

# **CREDIT CARD FRAUD DEDUCTION**

## **USING DATA SCIENCE**

### **PHASE-1**

#### **DATA SCIENCE**

Data science is an interdisciplinary field that encompasses a set of techniques, processes, and methods for extracting valuable insights, knowledge, and patterns from data. It combines elements of statistics, computer science, domain expertise, and data engineering to collect, analyze, and interpret large and complex datasets. The ultimate goal of data science is to use data to inform decision-making, solve problems, and drive improvements in various domains, including business, healthcare, finance, and more. Data scientists use a combination of data analysis, machine learning, data visualization, and domain-specific expertise to uncover meaningful information from data and provide valuable insights for organizations and individuals.

#### **Tools used**

##### **1. Python**

Python is often used as a support language for software developers, for build control and management, testing, and in many other ways.

##### **2. Python idle**

IDLE can be used to execute a single statement and create, modify, and execute Python scripts. IDLE provides a fully-featured text editor to create Python scripts that include features like syntax highlighting, autocompletion, and smart indent.

##### **3. PyCharm**

PyCharm is a dedicated Python Integrated Development Environment (IDE) providing a wide range of essential tools for Python developers, tightly integrated to create a convenient environment for productive Python, web, and data science development

##### **4. A Kaggle dataset**

It also known as a Kaggle Kernel is a set of data provided by companies, students, and alum. These data sets allow competitors to work through problems, or use them as a practice simulation.

##### **5. MATLAB**

In Data Science, MATLAB is used for simulating **neural networks** and fuzzy logic. Using the MATLAB graphics library, you can create powerful visualizations. MATLAB is also used in image and signal processing.

##### **6. Excel**

Excel is a powerful **analytical tool for Data Science**. While it has been the traditional tool for data analysis, Excel still packs a punch. Excel comes with various formulae, tables, filters, slicers, etc. You can also create your own custom functions and

formulae using Excel. While Excel is not for calculating the huge amount of Data, it is still an ideal choice for creating powerful data visualizations and spreadsheets.

#### 7. Jupyter

Project **Jupyter** is an open-source tool based on IPython for helping developers in making open-source software and experiences interactive computing. Jupyter supports multiple languages like Julia, **Python**, and R.

#### 8. Matplotlib

**Matplotlib is a plotting and visualization library** developed for Python. It is the most popular tool for generating graphs with the analyzed data. It is mainly used for plotting complex graphs using simple lines of code. Using this, one can generate bar plots, histograms, scatterplots etc.

#### 9. Scikit-learn

Scikit-learn is a library-based in Python that is used for implementing Machine Learning Algorithms. It is simple and easy to implement a tool that is widely used for analysis and data science. Scikit-learn makes it easy to use complex machine learning algorithms. It is therefore in situations that require rapid prototyping and is also an ideal platform to perform research requiring basic Machine Learning. It makes use of several underlying **libraries of Python** such as SciPy, Numpy, Matplotlib, etc.

## CREDIT CARD FRAUD DETECTION

Credit card fraud is a term that has been coined for unauthorized access of payment cards like credit cards or debit cards to pay for using services or goods. Hackers or fraudsters may obtain the confidential details of the card from unsecured websites. When a fraudster compromises an individual's credit/debit card, everyone involved in the process suffers, right from the individual whose confidential data has been leaked to the businesses (generally banks) who issue the credit card and the merchant who is finalizing the transaction with purchase. This makes it extremely essential to identify the fraudulent transactions at the onset. Financial institutions and businesses like e-commerce are taking firm steps to flag the fraudsters entering the system. Various advanced machine learning technologies are at play, assessing every transaction and stemming the fraud users in its nip using behavioral data and transaction patterns. The process of automatically differentiating between fraudulent and genuine users is known as "credit card fraud detection".



**Lost/Stolen cards:** People steal credit cards from the mail and use them illegally on behalf of the owner. The process of blocking credit cards that have been stolen and re-issuing them is a hassle for both customers and credit card companies. Some financial institutions keep the credit cards blocked until it is verified that the rightful owner has received the card.

**Card Abuse:** The customer buys goods and items on the credit card but has no intention to pay back the amount charged by the bank for the same. These customers stop answering the calls as the deadline to settle the dues approaches. Sometimes they even declare bankruptcy—this type of fraud results in losses of millions every year.

**Identity Theft:** The customers apply illegitimate information, and they might even steal the details of a genuine customer to apply for a credit card and then misuse it. In such cases, even card blocking can not stop the credit card from falling into the wrong hands.

**Merchant Abuse:** Some merchants show illegal transactions (that never occurred) for money laundering. For performing these illicit transactions, legal information of genuine credit card users is stolen to generate replicas of the cards and use it for illegal work.

## **CREDIT CARD FRAUD DETECTION USING DATA SCIENCE**

Data science is used for credit card fraud detection because it offers powerful tools and techniques to identify suspicious patterns and anomalies in large volumes of transaction data. Here's why data science is particularly well-suited for this task:

### **1. Handling Big Data:**

Credit card transactions generate massive amounts of data. Data science techniques, including big data processing frameworks like Hadoop and Spark, enable the efficient handling and analysis of this data.

### **2. Pattern Recognition:**

Data science algorithms excel at recognizing patterns within data. They can identify unusual behaviors or transactions that deviate from established patterns, which is crucial for fraud detection.

### **3. Real-Time Analysis:**

Many data science models can analyze transactions in real-time, allowing for **immediate** detection and response to suspicious activities as they occur.

### **4. Adaptability:**

Fraudsters constantly evolve their tactics, so fraud detection systems need to adapt. Data science models can be regularly retrained to stay up-to-date with emerging fraud patterns.

### **5. Anomaly Detection:**

Anomaly detection techniques in data science can flag transactions that are statistical outliers, helping to pinpoint potentially fraudulent activity.

### 6. False Positive Reduction:

Data science can help reduce false positives (legitimate transactions mistakenly identified as fraud), which is crucial to avoid inconveniencing genuine cardholders.

### 7. Ensemble Techniques:

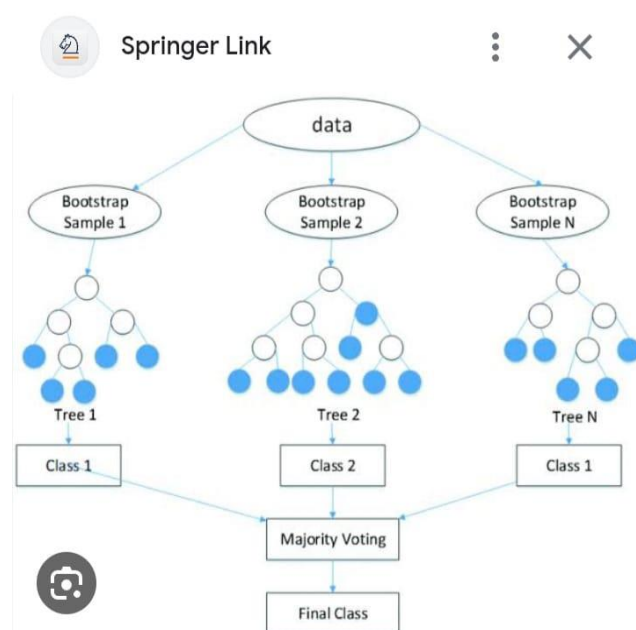
Combining multiple machine learning models using ensemble methods can improve fraud detection accuracy and reduce false alarms.

### 8. Predictive Analytics :

Data science can be used not only to detect ongoing fraud but also to predict future fraudulent activities based on historical data and trends.

## ALGORITHMS USED IN CREDIT CARD FRAUD DETECTION

1. Decision Tree
2. Predictive Analytics and Algorithms
3. Clustering Techniques
4. K-Nearest Neighbour Algorithm
5. Neural Networks
6. Naive Bayes Classifiers
7. Support Vector Machines (SVMs)



## **Problem Definition:**

The problem is to develop a machine learning-based system for real-time credit card fraud detection. The goal is to create a solution that can accurately identify fraudulent transactions while minimizing false positives. This project involves data preprocessing, feature engineering, model selection, training, and evaluation to create a robust fraud detection system.

## **Steps involved in design thinking:**

### **1.DATA SOURCE**

In credit card fraud detection using data science, the primary data source is typically a dataset containing transaction data. Here are the key data sources used in this type of project:

- ➤ **Transaction Data** : This is the most critical data source. It includes details of each credit card transaction, such as: - Transaction amount - Timestamp (date and time of the transaction) - Merchant information (merchant ID, location, category) - Card details (credit card number, expiration date)
- ➤ **Historical Data** : - Historical transaction data is used to train machine learning models. It provides a basis for the model to learn patterns of legitimate and fraudulent transactions. - The data sets are used to find the detection using kaggle website datasets

Link: <https://www.kaggle.com/datasets/mig-ulb/creditcardfraud>

- ➤ **Label Data** : This dataset includes labels or tags that specify whether each transaction is fraudulent or not. Labeling can be done manually by fraud analysts or through historical fraud data.
- ➤ **Customer Information** : While not always included, customer data can be useful in fraud detection. It may include details such as customer demographics, transaction history, and account status.
- ➤ **External Data** : Some projects incorporate external data sources to enhance fraud detection. This could include data from third-party fraud detection services, economic indicators, or geolocation data.
- ➤ **data model** : We use the logistics Regression Model to find the fraud detection using data set .

### **2.DATA PREPROCESSING :**

Data preprocessing is a crucial step in credit card fraud detection using data science. It helps ensure that the data you feed into your machine learning model is clean, standardized, and ready for analysis. Here are the common data preprocessing steps for this types :

**Handling Missing Values** : - Identify and assess missing values in the dataset. Missing values can occur in various fields like merchant information or card details. - Decide on an appropriate strategy to handle missing values. Common approaches include imputation

(filling in missing values with statistical measures like mean or median) or removing rows or columns with too many missing values. - In this use the method to find the Missing Values are NaN values is

`data.isnull().values.any() -> method used to handle missing values`

Row No.	Id	cluster	outlier	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V1
2660	2660	cluster_1	13.192	2177.000	-2.336	1.457	0.978	1.489	-0.841	2.219	-1.353	2.215	0.573	-0.637	-0.954	0.999	-0.1
542	542	cluster_0	12.559	405.000	-2.312	1.952	-1.810	3.998	-0.522	-1.427	-2.537	1.382	-2.770	-2.772	3.202	-2.906	-0.1
2042	2042	cluster_1	12.081	1580.000	-2.838	1.722	0.835	1.560	-1.704	0.790	-0.671	1.469	1.234	0.423	-1.579	0.452	-1.1
2027	2027	cluster_1	9.861	1572.000	-1.884	0.892	1.674	1.583	-0.697	0.732	-0.262	0.635	1.214	0.310	-1.338	0.141	-1.1
624	624	cluster_0	9.030	472.000	-3.044	-0.157	1.088	2.289	1.360	-1.965	0.326	-0.698	-0.271	-0.839	-0.415	-0.602	0.6
2791	2791	cluster_1	8.704	2336.000	-0.819	-0.372	2.087	-2.992	-0.918	0.033	0.146	0.028	-0.993	-1.195	-1.423	0.294	1.8
1633	1633	cluster_0	5.953	1284.000	-11.141	-0.613	-12.390	6.913	-32.992	21.393	34.303	-7.521	-1.926	-2.637	3.702	-1.843	2.4
2552	2552	cluster_1	5.513	2193.000	-1.379	0.052	1.548	-1.491	1.174	-1.525	0.811	-0.852	1.295	0.953	0.643	0.630	6.4
473	473	cluster_1	5.427	347.000	-1.531	1.400	-0.587	2.175	-2.138	-0.502	-1.215	0.957	-1.887	-2.311	2.770	-2.361	-0.1
1880	1880	cluster_1	5.338	1522.000	-0.893	0.135	0.952	-2.112	0.372	-0.832	1.185	-0.237	0.429	-1.813	-1.010	0.602	1.1
1504	1504	cluster_1	5.283	1170.000	-0.941	-0.446	2.188	-1.768	0.298	-0.590	0.422	-0.493	1.559	-1.235	-0.097	1.199	1.2
481	481	cluster_0	4.903	339.000	6.503	0.930	-0.858	2.943	-1.506	-1.300	-1.991	0.461	-1.124	-1.975	2.948	-2.150	-0.1
2028	2028	cluster_1	4.769	1573.000	-2.258	1.367	1.390	1.569	-1.087	0.700	-0.553	1.078	1.225	0.377	-1.454	0.301	-1.1
2703	2703	cluster_1	4.687	2242.000	-2.230	1.084	1.676	-2.676	-0.439	-0.427	0.419	-0.091	2.232	-0.839	-0.692	0.317	-0.1
2523	2523	cluster_1	4.569	2092.000	-0.279	0.647	1.220	-0.892	0.019	-0.986	0.829	-0.179	0.200	-1.158	0.390	0.836	0.6
2558	2558	cluster_1	4.144	2195.000	-2.290	-0.480	0.819	-1.705	0.822	-1.960	0.944	-0.542	1.323	-0.434	0.333	0.582	0.1
2548	2548	cluster_1	3.961	2191.000	-1.837	-0.382	1.647	-1.733	0.853	-1.597	0.950	-0.476	0.931	-0.918	0.349	0.833	0.5
2302	2302	cluster_1	3.950	1843.000	-4.719	1.259	-0.358	3.497	-3.571	1.584	-0.242	2.247	-0.556	0.360	-1.494	0.317	-0.1
2560	2560	cluster_1	3.888	2196.000	-1.319	-0.599	2.223	-2.301	0.069	-1.806	-0.211	0.062	1.062	-1.791	-0.417	0.197	-0.1
2964	2964	cluster_1	3.856	2597.000	-6.200	5.025	-2.742	-0.941	-6.656	5.432	-9.198	-22.589	-3.244	-5.453	-1.582	2.771	-0.1
996	996	cluster_1	3.815	751.000	-0.655	0.608	1.585	-3.009	0.038	-1.954	1.336	-0.613	0.690	-1.682	0.361	1.123	0.8
2717	2717	cluster_0	3.802	2257.000	1.238	-0.396	0.057	-0.474	-0.191	0.419	-0.466	0.258	0.387	-0.061	0.449	-0.913	-0.1
2958	2958	cluster_1	3.757	2691.000	-5.948	5.873	-2.932	-1.450	-2.522	2.642	-10.046	-22.746	-1.727	-5.962	-1.119	3.775	-1.1
2955	2955	cluster_1	3.725	2499.000	-5.390	5.141	-2.457	-1.434	-2.132	2.694	-9.704	-23.180	-1.747	-6.043	-0.995	3.607	-1.1
2952	2952	cluster_1	3.664	2497.000	-4.344	4.269	-1.580	0.273	-1.419	2.694	-11.155	-23.833	-1.840	-5.218	-0.374	3.607	-1.1
2652	2652	cluster_1	3.521	2185.000	-0.465	0.496	1.289	-2.193	0.358	-0.831	1.387	-0.888	1.491	-0.551	0.026	-0.102	-0.1
1326	1326	cluster_1	3.417	1038.000	-0.852	-0.246	1.352	-2.427	-0.549	-0.244	0.698	-0.677	1.213	-1.969	-1.154	0.992	1.6
2877	2877	cluster_1	3.417	2432.000	-1.102	-0.085	0.789	-2.019	1.073	-1.932	0.941	-0.325	0.914	-1.062	-0.472	0.398	0.3

**Data Cleaning :** - Check for and address data anomalies or errors, such as outliers that might indicate fraudulent transactions or incorrect entries. - Remove duplicates if they exist in the dataset. - Removing Missing values and NaN values .

**Data Transformation :** - Convert categorical variables into numerical format using techniques like one-hot encoding or label encoding, especially for features like merchant information or card details. - Normalize or scale numerical features to ensure they have similar scales. Common methods include Min-Max scaling or z-score standardization.

**Handling Imbalanced Data :** - Credit card fraud datasets often suffer from class imbalance, with a significantly higher number of non-fraudulent transactions compared to fraudulent ones. Address this issue using techniques like oversampling (creating more instances of the minority class), undersampling (reducing instances of the majority class), or using advanced methods like Synthetic Minority Over-sampling Technique (SMOTE).

**Data Splitting :** - Split the dataset into training and testing sets. The training set is used to train your machine learning model, while the testing set is used to evaluate its performance.

**Handling Time-based Data :** - If your dataset includes timestamps, consider converting them into useful features, **such** as time since the last transaction or time of day features

### **3.Feature Engineering : -**

Create relevant features that may aid in fraud detection. For example, calculate transaction frequency, average transaction amount, or time-based features like the hour of the day or the day of the week.

- **Data Splitting** : - Split the dataset into training and testing sets. The training set is used to train your machine learning model, while the testing set is used to evaluate its performance.
- **Handling Time-based Data** : - If your dataset includes timestamps, consider converting them into useful features, such as time since the last transaction or time of day features.
- **Data Scaling** : - Normalize numerical features to ensure that they have similar scales. This can help improve the performance of some machine learning algorithms.
- **Data Privacy and Security** : - Implement data privacy measures to protect sensitive information, such as credit card numbers. This may involve tokenization, encryption, or anonymization of data

### **4.MODEL SELECTION:**

Choose suitable machine learning algorithms

Example: logestic regression, random forest, gradient boosting

#### **LOGESTIC REGRESSION:**

Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables. Logistic regression predicts the output of a categorical dependent variable. To ensure high accuracy detection, two main method are used to clean the data. The mean-based method Deals with missing values, and the clustering-based Method deals with outliers. Extensive experiments are conducted to train and test The proposed classifier using a standard database. Import the relevant libraries

```

Import numpy as np
Import matplotlib.pyplot as plt
Import pandas as pd
Import seaborn as sns
From sklearn.linear_model
import LogisticRegression
From sklearn.model_selection
import _test_ssplrit

```



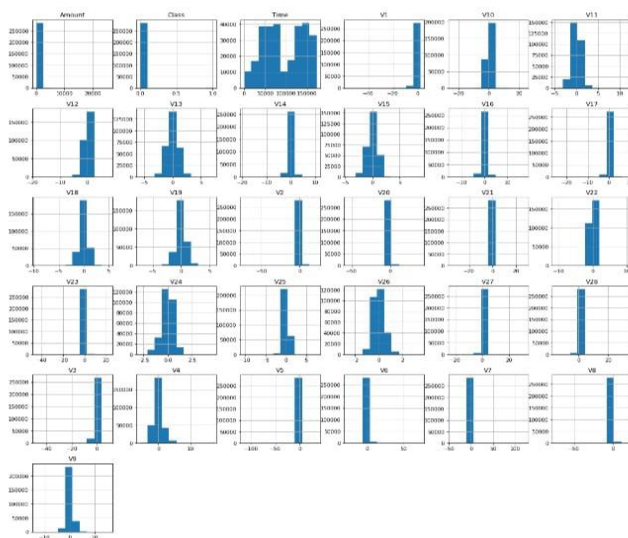
## RANDOM FOREST

Random forest is called a Random Forest because we use Random subsets of data and features and we end up building a Forest of decision trees (many trees). Random Forest is also a classic example of a bagging approach as we use different subsets of data in each model to make predictions.

## Credit Card Fraud Detection using Random Forest

Notebook Input Output Logs Co

```
plt.show()
```



## Implementation of Random Forest Algorithm

The procedure of the random forest algorithm execution is done there are several steps involved; first, there is a requirement together the information and to store the information. The gathered information are in the form of data set in an excel sheet. In data exploration, the entire data set checked and removed the unnecessary data that is present. However, the data which is further treated using a random forest algorithm in two ways by using train data set and then using the test data set.

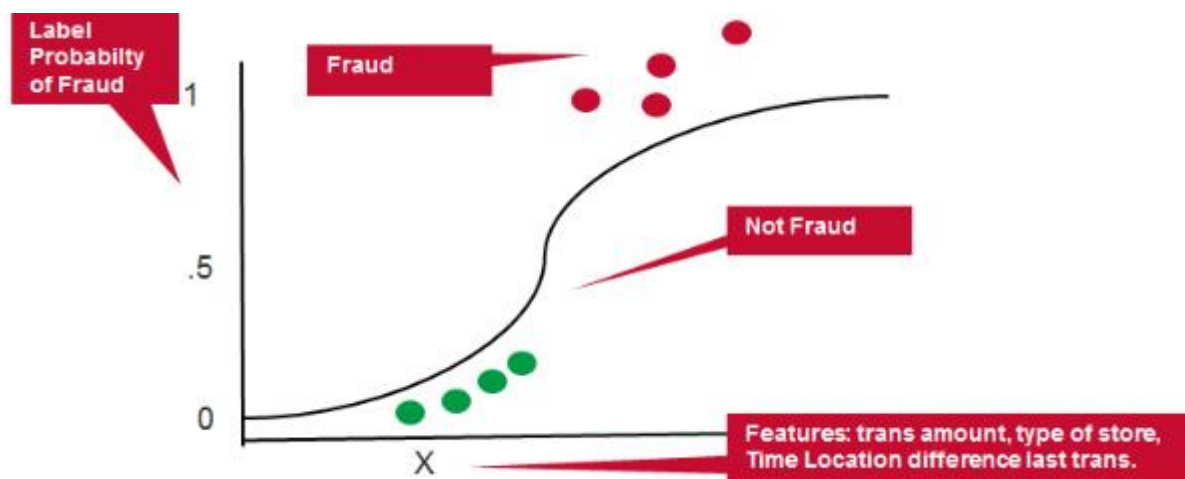
## 5.MODEL TRAINING

Data pre processing refers to the steps applied to make data more suitable for data mining. The steps used for Data Preprocessing usually fall into two categories: Selecting data objects and attributes for the analysis. Creating/changing the attributes. In this post I am going to walk through the implementation of Data Preprocessing methods using Python. It covers,



- Importing the libraries
- Importing the Dataset
- Handling of Missing Data
- Handling of Categorical Data
- Splitting the dataset into training and testing datasets
- Feature Scaling

For this Data Preprocessing script, I am going to use Anaconda Navigator and specifically Spyder to write the following code. If Spyder is not already installed when you open up Anaconda Navigator for the first time, then you can easily install it using the user interface.



## **6. EVALUATION**

Confusion matrix, also known as an error matrix, is a performance measurement for assessing classification models. Below is an example of a two-class confusion matrix.

Prediction Within the confusion matrix, there are some terms that you need to know, which can then be used to calculate various metrics:

Negative: Outcome where the model correctly predicts the negative class.

True Positive: Outcome where the model correctly predicts the positive class.

True Positive (Type 1 Error): Outcome where the model incorrectly predicts the positive class.

False Negative (Type 2 Error): Outcome where the model incorrectly predicts the negative class.

Accuracy: equal to the fraction of predictions that a model got right.

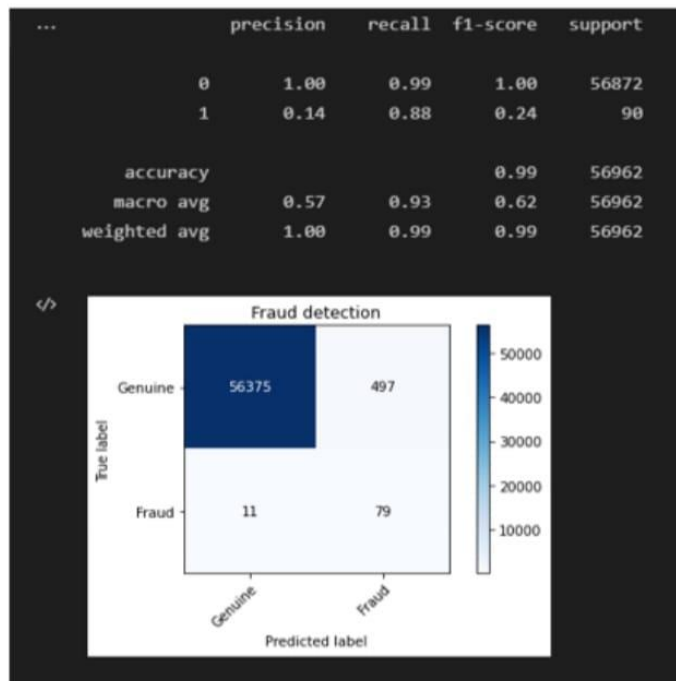


Figure: Confusion matrix for Naïve bayes

## Conclusion:

Credit card fraud is the unauthorized use of a credit or debit card to make purchases. Credit card companies have an obligation to protect their customers' finances and they employ fraud detection models to identify unusual financial activity and freeze a user's credit card if transaction activity is out of the ordinary for a given individual. The penalty for mislabeling a fraud transaction as legitimate is having a user's money stolen, which the credit card company typically reimburses. On the other hand, the penalty for mislabeling a legitimate transaction as fraud is having the user frozen out of their finances and unable to make payments. There is a very fine tradeoff between these two consequences and we will discuss how to handle this when training a model where data science has been used widely so far.