

# Fraud detection using the IEEE-CIS Fraud Detection

---

## **Problem Description :**

The goal of this project is to detect fraudulent transactions in financial transactions by using machine learning techniques. The dataset used in this project is the IEEE-CIS Fraud Detection dataset available on Kaggle. The dataset contains transactional data from various sources. The objective of this project is to predict the probability that a transaction is fraudulent or not.

## **Dataset Description**

The dataset consists of two files:

1. **train\_transaction.csv**: Contains the transaction data along with the target variable indicating whether the transaction is fraudulent or not.
2. **train\_identity.csv**: Contains identity information such as device type, operating system, browser, and IP address for each transaction.

The target variable is binary (0 or 1) indicating whether the transaction is fraudulent or not. There are a total of 49 features including the target variable.

## **Requirements**

1. Python 3.x
2. Jupyter Notebook
3. Libraries: pandas, numpy, sklearn, matplotlib, seaborn, xgboost

## **Deliverables**

1. A Jupyter notebook containing the data analysis, feature engineering, model selection, and evaluation.
2. A machine learning model that can predict fraudulent transactions with high accuracy.
3. A report summarizing the findings and recommendations.

## **Objectives**

1. Perform exploratory data analysis to understand the distribution of variables in the dataset.
2. Perform feature engineering to create new variables that can improve the performance of the machine learning model.
3. Train and evaluate different machine learning models using the dataset.
4. Select the best performing model and tune its hyperparameters to improve its performance.
5. Evaluate the final model on a holdout dataset to estimate its performance on unseen data.

**Background** Fraud detection is an important problem in the financial industry, where it is important to identify fraudulent transactions in order to prevent financial losses. Machine learning algorithms can be used to automatically detect fraudulent transactions based on historical data. The IEEE-CIS Fraud Detection dataset provides an opportunity to develop and test such algorithms.

1. Import Libraries: Import necessary libraries and load the dataset.
2. Exploratory Data Analysis (EDA): Perform basic EDA to understand the data and its distribution. Handle missing values, outliers, and other data inconsistencies.
3. Feature Engineering: Create new features, perform feature selection, and encode categorical variables.
4. Model Development: Build and train various machine learning models such as logistic regression, decision tree, random forest, XGBoost, and neural networks.
5. Model Evaluation: Evaluate the models using appropriate evaluation metrics such as ROC-AUC, F1-score, and accuracy. Identify the best performing model.
6. Hyperparameter Tuning: Fine-tune the best model using GridSearchCV or RandomizedSearchCV to obtain optimal hyperparameters.
7. Model Deployment: Deploy the best model to make predictions on new and unseen data.

Note: Repeat steps 4 to 7 as necessary to improve the model performance.

**Deliverables:**

- Exploratory Data Analysis (EDA) report
- Feature Engineering report
- Model Development report
- Model Evaluation report
- Hyperparameter Tuning report
- Final Model report
- Model Deployment report

## **Code Explanation :**

Here is the simple explanation for the code which is provided in the code.py file.

### **Section 1: Importing Libraries and Loading Data**

In this section, we import the necessary libraries for data processing, feature engineering, model building, and evaluation. We then load the dataset using Pandas and store it as a DataFrame object.

### **Section 2: Exploratory Data Analysis (EDA)**

In this section, we perform exploratory data analysis to gain insights into the dataset. We look at the distribution of target classes, the distribution of missing values, the distribution of categorical and numerical features, and the correlation between features.

### **Section 3: Feature Engineering**

In this section, we perform feature engineering to create new features that can improve model performance. We convert categorical features into numerical features using label encoding and one-hot encoding. We also fill missing values using imputation techniques and create new features based on existing features.

### **Section 4: Model Building**

In this section, we split the dataset into training and testing sets. We then build machine learning models using different algorithms such as Logistic Regression, Decision Trees, Random Forests, Gradient Boosting, and XGBoost. We evaluate the performance of each model using different evaluation metrics such as accuracy, precision, recall, F1 score, and ROC-AUC score.

### **Section 5: Hyperparameter Tuning**

In this section, we perform hyperparameter tuning to find the best set of hyperparameters for each model. We use techniques such as grid search and random search to find the optimal hyperparameters that can improve model performance.

## **Section 6: Model Evaluation**

In this section, we evaluate the performance of the models using different evaluation metrics on the testing set. We also compare the performance of different models to identify the best performing model. We then save the best model for future use.

To run the code, you need to have Python installed on your system along with the necessary libraries such as Pandas, NumPy, Scikit-learn, XGBoost, and LightGBM. You can install these libraries using pip or conda package manager. Once you have installed the required libraries, you can run the code in any Python environment such as Jupyter Notebook or Google Colab.

## **Future Work :**

**1. Feature Engineering** In the current project, we have only used a limited set of features for training our model. However, there are many other features available in the dataset that can be used for fraud detection. Further feature engineering can be done to extract more information from the dataset and improve the performance of the model.

**2. Hyperparameter Tuning** In the current project, we have used default hyperparameters for training our model. However, hyperparameters play a crucial role in determining the performance of the model. Therefore, a more thorough hyperparameter tuning process can be implemented to find the optimal set of hyperparameters for our model.

**3. Ensembling** Ensembling is a powerful technique that can be used to improve the performance of machine learning models. In this project, we can use different models and ensemble them to improve the performance of the fraud detection system.

**4. Data Augmentation** Data augmentation is a technique used to generate new training samples from the existing ones by applying random transformations such as flipping, rotating, scaling, etc. This technique can help to increase the diversity of the training data and prevent overfitting.

**5. Real-time Fraud Detection** Currently, the model is trained on historical data and used to detect fraud in batch mode. However, in real-world scenarios, fraud can occur in real-time. Therefore, a real-time fraud detection system can be implemented that continuously monitors the transactions and detects fraud in real-time.

### **Step-by-Step Implementation Guide**

1. Perform more extensive feature engineering to extract more information from the dataset.
2. Implement a hyperparameter tuning process to find the optimal set of hyperparameters for our model.
3. Ensemble different models to improve the performance of the fraud detection system.
4. Implement data augmentation to increase the diversity of the training data and prevent overfitting.
5. Develop a real-time fraud detection system that continuously monitors the transactions and detects fraud in real-time.

To implement these steps, we can use the existing code as a starting point and modify it accordingly. Additionally, we may require additional computational resources to implement these steps efficiently.

## **Exercise :**

**Try to answers the following questions by yourself to check your understanding for this project. If stuck, detailed answers for the questions are also provided.**

**1. What are some common challenges in fraud detection, and how did you address them in this project?**

Answer: Some common challenges in fraud detection include imbalanced datasets, a lack of labeled data, and evolving fraud tactics. To address these challenges, we used techniques such as oversampling the minority class, utilizing unsupervised learning for anomaly detection, and incorporating domain knowledge and expert input to identify and label potential fraud cases.

**2. How did you evaluate the performance of your fraud detection model?**

Answer: We evaluated the performance of our model using a combination of metrics such as accuracy, precision, recall, F1 score, and ROC AUC. We also used techniques such as cross-validation and hyperparameter tuning to ensure that our model was generalizable and robust.

**3. Can you explain how you used feature engineering in this project?**

Answer: Feature engineering is a crucial aspect of fraud detection, as it involves selecting and transforming the most relevant and informative features to improve model performance. In this project, we used a combination of domain knowledge and statistical analysis to select and engineer features such as transaction amount, time of day, and geographic location.

**4. How did you handle missing or incomplete data in the dataset?**

Answer: We utilized various techniques such as imputation, dropping columns or rows with high amounts of missing data, and creating new features to account for missing values. We also explored the patterns and distributions of missing data to identify potential sources of bias or errors in our analysis.

**5. Can you discuss any potential ethical or privacy concerns related to this project?**

Answer: Fraud detection often involves accessing and analyzing sensitive personal or financial information, which raises concerns around privacy and ethical considerations. To mitigate these concerns, we followed best practices such as anonymizing or pseudonymizing data, using secure and encrypted storage and transmission methods, and obtaining appropriate legal and ethical clearances for the collection and use of data.