# Detecting fraudulent credit card transactions

## Problem Description:

Credit card fraud is a growing problem that affects both individuals and businesses. According to the Nilson Report, global losses due to credit card fraud exceeded $27 billion in 2018. Fraudulent transactions can lead to financial losses for individuals and businesses, as well as damage to their credit scores and reputations.

The goal of this project is to build a machine learning model that can accurately detect fraudulent credit card transactions. The project will use a dataset from Kaggle, which contains a sample of credit card transactions made over a two-day period in September 2013 by European cardholders. The dataset contains a total of 284,807 transactions, of which 492 are fraudulent.

**Dataset Description:**

The dataset used in this project can be accessed through the following Kaggle link: https://www.kaggle.com/mlg-ulb/creditcardfraud

The dataset contains the following columns:

- Time: Number of seconds elapsed between this transaction and the first transaction in the dataset
- V1-V28: Features obtained through a principal component analysis (PCA) transformation for security purposes. Unfortunately, due to confidentiality concerns, we do not have access to the original features and more background information about the dataset.
- Amount: Transaction amount in Euros
- Class: 1 for fraudulent transactions, 0 for non-fraudulent transactions.

**Objectives:**

The primary objectives of this project are as follows:

- To build a machine learning model that can accurately detect fraudulent credit card transactions based on the provided dataset.
- To evaluate the performance of the model in terms of precision, recall, F1-score, and accuracy.
- To identify the most important features in the dataset that contribute to the detection of fraudulent transactions.
- To provide insights and recommendations for credit card issuers and payment processors to improve their fraud detection strategies.

**Deliverables:**

The following deliverables are expected at the end of this project:

- A well-documented Jupyter notebook or Python script that contains the code used to build and evaluate the machine learning model.
- A report that summarizes the findings and recommendations of the project, including the performance metrics of the model, the most important features, and the insights gained from the analysis.
- A presentation that communicates the key findings and recommendations of the project to stakeholders, including credit card issuers and payment processors.

# Possible Framework:

The following framework outlines the steps that will be taken to complete this project:

1.  **Data Exploration:**
*   Load the dataset into a pandas dataframe and explore the dataset to understand its structure and content.
*   Check for missing values, data types, and any other data quality issues that need to be addressed.
*   Visualize the distribution of the target variable (fraudulent vs non-fraudulent transactions) and the other features to gain insights into the dataset.
2.  **Data Preprocessing:**
*   Perform data preprocessing steps such as feature scaling and normalization, if necessary.
*   Split the dataset into training and testing sets.
*   Address class imbalance by using techniques such as oversampling or undersampling.
3.  **Model Selection:**
*   Select appropriate machine learning algorithms for the problem, such as logistic regression, decision trees, random forests, or neural networks.
*   Train the selected models on the training set and evaluate their performance using metrics such as precision, recall, F1-score, and accuracy.
*   Tune hyperparameters of the models using techniques such as grid search or randomized search.
4.  **Model Interpretation:**
*   Interpret the trained models to understand the most important features that contribute to the detection of fraudulent transactions.
*   Visualize the feature importance using techniques such as permutation importance or partial dependence plots.
5.  **Model Deployment:**
*   Deploy the final model to a production environment or cloud platform for real-time fraud detection.
*   Monitor the performance of the model over time and update it as necessary.

# Code Explanation :

**Problem Statement**

The objective of this project is to develop a model that can accurately detect fraudulent credit card transactions. The data used in this project is sourced from Kaggle, and consists of anonymized credit card transaction data.

**Section 1: Importing Required Libraries**

In this section, we import the necessary libraries that will be used throughout the project. The **pandas** library is used for data manipulation and analysis, **numpy** is used for numerical computations, **matplotlib** is used for data visualization, **seaborn** is used for enhanced data visualization, **sklearn** is used for machine learning models and metrics, and **warnings** is used to ignore any warning messages that may be displayed during the course of the project.

**Section 2: Loading and Exploring the Data**

In this section, we load the credit card transaction data into a Pandas dataframe and explore the data. We print out the first 5 rows of the dataset to get a sense of the columns and the data contained within them. We also use the **describe()** method to get some basic summary statistics for each numerical column in the dataset.

**Section 3: Exploratory Data Analysis (EDA)**

In this section, we perform exploratory data analysis (EDA) on the dataset to gain insights into the data and its characteristics. We create histograms to visualize the distribution of the two classes (fraudulent and non-fraudulent transactions) in the dataset, as well as scatter plots to visualize the relationships between different columns in the dataset.

**Section 4: Data Preprocessing**

In this section, we preprocess the data to prepare it for machine learning. We start by creating a new column called **Class_Label** that maps the original **Class** column (which contains the labels 0 for non-fraudulent and 1 for fraudulent transactions) to more meaningful class labels ('Non-Fraudulent' and 'Fraudulent'). We then split the dataset into training and testing sets, with 80% of the data being used for training and 20% for testing. We also standardize the numerical features in the dataset using the **StandardScaler** function from **sklearn.preprocessing**.

**Section 5: Handling Class Imbalance**

In this section, we handle the class imbalance in the dataset by oversampling the minority class (fraudulent transactions) using the **RandomOverSampler** function from **imblearn.over_sampling**. This is done to ensure that the machine learning model is not biased towards the majority class (non-fraudulent transactions) during training.

**Section 6: Model Training and Hyperparameter Tuning**

In this section, we train a random forest classifier on the oversampled training data. We perform hyperparameter tuning using a grid search to find the best hyperparameters for the model. We evaluate the performance of the model using cross-validation with 5 folds, and print out the mean accuracy score and standard deviation.

**Section 7: Model Evaluation**

In this section, we evaluate the performance of the trained model on the test data. We use the **predict()** method of the trained model to make predictions on the test data, and then print out the confusion matrix and classification report, which shows the precision, recall, and F1-score for each class in the dataset.

**Section 8: Model Interpretation**

In this section, we interpret the trained model by plotting the feature importances for the top 10 features in the dataset. We use the **feature_importances_** attribute of the trained model to calculate the feature importances, and then plot them as a horizontal bar chart.

**Conclusion**

In this project, we developed a machine learning model that can accurately detect fraudulent credit

# Future Work :

**Introduction**

In this section, we'll discuss some future work that can be done to improve the performance of the credit card fraud detection model. We'll discuss several steps that can be taken to enhance the model's accuracy and robustness.

**Section 1: Feature Engineering**

In this section, we can explore additional feature engineering techniques to extract more useful information from the dataset. This could include creating new features by combining existing features or by performing additional transformations on the data. Some potential techniques to explore include:

- Dimensionality reduction techniques such as PCA or t-SNE to identify the most important features in the dataset.
- Feature scaling and normalization techniques to improve the model's ability to distinguish between different features.
- Feature selection techniques to identify the most important features in the dataset.

**Section 2: Model Selection**

In this section, we can explore different machine learning models to determine which one performs best on the credit card fraud detection dataset. Some potential models to explore include:

- Logistic Regression
- K-Nearest Neighbors
- Support Vector Machines
- Gradient Boosting Machines
- Neural Networks

We can also try ensembling techniques such as stacking or bagging to combine multiple models and improve their performance.

**Section 3: Hyperparameter Tuning**

In this section, we can perform additional hyperparameter tuning to find the optimal hyperparameters for the chosen machine learning model. This can be done using a grid

search or a randomized search, depending on the size of the hyperparameter space. Some hyperparameters to tune include:

- Number of estimators in the model (for ensemble methods)
- Maximum depth of the decision trees (for tree-based models)
- Learning rate (for gradient boosting models)
- Regularization parameters (for logistic regression and support vector machines)

**Section 4: Imbalanced Class Handling**

In this section, we can explore additional techniques for handling the imbalanced class distribution in the dataset. Some potential techniques to explore include:

- Undersampling the majority class to balance the class distribution.
- Using more advanced oversampling techniques such as SMOTE or ADASYN to generate synthetic samples of the minority class.
- Using more advanced cost-sensitive methods to assign different misclassification costs to different classes.

**Section 5: Model Deployment**

In this section, we can deploy the trained machine learning model to a production environment. This can be done using a variety of tools and platforms, depending on the specific needs of the organization. Some potential techniques to explore include:

- Building a web-based dashboard to allow users to interact with the model and visualize its predictions.
- Integrating the model into an existing software application or API.
- Deploying the model to a cloud-based platform such as AWS or Google Cloud Platform.

**Conclusion**

In this future work, we explored several potential areas for improving the performance of the credit card fraud detection model. By exploring additional feature engineering techniques, model selection, hyperparameter tuning, and class balancing techniques, we can further improve the accuracy and robustness of the model. Additionally, by deploying the model to a production environment, we can make it available for real-world use and potentially save organizations millions of dollars in fraudulent transactions.

# Exercise Questions :

1. **What is the importance of imbalanced class handling in credit card fraud detection?**
The credit card fraud detection dataset typically has a highly imbalanced class distribution, where the number of non-fraudulent transactions is much higher than the number of fraudulent transactions. This can lead to the model being biased towards the majority class and having poor performance on the minority class. Therefore, it is important to handle the imbalanced class distribution by using techniques such as undersampling, oversampling, or cost-sensitive methods to ensure that the model can accurately detect fraudulent transactions.

2. **What are some potential feature engineering techniques that can be used in credit card fraud detection?**
Feature engineering is an important step in machine learning and can help to extract more useful information from the dataset. Some potential feature engineering techniques that can be used in credit card fraud detection include dimensionality reduction techniques such as PCA or t-SNE, feature scaling and normalization techniques, and feature selection techniques to identify the most important features in the dataset.

3. **What is the role of hyperparameter tuning in machine learning?**
Hyperparameter tuning is the process of selecting the optimal hyperparameters for a machine learning model. Hyperparameters are parameters that are set before training the model, such as the learning rate, regularization parameter, or maximum depth of the decision trees. By tuning the hyperparameters, we can improve the performance of the model and ensure that it is robust to different types of data.

4. **What are some potential drawbacks of using oversampling techniques in imbalanced class handling?**
Oversampling techniques such as SMOTE or ADASYN can generate synthetic samples of the minority class to balance the class distribution. However, oversampling can lead to overfitting if the model becomes too reliant on the synthetic samples. Additionally, oversampling can increase the computational complexity of the model, which can be problematic for large datasets.

5. **How can you deploy a machine learning model to a production environment?**
To deploy a machine learning model to a production environment, we can use a variety of tools and platforms, depending on the specific needs of the organization. Some potential techniques include building a web-based dashboard to allow users to interact with the model and visualize its predictions, integrating the model into an existing software application or API, or deploying the model to a cloud-based platform such as

AWS or Google Cloud Platform. We also need to ensure that the model is properly tested and validated before deploying it to a production environment to avoid any potential issues.

# Concept Explanation :

So, the algorithm we used in this project is called "Random Forest." No, it's not a forest that randomly grows trees, but it is a forest of decision trees that work together to make predictions.

Imagine you're lost in a forest and you need to find your way out. You come across a tree and it has a signpost pointing in two directions. One direction says "go left" and the other says "go right." You choose a direction and keep walking until you come across another tree with another signpost, and you make another decision based on the signpost.

This is basically what a decision tree does in machine learning. It asks a series of questions to make a decision or prediction. For example, in our credit card fraud detection project, a decision tree might ask questions like "Was the transaction amount greater than $500?" or "Did the transaction occur in a foreign country?" based on the features in the dataset.

But what if one decision tree isn't enough? What if there are multiple paths to follow and you don't know which one is the best? This is where the "Random Forest" comes in.

Think of it like a group of hikers in the forest. Each hiker has their own map and they all have different routes to get to the same destination. But by combining their maps and choosing the best parts of each route, they can find the most efficient path to their destination.

Similarly, in a Random Forest, multiple decision trees are trained on different parts of the dataset and their predictions are combined to make a final prediction. Each decision tree is like a hiker with their own map, and by combining their predictions, we can make a more accurate prediction.

For example, if one decision tree predicts that a transaction is fraudulent because the amount is unusually high, but another decision tree predicts that the transaction is non-fraudulent because it occurred in the same location as the cardholder's previous transactions, the Random Forest can combine these predictions and make a more informed decision.

And that's how the Random Forest algorithm works! It's a group of decision trees that work together to make more accurate predictions by combining their individual predictions. So the next time you're lost in a forest, remember that a group of decision trees might just be able to guide you out!