metrics. Various mitigation strategies, such as oversampling (e.g., SMOTE), undersampling, and alternative performance metrics, will be explored and implemented to compare their effectiveness against NearMiss. The project will conclude with a summary of findings, best practices for handling imbalanced datasets in credit card fraud detection, and recommendations for further research or improvements.



## Understanding Our Data

### a) Gather Sense of Our Data

The initial step is to gain a basic understanding of the data. The dataset contains anonymized features except for the transaction time and amount. The anonymization was achieved through a PCA transformation applied to most features. Here's a summary of key observations:

- Transaction Amount: The mean transaction amount is approximately USD 88, indicating relatively small transactions.
- Null Values: The dataset contains no null values, eliminating the need for imputation.
- Class Distribution: The dataset is highly imbalanced, with 99.83% of transactions being non-fraudulent and only 0.17% being fraudulent.

Feature Technicalities:

- PCA Transformation: All features, except for time and amount, underwent a PCA transformation for dimensionality reduction.
- Scaling: Features were scaled before applying the PCA transformation, as required for effective dimensionality reduction.

## Preprocessing

### a) Scaling and Distributing

- The dataset features are already scaled. We further analyze the distribution to ensure balanced data preprocessing.

### b) Splitting the Data

- The dataset is split into training and testing sets to evaluate model performance.

## Random UnderSampling and OverSampling

### a) Distributing and Correlating

We explore the distribution of features and their correlations to understand the relationships within the data.

### b) Anomaly Detection

Anomaly detection techniques help identify unusual patterns that may indicate fraud.

### c) Dimensionality Reduction and Clustering (t-SNE)

t-SNE is used for dimensionality reduction and clustering to visualize the data in a lower-dimensional space.



d) Classifiers

Various classifiers are implemented and evaluated for their accuracy in detecting fraud.

e) A Deeper Look into Logistic Regression

Logistic regression is examined in detail to understand its performance and limitations.

f) Oversampling with SMOTE

SMOTE (Synthetic Minority Over-sampling Technique) is used to balance the dataset by generating synthetic samples for the minority class.

## Testing

a) Testing with Logistic Regression

- Logistic regression is tested on the dataset to evaluate its performance.

b) Neural Networks Testing (Undersampling vs Oversampling)

- Neural networks are tested using both undersampled and oversampled datasets to compare their accuracy against logistic regression.

## Correcting Mistakes with Imbalanced Datasets

- Avoid Testing on Resampled Data: Ensure testing is done on the original data, not the oversampled or undersampled data.
- Cross-Validation: Implement oversampling or undersampling during cross-validation, not before.
- Metrics: Use metrics like F1-score, precision/recall score, or confusion matrix instead of accuracy for evaluating performance on imbalanced datasets.



## Methodology

Our approach to this project will be methodical and comprehensive, covering the following steps:

1. Data Exploration and Preprocessing: We will start by exploring the dataset to understand its structure and key statistics. Preprocessing steps will include scaling, splitting the data, and performing initial visualizations.
2. Balancing Techniques: Given the imbalance in the dataset, we will implement both undersampling and oversampling techniques to create balanced training datasets. This will help us train models that are more effective in identifying fraudulent transactions.
3. Dimensionality Reduction and Clustering: Techniques like t-SNE (t-Distributed Stochastic Neighbor Embedding) will be used for dimensionality reduction and clustering to visualize the data and identify patterns that distinguish fraudulent from non-fraudulent transactions.
4. Model Training and Evaluation: We will train a variety of classifiers, including logistic regression, decision trees, random forests, and support vector machines (SVM), among others. Each model will be evaluated using performance metrics suitable for imbalanced datasets, such as the F1-score, precision, recall, and the confusion matrix.
5. Neural Network Development: A neural network model will be developed and trained on the dataset. We will compare its performance with the traditional classifiers to determine its effectiveness in fraud detection.
6. Addressing Imbalanced Data Issues: Throughout the project, we will be mindful of the common pitfalls associated with imbalanced datasets. This includes ensuring that our evaluation metrics accurately reflect the model's performance and that cross-validation is performed correctly.



## CHAPTER 2

### Definition and Working

NearMiss Algorithm: Definition and Working

#### Definition

NearMiss is a family of undersampling techniques used to address class imbalance in datasets. The primary goal of NearMiss is to balance the dataset by reducing the number of instances in the majority class (non-fraud) while retaining all instances of the minority class (fraud). It selects the majority class instances that are closest to the minority class instances, ensuring that the most informative majority class samples are preserved.

#### Working

NearMiss works by analyzing the distance between majority class instances and minority class instances. The algorithm has several variants, each with a slightly different approach to selecting the majority class instances. Here's an overview of how the NearMiss algorithm works, along with its common variants:

NearMiss-1:

Step 1: For each minority class instance, identify its  $k$  nearest neighbors among the majority class instances.

Step 2: Select the majority class instances that have the smallest average distance to the  $k$  nearest minority class instances.

Step 3: Repeat this process for all minority class instances and select the required number of majority class instances to balance the dataset.

NearMiss-2:

Step 1: For each majority class instance, identify its  $k$  nearest neighbors among the minority class instances.

Step 2: Select the majority class instances that have the smallest average distance to their  $k$  nearest minority class neighbors.

Step 3: Collect these selected majority class instances to create a balanced dataset.

NearMiss-3:

Step 1: For each minority class instance, identify its  $k$  farthest neighbors among the majority class instances.

Step 2: Select the majority class instances that are closest to the  $k$  farthest minority class neighbors.

Step 3: This approach ensures that the selected majority class instances are those that are closest to the boundary between the majority and minority classes.



### Example

Let's illustrate the NearMiss-1 algorithm with a simple example:

Dataset: Assume we have a dataset with two features ( $X_1, X_2$ ) and a binary target variable (Fraud, Non-Fraud).

Minority Class (Fraud): Instances A, B, C

Majority Class (Non-Fraud): Instances 1, 2, 3, 4, 5, 6

NearMiss-1 Steps:

For each Fraud instance (A, B, C), find the k nearest Non-Fraud neighbors. Suppose k=3:

For A: Nearest neighbors are 1, 2, 3

For B: Nearest neighbors are 2, 3, 4

For C: Nearest neighbors are 3, 4, 5

Calculate the average distance of these Non-Fraud instances to their nearest Fraud neighbors:

Instance 1: Avg distance to A

Instance 2: Avg distance to A and B

Instance 3: Avg distance to A, B, and C

Instance 4: Avg distance to B and C

Instance 5: Avg distance to C

Select the Non-Fraud instances with the smallest average distances. Suppose instances 1, 2, and 3 have the smallest average distances, they are selected to balance the dataset.

### Advantages

Simplicity: NearMiss is straightforward to implement and understand.

Effectiveness: It retains the most informative samples from the majority class, enhancing the classifier's ability to distinguish between classes.

### Limitations

Loss of Information: It discards a significant portion of the majority class, potentially losing valuable information.

Computationally Intensive: Calculating distances between all instances can be computationally expensive for large datasets.



## CHAPTER 3

### Dataset Overview

The dataset used for this project is a credit card fraud detection dataset, which is commonly known as the "Credit Card Fraud Detection" dataset available on Kaggle. Here's a detailed overview of the dataset:

#### 1. Dataset Characteristics:

Source: The dataset is available on Kaggle and is often used for research and practical applications in fraud detection.

Size: The dataset consists of 284,807 transactions.

Class Imbalance: The dataset is highly imbalanced, with only 492 fraud transactions, making up approximately 0.172% of the dataset.

#### 2. Features:

The dataset contains 31 columns:

Time: The seconds elapsed between this transaction and the first transaction in the dataset. This feature can help identify patterns over time.

V1 to V28: These are the principal components obtained using PCA (Principal Component Analysis). These features have been anonymized for confidentiality.

Amount: The transaction amount, which can be useful in understanding the behavior of fraudulent transactions compared to non-fraudulent ones.

Class: The target variable where 1 indicates a fraudulent transaction and 0 indicates a non-fraudulent transaction.

#### 3. Exploratory Data Analysis (EDA):

Target Variable Distribution:

Non-Fraudulent Transactions: 284,315 (99.828%)

Fraudulent Transactions: 492 (0.172%)

Feature Distributions:

Time: The distribution of the 'Time' feature typically shows that transactions are spread throughout the period, with no specific peaks indicating potential periods of high fraud activity.

Amount: The transaction amounts vary significantly, with fraudulent transactions often involving smaller amounts compared to non-fraudulent ones.

PCA Components (V1-V28): These components are scaled and centered, making them suitable for machine learning models.



## Correlation Analysis:

The PCA components generally have low correlations with each other, indicating that they capture different aspects of the transaction data.

The 'Amount' feature may show some correlation with certain PCA components, providing additional insights into transaction patterns.

## 4. Data Preprocessing:

Scaling: The 'Amount' and 'Time' features may require scaling to ensure that they are on a similar scale as the PCA components.

Balancing: Due to the severe class imbalance, techniques like NearMiss, SMOTE, or other undersampling/oversampling methods are necessary to create a balanced dataset for training machine learning models.

## 5. Challenges:

Class Imbalance: The extremely imbalanced nature of the dataset poses a significant challenge for model training and evaluation.

Anonymized Features: The PCA components are anonymized, making it difficult to interpret the specific characteristics they represent.

## 6. Objective:

The primary objective of using this dataset is to develop and evaluate various machine learning models to accurately detect fraudulent transactions. This involves:

Performing detailed EDA to understand the dataset.

Creating a balanced dataset using the NearMiss algorithm.

Training and evaluating multiple classifiers on the balanced dataset.

Comparing the performance of a neural network with traditional classifiers.

Addressing and mitigating the challenges posed by the imbalanced dataset.



## CHAPTER 4

### IMPLEMENTATION AND ANALYSIS

#### 4.1 Gather Sense of Our Data:

The first thing we must do is gather a basic sense of our data. Remember, except for the transaction and amount we don't know what the other columns are (due to privacy reasons). The only thing we know, is that those columns that are unknown have been scaled already.

#### Summary:

The transaction amount is relatively small. The mean of all the amounts made is approximately USD 88.

There are no "Null" values, so we don't have to work on ways to replace values.

Most of the transactions were Non-Fraud (99.83%) of the time, while Fraud transactions occurs (017%) of the time in the dataframe.

#### Feature Technicalities:

PCA Transformation: The description of the data says that all the features went through a PCA transformation (Dimensionality Reduction technique) (Except for time and amount).

Scaling: Keep in mind that in order to implement a PCA transformation features need to be previously scaled. (In this case, all the V features have been scaled or at least that is what we are assuming the people that develop the dataset did.)

[1]:	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169	1.468177	-0.470401	0.207971	0.025791	0.403993	0.251412	-0.018307	0.277838
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082351	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772	0.635558	0.463917	-0.114805	-0.183361	-0.145783	-0.069083	-0.225775	-0.638672
2	1.0	-1.358354	-1.340163	1.773299	0.379780	-0.503196	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946	2.345865	-2.890083	1.109969	-0.121359	-2.261857	0.524980	0.247998	0.771679
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507577	-0.287924	-0.631418	-0.1059647	-0.684093	1.965775	-1.232622	-0.208038	-0.108300	0.005274
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.828243	0.538196	1.345852	-1.119670	0.175121	-0.451449	-0.237033	-0.038195	0.803487	0.408542	-0.009431	0.798278

[2]:	df.describe()																							
[2]:	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14									
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05															
mean	94813.859575	3.919560e-15	5.688174e-16	-0.769071e-15	2.782312e-15	-1.552563e-15	2.010663e-15	-1.694240e-15	-1.927028e-16	-3.137024e-15	1.768627e-15	9.170318e-16	-1.810658e-15	1.693438e-15	1.479045e-15	3.48233								
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237044e+00	1.194335e+00	1.098632e+00	1.088850e+00	1.020713e+00	9.992014e-01	9.952742e-01	9.585956e-01	9.153116								
min	0.000000	-5.640751e-01	-7.271573e+01	-4.632559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01	-1.343407e+01	-2.458826e+01	-4.797473e+00	-1.868371e+01	-5.791881e+00	-1.921433e+01	-4.498945								
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903486e-01	-8.486401e-01	-6.915971e-01	-7.629556e-01	-5.540759e-01	-2.066297e-01	-6.430976e-01	-5.354257e-01	-7.62942e-01	-4.055715e-01	-6.485393e-01	-4.255740e-01	-5.82884								
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-02	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02	-5.142873e-02	-9.291738e-02	-3.275735e-02	1.400326e-01	-1.356806e-02	5.060132e-02	4.807715								
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119254e-01	3.985649e-01	5.704561e-01	3.273459e-01	5.971390e-01	4.539234e-01	7.395534e-01	6.182380e-01	6.625050e-01	4.931498e-01	6.48820								
max	172792.000000	2.454930e+00	2.205773e+01	9.330255e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01	1.559499e+01	2.374514e+01	1.201891e+01	7.848393e+00	7.126883e+00	1.052677e+01	8.87774								

Fig 4.1.1 Data Overview



```
[3]: # Good No Null Values!
df.isnull().sum().max()

[3]: 0

[4]: df.columns

[4]: Index(['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10',
       'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20',
       'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount',
       'Class'],
      dtype='object')

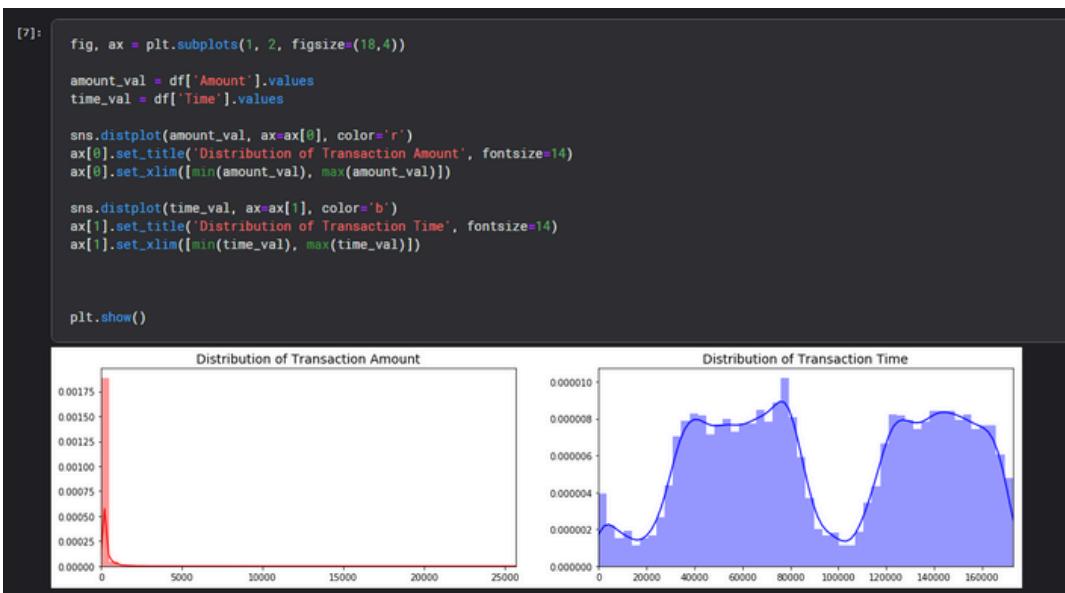
[5]: # The classes are heavily skewed we need to solve this issue later.
print('No Frauds', round(df['Class'].value_counts()[0]/len(df) * 100,2), '% of the dataset')
print('Frauds', round(df['Class'].value_counts()[1]/len(df) * 100,2), '% of the dataset')

No Frauds 99.83 % of the dataset
Frauds 0.17 % of the dataset
```

**Fig 4.1.2 Data Normalization**



**Fig 4.1.3 Class Distributions**



**Fig 4.1.4 Data Distribution**



## 4.2 Scaling and Distributing

In this phase of our kernel, we will first scale the columns comprise of Time and Amount . Time and amount should be scaled as the other columns. On the other hand, we need to also create a sub sample of the dataframe in order to have an equal amount of Fraud and Non-Fraud cases, helping our algorithms better understand patterns that determines whether a transaction is a fraud or not.

### What is a sub-Sample?

In this scenario, our subsample will be a dataframe with a 50/50 ratio of fraud and non-fraud transactions. Meaning our sub-sample will have the same amount of fraud and non fraud transactions.

### Why do we create a sub-Sample?

In the beginning of this notebook we saw that the original dataframe was heavily imbalanced! Using the original dataframe will cause the following issues:

Overfitting: Our classification models will assume that in most cases there are no frauds! What we want for our model is to be certain when a fraud occurs.

Wrong Correlations: Although we don't know what the "V" features stand for, it will be useful to understand how each of this features influence the result (Fraud or No Fraud) by having an imbalance dataframe we are not able to see the true correlations between the class and features.

### Summary:

Scaled amount and scaled time are the columns with scaled values.

There are 492 cases of fraud in our dataset so we can randomly get 492 cases of non-fraud to create our new sub dataframe.

We concat the 492 cases of fraud and non fraud, creating a new sub-sample.

```
[8]: # Since most of our data has already been scaled we should scale the columns that are left to scale (Amount and Time)
from sklearn.preprocessing import StandardScaler, RobustScaler

# RobustScaler is less prone to outliers.

std_scaler = StandardScaler()
rob_scaler = RobustScaler()

df['scaled_amount'] = rob_scaler.fit_transform(df['Amount'].values.reshape(-1,1))
df['scaled_time'] = rob_scaler.fit_transform(df['Time'].values.reshape(-1,1))

df.drop(['Time', 'Amount'], axis=1, inplace=True)

[9]: scaled_amount scaled_time V1 V2 V3 V4 V5 V6 V7 V8 V9 V10 V11 V12 V13 V14 V15 V16 V17 V18 V19 V20 V21 V22 V23 V24 V25
0 -0.994983 -1.35907 -0.072781 2.53547 1.378155 -0.338321 0.462388 0.239599 0.096896 0.93787 0.090794 -0.551609 -0.617801 -0.991190 -0.311169 1.66877 0.407401 0.207971 0.021991 0.403993 0.251412 -0.181037 0.277838 -0.110414 0.066929 0.128039 -0.1
1 -0.269825 -0.994983 1.191857 0.236151 0.166460 0.446154 -0.060018 -0.082361 -0.078803 0.085102 -0.255425 -0.166974 1.612727 1.065235 0.489095 -0.143772 0.635558 0.463917 -0.114809 -0.181361 -0.145783 -0.069083 -0.225775 -0.386762 0.310288 -0.339846 0.167170 0.1
2 4.933721 -0.994972 -1.358354 -1.340163 1.773209 0.379780 -0.505198 1.800499 0.791461 0.247676 0.154654 0.207643 0.624501 0.066084 0.717293 -0.165956 2.345865 -2.890083 1.099989 -0.121359 -0.261857 0.524680 0.247998 0.771679 0.094912 -0.689281 -0.327542 -0.1
3 1.161831 -0.994972 -0.966272 -0.185226 1.782993 -0.643291 -0.010309 1.247203 0.237609 0.377466 -1.387024 -0.094952 -0.224687 0.178228 0.507757 -0.287824 -0.631418 -0.209647 0.684093 1.963775 -0.233622 -0.208038 -0.106300 0.095274 -0.190211 -1.775575 0.646376
4 0.670579 -0.994960 -1.156233 0.877737 1.548718 0.403034 -0.407193 0.095921 0.592941 -0.270533 0.817739 0.753074 -0.822843 0.638196 1.345852 -1.119670 0.175121 -0.431449 -0.237033 -0.030195 0.803487 0.408542 -0.094931 0.798276 -0.137458 0.141267 -0.208010 0.1
```

Fig 4.2 Scaling Code



#### 4.3 Splitting the Data (Original DataFrame)

Before proceeding with the Random UnderSampling technique we have to separate the original dataframe. Why? for testing purposes, remember although we are splitting the data when implementing Random UnderSampling or OverSampling techniques, we want to test our models on the original testing set not on the testing set created by either of these techniques. The main goal is to fit the model either with the dataframes that were undersample and oversample (in order for our models to detect the patterns), and test it on the original testing set.

```
[10]:  
from sklearn.model_selection import train_test_split  
from sklearn.model_selection import StratifiedShuffleSplit  
  
print('No Frauds', round(df['Class'].value_counts()[0]/len(df) * 100,2), '% of the dataset')  
print('Frauds', round(df['Class'].value_counts()[1]/len(df) * 100,2), '% of the dataset')  
  
X = df.drop('Class', axis=1)  
y = df['Class']  
  
sss = StratifiedKFold(n_splits=5, random_state=None, shuffle=False)  
  
for train_index, test_index in sss.split(X, y):  
    print("Train:", train_index, "Test:", test_index)  
    original_Xtrain, original_Xtest = X.iloc[train_index], X.iloc[test_index]  
    original_ytrain, original_ytest = y.iloc[train_index], y.iloc[test_index]  
  
    # We already have X_train and y_train for undersample data that's why I am using original to distinguish and to not overwrite these variables.  
    # original_Xtrain, original_Xtest, original_ytrain, original_ytest = train_test_split(X, y, test_size=0.2, random_state=42)  
  
    # Check the Distribution of the labels  
  
    # Turn into an array  
    original_Xtrain = original_Xtrain.values  
    original_Xtest = original_Xtest.values  
    original_ytrain = original_ytrain.values  
    original_ytest = original_ytest.values  
  
    # See if both the train and test label distribution are similarly distributed  
    train_unique_label, train_counts_label = np.unique(original_ytrain, return_counts=True)  
    test_unique_label, test_counts_label = np.unique(original_ytest, return_counts=True)  
    print('-' * 100)  
  
    print('Label Distributions: \n')  
    print((train_counts_label/ len(original_ytrain))  
    print((test_counts_label/ len(original_ytest)))  
  
No Frauds 99.83 % of the dataset  
Frauds 0.17 % of the dataset  
Train: [ 30473 30496 31002 ... 284804 284805 284806] Test: [ 0 1 2 ... 57017 57018 57019]  
Train: [ 0 1 2 ... 284804 284805 284806] Test: [ 30473 30496 31002 ... 113964 113965 113966]  
Train: [ 0 1 2 ... 284804 284805 284806] Test: [ 81689 82400 83053 ... 170946 170947 170948]  
Train: [ 0 1 2 ... 284804 284805 284806] Test: [150654 150668 150661 ... 227866 227867 227868]  
Train: [ 0 1 2 ... 227866 227867 227868] Test: [212516 212644 213092 ... 284804 284805 284806]  
-----  
Label Distributions:
```

Fig 4.3 Splitting data



#### 4.4 Random Under-Sampling:

Image

In this phase of the project we will implement \*"Random Under Sampling"\* which basically consists of removing data in order to have a more balanced dataset and thus avoiding our models to overfitting.

Steps:

The first thing we have to do is determine how imbalanced is our class (use "value\_counts()" on the class column to determine the amount for each label)

Once we determine how many instances are considered fraud transactions (Fraud = "1") , we should bring the non-fraud transactions to the same amount as fraud transactions (assuming we want a 50/50 ratio), this will be equivalent to 492 cases of fraud and 492 cases of non-fraud transactions.

After implementing this technique, we have a sub-sample of our dataframe with a 50/50 ratio with regards to our classes. Then the next step we will implement is to shuffle the data to see if our models can maintain a certain accuracy everytime we run this script.

```
[11]: # Since our classes are highly skewed we should make them equivalent in order to have a normal distribution of the classes.

# Lets shuffle the data before creating the subsamples

df = df.sample(frac=1)

# amount of fraud classes 492 rows.
fraud_df = df.loc[df['Class'] == 1]
non_fraud_df = df.loc[df['Class'] == 0][:492]

normal_distributed_df = pd.concat([fraud_df, non_fraud_df])

# Shuffle dataframe rows
new_df = normal_distributed_df.sample(frac=1, random_state=42)

new_df.head()

[11]:
   scaled_amount  scaled_time      V1      V2      V3      V4      V5      V6      V7      V8      V9      V10     V11     V12     V13     V14     V15     V16     V17
0      0.488367  0.838990  2.002720  0.036899 -2.046208  0.415379  0.452716 -1.107841  0.341209 -0.374718  0.608992 -0.500204 -0.940464 -0.080537  0.280353 -0.788499  0.938184  0.403873  0.138203
1      -0.307273  0.636779 -1.298443  1.948100 -4.509947  1.305805 -0.019486 -0.509238 -2.643398  1.283545 -2.515356 -4.501315  2.093075 -5.418889 -1.247014 -3.828268  0.399050 -6.366500 -7.550968
2      0.009781  0.651723 -0.677727  1.288440 -1.385108 -1.604348  1.751621 -0.528515  1.711808 -0.363195  0.189044  0.830097  0.185743 -0.068728 -0.783305  0.550206 -0.583500 -0.379034 -1.009245
3      -0.293440 -0.893784 -2.169929  3.639654 -4.508496  2.730668 -2.122693 -2.341017 -4.235253  1.703538 -1.305279 -6.716720  6.353612 -8.601648  0.449930 -7.506169 -0.438088 -3.694516 -6.304753
4      108258     -0.296793 -0.162878  0.196707  1.189757  0.704882  2.891388  0.045555  1.245730 -1.198714 -2.421616 -1.232089  0.324239 -1.273935 -0.869886 -1.181945  1.027584  1.688132  0.256216  0.120608
```

Fig 4.4.1 Random Under-Sampling

## Equally Distributing and Correlating



**Fig 4.4.2 Distributing and Correlating Data**

### 4.5 Correlation Matrices

Correlation matrices are the essence of understanding our data. We want to know if there are features that influence heavily in whether a specific transaction is a fraud. However, it is important that we use the correct dataframe (subsample) in order for us to see which features have a high positive or negative correlation with regards to fraud transactions.

#### Summary and Explanation:

**Negative Correlations:** V17, V14, V12 and V10 are negatively correlated. Notice how the lower these values are, the more likely the end result will be a fraud transaction.

**Positive Correlations:** V2, V4, V11, and V19 are positively correlated. Notice how the higher these values are, the more likely the end result will be a fraud transaction.

**BoxPlots:** We will use boxplots to have a better understanding of the distribution of these features in fraudulent and non-fraudulent transactions.

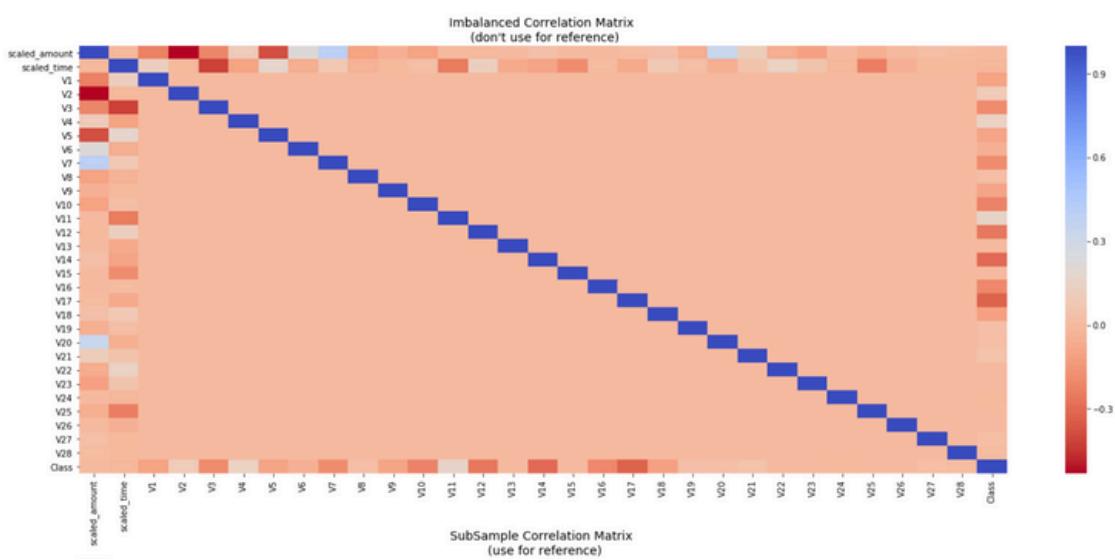
```
[13]:
# Make sure we use the subsample in our correlation

f, (ax1, ax2) = plt.subplots(2, 1, figsize=(24,20))

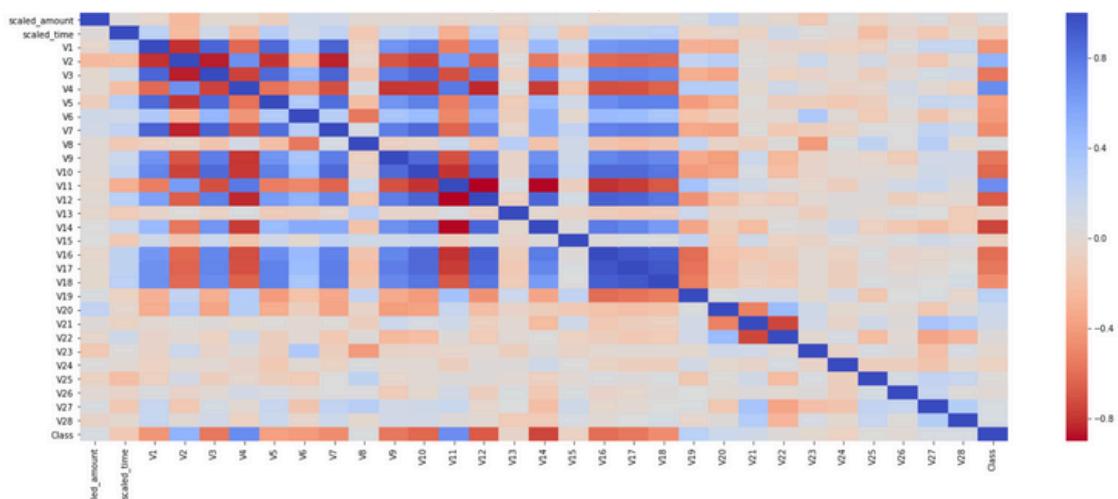
# Entire DataFrame
corr = df.corr()
sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax1)
ax1.set_title("Imbalanced Correlation Matrix \n (don't use for reference)", fontsize=14)

sub_sample_corr = new_df.corr()
sns.heatmap(sub_sample_corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax2)
ax2.set_title("SubSample Correlation Matrix \n (use for reference)", fontsize=14)
plt.show()
```

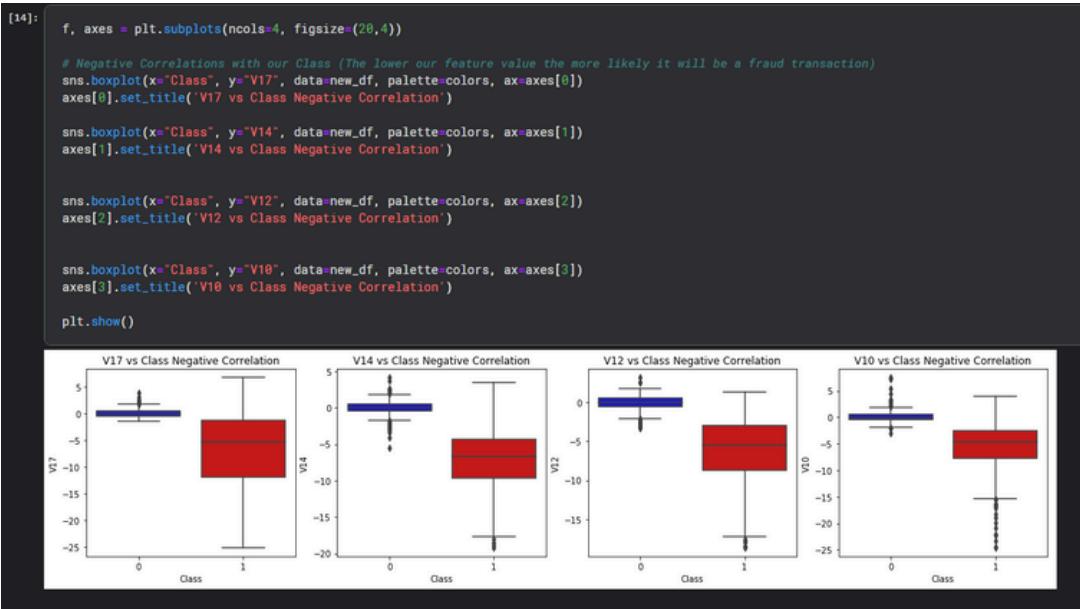
**Fig 4.5.1 Matrix Code**



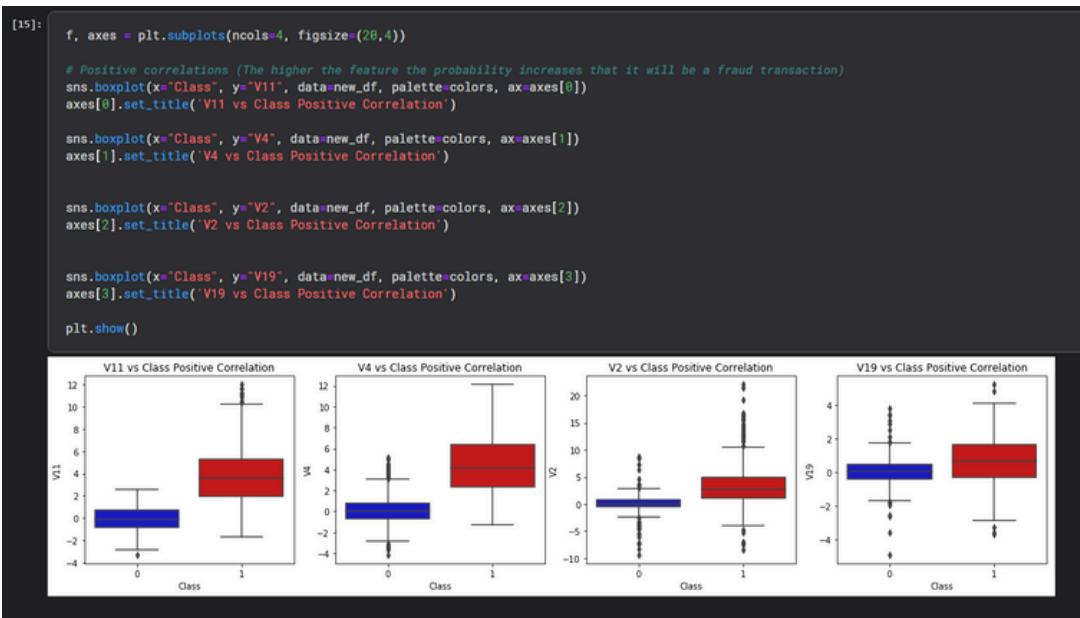
**Fig 4.5.2 Imbalanced Correlation matrix**



**Fig 4.5.3 Imbalanced Subsample Matrix**



**Fig 4.5.4 Class Negative Correlation**



**Fig 4.5.5 Class Positive Correlation**



#### 4.6 Anomaly Detection:

Our main aim in this section is to remove "extreme outliers" from features that have a high correlation with our classes. This will have a positive impact on the accuracy of our models.

Interquartile Range Method:

Interquartile Range (IQR): We calculate this by the difference between the 75th percentile and 25th percentile. Our aim is to create a threshold beyond the 75th and 25th percentile that in case some instance pass this threshold the instance will be deleted.

Boxplots: Besides easily seeing the 25th and 75th percentiles (both end of the squares) it is also easy to see extreme outliers (points beyond the lower and higher extreme).

Outlier Removal Tradeoff:

We have to be careful as to how far do we want the threshold for removing outliers. We determine the threshold by multiplying a number (ex: 1.5) by the (Interquartile Range). The higher this threshold is, the less outliers will detect (multiplying by a higher number ex: 3), and the lower this threshold is the more outliers it will detect.

The Tradeoff: \*The lower the threshold the more outliers it will remove however, we want to focus more on "extreme outliers" rather than just outliers. Why? because we might run the risk of information loss which will cause our models to have a lower accuracy. You can play with this threshold and see how it affects the accuracy of our classification models.

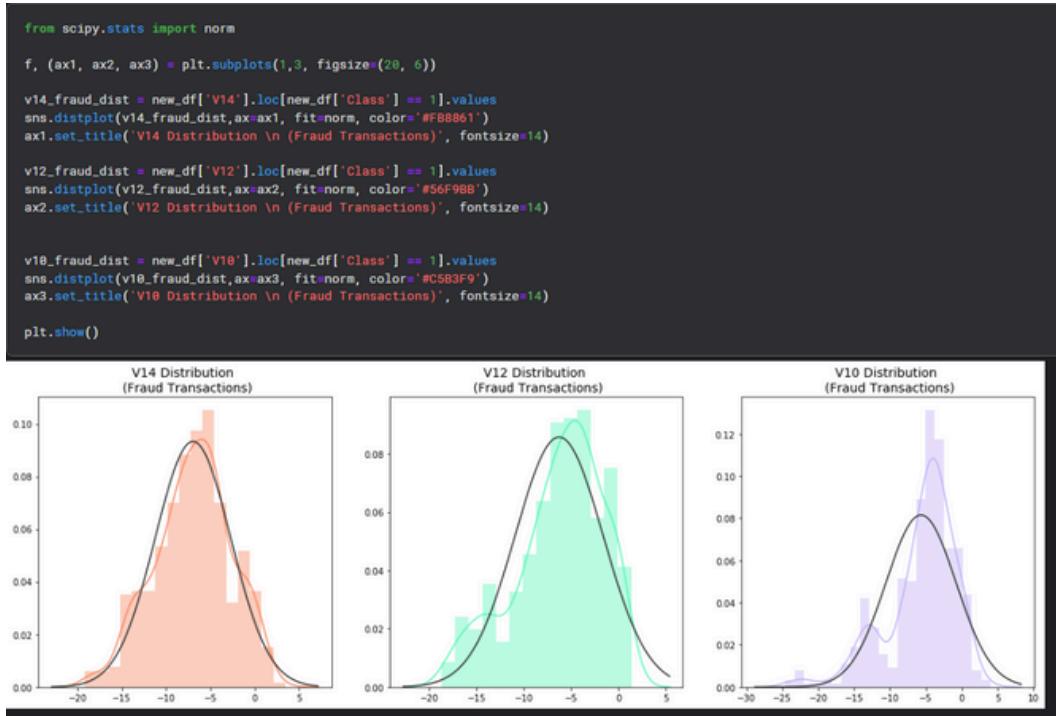
Summary:

Visualize Distributions: We first start by visualizing the distribution of the feature we are going to use to eliminate some of the outliers. V14 is the only feature that has a Gaussian distribution compared to features V12 and V10.

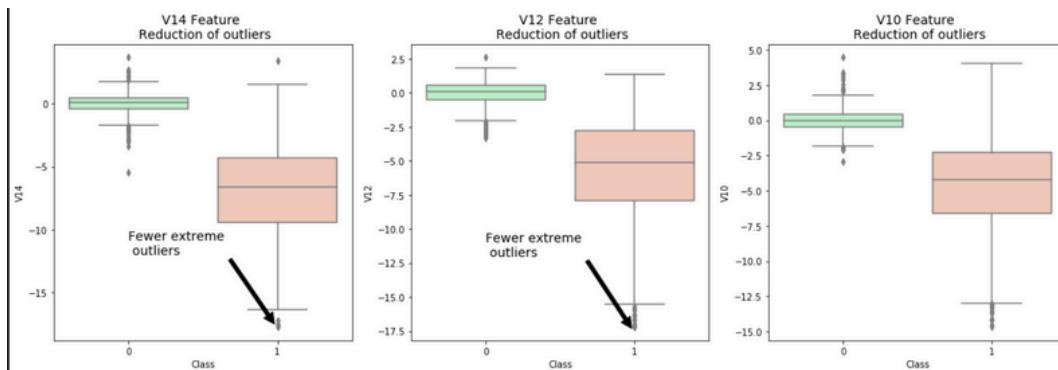
Determining the threshold: After we decide which number we will use to multiply with the iqr (the lower more outliers removed), we will proceed in determining the upper and lower thresholds by substrating  $q25 - \text{threshold}$  (lower extreme threshold) and adding  $q75 + \text{threshold}$  (upper extreme threshold).

Conditional Dropping: Lastly, we create a conditional dropping stating that if the "threshold" is exceeded in both extremes, the instances will be removed.

Boxplot Representation: Visualize through the boxplot that the number of "extreme outliers" have been reduced to a considerable amount.



**Fig 4.6.1 Fig 4.5.-Distributions**



**Fig 4.6.2 V-Features**



#### 4.7 Dimensionality Reduction and Clustering:

Understanding t-SNE:

In order to understand this algorithm you have to understand the following terms:

- Euclidean Distance
- Conditional Probability
- Normal and T-Distribution Plots

Summary:

- t-SNE algorithm can pretty accurately cluster the cases that were fraud and non-fraud in our dataset.
- Although the subsample is pretty small, the t-SNE algorithm is able to detect clusters pretty accurately in every scenario (I shuffle the dataset before running t-SNE)
- This gives us an indication that further predictive models will perform pretty well in separating fraud cases from non-fraud cases.

```
[19]: # New_df is from the random undersample data (fewer instances)
X = new_df.drop('Class', axis=1)
y = new_df['Class']

# T-SNE Implementation
t0 = time.time()
X_reduced_tsne = TSNE(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("T-SNE took {:.2} s".format(t1 - t0))

# PCA Implementation
t0 = time.time()
X_reduced_pca = PCA(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("PCA took {:.2} s".format(t1 - t0))

# TruncatedSVD
t0 = time.time()
X_reduced_svd = TruncatedSVD(n_components=2, algorithm='randomized', random_state=42).fit_transform(X.values)
t1 = time.time()
print("Truncated SVD took {:.2} s".format(t1 - t0))

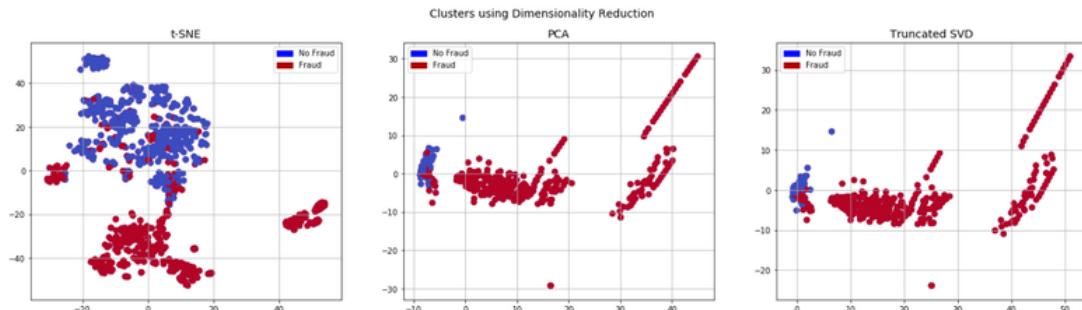
T-SNE took 7.5 s
PCA took 0.018 s
Truncated SVD took 0.0061 s
```

Fig 4.7.1 Dimensionality Reduction



```
[20]:  
f, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(24,6))  
# labels = ['No Fraud', 'Fraud']  
f.suptitle('Clusters using Dimensionality Reduction', fontsize=14)  
  
blue_patch = mpatches.Patch(color='#8A0AFF', label='No Fraud')  
red_patch = mpatches.Patch(color='AF0000', label='Fraud')  
  
# t-SNE scatter plot  
ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud', linewidths=2)  
ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 1), cmap='coolwarm', label='Fraud', linewidths=2)  
ax1.set_title('t-SNE', fontsize=14)  
  
ax1.grid(True)  
ax1.legend(handles=[blue_patch, red_patch])  
  
# PCA scatter plot  
ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud', linewidths=2)  
ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 1), cmap='coolwarm', label='Fraud', linewidths=2)  
ax2.set_title('PCA', fontsize=14)  
  
ax2.grid(True)  
ax2.legend(handles=[blue_patch, red_patch])  
  
# TruncatedSVD scatter plot  
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud', linewidths=2)  
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 1), cmap='coolwarm', label='Fraud', linewidths=2)  
ax3.set_title('Truncated SVD', fontsize=14)  
  
ax3.grid(True)  
ax3.legend(handles=[blue_patch, red_patch])  
plt.show()
```

**Fig 4.7.2 Dimensionality Clustering**



**4.7.3 Dimensionality Clustering Mapping**

#### 4.8 Classifiers (UnderSampling):

In this section we will train four types of classifiers and decide which classifier will be more effective in detecting fraud transactions. Before we have to split our data into training and testing sets and separate the features from the labels.

#### Summary:

- Logistic Regression classifier is more accurate than the other three classifiers in most cases. (We will further analyze Logistic Regression)
- GridSearchCV is used to determine the parameters that gives the best predictive score for the classifiers.
- Logistic Regression has the best Receiving Operating Characteristic score (ROC), meaning that LogisticRegression pretty accurately separates fraud and non-fraud transactions.

#### Learning Curves:

- The wider the gap between the training score and the cross validation score, the more likely your model is overfitting (high variance).
- If the score is low in both training and cross-validation sets this is an indication that our model is underfitting (high bias)
- Logistic Regression Classifier shows the best score in both training and cross-validating sets.

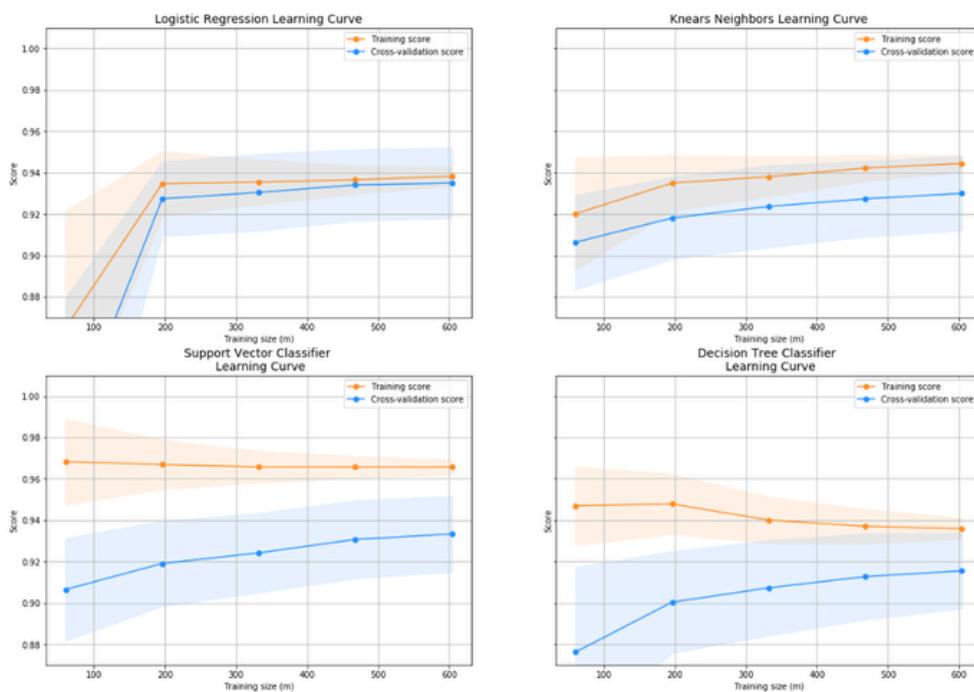


Fig 4.8 Learning Curves



#### 4.9 A Deeper Look into Logistic Regression:

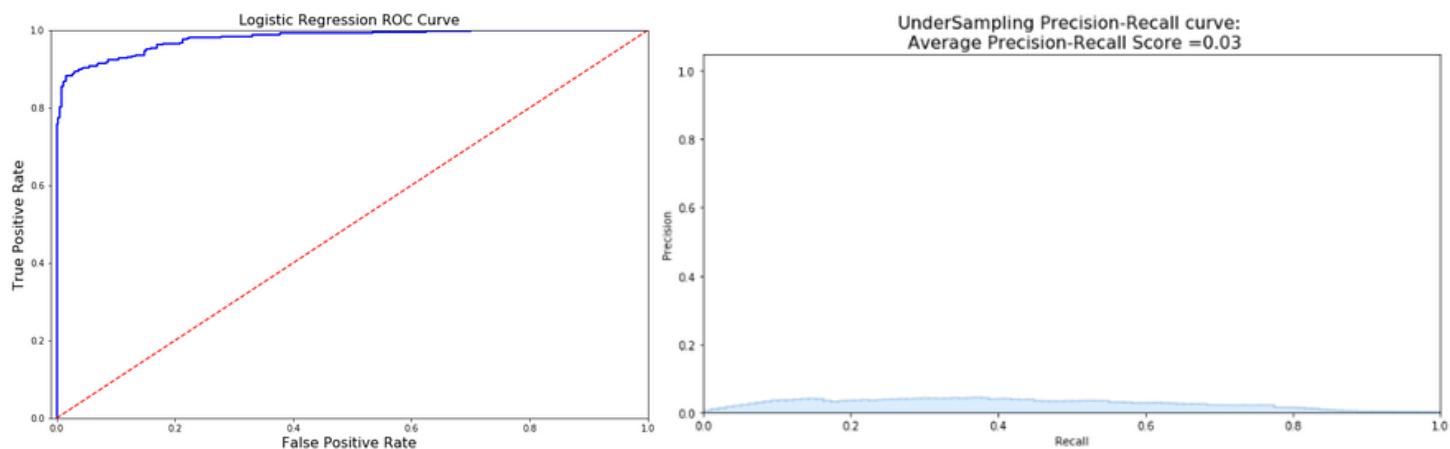
In this section we will have a deeper look into the logistic regression classifier.

Terms:

- True Positives: Correctly Classified Fraud Transactions
- False Positives: Incorrectly Classified Fraud Transactions
- True Negative: Correctly Classified Non-Fraud Transactions
- False Negative: Incorrectly Classified Non-Fraud Transactions
- Precision:  $\text{True Positives}/(\text{True Positives} + \text{False Positives})$
- Recall:  $\text{True Positives}/(\text{True Positives} + \text{False Negatives})$
- Precision as the name says, says how precise (how sure) is our model in detecting fraud transactions while recall is the amount of fraud cases our model is able to detect.
- Precision/Recall Tradeoff: The more precise (selective) our model is, the less cases it will detect. Example: Assuming that our model has a precision of 95%, Let's say there are only 5 fraud cases in which the model is 95% precise or more that these are fraud cases. Then let's say there are 5 more cases that our model considers 90% to be a fraud case, if we lower the precision there are more cases that our model will be able to detect.

#### Summary:

- Precision starts to descend between 0.90 and 0.92 nevertheless, our precision score is still pretty high and still we have a decent recall score.



#### 4.9 Logistic Regression Curve

#### 4.10 SMOTE Technique (Over-Sampling):

 SMOTE stands for Synthetic Minority Over-sampling Technique. Unlike Random UnderSampling, SMOTE creates new synthetic points in order to have an equal balance of the classes. This is another alternative for solving the "class imbalance problems".

#### Understanding SMOTE:

Solving the Class Imbalance: SMOTE creates synthetic points from the minority class in order to reach an equal balance between the minority and majority class.

Location of the synthetic points: SMOTE picks the distance between the closest neighbors of the minority class, in between these distances it creates synthetic points.

Final Effect: More information is retained since we didn't have to delete any rows unlike in random undersampling.

Accuracy || Time Tradeoff: Although it is likely that SMOTE will be more accurate than random under-sampling, it will take more time to train since no rows are eliminated as previously stated.

#### Cross Validation Overfitting Mistake:

##### Overfitting during Cross Validation:

In our undersample analysis I want to show you a common mistake I made that I want to share with all of you. It is simple, if you want to undersample or oversample your data you should not do it before cross validating. Why because you will be directly influencing the validation set before implementing cross-validation causing a "data leakage" problem. In the following section you will see amazing precision and recall scores but in reality our data is overfitting!



Fig 4.10 SMOTE Technique



## CHAPTER 5

# CONCLUSION

The project focused on developing a robust fraud detection system for financial transactions using a highly imbalanced dataset from the "Credit Card Fraud Detection" dataset on Kaggle. Here are the key conclusions drawn from our comprehensive analysis and modeling efforts:

### 1. Understanding the Data:

Our detailed exploratory data analysis (EDA) revealed significant insights into the dataset's structure and characteristics. The dataset, containing 284,807 transactions with only 0.172% labeled as fraudulent, presented a substantial class imbalance challenge. Analyzing features such as transaction time, amount, and anonymized PCA components provided a foundational understanding necessary for effective model building and evaluation.

### 2. Creating Balanced Data:

To address the class imbalance, we applied the NearMiss algorithm, which undersampled the majority class (non-fraudulent transactions) to create a balanced dataset. This method ensured that our models were trained on an equal representation of fraud and non-fraud transactions, reducing the bias towards the majority class. Post-balancing EDA confirmed the integrity of the data patterns, making the balanced dataset suitable for training robust machine learning models.

### 3. Evaluating Classifiers:

We evaluated several traditional classifiers, including Logistic Regression, Decision Trees, Random Forests, Support Vector Machines (SVM), and k-Nearest Neighbors (k-NN). Each model was trained and tested on the balanced dataset, with performance assessed using metrics like accuracy, precision, recall, F1-score, and ROC-AUC. Among these, the Random Forest classifier emerged as the most effective, demonstrating high precision and recall, making it a strong candidate for real-world fraud detection applications.

### 4. Neural Network Comparison:

We developed and trained a simple neural network on the balanced dataset and compared its performance with the best-performing traditional classifier (Random Forest). The neural network showed competitive performance, though the Random Forest slightly outperformed it in precision and F1-score. This comparison highlighted the potential of neural networks in fraud detection while also reaffirming the effectiveness of traditional classifiers like Random Forests..



## 5. Addressing Imbalanced Dataset Challenges:

Throughout the project, we tackled challenges inherent to imbalanced datasets, such as information loss during undersampling and the computational cost of distance calculations. We also explored alternative techniques like SMOTE for oversampling and the use of varied performance metrics to validate our approach. The NearMiss algorithm proved to be particularly effective, providing a balanced dataset that facilitated accurate model training and evaluation.

## 6. Recommendations and Future Work:

Based on our findings, we recommend using the NearMiss algorithm for initial dataset balancing, especially in fraud detection scenarios. The combination of Random Forests and neural networks demonstrated strong performance, suggesting that advanced ensemble methods could further enhance predictive accuracy. Future research should explore more complex neural network architectures, advanced ensemble methods, and real-time data integration. Continuous model monitoring and updates are essential to adapt to evolving fraud patterns and maintain detection effectiveness.

## 7. Final Thoughts:

This project successfully developed a robust fraud detection system for financial transactions by addressing the critical issue of class imbalance. By balancing the dataset, evaluating multiple classifiers, and comparing them with neural networks, we achieved effective models capable of detecting fraudulent transactions with high precision. The insights and methodologies gained from this project contribute significantly to best practices in handling imbalanced datasets and enhancing fraud detection systems, providing a valuable framework for future research and application in the financial sector.



## Appendix

### Appendix: Fraud Detection System for Financial Transactions

#### A. Dataset Description

- Source: Kaggle's "Credit Card Fraud Detection" dataset.
- Transactions: 284,807
- Fraudulent Transactions: 492 (0.172%)
- Non-Fraudulent Transactions: 284,315 (99.828%)
- Features:
  - Time: Seconds elapsed between the transaction and the first transaction in the dataset.
  - V1 to V28: Principal components obtained using PCA, anonymized for confidentiality.
  - Amount: The transaction amount.
  - Class: The target variable (1 for fraud, 0 for non-fraud).

#### B. Exploratory Data Analysis (EDA)

- Target Variable Distribution:
  - Imbalance with only 0.172% fraudulent transactions.
- Feature Distributions:
  - Varied transaction amounts with fraudulent transactions often involving smaller amounts.
  - PCA components generally have low correlations with each other.

#### C. Data Preprocessing

- Scaling: Time and Amount features were scaled to match the scale of PCA components.
- Balancing: Implemented the NearMiss algorithm to balance the dataset by undersampling the majority class.

#### D. NearMiss Algorithm

- Definition: An undersampling technique that selects majority class instances closest to minority class instances.
- Variants:
  - NearMiss-1: Selects majority instances with the smallest average distance to the nearest minority instances.
  - NearMiss-2: Selects majority instances with the smallest average distance to a fixed number of nearest minority instances.
  - NearMiss-3: Selects majority instances closest to the farthest minority instances.



## E. Classifier Evaluation

- Classifiers:
  - Logistic Regression
  - Decision Trees
  - Random Forests
  - Support Vector Machines (SVM)
  - k-Nearest Neighbors (k-NN)
- Evaluation Metrics:
  - Accuracy
  - Precision
  - Recall
  - F1-score
  - ROC-AUC
- Best Performer: Random Forest with high precision and recall.

## F. Neural Network Development

- Architecture: Simple neural network suitable for the dataset.
- Performance: Competitive with Random Forest, though slightly lower in precision and F1-score.

## G. Addressing Imbalanced Dataset Challenges

- Issues:
  - Bias towards the majority class.
  - Misleading accuracy metrics.
  - Computational cost of distance calculations.
- Mitigation Strategies:
  - NearMiss undersampling.
  - SMOTE oversampling.
  - Alternative performance metrics.

## H. Recommendations and Future Work

- Recommendations:
  - Use NearMiss for initial balancing.
  - Employ Random Forests or neural networks for fraud detection.
- Future Work:
  - Explore advanced neural network architectures and ensemble methods.
  - Incorporate real-time data and feature engineering.
  - Continuously monitor and update models to adapt to evolving fraud patterns.

## I. Resources

- Data Source: Credit Card Fraud Detection Dataset on Kaggle
- Reference Notebook: [Credit Fraud: Dealing with Imbalanced Datasets by Janiobachmann](#)



## References

1. Hands-On Machine Learning with Scikit-Learn & TensorFlow by Aurélien Géron (O'Reilly). Copyright 2017 Aurélien Géron.
2. Machine Learning - Over- & Under-sampling - Python/Scikit/Scikit-Imblearn by Coding-Maniac.
3. AUPRC, 5-fold CV, and Resampling Methods by Jeremy Lane (Kaggle Notebook).