# Credit Card Fraud Detection using Machine Learning Techniques

Abhijit Gokhale, Shubham Sharma

*Abstract*— **Building Classification and Regression models to predict if a credit card transaction is fraudulent or legitimate**

## I. INTRODUCTION

In this modern era, there have always been people who will find new ways to access someone's finances illegally. As all transactions can easily be completed online with a credit card, consumers, banks, and merchants are put at risk when a data breach leads to monetary theft and ultimately the loss of customers' loyalty along with the company's reputation. Even after numerous mechanisms to stop fraud, fraudsters are continuously trying to find new ways and tricks to commit fraud.

## II. PROJECT

### A. Why fraud detection ?

Consumers in 2020 reported losing over $3.3 billion to **fraud**, an increase of $1.5 billion over the amount reported in 2019. Worldwide, businesses are expected to lose $75 billion to e-commerce **fraud** from 2019 to 2023. In 2019, the U.S. accounted for 33.57 percent of all gross card **fraud** losses worldwide.

Credit card fraud represented the second most commonly reported type of identity theft in 2020. In 2020, credit cards were most frequently identified as the payment method in fraud reports. According to Federal Trade Commission, there were over 390k reports of credit card fraud in 2020 and 149M dollars were lost only in the United States. In 2020 the rate of new account credit card fraud attempts rose 48% as compared to 2019.

Thus, in order to stop these frauds we need a powerful fraud detection system that not only detects the fraud but also detects it before it takes place and in an accurate manner. Detecting real-time fraud transactions is dreadfully difficult.

Given the demographics data, geolocation data and transaction data, our model will try to classify if a transaction is fraudulent. If our model were to achieve a high detection rate for fraudulent transactions, it could help the credit card merchants to successfully mitigate the frauds committed by hackers and save millions of dollars.

### B. Applications of ML and Big Data

Most of the fraudulent transactions follow a similar pattern. Using data mining and pattern recognition techniques, we can segment and classify data to search millions of transactions to find patterns and detect fraud. For our model to be of high quality, we need to ensure the amount and quality of data is significant.

### C. Project Goals

The project goal is to utilize existing credit card fraud data from January 2019 to December 2020 and develop a model that will provide the best results in revealing and preventing fraudulent transactions. For this objective the project will be supported by results of data visualizations to explore the data and feature engineering and dimensionality reduction to incorporate important features. These steps will get the data ready for our analyses. Next, we will perform time series analysis to identify the trends, and regression, classification and unsupervised machine learning techniques like clustering and association rules mining to identify whether a transaction is fraudulent or legitimate.

### D. Dataset

Our Data containing geographical, demographics and transaction features is a simulated credit card transaction dataset containing legitimate and fraud transactions from the duration 1st Jan 2019 - 31st Dec 2020. It was generated using Sparkov Data Generation tool and covers credit cards of 1000 customers doing transactions with a pool of 800 merchants. It is already split into training and test data, where training data contains transactions from 1st Jan 2019 to 21st June 2020 and test data contains transactions from 21st June 2020 to 31st Dec 2020. We will create a baseline model after data preprocessing but that might be underfitted or overfitted, which could be mitigated if a model is trained and tested by a randomly sampled dataset to create a robust model. This is an important step to detect and predict fraudulent or legitimate transactions and not to deteriorate the detection rate and false alarm rate.

We plan to build machine learning models like Support Vector Machines, Logistic Regression, k-Nearest Neighbours, Artificial Neural Networks, XGBoost and Random Forest and perform a comparative analysis of the model using quantitative measures like accuracy, sensitivity, specificity, detection rate and false alarm rate. For implementing the models, we will use both MLlib library in PySpark and scikit-learn package. Our ultimate intent is to classify and distinguish fraud transactions that can be used to lower the number of credit card frauds.