# Bangalore Institute of Technology
## Department of Computer Science and Engineering
## K R Road, V V Puram, Bengaluru-560004

Mini Project Synopsis on
**Credit Card Fraud Detection**

Submitted as the mini project for the subject
Data Mining and Data Warehousing (18CS641)

**Submitted by**

| | |
|---|---|
| **Samarth S Hiremath** | **1BI20CS150** |
| **Vamshik R** | **1BI20CS144** |
| **Rakshith Rajesh N B** | **1BI21CS408** |

For academic year 2022-23

**Under the guidance of**
**Dr. Bhargavi M S**
**Associate Professor**

# Introduction to the problem:

With the increase of people using credit cards in their daily lives, credit card companies should take special care in the security and safety of the customers. According to (Credit card statistics 2021) the number of people using credit cards around the world was 2.8 billion in 2019, in addition 70% of those users own a single card at least.

Reports of Credit card fraud in the US rose by 44.7% from 271,927 in 2019 to 393,207 reports in 2020. There are two kinds of credit card fraud, the first one is by having a credit card account opened under your name by an identity thief, reports of this fraudulent behavior increased 48% from 2019 to 2020. The second type is by an identity thief uses an existing account that you created, and it's usually done by stealing the information of the credit card, reports on this type of fraud increased 9% from 2019 to 2020 (Daly, 2021). Those statistics caught my attention as the numbers are increasing drastically and rapidly throughout the years, which gave me the motive to try to resolve the issue analytically by using different machine learning methods to detect the credit card fraudulent transactions within numerous transactions.

Nowadays as we can see that there is a huge increase online payment, and the payment is mostly done with the help of credit cards. It becomes a big problem for marketing company to overcome with the credit card fraudulent activities. Fraudulent can be done in many ways such as tax return in any other account, taking loans with wrong information etc. Therefore, we need an efficient fraudulent detection model to minimize fraudulent activity and to minimize their losses.

The primary objectives of this project are as follows:

**Dataset Selection**: We will carefully select diverse datasets from various domains, ensuring they represent real-world scenarios and challenges. These datasets will encompass different dimensions and characteristics, allowing us to evaluate the effectiveness of **Random Forest classifier** across various contexts.

**Random Forest Algorithm**: Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, **"**Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset**."** Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

**Feature Engineering and Preprocessing**: Prior to applying the Random Forest algorithm, we will conduct comprehensive feature engineering and preprocessing tasks to ensure the data's quality and suitability for clustering. This may involve handling missing values,

normalizing features, and addressing outliers to enhance the accuracy and robustness of the clustering results.
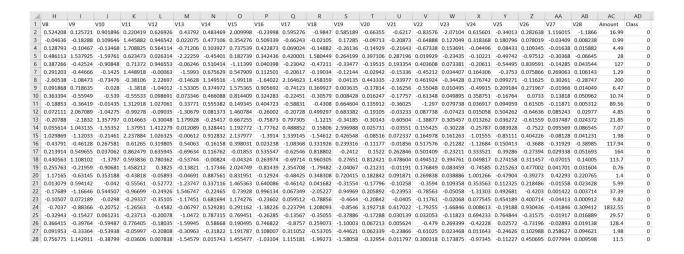
**Evaluation and Interpretation**: We will evaluate the effectiveness of the Random Forest classifier results using internal evaluation metrics like the silhouette score or external metrics if ground truth labels are available. Additionally, we will interpret the resulting clusters by analyzing their characteristics and visualizing the clusters to gain insights and facilitate decision-making.

By conducting this project, we aim to showcase the potential of the k-means clustering algorithm as a valuable tool for pattern discovery and data segmentation. The outcomes of our analysis will contribute to identifying distinct groups, understanding relationships between data points, and informing decision-making processes across various domains.

## Dataset Description:

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 30 attributes, 11558 rows. 28 attributes are numeric variables that due to confidentiality and privacy of the customers have been transformed using PCA transformation, the two remaining attributes are "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.

## Snapshot

## Preprocessing Techniques to be used:

As there are no NAs nor duplicated variables, the preparation of the dataset was simple the first alteration that was made to be able to open the dataset on Weka program is changing the type of the class attribute from Numeric to Class and identify the class as {1,0} using the program Sublime Text. Another alteration was made on the type as well on the R program to be able to create the model and the visualization.

## Classifier technique to be applied:

### Random Forest Classifier algorithm: -

Random Forest works in two-phase first is to create the random forest by combining N decision tree, and second is to make predictions for each tree created in the first phase.

The Working process can be explained in the below steps and diagram:

**Step-1:** Select random K data points from the training set.

**Step-2:** Build the decision trees associated with the selected data points (Subsets).

**Step-3:** Choose the number N for decision trees that you want to build.

**Step-4:** Repeat Step 1 & 2.

**Step-5:** For new data points, find the predictions of each decision tree, and assign the new data points to the category that wins the majority votes.

## Tools to be used:

1. Google Colaboratory