**Project: Fraud Detection Analysis**
**Group Members:** Aaron Marshall, Ashleigh Clark
**Professor Monadjemi**
**December 3rd, 2024**

## Abstract

Fraud detection in digital wallet transactions is a growing concern as digital payment systems become increasingly widespread. This project leverages a synthetic dataset to explore fraud detection methodologies, focusing on creating and analyzing fraud indicators. We engineered features to identify potentially fraudulent transactions based on transaction amount, geographic anomalies, payment method irregularities, and frequency patterns. Our analysis uncovered distinct trends, such as rapid transaction bursts and cross-regional usage, that align with known fraud behaviors. This report summarizes our approach, findings, and recommendations for developing effective fraud detection systems.

## Introduction

Digital payment platforms have revolutionized financial transactions but are increasingly targeted by fraudsters. Fraudulent activities not only erode consumer trust but also result in significant financial losses for businesses. This project aims to develop a data-driven framework to identify suspicious transactions and provide actionable insights for mitigating fraud risks. Using a synthetic digital wallet transaction dataset, we designed five fraud indicators, each targeting a specific fraud behavior. Exploratory data analysis highlighted significant patterns, such as anomalies in transaction locations and excessive use of multiple payment methods.

These indicators were crafted based on the specific qualities of this dataset but can be extended to similar datasets including real-world datasets. The design of each indicator relies on a general idea of potential evidence of fraud (transactions that have very high values, transactions from atypical locations, frequent transactions made in a very short time, etc.). Moreover, given that the general principles used in the creation of indicators are relevant and evident in various digital transactions datasets, the framework of our indicators remains effective with those other datasets.

Through this project, we aim to empower businesses with tools to preemptively address fraud, ensuring secure and reliable digital transactions. Our findings also provide a foundation for refining fraud detection systems based on real-world scenarios. Combining various tools and methods including machine learning, these fraud indicators can be adjusted to improve effectiveness and reliability. Consequently, future work can include the application of our indicators to real-world datasets and the implementation of machine learning models or algorithms to enhance the indicators improving the accuracy and reliability of the indicators.

## Background

Fraud detection involves identifying transactions that deviate from typical behavioral patterns. In digital payments, fraud often manifests through:

- **Transaction anomalies**: Extremely high-value transactions relative to the norm.
- **Geographic inconsistencies**: Purchases from atypical locations.
- **Behavioral irregularities**: Rapid-fire transactions or unusual payment methods.

Our dataset includes detailed transaction records, such as timestamps, locations, and payment methods, allowing for a comprehensive fraud analysis. These domain-specific features are crucial in distinguishing legitimate transactions from fraudulent ones.

## Analytic Goals

This project aims to answer the following questions:
1. What transaction patterns are indicative of fraudulent behavior?
2. Can transaction-level data accurately flag high-risk transactions?
3. What insights can be drawn to refine fraud detection algorithms?

By addressing these questions, we aim to propose actionable steps for fraud mitigation.

## Data

**Source:** Synthetic dataset from Kaggle, comprising digital wallet transactions.
**Size:** Approximately 250,000 records with attributes such as transaction amount, timestamps, locations, and payment methods.

Key preprocessing steps included:

- Handling missing values and duplicates.
- Normalizing categorical data for consistency.
- Engineering features such as high transaction flags and unusual location flags.

## Methods

1. **Data Cleaning:** Addressed null values, removed duplicates, and standardized formats.
2. **Feature Engineering:** Created five fraud indicators:
   a. High Transaction Amounts
   b. Frequent Transactions in Short Time Periods
   c. Unusual Locations
   d. Irregular Product Categories
   e. Multiple Payment Methods
3. **Aggregation:** Consolidated user-level data to derive behavioral baselines.
4. **Analysis:** Applied statistical thresholds to flag anomalies and identify trends.
5. **Evaluation:** Iteratively refined fraud indicators based on exploratory analysis results.
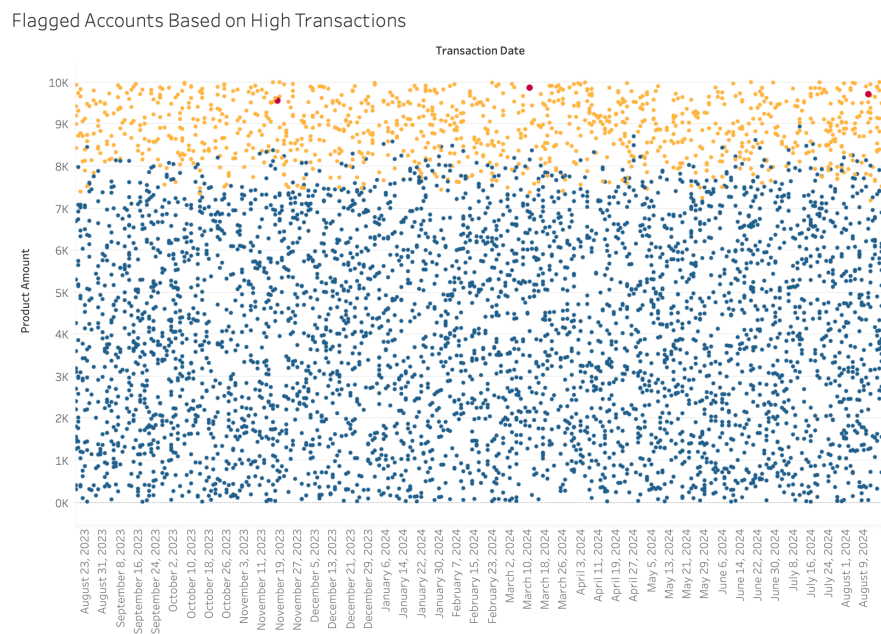
## Findings

1. **High Transaction Amounts:**
   a. Transactions exceeding three standard deviations from the mean flagged a significant portion of high-risk entries.
   b. Adjusting the mean and standard deviation to account for averages influenced by product type, the same 3-sigma rule produced 0 flagged transactions.
   c. 2 and 1 standard deviation thresholds given the adjusted mean and standard deviations were more descriptive.
2. **Frequent Transactions:** Several users exhibited rapid transaction bursts, indicative of potential automation.
3. **Unusual Locations:** Cross-regional transactions often aligned with flagged high-value purchases.

4. **Multiple Payment Methods:** Users utilizing three or more methods showed a higher likelihood of anomalous activity.
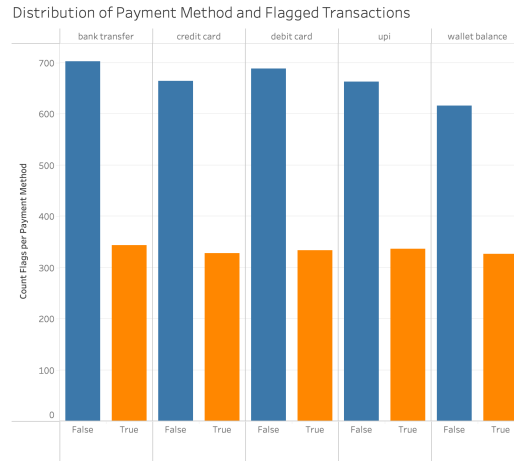
**Visualizations and Key Statistics:**

● High transaction flaggs using adjusted mean and standard deviations for different sigma-values (thresholds)
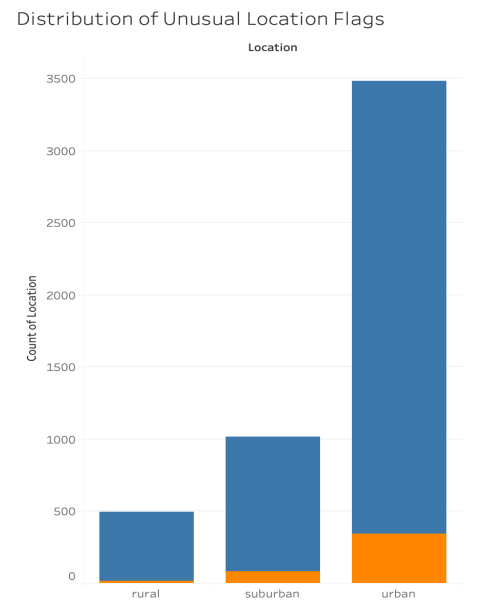


Flagged Accounts Based on High Transactions

This scatter plot shows the distribution of which of the transactions were flagged with which threshold. The red data points represent the 3 transactions flagged with a 2-sigma threshold, the yellow represents the 1032 flagged with a 1-sigma threshold, and the blue represents the remaining unflagged transactions. Given the adjusted mean and standard deviations, the 3 flagged with the highest threshold were all above 9,000 in price, but were not the most expensive transactions in the dataset.

● Distribution of the number of transactions flagged using multiple payment method indicators based on each payment method

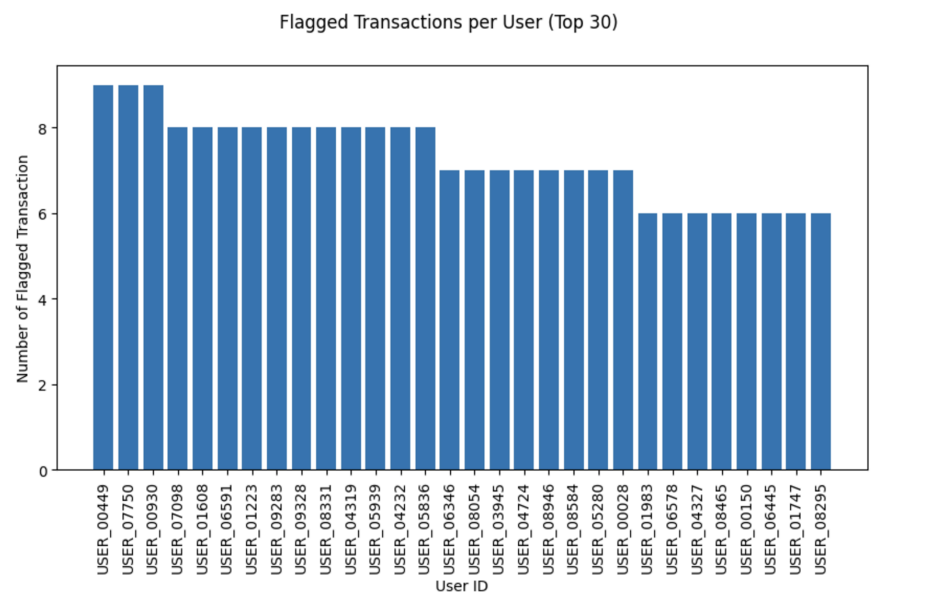Distribution of Payment Method and Flagged Transactions

The bar charts represented the total number of transactions with each payment method (the separation of the colors represents the distinction between transactions that were flagged or not). Specifically, orange represents the flagged transactions, and the orange transactions' distribution shows that it was fairly consistent across the payment method categories. Furthermore, the sum of the orange is precisely 1666 corresponding to the 1666 transactions flagged with this indicator.

● Distribution of the number of transactions flagged using the unusual location indicator separated based on different location options



Distribution of Unusual Location Flags

These bar charts show the distribution of flagged transactions based on location categories. The height of the bars demonstrates the trend of the distribution of transactions over the location categories is similar to the trend of the flagged transactions. This suggests that the specific location was not necessarily more likely to be flagged than others. Ultimately, however, this indicator flagged 447 transactions for unusual locations representing the sum of the height of the orange blocks in the visualization.

- Distribution of the users who were most often flagged



Flagged Transactions per User (Top 30)

This bar chart represents the count of how many times each user was flagged based on our 5 indicators. Since there were thousands of users, the visualization is limited to the top 30 to determine which users were most commonly flagged. Ultimately, given the overall since of the data frame and the size of each user's transaction history, the variability between users was not too high. Therefore it is hard to say that there were any significant outliers that would suggest a stronger link to fraudulent activity.

## Discussion and Conclusions

**Discussion:**

Our findings support the efficacy of data-driven fraud indicators. By automating fraud detection

processes, businesses can proactively safeguard against financial risks. However, the reliance on synthetic data limits the generalizability of results. Future work should include real-world datasets and machine-learning models to improve precision and scalability.

To elaborate, the effectiveness of some of the indicators was limited to the size and reliability of the data. In particular, the indicators that focused on each user's transaction history were more limited since the size of this dataset did not provide as extensive of history as a real digital transaction dataset may provide over the same time frame. Furthermore, the structure of some of the columns such as location did not always reflect what would be seen in a realistic dataset (potential longitude and latitude coordinates instead of suburban/urban/rural) confining the possibilities of the analysis of the location indicator. However, our indicator based on unusual locations can easily be extended and adjusted to account for more detailed location data by establishing a typical location area and determining when users made transactions outside of their typical location area.

Consequently, our indicators offer a nice potential framework and foundation for fraud analysis, yet most likely would need to be adjusted to adapt to the different characteristics and qualities of other transaction datasets. Ultimately, context and background information can be instrumental in the construction of similar indicators for different real-world datasets. Additionally, machine-learning models can also help improve the accuracy and efficacy of the indicators.

**Conclusions:**
This project demonstrates the utility of engineered features in fraud detection. By identifying high-risk behaviors, our framework provides a robust starting point for developing comprehensive fraud mitigation systems. Furthermore, the structure of many of the indicators allows for customization based on specific characteristics or qualities of various datasets allowing easier adaption to real-world datasets.