

Data Quality Report

Sheetal Srivastava

1. Data Description

The dataset is **Card Transactions Data**, which contains **business transaction** information of credit cards for an organization. The data consists of real card transactions coming from a U.S. government organization. The data spans over the time period of a year, from **1st January, 2010** to **31st December, 2010**. There are **10 fields** and **96753 records**.

2. Summary Tables

Numeric Fields Table

Field Name	# Records With Values	% Populated	# Zeros	Min	Max	Mean	Most Common	Stdev
Date	96753	100.00%	0	2010-01-01	2010-12-31	2010-06-25	2010-02-28	98 days
Amount	96753	100.00%	0	0.01	3102045.53	427.89	3.62	10006.14

Categorical Fields Tables

Field Name	# Records With Values	% Populated	# Zeros	# Unique Values	Most Common
Recnum	96753	100.00%	0	96753	1
Cardnum	96753	100.00%	0	1645	5142148452
Merchnum	93378	96.50%	0	13091	930090121224
Merch description	96753	100.00%	0	13126	GSA-FSS-ADV
Merch state	95558	98.80%	0	227	TN
Merch zip	92097	95.20%	0	4567	38118
Transtype	96753	100.00%	0	4	P
Fraud	96753	100.00%	95694	2	0

3. Visualization of Each Field

a. Field Name : **Recnum**

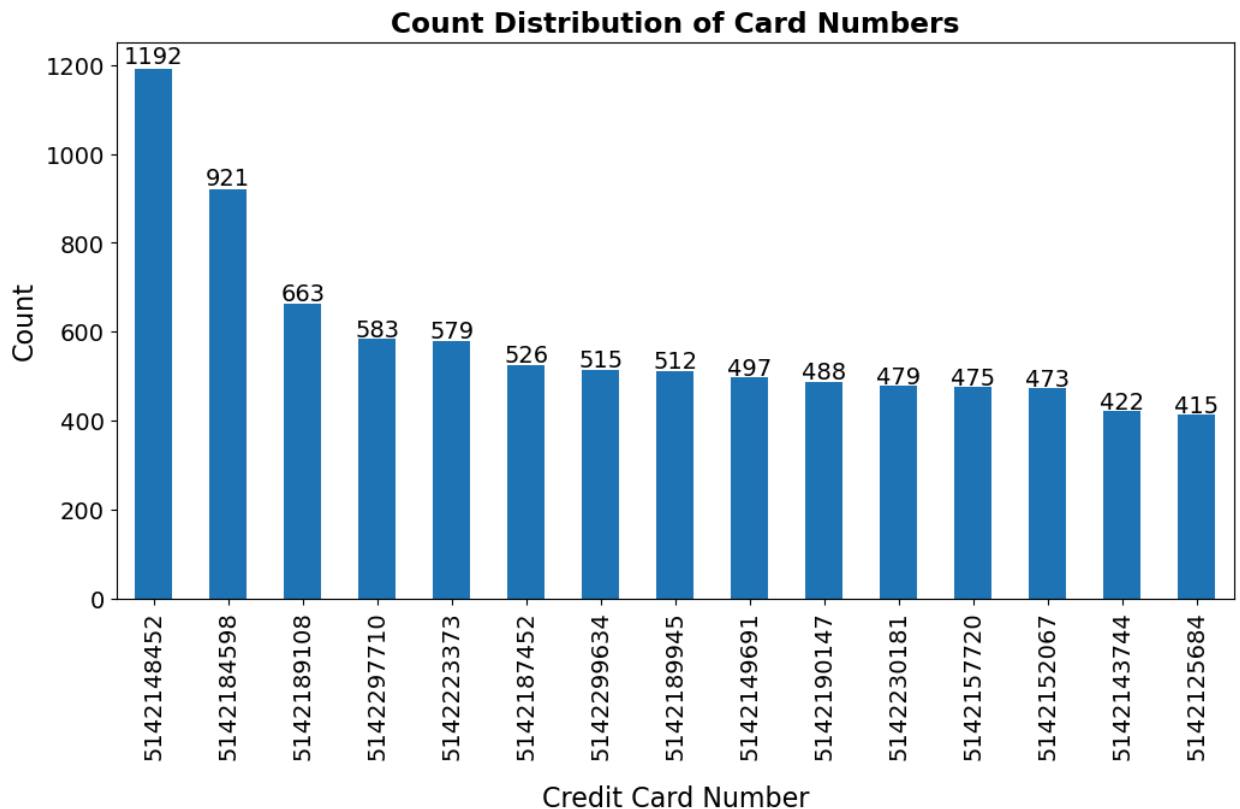
Description : Ordinal unique positive integer for each credit card transaction record, from 1 to 96753.

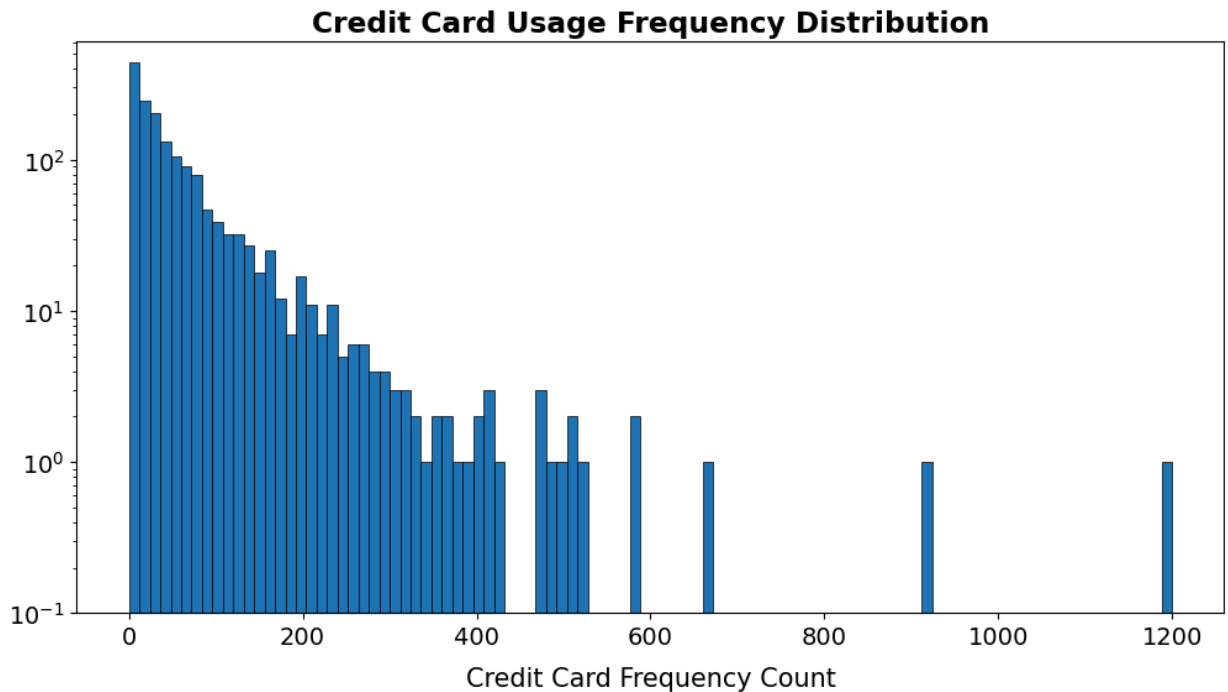
b. Field Name : **Cardnum**

Description : Credit card number. This field has 1645 unique values.

The first graph shows the top 15 credit cards with maximum occurrences in the dataset. The credit card with the card number "5142148452" is the most common card used, with 1192 total transactions.

Observing the frequency distribution graph (second figure) we can see that only 1 card has been used more than 1000 times (1192 times, to be precise). Majority of the credit cards have been used 1 to 200 times only.



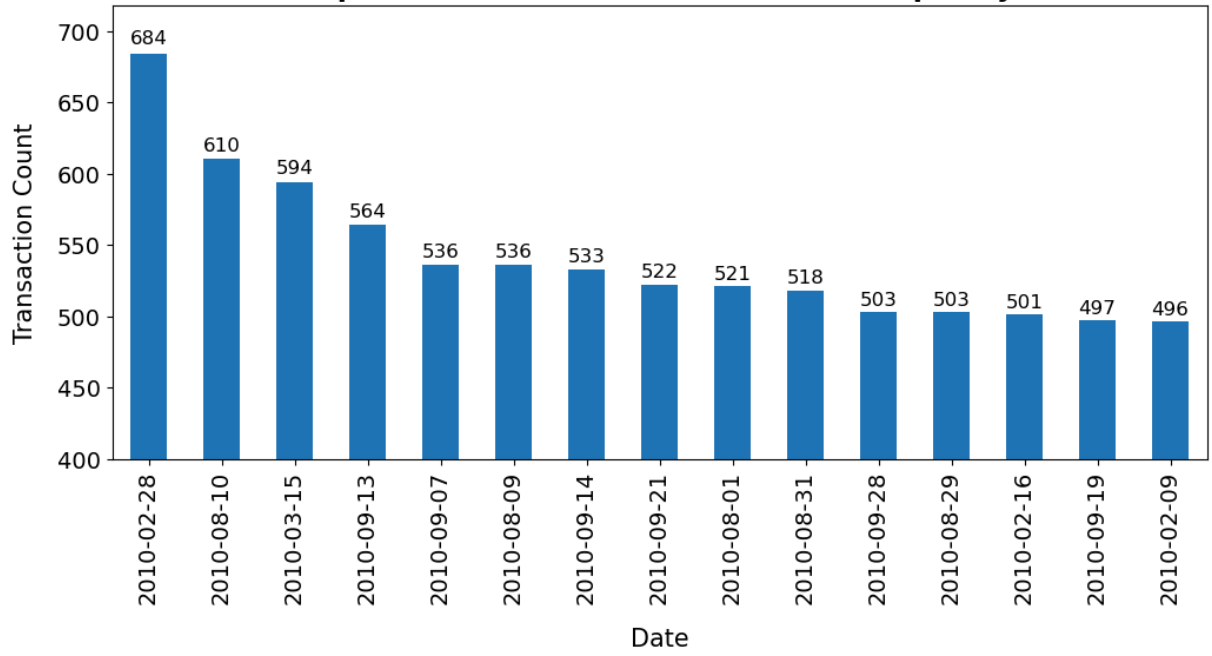


c. Field Name : Date

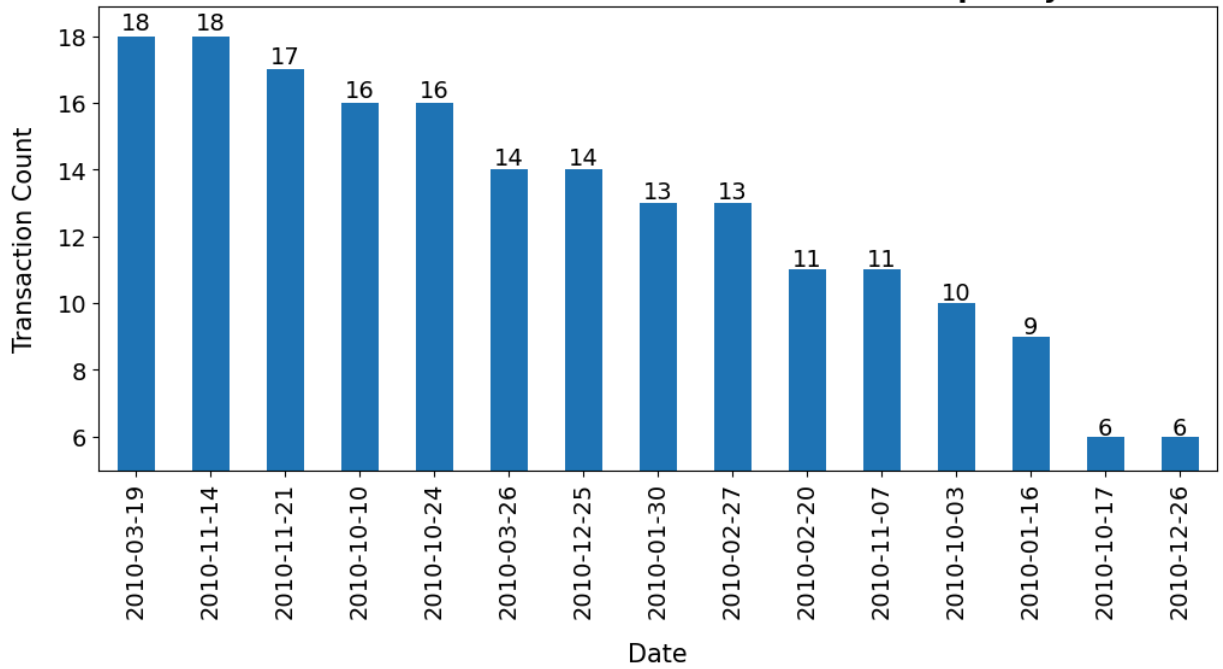
Description : Transaction date. This field has 365 unique values. We have exactly a year's worth of data, for the year 2010 (Jan 1st to Dec 31st). Based on the two Transaction Frequency graphs, we can see that the minimum number of transactions on any given date were 6, and the maximum number of transactions made on a single day were 684 (on 28th February, 2010).

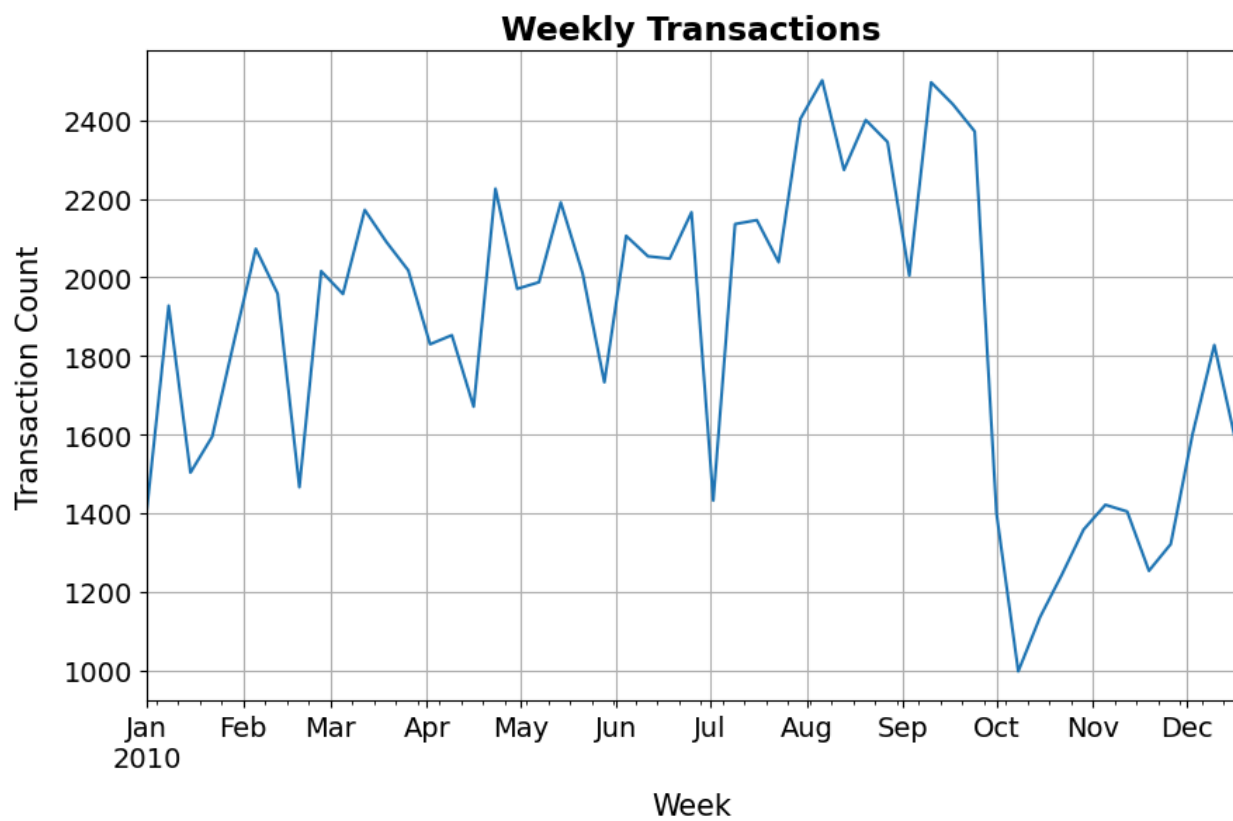
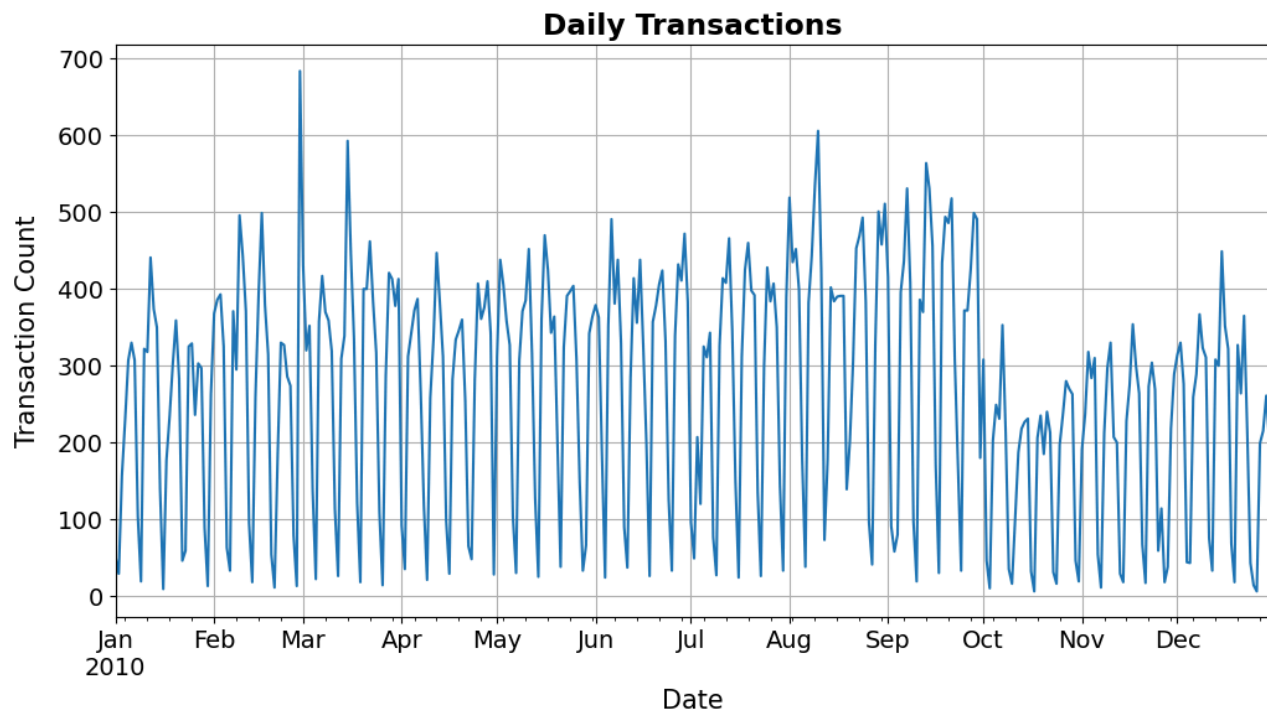
The remaining three plots show the daily, weekly and monthly transaction frequency distributions over time. We can see that the month of August has recorded the highest number of transactions, in 2010.

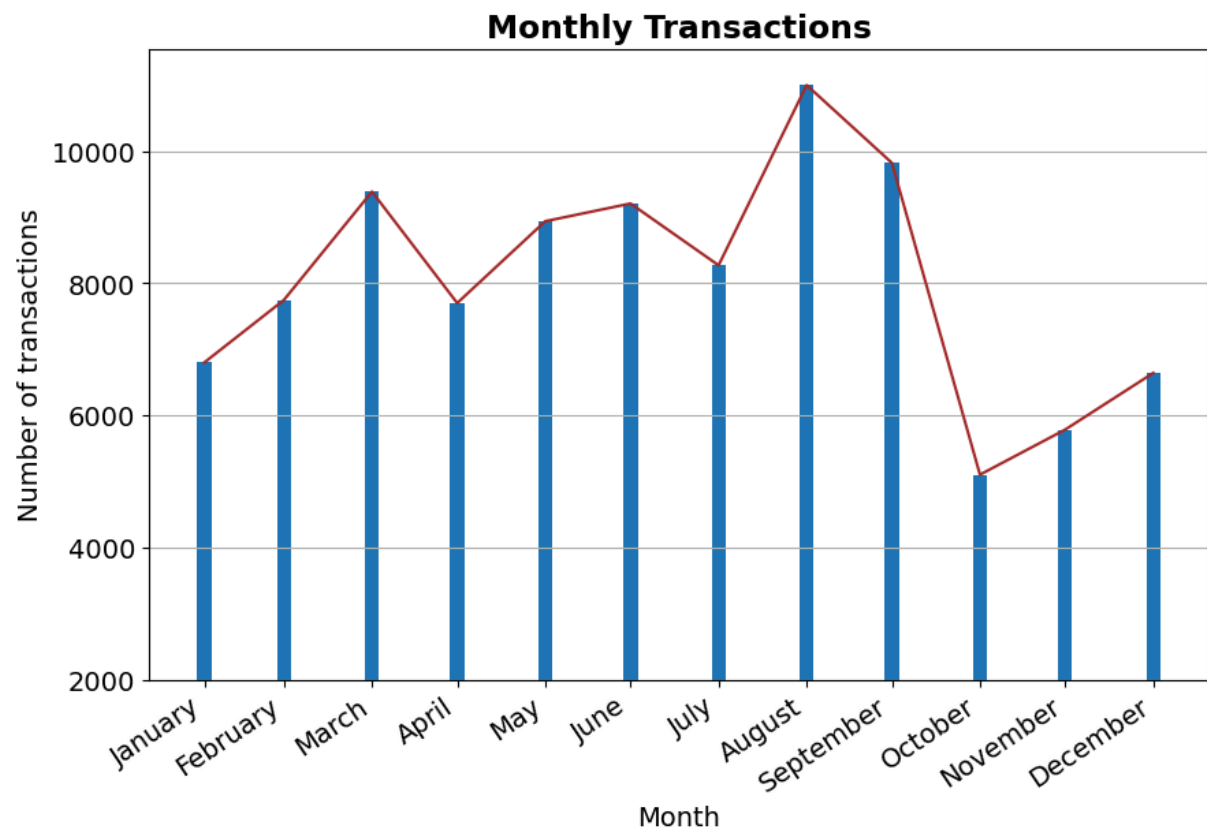
Top 15 Dates based on Transaction Frequency



Bottom 15 Dates with Least Transaction Frequency

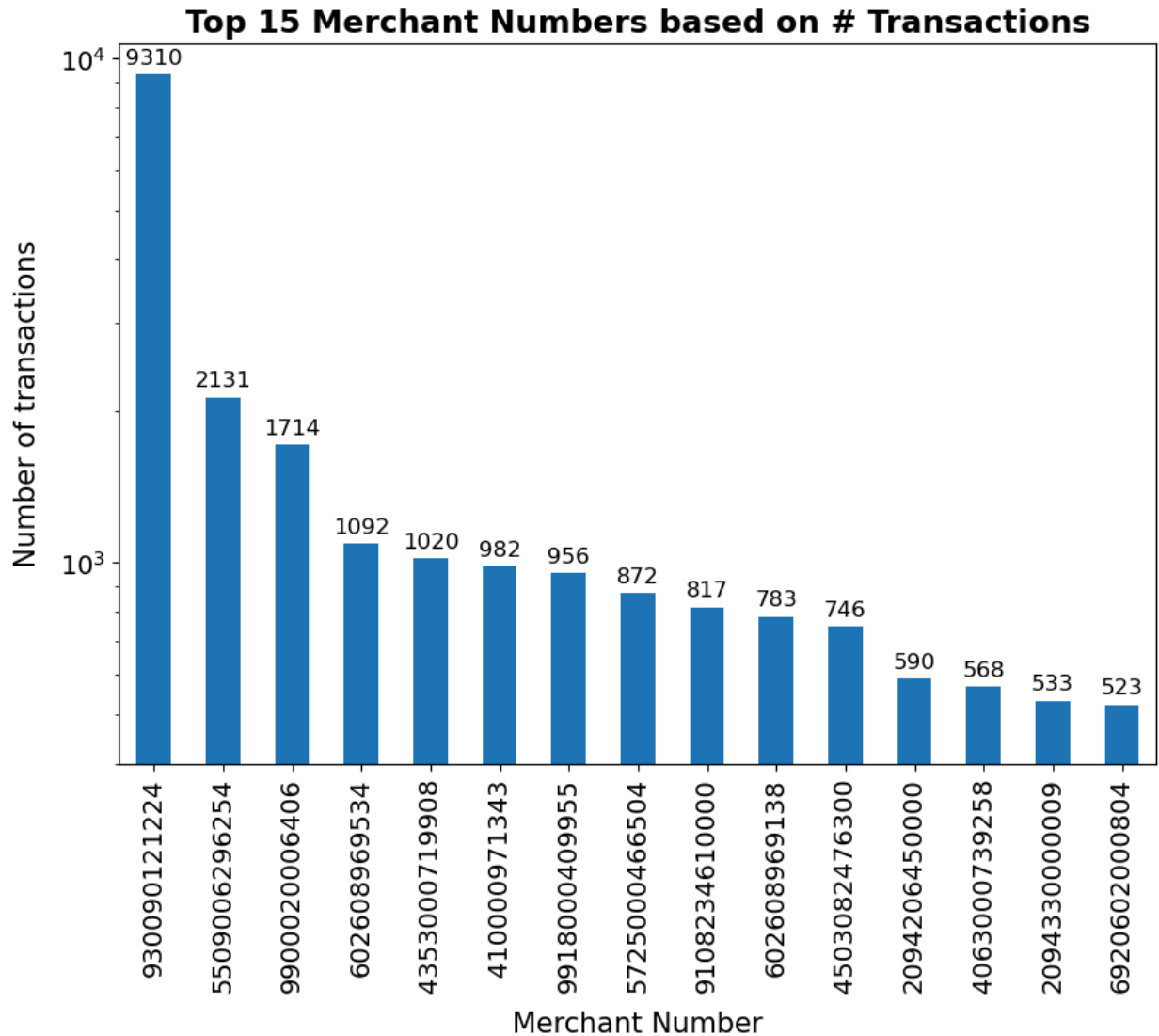






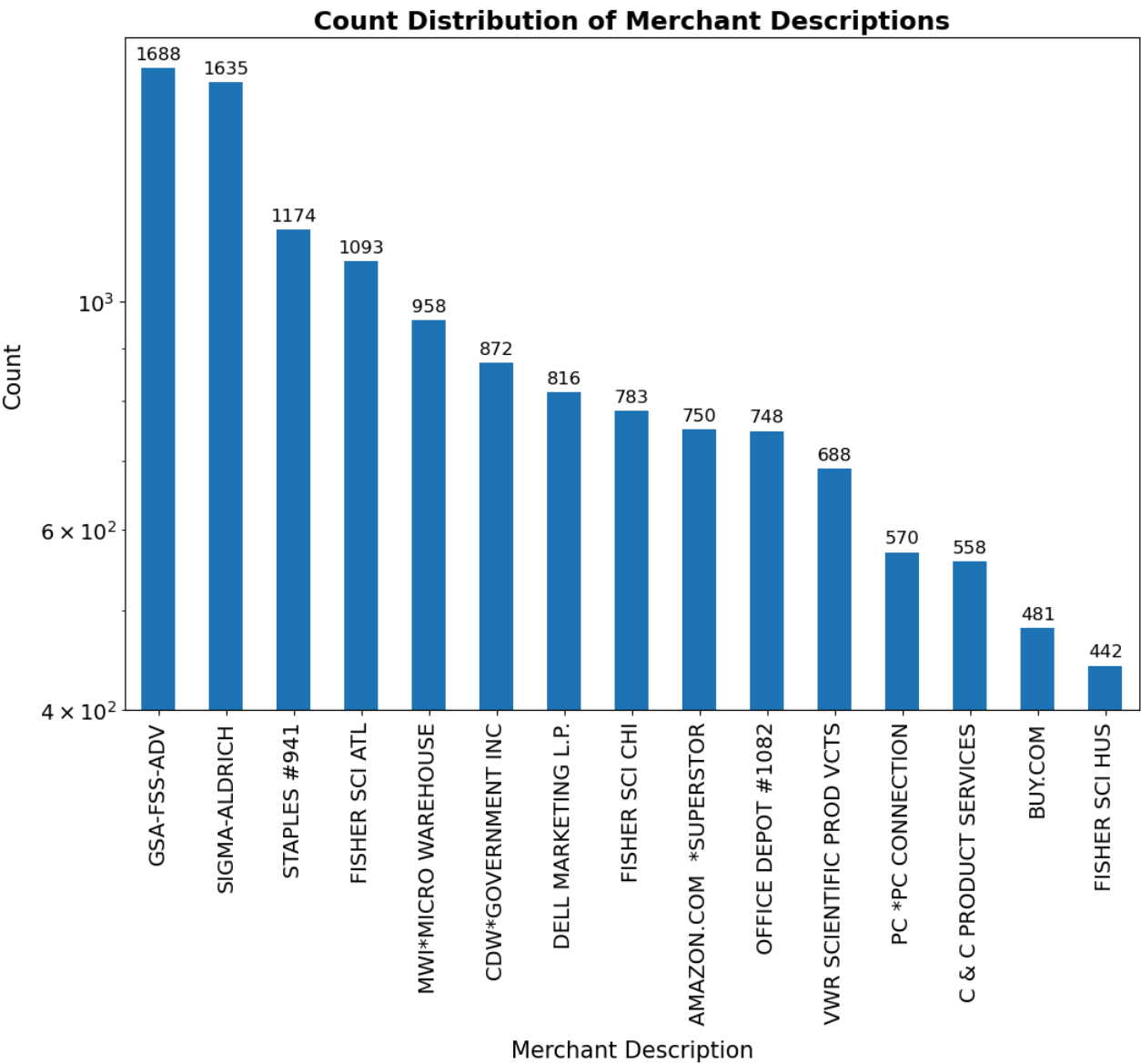
d. Field Name : **Merchnum**

Description : Merchant Number. This is a categorical field with 13092 distinct values. The most common merchant number is “930090121224”, with a total transaction count of 9310.



e. Field Name : **Merch description**

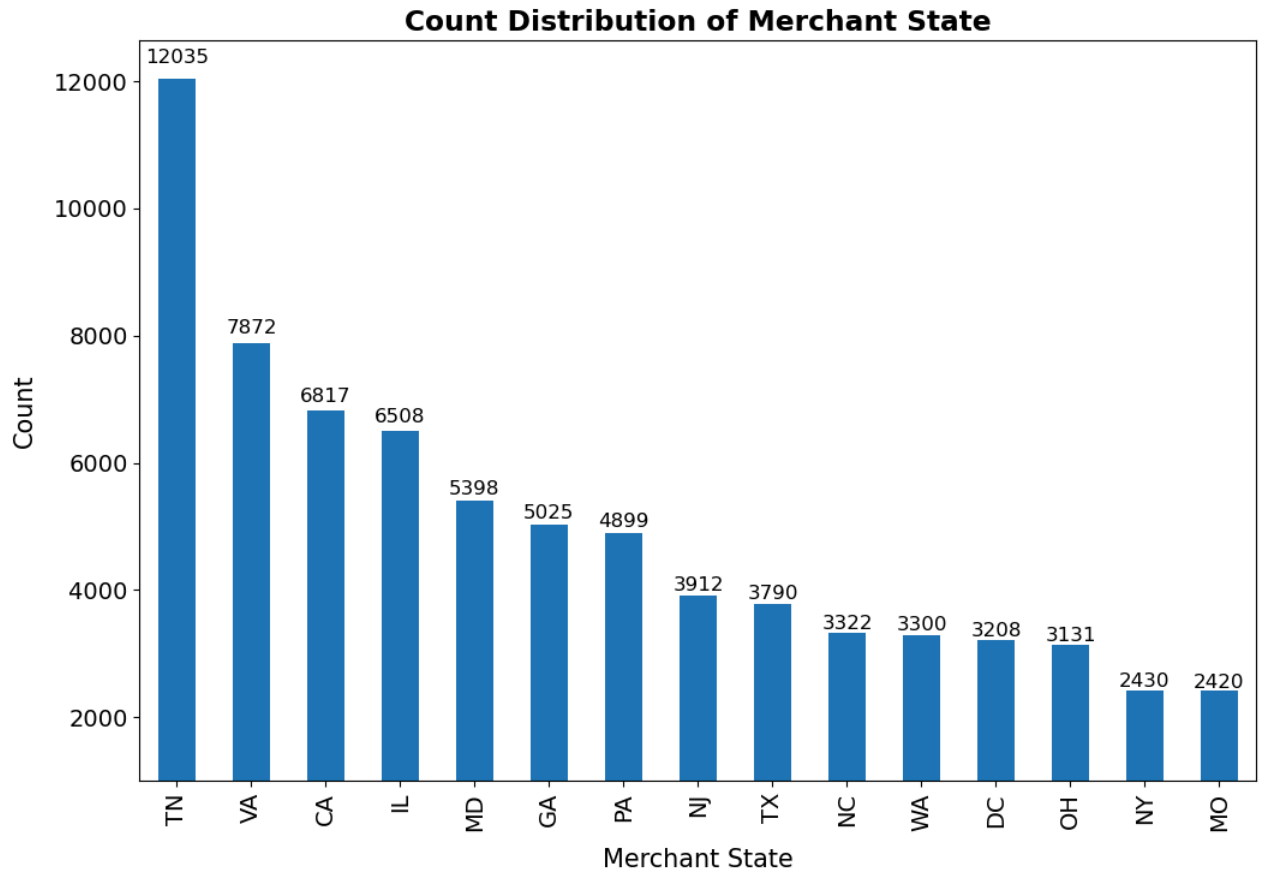
Description : Merchant Description. This a text field, giving a brief description of the merchant. This field has 13126 distinct values. The most common merchant description is “GSA-FSS-ADV” with a total count of 1688.



f. Field Name : **Merch state**

Description : Merchant State. This is a categorical field that contains the two letter state code of the state to which the merchant belongs. This field has 228 unique values.

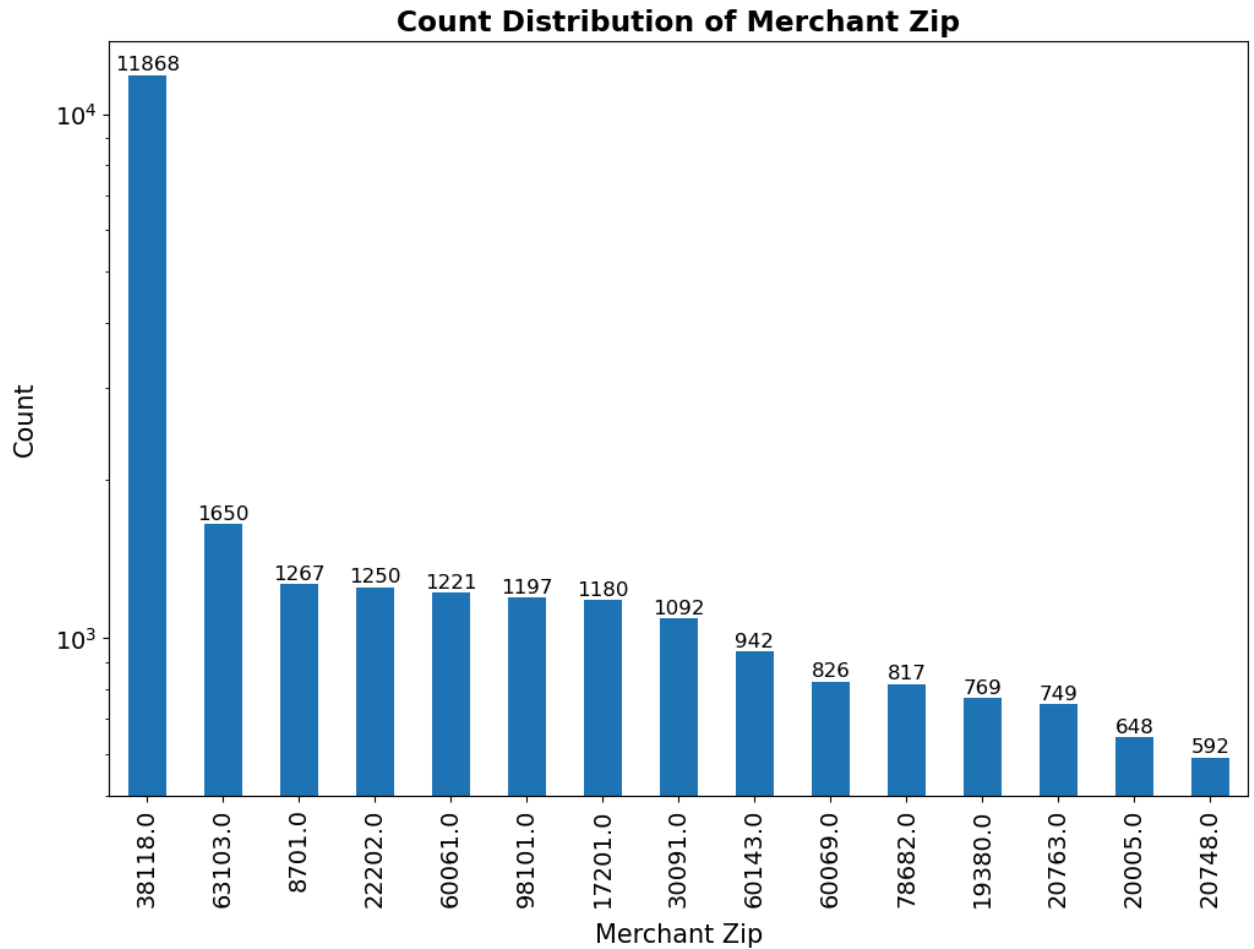
The figure below shows the top 15 merchant states based on the number of occurrences. The most common merchant state is “TN” with a total count of 12035 transactions.



g. Field Name : Merch zip

Description : Merchant zip code. This is a categorical field consisting of the zip code values of the merchant. This field has 4568 distinct values.

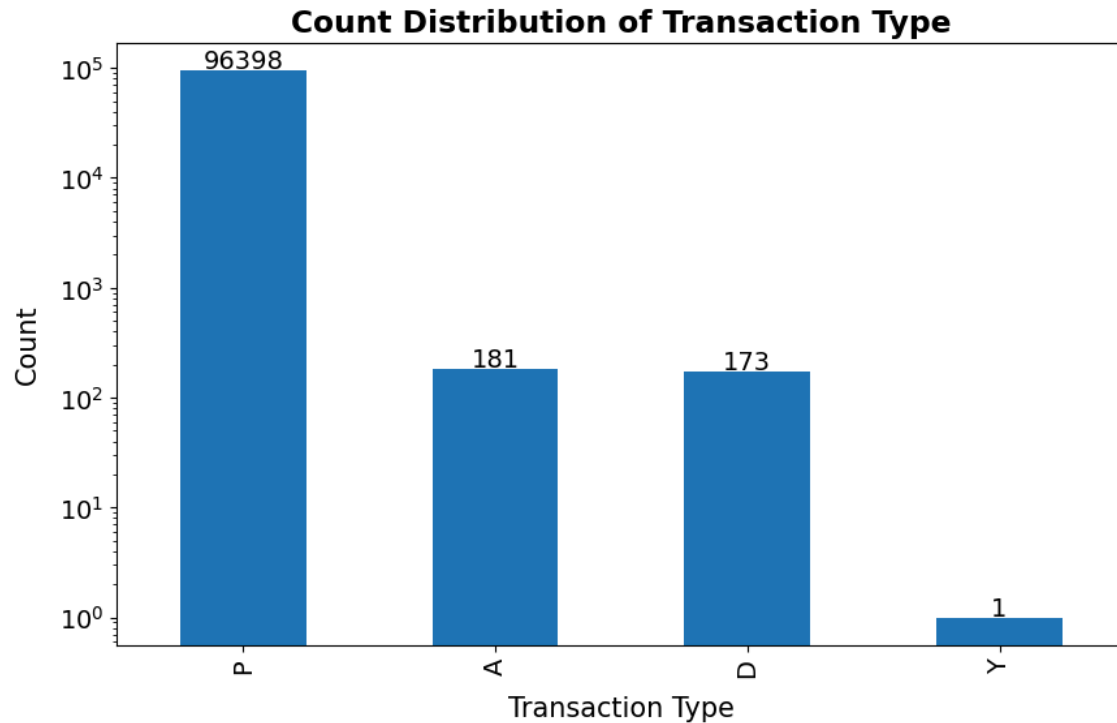
The most common Merchant Zip is “38118” with 11868 total occurrences.



h. Field Name : Transtype

Description : Transaction type. This is a categorical field classifying the type of transaction. This field has 3 distinct values which are - P, A, D, and Y.

The most common transaction type is “P” (purchase) with a total of 96398 occurrences out of a total of 96753 records.

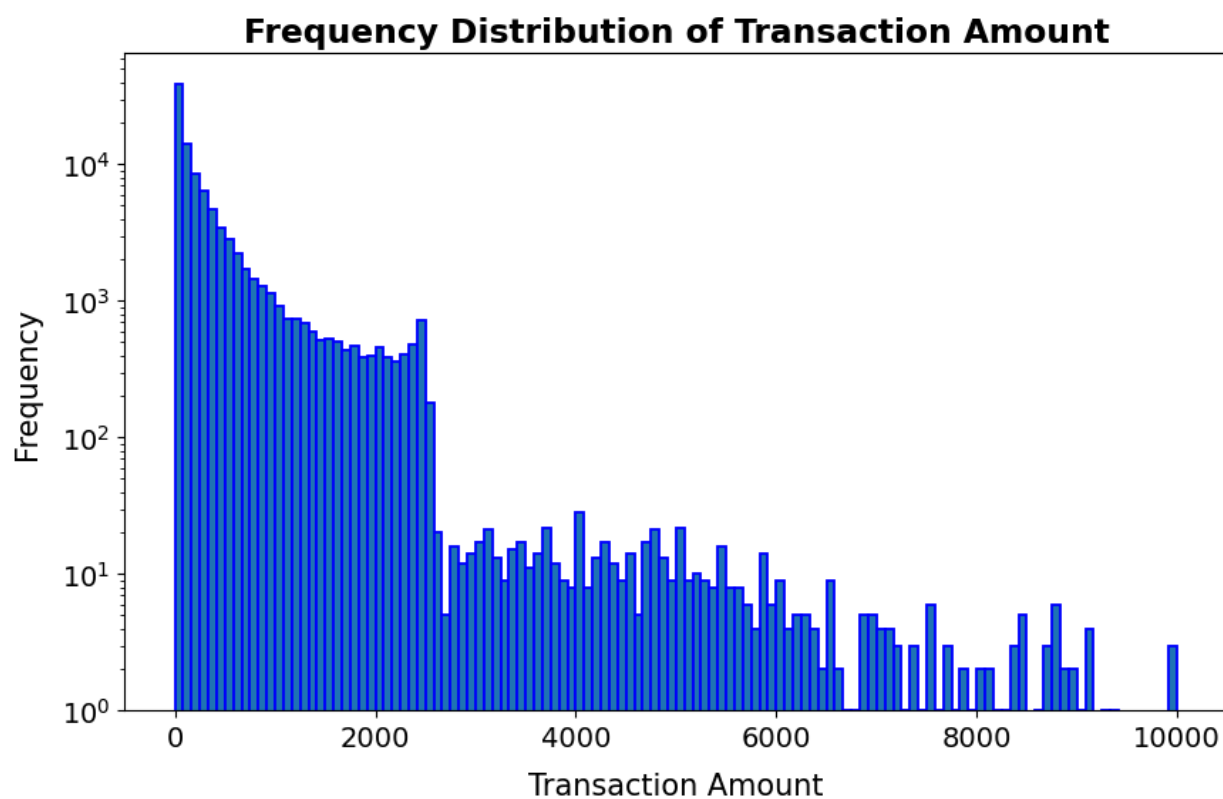
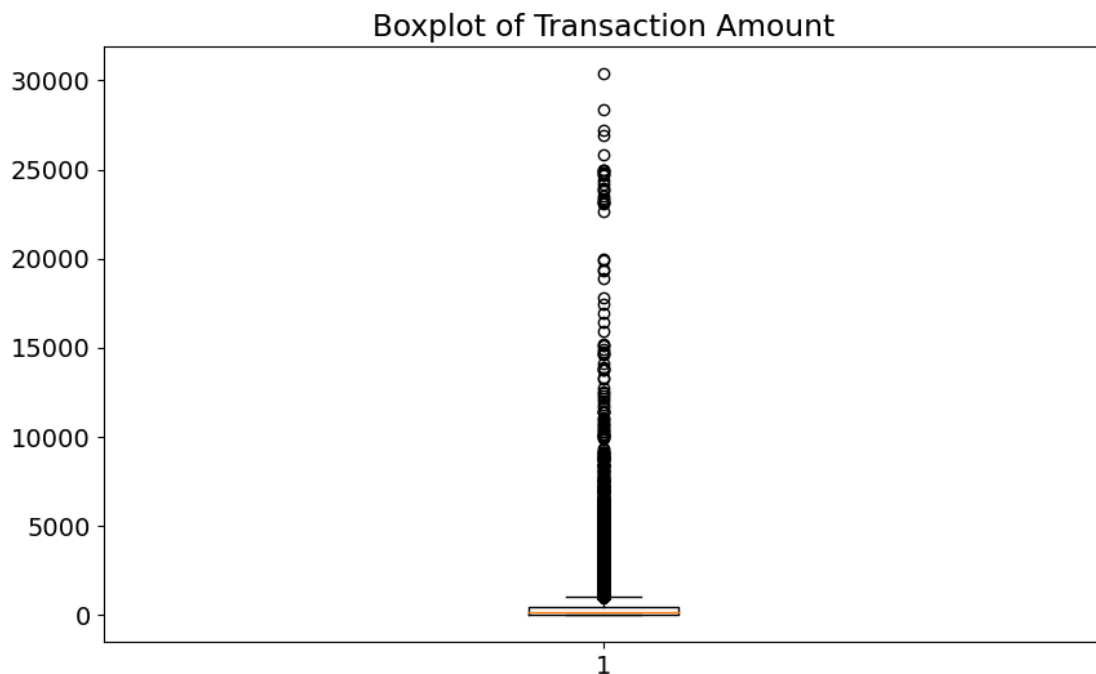


i. Field Name : Amount

Description : This numerical field contains the value of the transaction amount. The amount value ranges from 0.01 to 3102045.53 with a mean value of 427.89 and mode value of 3.62.

Examination of this data shows that the max value (3102045.53) is an outlier. Removing this outlier value, we get the first figure, a boxplot, showing that the majority of the transactions fall within the range of 0.01 to 10,000.

Therefore, we plot the second graph, which gives us the distribution of the transaction amount, keeping the amount limit as 10,000. This plot shows us that the bulk of the transactions are of an amount ranging from 0.01 to 2700.



j. Field Name : **Fraud**

Description : This is a binary categorical field that specifies the label of the records, classifying it as fraud if the value = 1, and not a fraud transaction if value = 0. Since it's a binary field, it only takes 2 values : 0 and 1.

The total count of fraudulent transactions = 1,059.

The total count non-fraudulent transactions = 95,694.

