# Credit Card Fraud Detection-Innovation

## Abstract

The purpose of this project is to detect the fraudulent transactions made by credit cards by the use of machine learning techniques, to stop fraudsters from the unauthorized usage of customers' accounts. The increase of credit card fraud is growing rapidly worldwide, which is the reason actions should be taken to stop fraudsters. Putting a limit for those actions would have a positive impact on the customers as their money would be recovered and retrieved back into their accounts and they won't be charged for items or services that were not purchased by them which is the main goal of the project. Detection of the fraudulent transactions will be made by using three machine learning techniques KNN, SVM and Logistic Regression, those models will be used on a credit card transaction dataset.

## Introduction

We are using credit card daily for our expenses. In a physical-card based purchase, the card holder presents his card physically to a merchant for making a payment. If the card holder does not realize the loss of card, it can lead to a substantial financial loss to the credit card company. Credit Card Misrepresentation is one of the greatest dangers to business and business foundations today. Just, Master card Misrepresentation is characterized as, "when an individual uses another individuals". Fraud detection involves monitoring the activities of populations of users in order to estimate, perceive or avoid objectionable behaviour, which consist of fraud, intrusion, and defaulting. This problem is particularly challenging from the perspective of learning, as it is characterized by various factors such as class imbalance. The number of valid transactions far outnumber fraudulent ones. Also, the transaction patterns often change their statistical properties over the course of time. In order to minimize costs of detection it is important to use expert rules and statistical based models to make a first screen between genuine and potential fraud and ask the investigators to review only the cases with high risk. When a fraud cannot be prevented, it is desirable to identify. The number of domain constraints and characteristics exaggerate the problem of detection and prevention. Customer irritation is to be avoided. Most banks considers huge

transactions, among which very few is fraudulent, often less. Also, only a limited number of transactions can be checked by fraud investigators, i.e. we cannot ask a human person to check all transactions one by one if it is fraudulent or not.

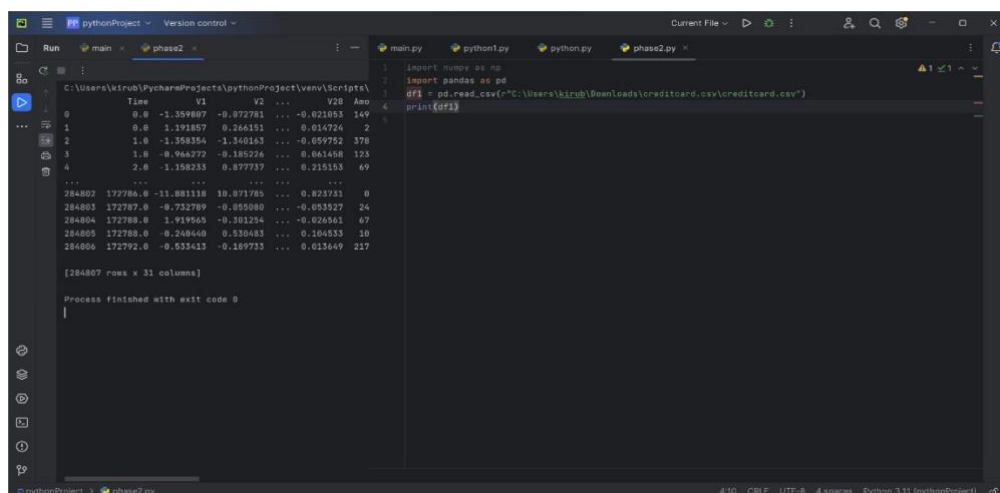# Implementation

## ➢ Python Libraries

Python libraries are collections of modules that contain useful codes and functions, eliminating the need to write them from scratch. There are tens of thousands of Python libraries that help machine learning developers, as well as professionals working in data science, data visualization, and more.

- ❖ Pandas
- ❖ Numpy
- ❖ Mtplotlib
- ❖ Sci-Kit Learn

## ➢ Data Reading and Activities

The dataset was retrieved from an open-source website, Kaggle.com. it contains data of transactions that were made in 2013 by credit card users in Europe, in two days only. The dataset consists of 31 attributes, 284,808 rows. 28 attributes are numeric variables that due to confidentiality and privacy of the customers have been transformed using PCA transformation, the three remaining attributes are "Time" which contains the elapsed seconds between the first and other transactions of each attribute, "Amount" is the amount of each transaction, and the final attribute "Class" which contains binary variables where "1" is a case of fraudulent transaction, and "0" is not as case of fraudulent transaction.

Dataset Link: https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud

> ## Activity Steps:

1) Data Collection

2) Data Preparation.

3) Choose a Model.

4) Train the Model.

5) Evaluate the Model.

6) Parameter Tuning.

7) Make Predictions.



> ## Data Cleaning:

The preparation of the dataset includes selecting the wanted attributes or variables, cleaning it by excluding Null rows, deleting duplicated variables, treating outlier if necessary, in addition to transforming data types to the wanted type, data merging can be performed as well where two or more attributes get merged. All those alterations lead to the wanted result which is to make the data ready to be modeled.

## EDA(Exploratory Data Analysis)

Exploratory Data Analysis (EDA) refers to the method of studying and exploring record sets to apprehend their predominant traits, discover patterns, locate outliers, and identify relationships between variables.

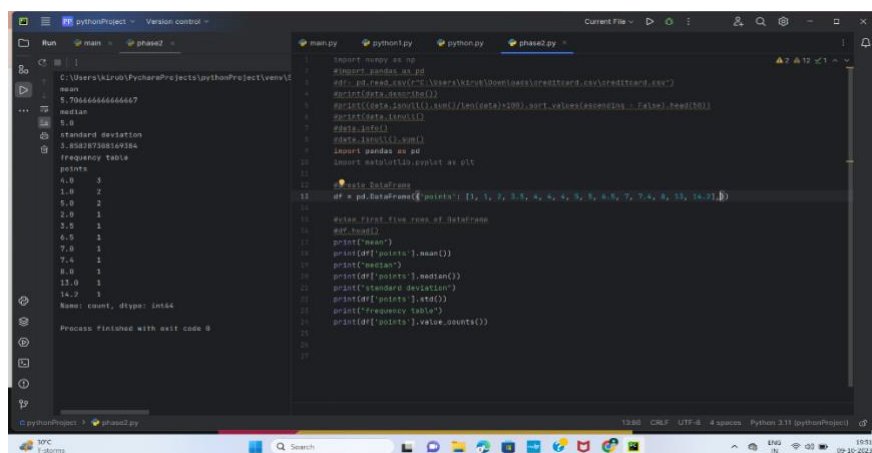Steps of EDA

1)variable Identification

2)Univariant Analysis

3)Bivariant Analysis

4)Missing value

5)Outlier Handling

## ➤ Univriant Analysis

Univriant Analysis refers to the analysis of one variable.



## ➤ Bivariate analysis

Bivariate analysis refers to the analysis of two variables to perform bivariate analysis Scatterplots:

➤ Outliers

Outliers are those data points that are significantly different from the rest of the dataset.



➤ Data testing and training

train and test datasets are the two key concepts of machine learning, where the training dataset is used to fit the model, and the test dataset is used to evaluate the model.

❖ x_train: It is used to represent features for the training data
❖ x_test: It is used to represent features for testing data
❖ y_train: It is used to represent dependent variables for training data
❖ y_test: It is used to represent independent variable for testing data

# METHODOLOGY

We are using supervised learning techniques to extract the information as a part of the fraud analysis. The dataset used is a binary classification. Fraud detection is a binary classification task in which any transaction will be predicted and labelled as a fraud or legit.classification techniques were tried for this task and their performances were compared. The following subsections briefly explain these classification techniques, data set and metrics used for performance measure.

## ➢ Support vector machines (SVM)

SVM is introduced in 1992 to solve binary classification problems and then they are extended to nonlinear regression problems. SVMs are based on structural risk minimization unlike ANNs which is based on empirical risk minimization. SVM map the data to a predetermined very high- dimensional space via a kernel function and finds the hyper plane that maximizes the margin between NOVATEUR PUBLICATIONS INTERNATIONAL JOURNAL OF INNOVATIONS IN ENGINEERING RESEARCH AND TECHNOLOGY. The solution is based only on those data points, which are at the margin. These points are called support vectors.

```
Correctly Classified Instances         85391               99.9403 %
Incorrectly Classified Instances        51                 0.0597 %
Kappa statistic                          0.8388
Mean absolute error                      0.0006
Root mean squared error                  0.0244
Relative absolute error                 16.4433 %
Root relative squared error             54.1917 %
Total Number of Instances               85442

=== Detailed Accuracy By Class ===

               TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
               0.764    0.000    0.930      0.764   0.839      0.843  0.882     0.711     1
               1.000    0.236    1.000      1.000   1.000      0.843  0.882     1.000     0
Weighted Avg.  0.999    0.235    0.999      0.999   0.999      0.843  0.882     0.999

=== Confusion Matrix ===

    a      b    <-- classified as
   133    41 |     a = 1
    10 85258 |     b = 0
```

## ➢ Naive Bayes Algorithm

Naive Bayes is based on two assumptions. Firstly, all features in an entry that needs to be classified are contributing evenly in the decision Secondly, all attributes are statistically independent, meaning that, knowing an attribute's value does not indicate anything about other attributes' values which is not always true in practice. The process of classifying an instance is done by applying the Bayes rule for each class given the instance. In the fraud detection task, the following formula is calculated for each of the two classes and the class associated with the higher probability is the predicted class for the instance.

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Steps:
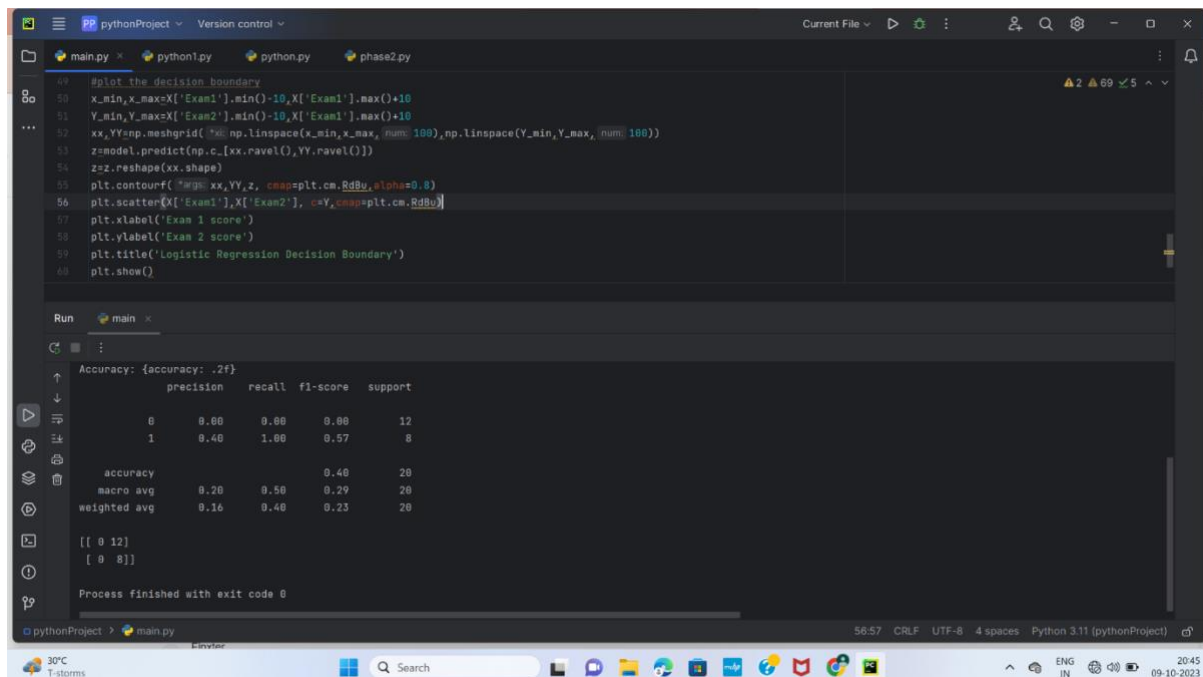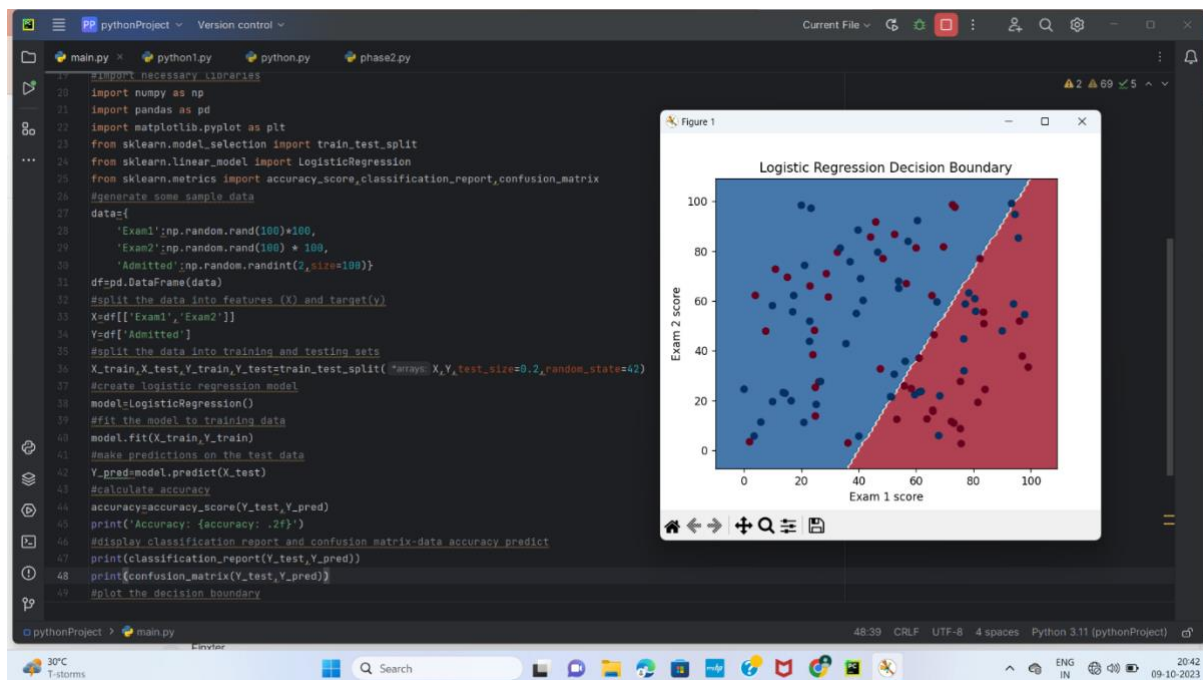
1: Calculate the prior probability for given class labels

2: Find Likelihood probability with each attribute for each class

3: Put these value in Bayes Formula and calculate posterior probability.

4: See which class has a higher probability, given the input belongs to the higher probability class.

```
Correctly Classified Instances        83504              97.7318 %
Incorrectly Classified Instances       1938               2.2682 %
Kappa statistic                        0.1292
Mean absolute error                    0.0227
Root mean squared error                0.1491
Relative absolute error              626.539  %
Root relative squared error          330.6127 %
Total Number of Instances             85442

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                0.851    0.022    0.072      0.851   0.132      0.243  0.968     0.091     1
                0.978    0.149    1.000      0.978   0.989      0.243  0.964     1.000     0
Weighted Avg.   0.977    0.149    0.998      0.977   0.987      0.243  0.964     0.998

=== Confusion Matrix ===

    a     b    <-- classified as
  148    26 |    a = 1
 1912 83356 |    b = 0
```

Confusion Matrix and Statistics

```
                         Reference
Prediction        Not Fraudulent  Fraudulent
    Not Fraudulent           89684          33
    Fraudulent                2018         139

            Accuracy : 0.9777
```

## ➢ Logistic Regression

Logistic regression also does not require independent variables to be linearly related, nor does it require equal variance within each group, which also makes it a less stringent procedure for statistical analysis. As a result, logistic regression was used to predict the probability of fraudulent credit cards. Assumptions and Limitations of Logistic Regression. Logistic regression analysis uses maximum likelihood estimation to predict group membership. However, to interpret the results of the prediction of group membership with precision and accuracy, a preliminary analysis of the cleaned dataset was conducted to observe if the assumptions of logistic regression were met. The above techniques are used for the detection of fraud transactions involved with the banks. After the classification, the intensity of the individual cards is inquired and calculated, which prepares and classifies the credit card trend and spending based anomalies. The regression models are used to perform the operation on the given data stream obtained from the credit card company for the detection of the credit card frauds by analysing the spending behaviour of the customers.

> ➤ Linear Regression:

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables. Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (y) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

$$y = a0 + a1x + \varepsilon$$

Where, Y= Dependent Variable (Target Variable)

X= Independent Variable (predictor Variable)

a0= intercept of the line (Gives an additional degree of freedom)

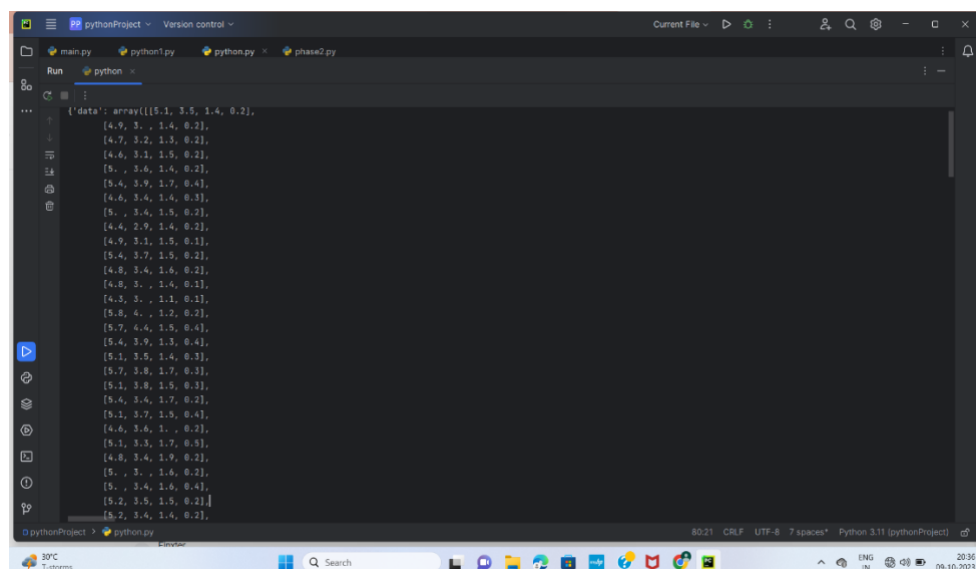a1 = Linear regression coefficient (scale factor to each input value).

ε = random error

➤ Accuracy and Other Metrics for Prediction evaluation

The last stage of the CRISP-DM model is the evaluation and deployment stage, as presented in table 2 below all models are being compared to each other to figure the best model in identifying fraudulent credit card transactions. Accuracy is the overall number of instances that are predicted correctly, accuracies are represented by confusion matrix where it showed the True Positive (TP), True Negative (TN), False Positive (FP) and False Negative (FN). True Positive represents the transactions that are fraudulent and was correctly classified by the model as fraudulent. True Negative represents the not fraudulent transactions that were correctly predicted by the

model as Not fraudulent. The third rating is False positive which represents the transaction that are fraudulent but was misclassified as not fraudulent. And finally False Negative which are the not fraudulent transactions that were identified as fraudulent, table 1 below shows the confusion matrix.

| Actual/Predicted | Positive | Negative |
|---|---|---|
| Positive | TP | FN |
| Negative | FP | TN |

| Model | | Accuracy |
|---|---|---|
| KNN | K = 3 | 99.89% |
| | K = 3 | |
| | K = 7 | 99.88% |
| | K = 7 | |
| Naïve Bayes | Naïve Bayes | 97.76% |
| | Naïve Bayes | |
| Logistic Regression | Logistic Regression | 99.92% |
| | Logistic Regression | |
| Support Vector Machine | SVM | 99.94% |