# Credit Card Fraud Detection Project Report

Prepared by: Neha Jhakra

Date: Sunday, August 24, 2025, 11:12 AM IST

# Executive Summary

This project presents an end-to-end analysis of credit card transaction data to detect and mitigate fraud risks. Utilizing a dataset comprising transaction details such as Transaction ID, Customer Name, Merchant Name, Transaction Date, Transaction Amount (INR), Fraud Risk, Fraud Type, State, Card Type, Bank, IsFraud, Fraud Score, Transaction Category, and Merchant Location, the analysis employs SQL for data querying, Python for data cleaning, exploratory data analysis (EDA), and machine learning (ML) modeling, and Power BI for interactive dashboard visualization.

Key findings include:

- Rajasthan records the highest fraud incidents (42 medium-risk cases), followed by West Bengal (25 high-risk cases) and Kerala (38 medium-risk incidents).

- Transaction volumes peaked at 1.13M in May, dipped to 625K in September, and recovered to 1.28M in December.

- Electronics and E-commerce categories show high fraud associations (52 and 36 cases, respectively).

- Models like Random Forest achieved high ROC-AUC scores ($>0.90$), with emphasis on recall to prioritize fraud detection.

The project demonstrates practical applications in banking for fraud prevention, with recommendations for real-time monitoring and targeted interventions. The dashboard enables stakeholders to explore fraud hotspots, trends, and risk profiles interactively.

# 1 Introduction

## 1.1 Project Background

Credit card fraud is a pervasive issue in the financial sector, causing significant losses and eroding customer trust. This project analyzes a dataset of credit card transactions (adapted from sources like Kaggle and DataCamp, focused on Indian states despite initial western U.S. mentions) to uncover patterns, assess risks, and develop predictive insights. The objective is to build a fraud detection system that errs on the side of cautionprioritizing the identification of fraudulent transactions (high recall) even if it increases false positives.

## 1.2 Objectives

- Perform data exploration and cleaning to ensure quality.

- Conduct aggregate and time-based analyses to identify fraud trends.

- Develop ML models for fraud prediction.

- Create an interactive Power BI dashboard for visualization and decision-making.

- Provide actionable recommendations for fraud mitigation.

## 1.3   Scope and Limitations

The analysis covers transaction data from January to December (year unspecified, but trends are monthly). Limitations include potential class imbalance in fraud data, no real-time data access, and reliance on provided features without external enrichments like IP or device info. The project assumes the dataset is representative but notes possible biases in labeling.

# 2   Methodology

## 2.1   Data Source and Structure

The dataset is a CSV file ("Credit Card Fraud Risk Analysis.csv") with the following fields:

- **Transaction ID**: Unique identifier (BIGINT).

- **Customer Name**: Customer's name (VARCHAR).

- **Merchant Name**: Merchant involved (VARCHAR, e.g., Flipkart, Zomato).

- **Transaction Date**: Date and time (DATETIME).

- **Transaction Amount (INR)**: Amount in Indian Rupees (DECIMAL).

- **Fraud Risk**: Risk level (VARCHAR, e.g., Low, Medium, High, Critical).

- **Fraud Type**: Type of fraud (VARCHAR, e.g., Card Skimming, Phishing).

- **State**: Indian state (VARCHAR, e.g., Rajasthan, West Bengal).

- **Card Type**: Card type (VARCHAR, e.g., Visa, Mastercard).

- **Bank**: Issuing bank (VARCHAR).

- **IsFraud**: Fraud flag (TINYINT, 0=Genuine, 1=Fraud).

- **Fraud Score**: Score indicating fraud likelihood (INT).

- **Transaction Category**: Category (VARCHAR, e.g., Electronics, Apparel).

- **Merchant Location**: Location (VARCHAR).

Data was loaded into a MySQL table for querying.

## 2.2   Tools and Technologies

- **SQL (MySQL)**: For data ingestion, exploration, and advanced queries (e.g., window functions for rolling averages).

- **Python**: For cleaning (Pandas), EDA (Seaborn, Matplotlib), feature engineering, and ML (Scikit-learn: Logistic Regression, Random Forest).

- **Power BI**: For dashboard creation, including KPIs, charts, and slicers.

- Environment: Python 3.x with libraries like NumPy, Pandas, Scikit-learn.

## 2.3 Data Preparation

1. **SQL Table Creation and Loading**:

   - Created a `transactions` table.
   - Loaded data using `LOAD DATA INFILE`, handling Windows line terminators (`\r\n`).
   - Checked for duplicates and nulls (e.g., no nulls in key fields like amounts).

2. **Python Cleaning**:

   - Loaded CSV into Pandas DataFrame.
   - Converted dates to datetime, amounts and scores to numeric.
   - Filled missing values (e.g., median for amounts/scores, 0 for IsFraud).
   - Dropped duplicates by Transaction ID.
   - Added features: Log-transformed amounts, hour/day/month/weekday from dates, high-value flag (>100,000), amount-score interaction.
   - Encoded categoricals using LabelEncoder.

3. **Data Quality Checks**:

   - Total transactions: Queried via SQL.
   - Unique counts: Merchants, categories, states, etc.
   - Distribution of Fraud Score and duplicates check.

## 2.4 Analysis Approach

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

   - Visualized fraud vs. non-fraud counts, boxplots for amounts by fraud status, histograms for distributions.
   - Grouped fraud rates by category, state, card type, bank.
   - Time-based: Fraud by hour, day of week.

2. **Aggregate Analysis**:

   - Fraud counts and percentages by category, state, merchant.
   - Top 10 merchants/states with highest fraud.

3. **Advanced SQL**:

   - Rolling averages per state.
   - Flagged high-amount transactions (>3x state average).

4. **Machine Learning**:

   - Split data (80/20, stratified).

- Models: Logistic Regression (baseline, balanced weights), Random Forest ($n_estimators = 200, balanced$). $Metrics: Confusion matrix, classification report, ROC - AUC, Precision - Recall curve.$

- Handling imbalance: Class weights; optional SMOTE.
- Feature importance from Random Forest.
- Threshold tuning to maximize recall (e.g., min precision 0.10).

  5. **Power BI Dashboard**:

     - **Pages**:
       - **Overview**: KPIs (total txns, frauds, avg amount), monthly trends line chart.
       - **Geo & Demographics**: State fraud distribution bar/map, fraud by age (if derived).
       - **Investigation**: Table of flagged transactions with conditional formatting.
     - **Visuals**: Pie for merchant share, bar for fraud type vs. category.
     - **DAX Measures**: Total Txns, Fraud Rate, Precision/Recall (if scored data imported).
     - Exports: CSVs for fraud by state/category, summary stats.

# 3   Results and Findings

## 3.1   Key Insights from Dashboard and Analysis

- **Transaction Trends**: Volumes dipped in February, peaked at 1.13M in May, sharp decline to 625K in September, recovery to 1.28M in December (highest yearly).

- **Fraud by State** (from bar chart):

| State | Medium | Low | Critical | High | Total Incidents |
|---|---|---|---|---|---|
| Rajasthan | 42 | 44 | 14 | 13 | 113 |
| Maharashtra | 30 | 39 | 23 | 15 | 107 |
| West Bengal | 28 | 45 | 25 | 11 | 109 |
| Tamil Nadu | 34 | 45 | 18 | 22 | 119 |
| Uttar Pradesh | 26 | 44 | 19 | 13 | 102 |
| Telangana | 28 | 36 | 15 | 19 | 98 |
| Karnataka | 25 | 29 | 22 | 13 | 89 |
| Kerala | 38 | 36 | 14 | 19 | 107 |

Rajasthan leads in total frauds, West Bengal in high-risk, Kerala in medium-risk.

- **Fraud Type vs. Category**:

  - Card Not Present: 40 in Transportation.
  - Card Skimming: 52 in Electronics.

  - – Identity Theft: 35 in E-commerce.
  - – Account Takeover: 32 in Food Delivery.
  - – Phishing: 31 in Apparel, 23 in Groceries.

- **Merchant Share** (Pie Chart):

  - – Tata Cliq: 166
  - – Flipkart: 105
  - – Big Bazaar: 105
  - – Zomato: 107
  - – Myntra: 97
  - – Lifestyle: 121
  - – Other: 366

- **EDA Visuals**:

  - – Fraud percentage: Typically low ( 1-5%), but imbalanced.
  - – Higher amounts correlate with fraud (boxplot shows outliers).
  - – Fraud ratios highest in Electronics and E-commerce categories.
  - – Peaks in fraud during late hours (SQL query: highest by hour).

- **ML Performance**:

  - – Logistic Regression: ROC-AUC  0.85, balanced for recall.
  - – Random Forest: ROC-AUC >0.90, better precision-recall.
  - – Top Features: Fraud Score, Amount, Hour, Merchant Name (from importance plot).
  - – Tuned threshold: Achieved  90% recall with  15-20% precision.

- **Anomaly Detection**:

  - – Flagged  10-20% of transactions as high-risk using SQL (e.g., >3x avg amount).
  - – Cumulative frauds over time show seasonal spikes (October-November recovery).

## 3.2   Anomaly Detection

Flagged  10-20% of transactions as high-risk using SQL (e.g., >3x avg amount). Cumulative frauds over time show seasonal spikes (October-November recovery).

# 4   Discussion

The analysis reveals geographic hotspots (Rajasthan, West Bengal) and category vulnerabilities (Electronics), likely due to online/high-value nature. ML models perform well but could improve with more features (e.g., velocity). The "err on caution" approach aligns with business needs, accepting more alerts to catch frauds.

Challenges: Imbalanced data (addressed via weights), potential overfitting (mitigated by cross-validation, CV ROC-AUC  0.88).

# 5   Conclusion

This project successfully builds a fraud risk analysis framework, from data ingestion to predictive dashboarding. It highlights actionable insights for reducing fraud through targeted monitoring.

# 6   Recommendations

- Implement real-time alerts for high-score transactions.

- Enforce stricter rules in high-risk states/categories (e.g., OTP for Electronics >50K).

- Integrate API for ongoing model updates.

- Pilot A/B testing for threshold impacts on operations.

- Expand dataset with customer demographics for better age-based analysis.

# 7   Appendices

- SQL Queries: Included in methodology.

- Python Code: Full script for cleaning, EDA, ML.

- Power BI Layout: As described, with suggested visuals.