# CREDIT CARD FRAUD DEDUCTION USING APPLIED DATA SCIENCE

# PHASE-02(INNOVATION)

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.

## 1.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.
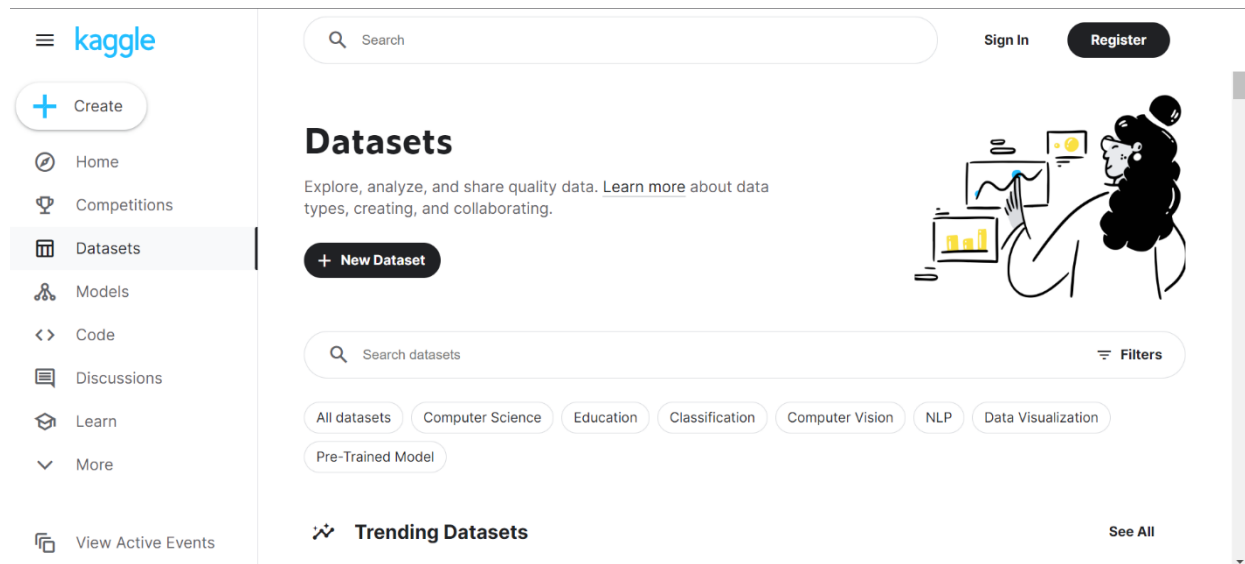
### 8. Predictive Analytics :

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

## 2.SOURCE OF THE DATASET

This dataset is taken from the website name "**kaggle**" and the link of the website is given below

**Link: www.kaggle.com/data**



In this page search for your desired dataset about the credit card fraud deduction

The link for this project:

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

## 3.DETAILS ABOUT DATASET

The dataset consists about the data's of the credit cards and the time in which the fraud has been taken place with the attributes

- In the given dataset for Credit card fraud detection,the rows are aligned from V1 to V28 and with amount and class.
- The columns are aligned with integers from 0 to 114.
- The values in dataset are both positive and negative.
- In the amount column, the highest value is 3828 and the least value is 0.75.
- In the class column,the values are almost equal to zero's and one's(Boolean values)
- There are almost 284808 tuples present

## 4.Libraries used in credit card fraud detection

We use the following libraries and frameworks in credit card fraud detection project.

1. Python – 3.x
2. Numpy – 1.19.2
3. Scikit-learn – 0.24.1
4. Matplotlib – 3.3.4
5. Imblearn – 0.8.0

## How to download the libraries:

### Download python 3.x

This document describes how to install Python 3.6 or 3.8 on Ubuntu Linux machines.To see which version of Python 3 you have installed, open a command prompt and run

```
$ python3 –version
```

If you are using Ubuntu 16.10 or newer, then you can easily install Python 3.6 with the following

commands:

```
$ sudo apt-get update
$ sudo apt-get install python3.6
```

If you're using another version of Ubuntu (e.g. the latest LTS release) or you want to use a more current

Python, we recommend using the deadsnakes PPA to install Python 3.8:

$ sudo apt-get install software-properties-common

$ sudo add-apt-repository ppa:deadsnakes/ppa

$ sudo apt-get update

$ sudo apt-get install python3.8

If you are using other Linux distribution, chances are you already have Python 3 pre-installed as well. If not, use your distribution's package manager. For example on Fedora, you would use dnf:

$ sudo dnf install python3

Download numpy 1.19.2

Installed Pythons found by C:\WINDOWS\py.exe Launcher for Windows

I use Python 3.6 for ESA SNAP snappy, so

PS C:\> cd $env:USERPROFILE\.snap\snap-python\

PS C:\Users\XXXXX\.snap\snap-python> py -3.6

Python 3.6.8 (tags/v3.6.8:3c6b436a57, Dec 24 2018, 00:16:47) [MSC v.1916 64 bit (AMD64)] on win32

Type "help", "copyright", "credits" or "license" for more information.

```
>>> import snappy
```

**Numpy installation:**

To install numpy on a Python.org 2 version for use with snappy:

```
PS C:\> py -3.6 -m pip install numpy
```

Collecting numpy

Downloading

51986f6d0c5c2/numpy-1.19.5-cp36-cp36m-win_amd64.whl (13.2MB)

100% |████████████████████████████████| 13.2MB 3.9MB/s

Installing collected packages: numpy

 The script f2py.exe is installed in 'D:\Python36\Scripts' which is not on PATH.

 Consider adding this directory to PATH or, if you prefer to suppress this warning, use –no-warn-script  location.

Successfully installed numpy-1.19.5

You are using pip version 18.1, however version 21.0.1 is available.

You should consider upgrading via the 'python -m pip install –upgrade pip' command.

Note that using py -M.N doesn't adjust the PATH

**Download matplotlib**

Matplotlib releases are available as wheel packages for macOS, Windows and Linux on PyPI. Install it

**using pip:**

```
Python -m pip install -U pip

Python -m pip install -U matplotlib
```

If this command results in Matplotlib being compiled from source and there's trouble with the compilation, you can add –prefer-binary to select the newest version of Matplotlib for which there is a precompiled wheel for that.

Install an official release

Matplotlib releases are available as wheel packages for macOS, Windows and Linux on PyPI. Install it

**using pip:**

```
Python -m pip install -U pip

Python -m pip install -U matplotlib
```

If this command results in Matplotlib being compiled from source and there's trouble with the

compilation, you can add –prefer-binary to select the newest version of Matplotlib for which there is a

precompiled wheel for your OS and Python.

6

Third-party distributions

Various third-parties provide Matplotlib for their environments.

Conda packages

Matplotlib is available both via the anaconda main channel

Conda install matplotlib

As well as via the conda-forge community channel

Conda install -c conda-forge matplotlib

**Downloading Imblearn 0.8.0**

Imbalanced-learn

Imbalanced-learn is a python package offering a number of re-sampling techniques commonly used in datasets showing strong between-class imbalance. It is compatible with scikit-learn and is part of scikit learn-contrib projects.

Documentation

Installation documentation, API documentation, and examples can be found on the documentation.

**Dependencies**

Imbalanced-learn is tested to work under Python 3.6+. The dependency requirements are based on the

last scikit-learn release:

Scipy(>=0.19.1)

Numpy(>=1.13.3)

Scikit-learn(>=0.23)

Joblib(>=0.11)

Keras 2 (optional)

Tensorflow (optional)

Additionally, to run the examples, you need matplotlib(>=2.0.0) and pandas(>=0.22).

# 5.How to  Train/Test

Train/Test is a method to measure the accuracy of your model .It is called Train/Test because you split the data set into two sets: a training set and a testing set.

> 80% for training, and 20% for testing.

You train the model using the training set.

You test the model using the testing set.

Train the model means create the model.

Test the model means test the accuracy of the model.

**Split Into Train/Test**

The training set should be a random selection of 80% of the original data.

The testing set should be the remaining 20%.

train_x = x[:80]

train_y = y[:80]

test_x = x[80:]

test_y = y[80:]

**Display the Training Set**

Display the same scatter plot with the training set.

**Example**

plt.scatter (train_x, train_y)

plt.show()

# 6.EXTRA EXPLANATION

ALGORITHM USED IN CREDIT CARD FRAUD DEDUCTION

 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)


# 7.METRICS USED FOR CREDIT CARD FRAUD DETECTION

The most frequent metric used is 'Accuracy'. High Accuracy doesn't mean that the model is performing better in all the situations. It is not always considered to be accurate as it sometimes is misleading in some situations like imbalanced class datasets.

## ACCURACY METRICS

Accuracy is a metric that generally describes how the model performs across all classes. It is useful when all classes are of equal importance. It is calculated as the ratio between the number of correct predictions to the total number of predictions

## CODE FOR ACCURACY OF A CLASSIFER

Importing All necessary libraries

from sklearn.metrics import accuracy_score

Calculating the accuracy of classifier

```
print(Accuracy of the classifier is: {accuracy_score(y_test, predictions)}")
```

## IMPLEMENTATION OF ACCURACY METRICS

The Accuracy score is calculated by dividing the number of correct predictions by the total

prediction number.

$$Accuracy = \frac{True\,Negatives + True\,Positive}{True\,Positive + False\,Positive + True\,Negative + False\,Negative}$$

Another formula:

$$Accuracy = \frac{Number\ of\ correct\ predictions}{Total\ number\ of\ predictions}$$