

StudentName:KAZIMAHMED .T.M

Register Number:511923205025

Institution: PEC

Department:B TECH-IT

Date of Submission:05\05\2025

GitHub Repository Link:<https://github.com/kazimahmed4u/team-5-project.git>

1. PROBLEM STATEMENT

Dynamically adapt to new and evolving fraud techniques, ensuring continuous protection against emerging threats that traditional systems fail to recognize.

- **Provide a holistic fraud management solution** encompassing real-time transaction analysis, risk scoring, immediate alerts for suspicious activities, and proactive prevention mechanisms.
- **Minimize disruptions to legitimate cardholders** by drastically reducing false positives and ensuring a seamless and trustworthy transaction experience.
- **Effectively address the inherent class imbalance** in transaction data, where fraudulent transactions are a small minority compared to legitimate ones, without compromising detection of the minority class.
-

2. PROJECT OBJECTIVES

i. Develop a highly accurate AI-powered fraud detection model:

- Achieve a significant improvement in fraud detection accuracy compared to traditional rule-based systems, aiming for a substantial reduction in both false positive and false negative rates.
- (e.g., precision, recall, F1-score, AUC)

ii. Build an adaptive system capable of learning and evolving:

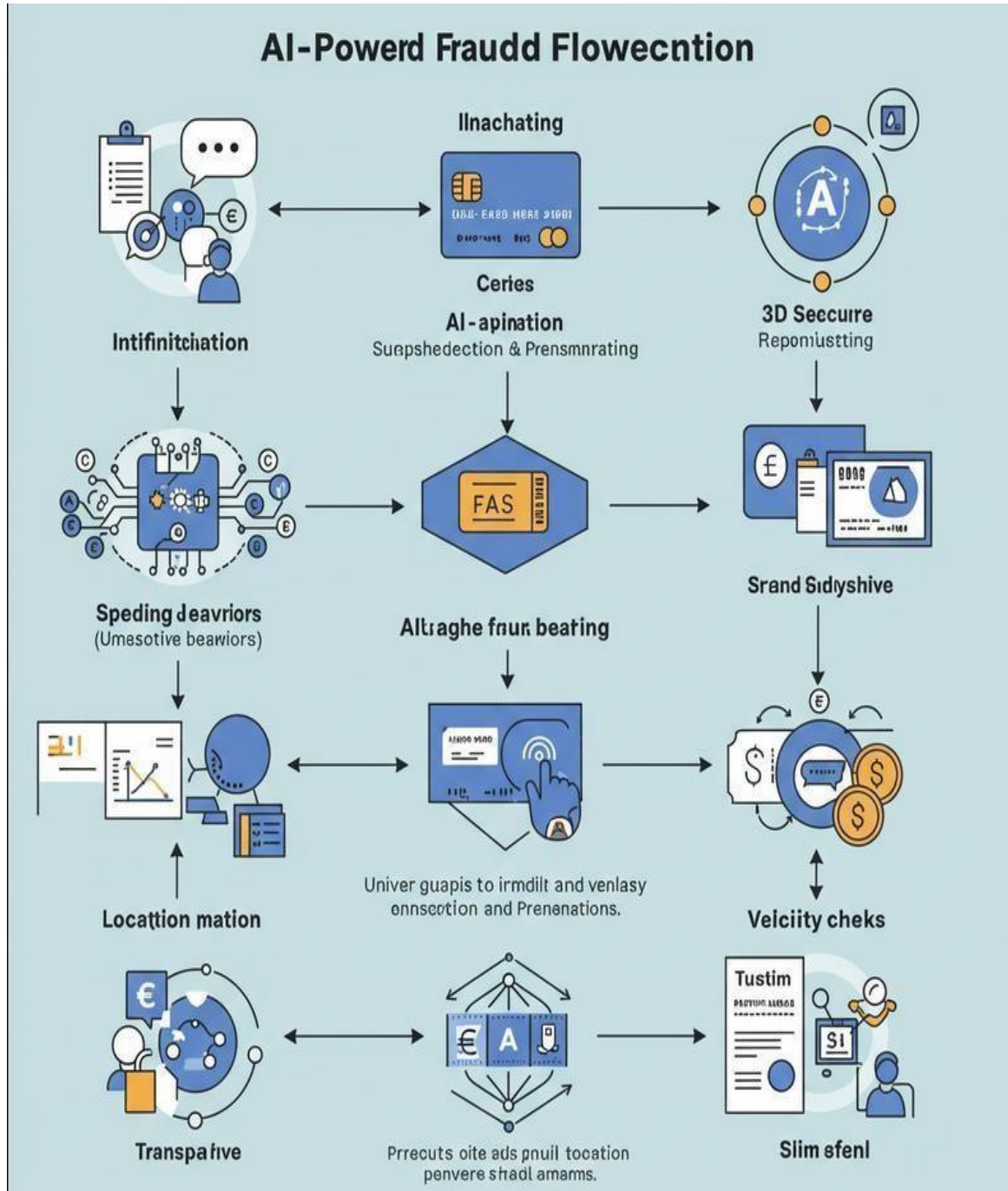
- Implement mechanisms for continuous learning and model retraining using new transaction data to ensure the system adapts to emerging fraud patterns and techniques.
- Incorporate feedback loops to refine the AI models based on the outcomes of transaction investigations and fraud confirmations.

iii. Achieve scalability and high performance:

- Design the system architecture to handle the increasing volume of credit card transactions efficiently and with low latency.
- Ensure the system can scale to accommodate future growth in transaction volume.

3. FLOWCHART OF THE WORKFLOW

4. DATA



DESCRIPTION

Transaction Data: This is the most crucial data source and includes details of every credit card transaction, such as:

- **Amount:** The monetary value of the transaction. For example, a sudden increase in transaction amount compared to the cardholder's usual spending habits can be a red flag.

- **Timestamp:** The exact time and date of the transaction. Multiple transactions occurring within a very short time frame but from geographically distant locations can be suspicious.
- **Merchant Information:** Details about the seller or service provider. Transactions with high-risk or blacklisted merchants can be flagged.

II. Data Processing and Feature Engineering:

- **Big Data Processing Frameworks:**
 - **Apache Spark:** A powerful and widely used distributed computing framework for large-scale data processing, feature engineering, and machine learning.
 - **Apache Hadoop (HDFS, MapReduce, YARN):** While Spark is often preferred for its speed, Hadoop remains relevant for distributed storage and batch processing of massive datasets.
 - **Apache Flink:** A stream processing framework ideal for real-time feature engineering and analysis of streaming transaction data.

The dataset used in this project comprises historical transaction records, including both legitimate and fraudulent instances. Key features include:

- **Transaction ID:** Unique identifier for each transaction
- **Timestamp:** Date and time of the transaction
- **Customer ID:** Anonymized customer identifier
- **Credit Details:** Information such as transaction amount, credit score, credit limit, and outstanding balance
- **Product Information:** Type of product purchased, product ID, category, price, and merchant details
- **Location Data:** Transaction origin, such as city, region, or IP address
- **Device and Channel Info:** Device used, browser type, platform (e.g., mobile, web), and payment method

- **Fraud Label:** Binary indicator (1 = fraudulent, 0 = legitimate)

This data enables the training and testing of AI models to accurately detect anomalies and guard against fraud by learning patterns associated with high-risk transactions across both credit-related and product-level behaviors.

5. DATAPREPROCESSING

1. Data Collection and Ingestion:

- **Real-time Streaming:** Transaction data is captured in real-time as it occurs through various channels (e.g., online purchases, in-store transactions, ATM withdrawals). This often involves streaming data platforms like Apache Kafka or Amazon Kinesis for efficient and low-latency data transfer

2. Data Preprocessing:

- **Data Cleaning:** Identifying and handling missing values (imputation using mean, median, or more advanced techniques), correcting errors, and removing irrelevant or duplicate data.

3. Data Analysis and Modeling:

- **Exploratory Data Analysis (EDA):** Analyzing the preprocessed data to understand its characteristics, identify patterns, and gain insights into potential fraud indicators. This involves visualizations (histograms, scatter plots), statistical summaries, and correlation analysis.

6. EXPLORATORY DATA ANALYSIS(EDA)

Exploratory Data Analysis (EDA) is the foundational step in building effective AI-powered credit card fraud detection systems.

High-Level Overview: Examine the dataset's structure:

- Number of rows (transactions) and columns (features).
- Data types of each column (e.g., integer, float, object/string, datetime). This helps identify potential data type issues that need to be addressed.
- Initial rows of the data to get a feel for the content and format.

Geospatial Analysis (If Location Data is Available):

- **Mapping Transactions:** Visualize transaction locations on a map. Unusual transaction locations for a cardholder can be a strong indicator of fraud.
- **Identify Potential Fraud Indicators:** Discover features or combinations of features that seem to differ significantly between legitimate and fraudulent transactions.
- **Formulate Hypotheses:** Develop initial ideas about which features might be most predictive of fraud.

To understand the underlying patterns and anomalies within the dataset, we conducted an in-depth exploratory data analysis focused on transaction behavior, credit attributes, and product-related insights:

Missing Values and Data Types:

- Checked for null or inconsistent entries in critical fields like `credit_score`, `product_category`, and `transaction_amount`.
- Ensured appropriate data types for numerical (e.g., amount, credit limit) and categorical (e.g., product type, device) features.

2. Fraud Distribution:

- Examined class imbalance in the `fraud_label` column.
- Found that fraudulent transactions account for a small fraction of the data, requiring strategies like oversampling or SMOTE for model training.

3. Transaction Amount Patterns:

- Visualized transaction amounts by fraud status.
- Detected higher variability and unusual spikes in amounts associated with fraudulent cases.

4. Credit Behavior Insights:

- Analyzed the relationship between credit score, credit limit, and fraud probability.
- Observed that lower credit scores and near-limit usage often correlated with fraudulent transactions.

5. Product and Category Trends:

-
-
- Identified specific product categories (e.g., electronics or high-value items) with a higher fraud ratio.
- Explored merchant behavior and frequency of returns or cancellations.

6. Temporal and Geographic Trends:

- Time-series plots revealed fraud peaks during non-business hours.
- Geographic heatmaps indicated clusters of fraudulent activity in specific regions or IP zones.

7. Device and Access Channel Analysis:

- Compared fraud rates across devices (mobile vs. desktop) and browsers.
- Noted that unfamiliar or spoofed devices contributed disproportionately to fraud.

8. Correlation Matrix:

- Computed pairwise correlations to identify significant relationships between features like `credit_utilization` and `transaction_risk_score`.

7. FEATURE ENGINEERING

Transaction-Based Features: These features are derived directly from the attributes of individual transactions.

- **Basic Transaction Attributes:**

- **Transaction Amount:** The raw amount of the transaction. Often, fraudulent transactions have different amount distributions (e.g., tend to be higher or lower than typical)
- **Cardholder-Based Feature**Geographic distribution of fraud.
- **s:** these features aggregate information about the cardholder's past behavior.

- **Historical Spending Patterns:**

- **Average Transaction Amount (over different time windows):** Captures the typical spending habit.
- **Standard Deviation of Transaction Amounts:** Measures the variability in spending.
- **Most Frequent Transaction Amounts/Merchants:** Establishes the "normal" behavior.

- **Recency of Activity:**

- **Time Since Last Login to Online Banking:** Unusual inactivity followed by transactions could be suspicious.
- **Time Since Card Activation:** Newly activated cards might be more vulnerable.

8. MODEL BUILDING

I. Model Selection

Supervised Learning Models

Logistic Regression

Random Forests

Decision Trees

II. Data Preparation for Model Training

Training set

Validation set

Testing set

II. Data Processing and Feature Engineering:

- **Big Data Processing Frameworks:**

- **Apache Spark:** A powerful and widely used distributed computing framework for large-scale data processing, feature engineering, and machine learning.
-
- **Apache Hadoop (HDFS, MapReduce, YARN):** While Spark is often preferred for its speed, Hadoop remains relevant for distributed storage and batch processing of massive datasets.
- **Apache Flink:** A stream processing framework ideal for real-time feature engineering and analysis of streaming transaction data.

III. Model Evaluation

Confusion Matrix

Precision

Recall (Sensitivity)

Average Precision (AP)

9. VISUALIZATION OF RESULTS & MODEL INSIGHTS

I. Visualizing Model Performance:

Heatmaps: Display the counts or percentages of True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN) in a visually intuitive way. Different color intensities can represent the magnitude of each category. This helps quickly assess the types of errors the model is making.

II. Visualizing Model Insights (Understanding What the Model Learned):

Bar Charts or Dot Plots: Show the relative importance of different features as determined by the model (e.g., using techniques like Gini importance for tree-based models, coefficient magnitudes for linear models, or permutation importance). This helps understand which factors the AI considers most influential in predicting fraud

III. Interactive Visualizations and Dashboards:

Geographic distribution of fraud.

- False positive rate.
- Number of blocked fraudulent transactions.
- Trends in fraud patterns over time.

10. TOOLS AND TECHNOLOGIES USED

I. Data Storage and Management:

- **Databases:**

- **Relational Databases (SQL):** PostgreSQL, MySQL, Oracle, Microsoft SQL Server are used for structured transaction data, cardholder information, and historical records.
- **NoSQL Databases:** Cassandra, MongoDB, HBase are employed for handling large volumes of unstructured or semi-structured data, real-time data streams, and scalability.
- **Data Warehouses:** Amazon Redshift, Google BigQuery, Snowflake are used for storing and analyzing large historical datasets for model training and batch processing.

II. Data Processing and Feature Engineering:

- **Big Data Processing Frameworks:**

- **Apache Spark:** A powerful and widely used distributed computing framework for large-scale data processing, feature engineering, and machine learning.
- **Apache Hadoop (HDFS, MapReduce, YARN):** While Spark is often preferred for its speed, Hadoop remains relevant for distributed storage and batch processing of massive datasets.
- **Apache Flink:** A stream processing framework ideal for real-time feature engineering and analysis of streaming transaction data.

1. Data Cleaning and Preprocessing

- **Handling Missing Values:**

Imputed missing `credit_score` and `product_category` using median/mode imputation or predictive modeling.

- **Data Type Conversion:**

Converted timestamps to datetime objects and encoded categorical features like `device_type`, `product_category`, and `payment_method`.

- **Outlier Detection:**

Applied IQR and z-score methods to detect and optionally cap extreme values in `transaction_amount`, `credit_utilization`, etc.

2. Feature Engineering

Credit-Based Features

- **Credit Utilization Ratio:**
 $\text{'creditutilization'} = \text{'transactionamount'} / \text{'creditlimit'}$
High utilization often signals risky behavior.
- **Credit Score Binning:**
Grouped credit_score into ranges (e.g., poor, average, good) to simplify analysis.

Product-Based Features

- **High-Risk Product Flag:**
Created a binary feature to flag purchases in categories frequently linked to fraud (e.g., electronics, luxury items).
- **Product Price Deviation:**
Compared current transaction price to a rolling mean of previous purchases by the same user to detect anomalies.

Transaction Behavior Features

- **Transaction Frequency:**
Number of transactions by a user within a certain time window (e.g., last hour, day).
- **Time-Based Features:**
Extracted hour, day, and weekday from timestamp to capture off-hour or holiday activity.
- **Geo-Distance:**
Calculated distance between current and previous transaction locations for the same user to flag unusual access patterns.

Device and Channel Features

- **New Device or Location Indicator:**
Boolean flags for whether the transaction came from a new device or location.
- **Channel Risk Score:**
Assigned weights to transaction channels (e.g., web, mobile) based on historical fraud likelihood.

11. TEAM MEMBERS AND CONTRIBUTION

Data Cleaning : Mokesh S

EDA: kazhim ahmed A

MODEL development : Lakshmanan

Documentation: KAZIM AHMED