# Credit Card Fraud Detection Report

## 1. Introduction

Credit card fraud detection is an essential challenge for financial institutions due to the increasing volume of online transactions. Fraudulent activities result in **financial losses, reputational damage, and regulatory risks**. Machine learning offers an effective way to **detect and prevent fraudulent transactions** by identifying suspicious patterns.

This report describes a **machine learning pipeline** designed for fraud detection, covering **data preprocessing, model selection, performance evaluation, and future improvements**.

---

## 2. Design Choices

Developing a fraud detection system requires a structured approach, including **data preprocessing, feature selection, model evaluation, and handling class imbalance**.

### 2.1 Data Exploration and Preprocessing

**Step 1: Exploratory Data Analysis (EDA)**

- Checked for **missing values** and outliers.
- Examined the **distribution of legitimate vs. fraudulent transactions** (highly imbalanced dataset).
- Visualized key features to **detect anomalies** in transaction behavior.

**Step 2: Feature Transformation**

- **Time Feature:** Converted into a more useful format, such as **hour of the day** to detect time-based fraud patterns.
- **Amount Feature:** Standardized to ensure numerical consistency across different transactions.

**Step 3: Handling Class Imbalance**

- Fraud cases are significantly lower than legitimate transactions.
- Used **SMOTE (Synthetic Minority Oversampling Technique)** to **oversample fraudulent transactions** and balance the dataset.
- Ensured that the **model does not favor the majority class (legitimate transactions)**.

---

**2.2 Feature Selection**

To improve the **efficiency and accuracy** of the model, we selected the most relevant features:

- **Highly correlated or redundant features were removed** to reduce noise.
- **Recursive Feature Elimination (RFE)** and **Random Forest Feature Importance** were used to identify the most influential predictors.

---

**2.3 Model Selection**

Several machine learning models were tested to find the **most effective fraud detection approach**:

- **Logistic Regression** – A simple, interpretable model but struggles with non-linear relationships.
- **Random Forest** – Uses multiple decision trees to capture complex patterns.
- **XGBoost** – A high-performance gradient boosting algorithm that works well with imbalanced data.
- **Support Vector Machine (SVM)** – Effective for classification but computationally expensive for large datasets.

---

**2.4 Model Evaluation Metrics**

Fraud detection requires balancing **precision and recall** to reduce false positives (blocking legitimate users) and false negatives (missing fraud cases).

- **Accuracy:** Measures overall correctness but can be misleading in imbalanced datasets.
- **Precision:** The proportion of predicted fraud cases that are actually fraudulent.
- **Recall:** The proportion of actual fraud cases that were correctly detected.
- **F1 Score:** A balance between precision and recall, useful for imbalanced datasets.

---

# 3. Performance Evaluation

### 3.1 Model Training

- The dataset was split into **80% training** and **20% testing** to evaluate model performance.
- **Stratified sampling** was used to ensure the fraud-to-legitimate ratio remained consistent.

---

### 3.2 Model Performance Results

The models were tested using the **test dataset**, and their performance was evaluated using the defined metrics.

| Model | Precision | Recall | F1 Score |
|---|---|---|---|
| **Logistic Regression** | 0.98 | 0.93 | 0.95 |
| **Random Forest** | 0.99 | 1.00 | 0.99 |
| **XGBoost** | 0.99 | 1.00 | 0.99 |

**Key Takeaways:**

- ➢ **XGBoost achieved the best performance** and was selected for deployment.
- ➢ **Random Forest performed similarly well**, making it a viable backup model
- ➢ **Logistic Regression, while interpretable, showed slightly lower recall**, meaning it missed more fraud cases than XGBoost.

---

# 4. Future Work

Although the model performs well, several improvements can be explored:

### 4.1 Enhanced Feature Engineering

- Incorporate **transaction metadata** (e.g., device type, location).
- Analyze **customer spending behavior** for better fraud detection.

### 4.2 Anomaly Detection Techniques

- Implement **Autoencoders and Isolation Forests** for semi-supervised learning.
- Detect emerging fraud patterns **not seen in the training dataset**.

### 4.3 Real-Time Deployment

- Deploy the model using **AWS SageMaker, Google Cloud AI, or Microsoft Azure**.
- Monitor **fraud trends in real-time** and update the model periodically.

**4.4 Model Explainability**

- Use **SHAP (SHapley Additive Explanations)** to make the model's predictions **transparent**.
- **LIME (Local Interpretable Model-Agnostic Explanations)** can help regulators and financial analysts understand **why a transaction was flagged as fraudulent**.

---

# 5. Conclusion

This report outlines a **highly effective fraud detection pipeline** using **machine learning**. The **XGBoost model demonstrated the best performance**, and future enhancements can improve its **scalability and interpretability**.

**Key Achievements:**
✔ Successfully **handled imbalanced data** using SMOTE.
✔ Used **advanced models** (XGBoost, Random Forest) for **high fraud detection accuracy**.
✔ Identified areas for **future improvements** to make the system even more robust.

---

# 6. References

- **Dataset Source:** [Mention dataset origin]
- **Libraries Used:** Scikit-learn, XGBoost, Pandas, NumPy
- **SMOTE Paper:** Chawla et al., 2002 – "SMOTE: Synthetic Minority Over-sampling Technique"