

# 1. Introduction

## 1.1 Problem Statement

This project aims to analyze how different transaction factors influence the likelihood of fraud.

## 1.2 Dataset

### 1.2.1 Dataset Source:

<https://github.com/VineetKumar711/creditcardfraud/tree/main>

### 1.2.2 Dataset Description:

The Credit Card Fraud Dataset is a publicly available dataset containing transactional records spanning seven months, from June to December. It consists of multiple CSV files representing monthly transaction data, along with a separate file identifying fraudulent transactions. The dataset captures a wide range of attributes that provide valuable insights into financial activities, including transaction details (timestamp, transaction amount, merchant information, transaction category, and unique transaction identifiers), cardholder demographics (name, gender, address, city, state, and date of birth), and geospatial data (latitude and longitude of both the cardholder and merchant). Additionally, it includes city population and cardholder job title, which may contribute to fraud risk assessment. The fraud dataset specifically contains an `is_fraud` column, indicating whether a transaction is fraudulent (1) or legitimate (0). By merging the transaction and fraud datasets, we obtain a comprehensive dataset with 555,719 transactions and 23 attributes, enabling an in-depth analysis of transaction patterns and fraud detection.

## 1.3 Tools

To analyze our dataset and generate our report, we utilized the following tools and techniques:

- Unzip tools were used to extract one of the transaction data from a compressed .zip file.
- Python was used for data cleaning, exploratory data analysis, data preprocessing and feature engineering.

## 1.4 Why Analyze this Dataset

Analyzing this dataset is essential for understanding transaction trends and identifying key factors associated with fraudulent activity. Through data cleaning, exploratory data analysis (EDA), data preprocessing, and feature engineering, we can uncover patterns in transaction behavior and detect significant trends. EDA allows us to visualize transaction distributions, spot anomalies, and identify relationships between variables, while feature engineering helps refine the dataset for further modeling. By focusing on these steps, we aim to extract meaningful insights that can guide future fraud detection efforts and enhance data-driven decision-making.

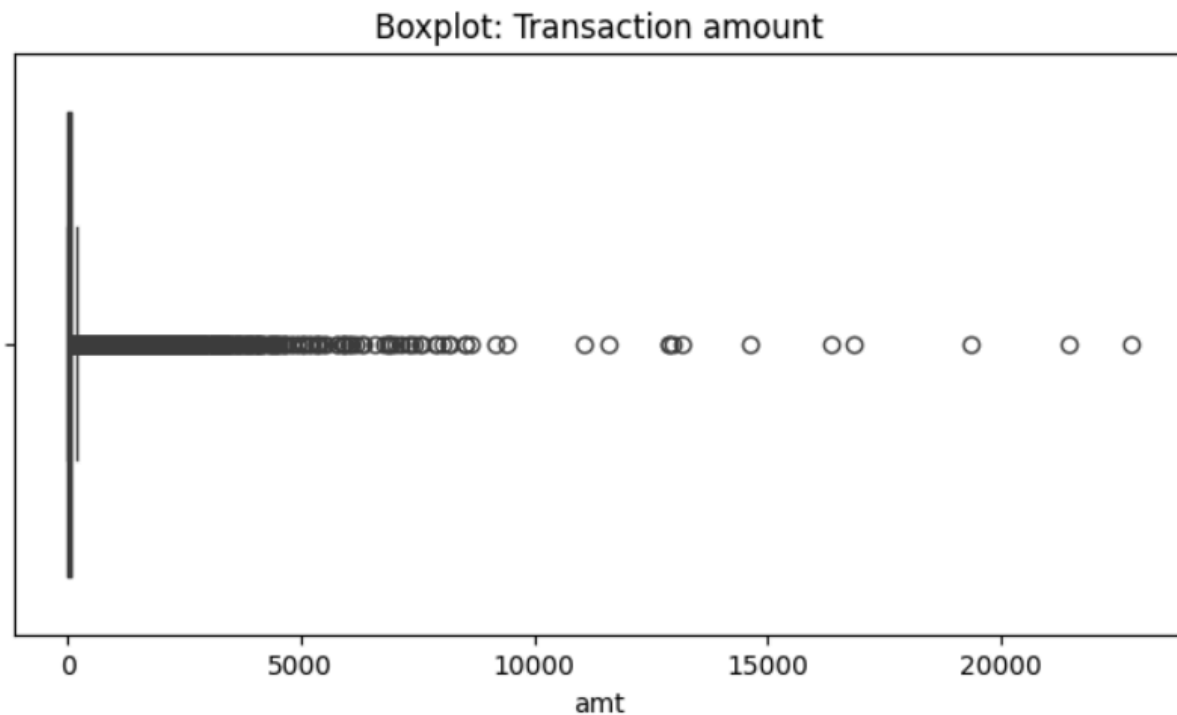
# 2. Objectives

## 2.1 Data Cleaning and Handling Inconsistencies

Firstly, a fraud dataset and six individual monthly transaction datasets from June to December are loaded individually, with each dataset containing transactional records. These datasets are then combined into a single unified dataset through concatenation. A merging operation is performed with the fraud dataset using a left join on the transaction ID, incorporating an indicator variable that designates whether each transaction is fraudulent or legitimate. Then, the shape of the merged data frame is (555719, 23) which means that there are 555719 rows and 23 columns.

A series of data cleaning steps were performed to ensure data integrity and consistency. First, duplicate records were checked and removed, but no duplicates were found. Next, missing values across multiple columns, including merchant, category, last name, street, and job, were identified and subsequently dropped to maintain data completeness. After dropping the missing value, the shape of the data frame becomes (555624,23). The dataset was then reset to ensure proper indexing. The number of unique values per column was analyzed to understand variability in attributes such as transaction ID, merchant, category, and customer demographics. State values were converted to uppercase, reducing the number of unique state entries from 96 to 50, ensuring consistency. Gender labels were standardized by replacing "Male" and "Female" with "M" and "F".

A box plot was generated to visualize transaction amounts, revealing the presence of outliers.



The box plot of transaction amounts reveals a right-skewed distribution with a significant number of high-value outliers, indicating that most transactions are of relatively small amounts while a few extend beyond \$5,000, \$10,000, and even \$20,000. The presence of numerous outliers suggests that large transactions are uncommon and may warrant further investigation, particularly in the context of fraud detection. Since fraudulent activities often involve unusually high transaction amounts, these outliers could serve as potential fraud indicators. Additionally, the interquartile range (IQR) method identified an upper bound of approximately \$266.96, suggesting that most transactions fall within this range, while those exceeding this threshold might be considered suspicious.

Therefore, given the unique characteristics of fraud detection, where fraudulent transactions often involve unusually high amounts, these outliers hold valuable information. Removing them could lead to the loss of critical fraud indicators, potentially reducing the effectiveness of exploratory data analysis and future fraud detection models. Thus, the decision was made not to remove these outliers, allowing for a more comprehensive analysis of transaction patterns and anomalies.

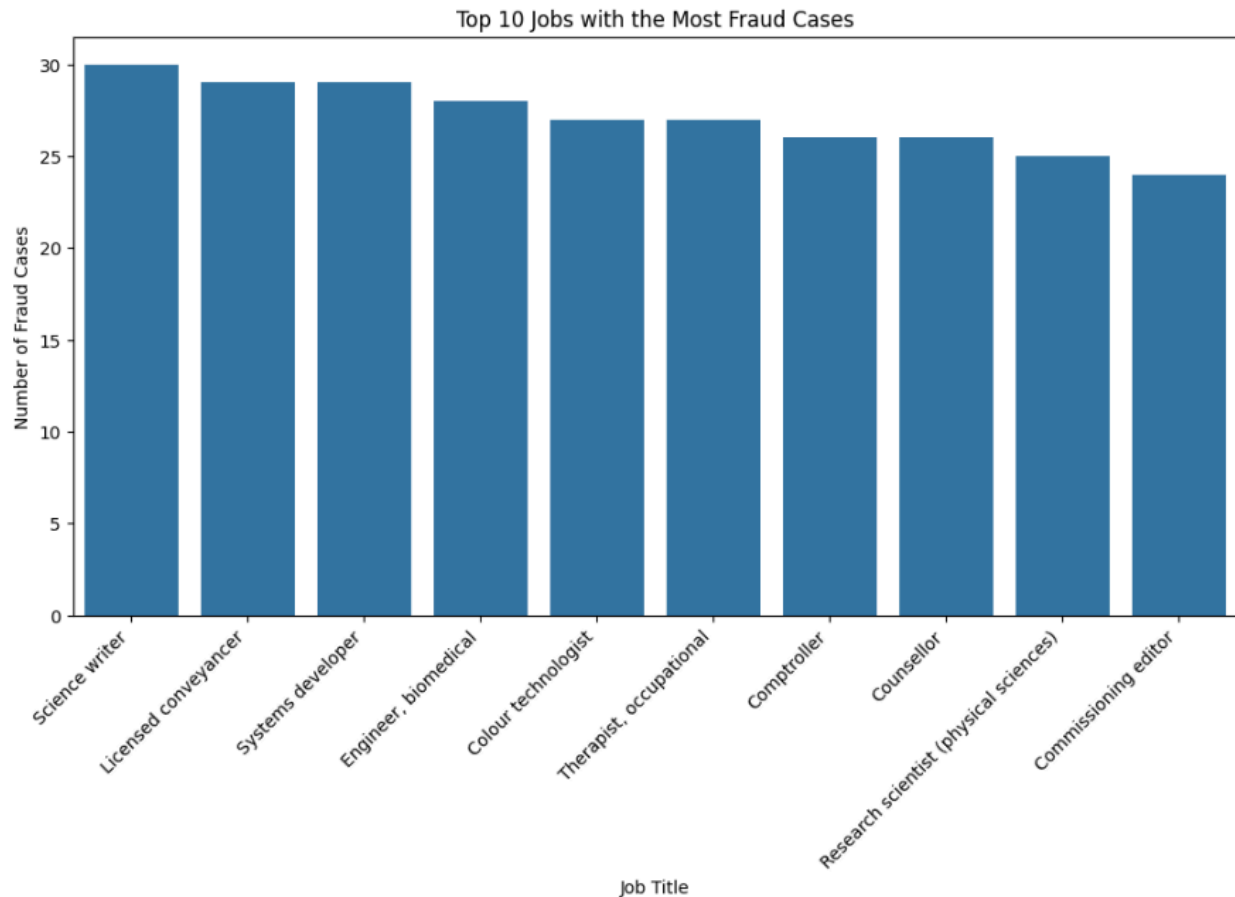
## **2.2 Exploratory Data Analysis**

Important details about the dataset's numerical characteristics are disclosed by the summary statistics. There is substantial variation in the transaction values, as seen by the transaction amount's (amt) large standard deviation of \$156.75 and mean of \$69.39. The existence of high-value transactions is further supported by the maximum transaction amount of \$22,768.11. Both cardholders' and merchants' latitude and longitude data indicate that transactions take place over a large geographic region. The wide range of city populations—from 23 to over 2.9 million—indicates that transactions take place in both urban and rural areas.

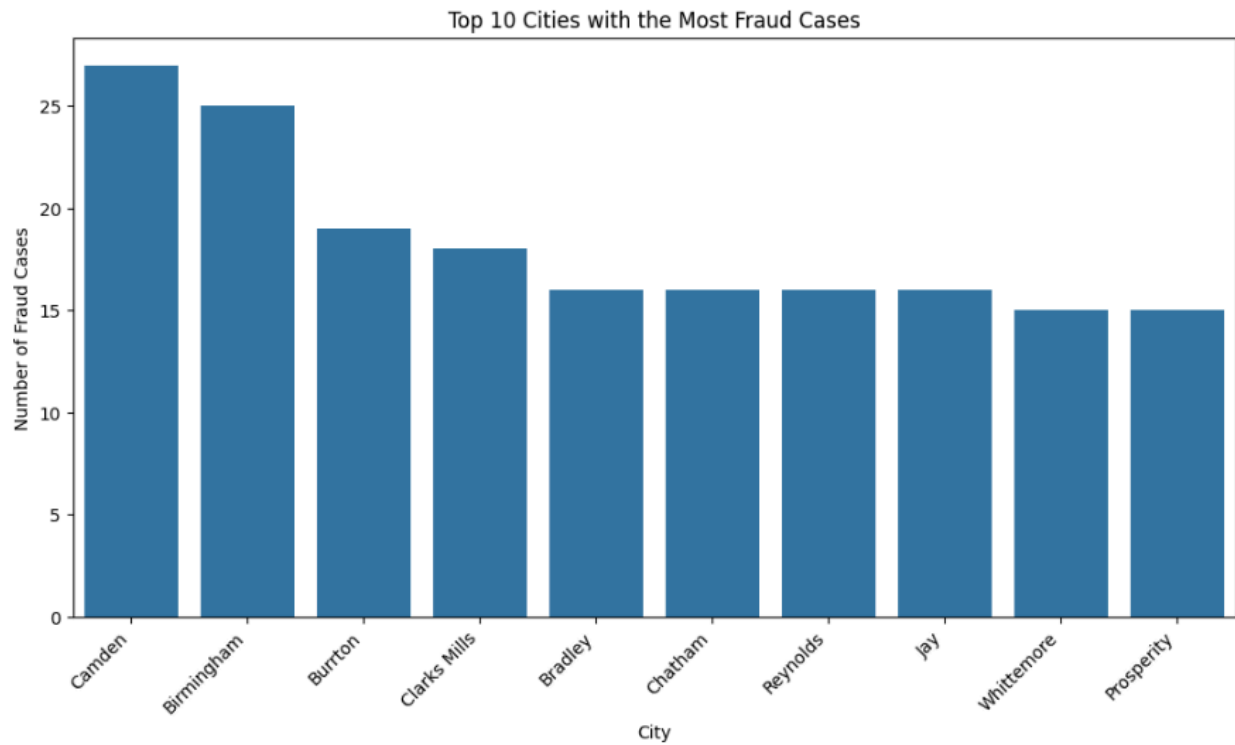
With 99.61% of transactions being genuine and just 0.39% being categorized as fraudulent, the class distribution of fraudulent transactions (is\_fraud) shows a significant class imbalance. This disparity implies that in order to enhance fraud classification performance, fraud detection models would need to employ strategies like oversampling, undersampling, or synthetic data synthesis. Reducing false negatives and guaranteeing efficient fraud detection depend on addressing this mismatch.



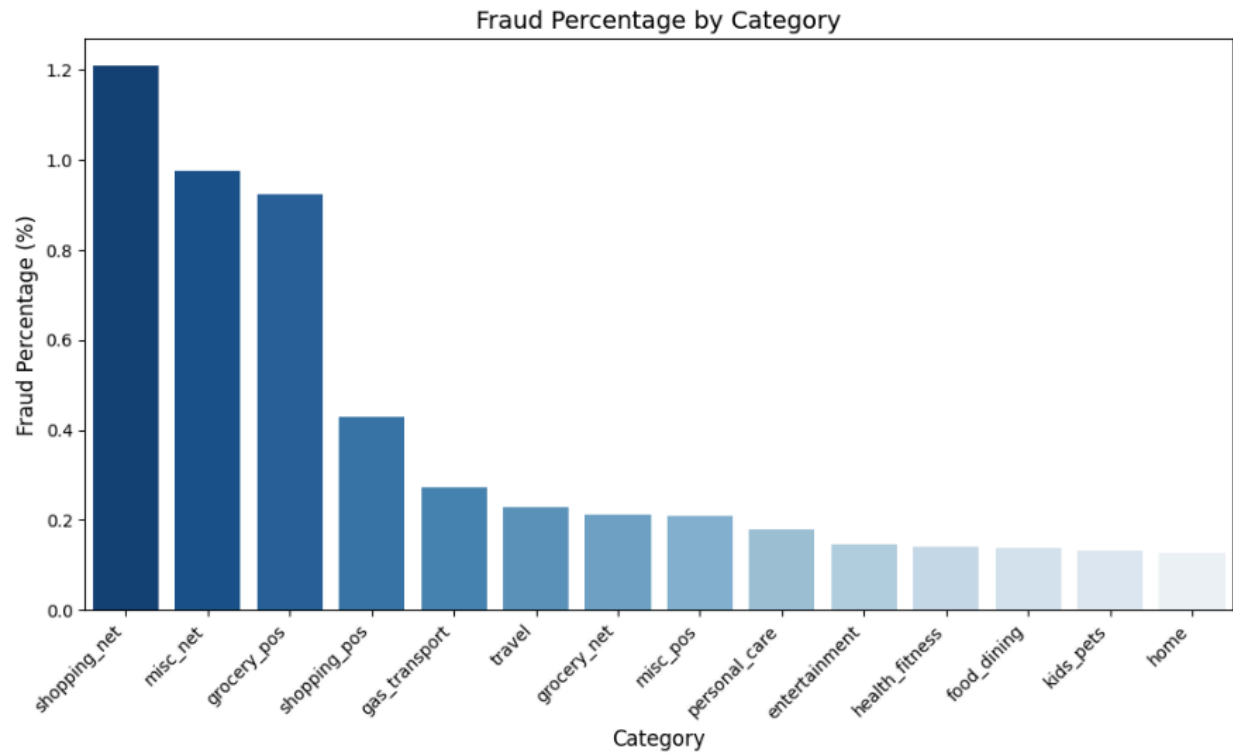
The pie charts illustrate the distribution of transactions by gender for both fraudulent and non-fraudulent transactions. In both cases, the proportions of female and male cardholders are relatively similar, with females accounting for approximately 54.9% of non-fraudulent transactions and 54.3% of fraudulent transactions, while males represent 45.1% and 45.7%, respectively. This suggests that fraudulent transactions do not appear to be strongly associated with gender, as the distribution remains consistent between fraudulent and legitimate transactions. Therefore, gender may not be a significant distinguishing factor for fraud detection. However, further analysis is needed to assess whether other factors, such as transaction amount or location, play a more critical role in fraudulent activity.



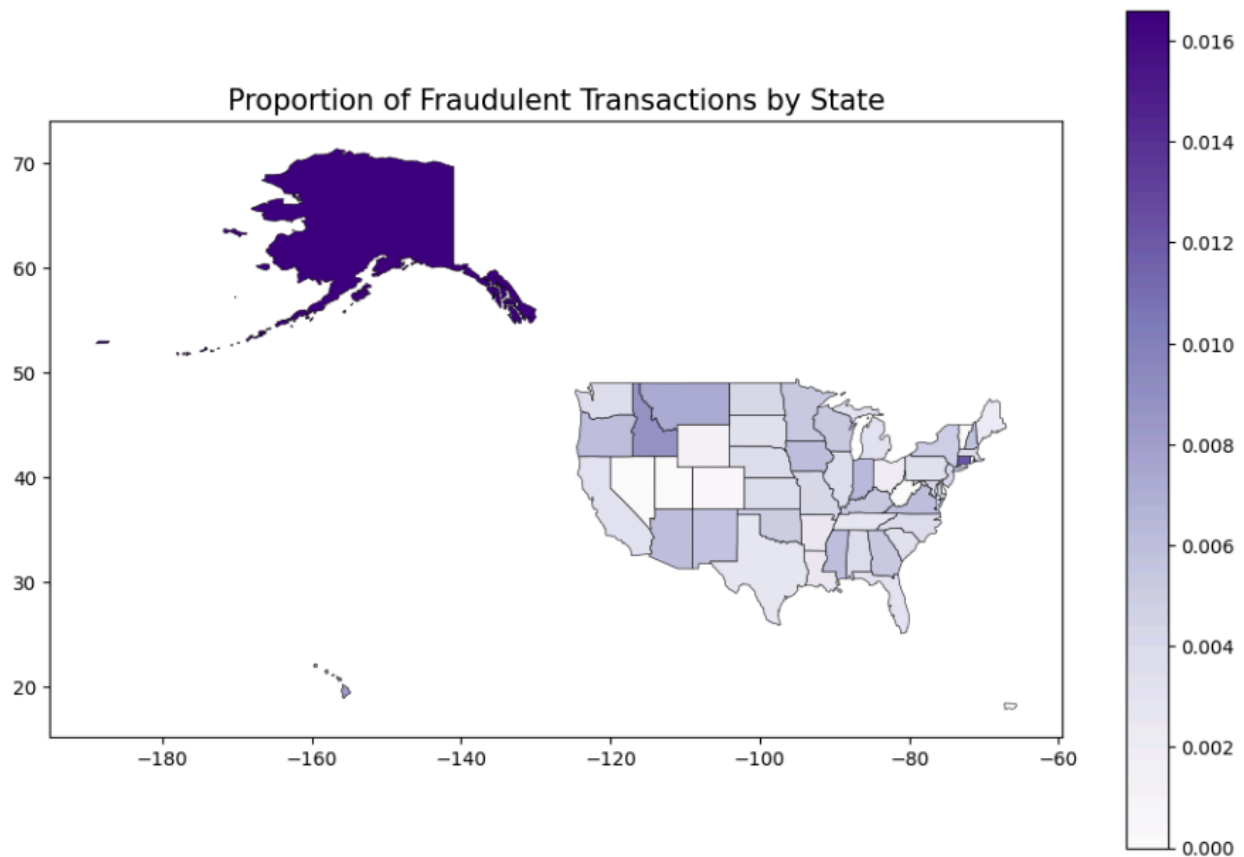
The examination of fraudulent transactions by job title reveals noticeable outliers, with certain occupations experiencing significantly higher fraud rates. Jobs such as Science Writer (30 cases), Licensed Conveyancer (29 cases), Systems Developer (29 cases), and Biomedical Engineer (28 cases) have a disproportionately high number of fraudulent transactions compared to other professions. Identified as outliers using the Interquartile Range (IQR) method, these job titles suggest that individuals in these roles may be targeted more frequently or exhibit transaction patterns that increase their likelihood of being flagged as fraudulent. The bar chart further illustrates the top 10 job titles associated with fraud, reinforcing the prominence of these occupations in fraudulent transactions. This trend raises the question of whether certain job titles are inherently more prone to fraud or if external factors contribute to their overrepresentation. Further analysis is needed to assess whether job titles serve as a meaningful predictor of fraudulent activity or if the observed pattern is due to data distribution biases.



The analysis examines the distribution of fraudulent transactions across different cities by identifying outlier cities using the Interquartile Range (IQR) method. The results indicate that Camden (27 cases), Birmingham (25 cases), and Burrton (19 cases) exhibit significantly higher fraud occurrences compared to other cities, classifying them as fraud hotspots. A bar chart displaying the top 10 cities with the most fraud cases further highlights these locations, suggesting that certain geographic areas may be more prone to fraudulent activities. Further investigation is necessary to determine whether these cities are inherently high-risk for fraud.



The bar chart presents the fraud percentage by transaction category, calculated as the proportion of fraudulent transactions within each category. The highest fraud percentages are observed in shopping\_net, misc\_net, and grocery\_pos, each exceeding 1% fraud occurrence relative to total transactions in their respective categories. Shopping\_pos, gas\_transport, and travel also show noticeable fraud rates, albeit lower than the top categories. Other categories, such as personal\_care, entertainment, and health\_fitness, have comparatively lower fraud percentages.



The choropleth map displays the proportion of fraudulent transactions by state, illustrating the relative fraud rates rather than absolute transaction counts. The darker shades indicate a higher proportion of fraudulent transactions relative to total transactions within each state. Alaska exhibits the highest fraud proportion, significantly darker than other states, suggesting that a larger fraction of transactions in Alaska are fraudulent. Other states with higher fraud proportions include Connecticut, Idaho, Montana and Oregon while several states remain light-colored, indicating a lower proportion of fraudulent transactions. This visualization helps identify regions where fraudulent transactions make up a greater share of overall transactions.



## 2.3 Data Preprocessing and Feature Engineering

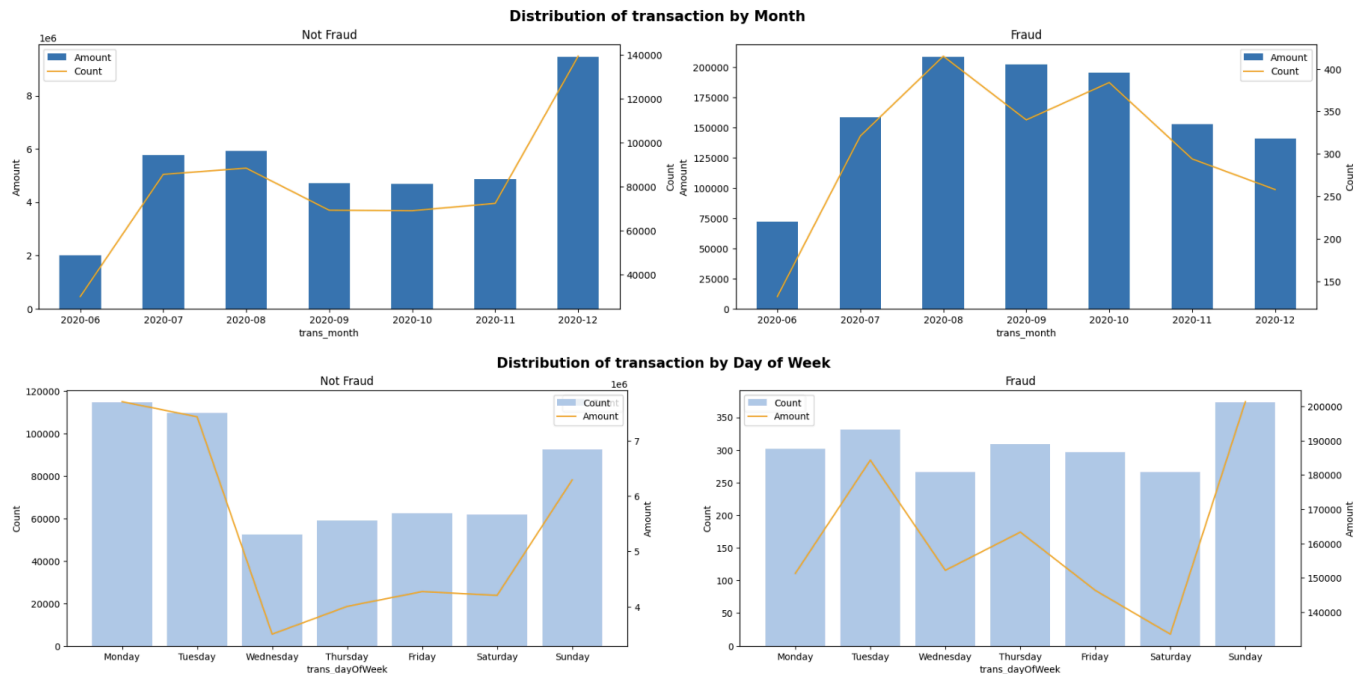
Through data preprocessing and feature engineering, we've transformed the raw dataset into a clean, structured format ready for analysis and modeling. The dataset includes details like transaction times, amounts, merchant information, and customer demographics. Our goal is to identify fraudulent transactions, which are labeled in the data. To prepare the dataset for analysis and modeling, we performed several preprocessing tasks, such as handling datetime features, creating new features, encoding categorical variables, and balancing the dataset. We also engineered features to extract meaningful insights that can help improve the performance of ML models.

### 2.3.1 Data Preprocessing

#### 2.3.1.1 Handling Datetime Features

The dataset includes a column called **trans\_date\_trans\_time**, which originally stored dates and times as strings. To make this data easier to work with, we converted it into a datetime format. We then split this column into two separate columns: **trans\_date** and **trans\_time**. This allowed us to analyze the data more granularly. Additionally, we created several new features:

- **trans\_time\_group**: This captures the hour of the transaction, helping us analyze patterns based on the time of day.
- **trans\_month**: This identifies the month of the transaction, which can help identify seasonal trends.
- **trans\_dayOfWeek**: This identifies the day of the week, helps us to study weekly patterns.



### 2.3.1.2 Handling Customer Demographics

The dataset includes a **dob** (date of birth) column, which we converted into a datetime format. From this, we calculated each customer's age. To make the age data more useful, we grouped customers into categories like "**Minors**," "**Young Adults**," "**Adults**," "**Mature Adults**," and "**Seniors**." This helps us analyze how fraud might vary across different age groups.



### 2.3.1.3 Creating New Features

We created several new features to better understand the data:

- **full\_name**: We combined the first and last name columns to create a full name for each customer.
- **time\_diff**: This calculates the time difference between consecutive transactions for each customer, which can help identify unusual patterns.

- **time\_window**: We categorized the time differences into predefined windows (e.g., "10 min," "1 hour," "24 hours") to analyze how frequently transactions occur within certain timeframes.

#### 2.3.1.4 Dropping Irrelevant Columns

To streamline the dataset, we removed columns that weren't relevant to our analysis. These included zip, trans\_num, unix\_time, merch\_lat, merch\_long, and others. This step helps reduce noise and focus on the most important features.

---

### 2.3.2 Feature Engineering

#### 2.3.2.1 Balancing the Dataset

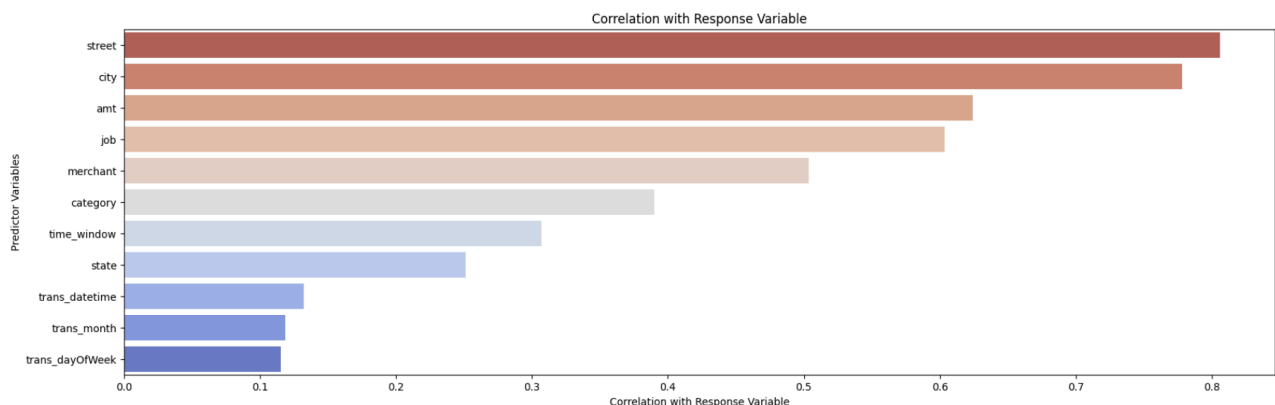
The dataset was highly imbalanced, with far fewer fraudulent transactions than legitimate ones. To address this, we created a balanced subsample by combining all fraudulent transactions with an equal number of randomly selected non-fraudulent transactions. This ensures our model isn't biased toward the majority class.

#### 2.3.2.2 Encoding Categorical Variables

Categorical variables like gender, job, state, street, merchant, city, category, trans\_dayOfWeek, and time\_window were encoded to make them suitable for machine learning models. We used one-hot encoding for the gender column and target encoding for the other categorical variables.

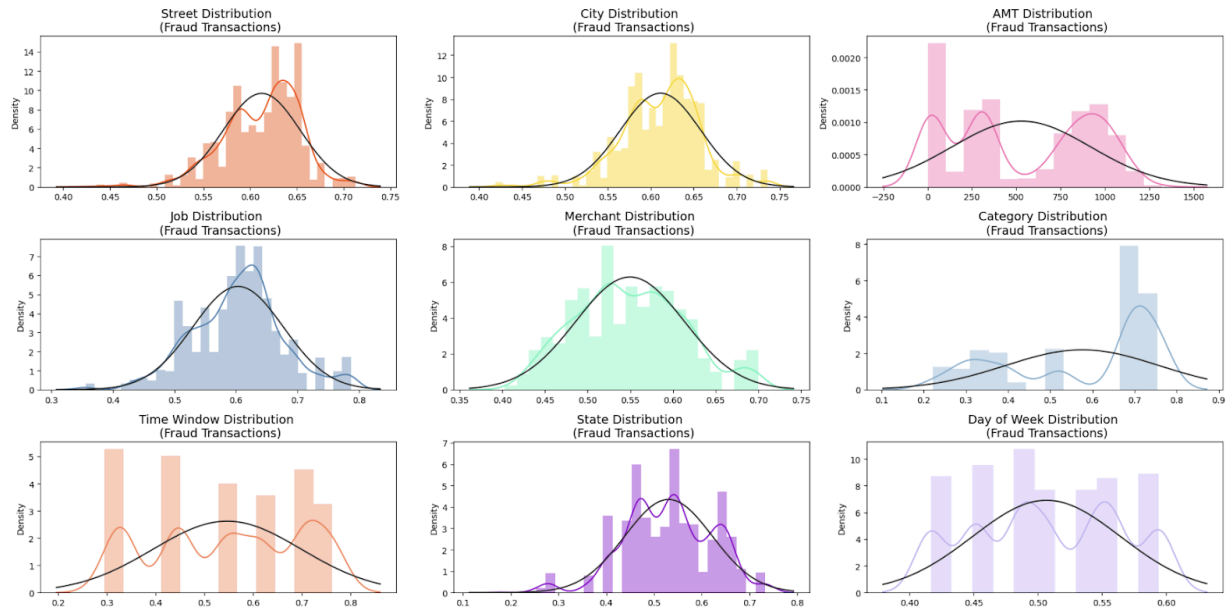
#### 2.3.2.3 Correlation Analysis

We performed a correlation analysis to identify which features have a strong relationship with the target variable (**is\_fraud**). Features with a correlation greater than 0.1 were considered relevant for modeling.



### 2.3.2.4 Visualizing Feature Distributions

To better understand the data, we visualized the distributions of key features like **street**, **city**, **amt**, **job**, **merchant**, **category**, **time\_window**, **state**, and **trans\_dayOfWeek** for fraudulent transactions. These visualizations help us spot patterns and anomalies that could be useful for fraud detection.



## 3. Summary of Key Findings

## 4. Challenges Faced and Future Recommendations

### 4.1 Challenges Faced

#### 1. Imbalanced Dataset:

- The dataset was highly imbalanced, with fraudulent transactions making up a very small percentage of the total transactions. This imbalance can lead to models that are biased toward the majority class (non-fraudulent transactions), making it difficult to accurately detect fraud.

- **Solution:** We addressed this by creating a balanced subsample.
- 2 Handling Datetime Features:
  - The **trans\_date\_trans\_time** column required careful handling to extract meaningful features like **trans\_time\_group**, **trans\_month**, and **trans\_dayOfWeek**. Ensuring consistency in datetime formatting and handling time zone differences was a challenge.
  - **Solution:** We standardized the datetime format and derived relevant features.
- 3 High Dimensionality After Encoding:
  - Encoding categorical variables like job, state, city, and category using one-hot encoding and target encoding increased the dimensionality of the dataset. This can lead to computational inefficiency and overfitting.
  - **Solution:** We used target encoding to reduce dimensionality.
- 4 Identifying Relevant Features:
  - With a large number of features, identifying the most relevant ones for fraud detection was challenging. Some features had weak correlations with the target variable, making it difficult to determine their importance.
  - **Solution:** We used correlation analysis to filter features with a correlation greater than 0.1.
- 5 Missing or Inconsistent Data:
  - The dataset had missing or inconsistent values in some columns (e.g., dob, street, city), which required careful handling to avoid introducing bias or errors into the analysis.
  - **Solution:** We dropped irrelevant columns and handled missing values.

## 4.2 Future Recommendations

### Deeper Feature Exploration:

It would be beneficial to further explore the nuances of customer behavior. For instance, looking at patterns over longer time periods or combining transaction details in new ways might reveal hidden signals that differentiate normal activity from fraudulent behavior.

### Enhanced Data Balancing Methods:

While our current approach balances the dataset by undersampling non-fraudulent transactions, exploring techniques like SMOTE (Synthetic Minority Over-sampling Technique) might help us retain more information from the non-fraudulent side while still addressing the imbalance effectively.

## **6. Each member's contribution to the project**

Zhuoxuan Li <zl3429> – Data Cleaning and Handling Inconsistencies, Report Preparation

Ziyue Gao <zg2520> – Exploratory Data Analysis (EDA), Report Preparation

Dailin Song <ds4354> – Data Preprocessing and Feature Engineering, Report Preparation

Fatih Uysal <fu2137> – Data Preprocessing and Feature Engineering, Report Preparation

For Data Acquisition, each member searched for different datasets, and then, as a team, we decided on one of them.