# Fraud Detection in Credit Card Transactions

**Position:** Data Science Intern
**Company:** XPace Technologies Pvt Ltd

---

## 1. Introduction

Fraud detection in credit card transactions is a critical challenge for financial institutions, requiring sophisticated models to identify fraudulent activities while minimizing false positives. At XPace Technologies, we developed a fraud detection system leveraging machine learning techniques to improve the accuracy of predicting fraudulent transactions. This project aims to enhance the institution's ability to detect fraudulent behavior and ensure secure transactions.

---

## 2. Data Preprocessing

The dataset comprised several variables related to transaction details, including transaction ID, customer ID, transaction date, amount, merchant, location, and fraud labels. To ensure the dataset was ready for model development, several preprocessing steps were undertaken:

- **Handling Missing Values:** Missing values were handled by filling numerical columns with median values, ensuring data integrity without introducing bias.
- **Categorical Data Encoding:** Categorical variables such as "Merchant," "Location," "Transaction Type," and "Card Type" were converted into numerical form using label encoding to make them compatible with machine learning models.
- **Target Variable Encoding:** The target variable "Is Fraudulent" was encoded into binary form, where fraudulent transactions were labeled as 1 and legitimate transactions as 0.

By thoroughly cleaning and transforming the dataset, we prepared it for robust model training and evaluation.

---

## 3. Exploratory Data Analysis (EDA)

Exploratory Data Analysis provided crucial insights into the dataset's structure and patterns. Key findings from the analysis include:

- **Balanced Dataset:** The transactions were legitimate, with fraudulent cases representing a equal portion of the dataset. This balance highlighted the no need for techniques to handle class disparity.

- **Transaction Trends:** Patterns in transaction amounts, merchant types, and geographic locations were examined to identify potential red flags indicative of fraud.

This initial analysis was crucial for understanding the data and guiding the model-building process.

---

## 4. Feature Engineering

To enhance model performance, we engineered additional features that could reveal hidden patterns in the data:

- **Transaction Hour:** Extracting the hour from the transaction timestamp helped uncover potential time-based fraud trends, such as unusual activity during specific periods.
- **Transaction Frequency:** The number of transactions per customer was calculated to identify abnormal transaction behaviors that may indicate fraud.

By incorporating these engineered features, we aimed to improve the model's predictive power and its ability to detect fraud.

---

## 5. Model Development

For this project, we selected a **Random Forest Classifier** due to its ability to handle imbalanced datasets and its robustness in classification tasks. The dataset was split into training and testing sets (70% training, 30% testing) to evaluate the model's performance.

- **Training the Model:** The Random Forest Classifier was trained on the processed dataset, learning from both legitimate and fraudulent transactions.
- **Handling Imbalance:** Techniques like **class weighting** were employed to balance the impact of minority (fraud) and majority (legitimate) classes.

---

## 6. Model Evaluation

The model's performance was evaluated using several metrics to ensure its effectiveness:

- **Accuracy:** The model exhibited a high accuracy rate, correctly classifying most transactions as either legitimate or fraudulent.
- **Precision & Recall:** High precision indicated that most predicted fraud cases were true frauds, while satisfactory recall ensured that the model detected a significant portion of actual fraud cases.
- **Confusion Matrix:** The confusion matrix provided insight into the model's classification, showing the balance between true positives, true negatives, false positives, and false negatives.

While the model performed well overall, some fraudulent cases were still misclassified, indicating room for further refinement.

---

## 7. Recommendations

Based on the model's performance and our observations, the following recommendations were made to further improve fraud detection capabilities:

- **Real-Time Monitoring:** Implementing real-time transaction monitoring to quickly flag suspicious transactions.
- **Regular Model Retraining:** To stay ahead of evolving fraud patterns, retraining the model periodically with new data is essential.
- **Integration with External Data:** Incorporating additional external data sources (e.g., customer behavior, geographical trends) can enhance detection accuracy.

---

## 8. Conclusion

In conclusion, this project successfully developed a fraud detection model that leverages Random Forest techniques to identify fraudulent credit card transactions with a high level of accuracy. With continuous improvement, real-time integration, and the adoption of advanced fraud-detection mechanisms, this solution can significantly strengthen a financial institution's fraud prevention efforts, providing safer and more reliable transactions for customers.