

# Introduction

This automated data analysis report provides an initial overview of the 'financial\_fraud\_detection\_dataset.csv' file, which contains 5,000,000 rows of transactional data across 18 distinct columns. The dataset appears to be comprehensive, capturing various aspects of financial transactions, including transaction details, sender and receiver information, merchant category, location, device used, and indicators of potential fraud. The columns included in the dataset are: transactionid, timestamp, senderaccount, receiveraccount, amount, transactiontype, merchantcategory, location, deviceused, isfraud, fraudtype, timesincelasttransaction, spendingdeviationscore, velocityscore, geoanomalyscore, paymentchannel, ipaddress, and devicehash. The purpose of this report is to undertake an exploratory analysis of the dataset, highlighting key characteristics, identifying potential data quality issues, and providing basic insights into the nature of the data. This initial assessment will lay the foundation for further, more in-depth analysis, and will help to inform strategies for data preprocessing, feature engineering, and model development. By examining the distribution of values within each column, relationships between columns, and patterns within the data, we can gain a better understanding of the dataset's structure and content. Through this analysis, we aim to identify areas of interest, such as correlations between transaction characteristics and indicators of fraud, patterns of behavior among senders and receivers, and geographic or temporal trends in transactional activity. Additionally, we will be on the lookout for potential data quality issues, including missing or anomalous values, inconsistencies in data formatting, and other factors that could impact the accuracy and reliability of our analysis. By providing a thorough and informative overview of the dataset, this report will serve as a valuable resource for anyone seeking to work with the 'financial\_fraud\_detection\_dataset.csv' file, and will help to ensure that subsequent analyses are well-informed and effective.

## Descriptive Statistics

This section provides a statistical summary of the dataset.

Summary for Numerical Features

Feature	count	mean	std	min	25%	50%	75%	max
amount	5000000.0	358.93	469.93	0.01	26.57	138.67	503.89	3520.57
timesincelasttransaction	4103487.0	1.53	3576.57	-8777.81	-2562.38	0.84	2568.34	8757.76
spendingdeviationscore	5000000.0	-0.0	1.0	-5.26	-0.68	0.0	0.67	5.02
velocityscore	5000000.0	10.5	5.77	1.0	5.0	11.0	16.0	20.0
geoanomalyscore	5000000.0	0.5	0.29	0.0	0.25	0.5	0.75	1.0

Summary for Categorical Features

Feature	count	unique	top	freq
transactionid	5000000	5000000	T100000	1
timestamp	5000000	4999998	2023-12-14T01:56:37.401698	2
senderaccount	5000000	896513	ACC983922	20
receiveraccount	5000000	896639	ACC400278	23
transactiontype	5000000	4	deposit	1250593
merchantcategory	5000000	8	retail	626319
location	5000000	8	Tokyo	625994

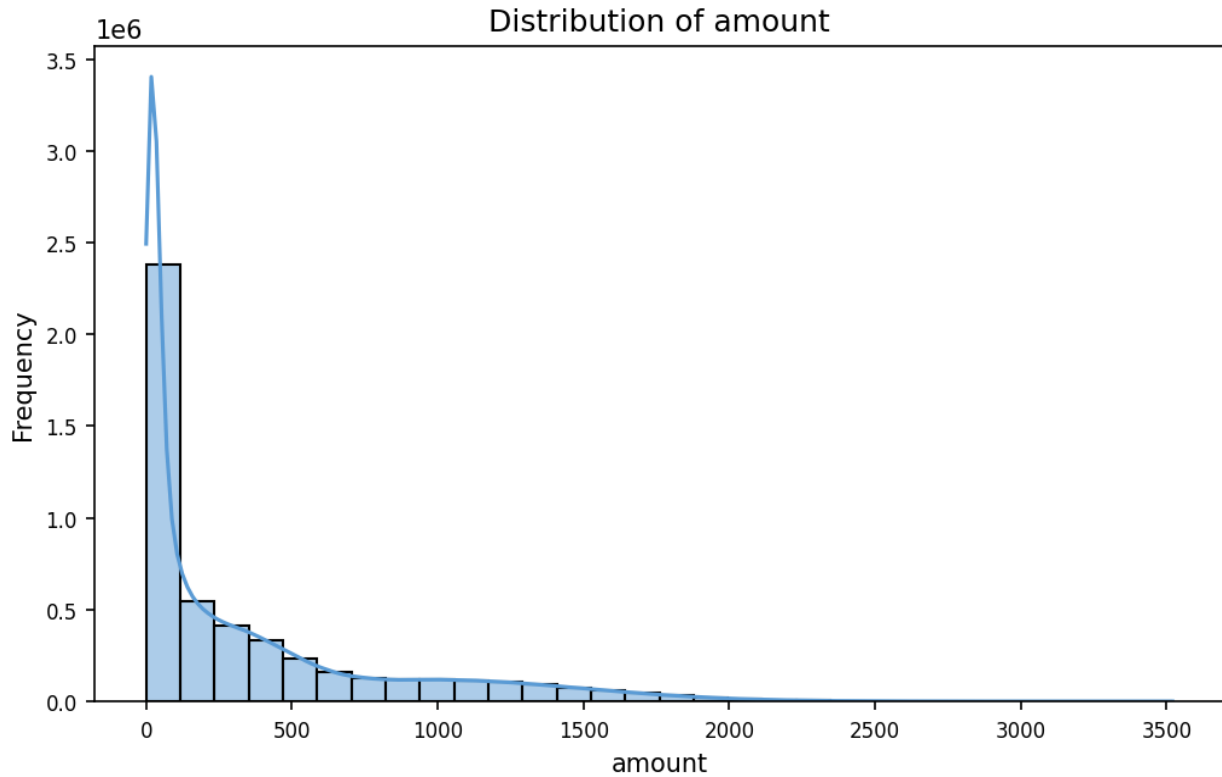
Feature	count	unique	top	freq
deviceused	5000000	4	mobile	1251131
fraudtype	179553	1	card_not_present	179553
paymentchannel	5000000	4	wire_transfer	1251219
ipaddress	5000000	4997068	135.227.29.146	2
devicehash	5000000	3835723	D7441961	9

Interpretation: The provided tables offer a summary of a dataset, which appears to be related to financial transactions. The numerical features table presents an overview of the quantitative data, including the count of observations, mean, standard deviation, and percentiles for each feature. On the other hand, the categorical features table provides an insight into the qualitative data, showcasing the count of observations, unique values, and the most frequent categories for each feature. Upon examining the numerical features, several striking observations emerge. For instance, the "timesincelasttransaction" feature has a remarkably wide range, with values spanning from -8777.81 to 8757.76, and a standard deviation of 3576.57, which is significantly larger than its mean of 1.53. This suggests the presence of potential outliers, which may warrant further investigation. Similarly, the "amount" feature has a substantial range, with values between 0.01 and 3520.57, and a standard deviation of 469.93, which is roughly 1.3 times its mean of 358.93. In contrast, features like "velocityscore" and "geoanomalyscore" exhibit relatively narrower ranges and smaller standard deviations, indicating more consistent values. The categorical features table reveals some notable patterns as well. Several features, such as "transactionid", "timestamp", "ipaddress", and "devicehash", have an extremely high number of unique values, which is not surprising given their nature (e.g., unique identifiers, timestamps, and IP addresses). On the other hand, features like "transactiontype", "merchantcategory", "location", and "deviceused" have a relatively small number of unique values, with some categories dominating the distribution (e.g., "deposit" in "transactiontype" and "retail" in "merchantcategory"). These observations can inform further analysis and data processing, such as data transformation, feature engineering, and potentially, the development of predictive models to identify patterns and anomalies in financial transactions.

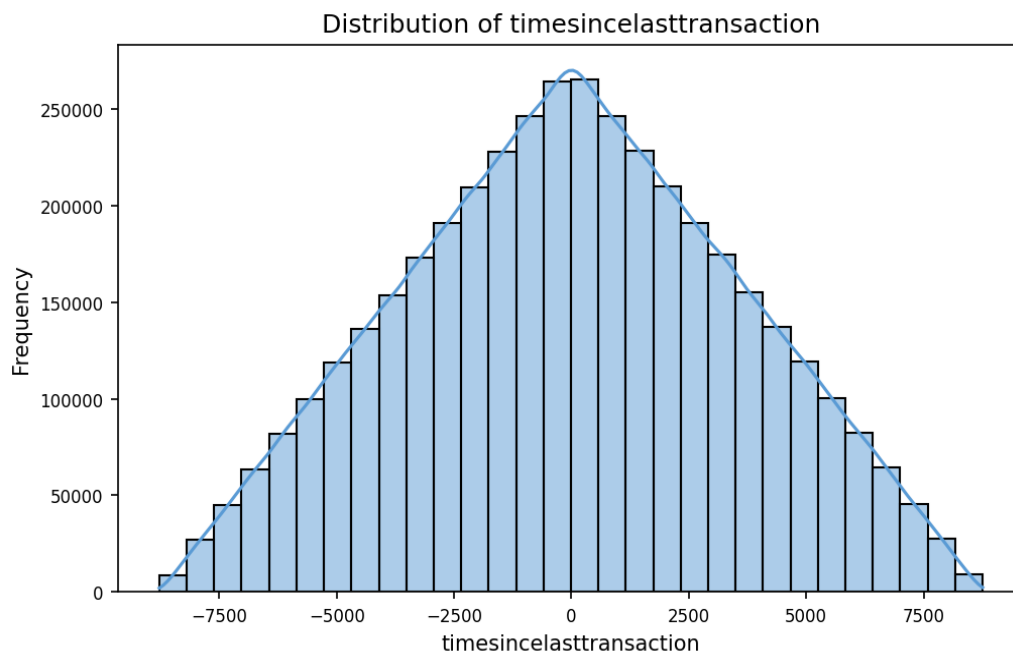
# Data Visualizations

Visual exploration can reveal patterns, trends, and outliers that are not obvious from summary statistics alone.

## Distribution of Numerical Features



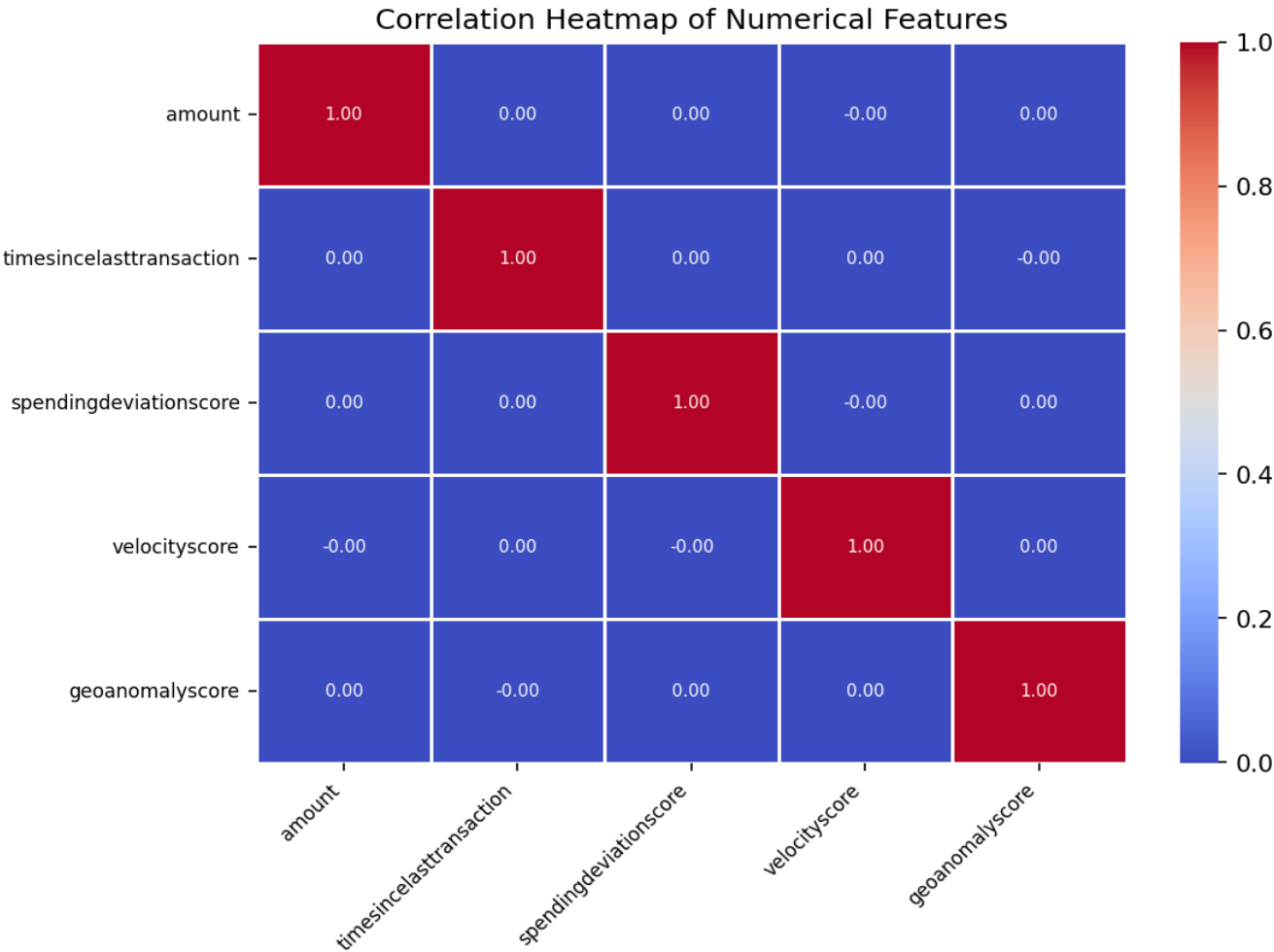
Histogram showing the distribution of 'amount'. This helps understand its central tendency, spread, and skewness.



Histogram showing the distribution of 'timesincelasttransaction'. This helps understand its central tendency, spread, and skewness.

## ***Distribution of Categorical Features***

# Correlation Analysis



Heatmap visualizing linear correlations between numerical features. Values range from -1 (strong negative) to +1 (strong positive).

Interpretation: A correlation matrix is a table that displays the correlation coefficients between different numerical features or variables. The values in the matrix range from -1 to 1, where 1 indicates a perfect positive correlation, -1 indicates a perfect negative correlation, and 0 indicates no correlation. By examining the correlation matrix, one can identify which features tend to move together or in opposite directions, providing insights into the relationships between the variables. Upon reviewing the given correlation matrix, it appears that there are no strong correlations ( $|value| > 0.6$ ) between the features. Most of the correlation coefficients are close to 0, indicating weak or negligible relationships between the variables. For example, the strongest correlation observed is between "amount" and "spendingdeviationscore" with a coefficient of 0.000799, which is still extremely weak. Similarly, the correlation between "spendingdeviationscore" and "geoanomalyscore" is 0.000472, also a very weak relationship. The lack of strong correlations suggests that these features may be relatively independent of each other, and there may not be significant patterns or relationships between them that can be exploited for further analysis or modeling.

# Data Quality Assessment

## Missing Values

Features with Missing Values

Feature	Missing Count	Percentage Missing
fraudtype	4820447	96.41
timesincelasttransaction	896513	17.93

## Duplicate Records

No duplicate rows were found in the dataset.

Overall Data Quality Remarks: Based on the findings of the data quality assessment, the overall quality of the data appears to be compromised. The presence of a significant percentage of missing values in two key features, 'fraudtype' and 'timesincelasttransaction', is a major concern. The 'fraudtype' feature has a staggering 96.41% missing values, which suggests that this feature may not be reliable for analysis. The 'timesincelasttransaction' feature has a relatively lower 17.93% missing values, but it is still a substantial proportion. The lack of duplicate rows is a positive finding, but it is overshadowed by the significant missing value issue. The potential implications of these findings for further analysis are significant. The high percentage of missing values in the 'fraudtype' feature may limit the ability to perform meaningful analysis, particularly if this feature is a key predictor or outcome variable. The missing values in 'timesincelasttransaction' may also impact the accuracy of models or analysis that rely on this feature. Overall, these findings suggest that data imputation or feature engineering techniques may be necessary to address the missing value issue before proceeding with further analysis. Alternatively, the analysis may need to be restricted to the subsets of data where these features are not missing, which could impact the generalizability of the findings.

## Conclusion and Potential Next Steps

The automated analysis of the dataset with 5000000 rows and 18 columns has provided a foundational understanding of the data, offering insights into the distribution of features, correlations between variables, and data quality. This initial report serves as a crucial starting point for further exploration, highlighting key characteristics of the data and identifying potential areas of interest for more in-depth analysis. Building on this foundation, several potential next steps can be explored to gain a deeper understanding of the data. Conducting hypothesis tests to validate observed relationships could help confirm initial findings and provide a more robust understanding of the data. Additionally, segmenting the data based on key categorical features, such as transaction type or merchant category, could facilitate comparative analysis and reveal subtle patterns or trends that may not be immediately apparent. Further analysis could also involve performing feature engineering to create new insightful variables, such as calculating the average transaction amount by sender account or receiver account, or identifying the most common transaction types by location. By pursuing these next steps, a more comprehensive and nuanced understanding of the data can be developed, ultimately informing business decisions or guiding further investigation.