

Bilenkin550Week8_Term_Project_Milestones

May 1, 2025

1 DSC 550 Term Project - Milestone 1

1.0.1 Exploratory Data Analysis (EDA)

Finding Patterns in Credit Card Transactions to Detect Fraud

Today, credit cards have become an essential part of our daily lives, offering a convenient way to purchase goods and services. Instead of carrying large amounts of cash, consumers can simply use a single credit card to make transactions. Whether by swiping, inserting, or tapping against a contactless reader, making a payment has never been easier.

However, with the growing adoption of credit cards, fraudulent activities have also increased. Fraudsters can use stolen credit card information to make unauthorized purchases, often without the cardholder realizing it until it's too late. These fraudulent transactions lead to significant financial losses for banks and credit card issuers.

By leveraging big data from credit card transaction histories, banks can identify patterns and correlations associated with fraudulent activities. Machine learning models can analyze various features, such as transaction location, merchant category, purchase amounts, and unusual spending behaviors, to detect potential fraud. More importantly, predictive analytics can help banks assess the risk of issuing credit cards to individuals who may have a high likelihood of committing fraud.

Implementing machine learning-based fraud detection not only helps banks minimize financial losses but also enhances security for both consumers and merchants. A robust fraud prevention system ensures that legitimate transactions are processed smoothly while fraudulent attempts are blocked in real-time. As a result, customers can feel more secure using their credit cards, and merchants can conduct transactions with greater confidence.

Banks can take preventive actions by freezing suspicious transactions, sending alert messages to cardholders, or implementing stricter verification processes.

```
[150]: import pandas as pd

# Loading dataset
file_path = r"C:\Users\maxim\OneDrive\Desktop\BU\DSC 550\Term_
↳Project\credit_card_transactions.csv"
df = pd.read_csv(file_path)
```

1. Fraud vs. Non-Fraud Transactions Count (Bar Chart)

```

[151]: import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib.ticker import FuncFormatter

# Counting fraud vs. non-fraud transactions
fraud_counts = df['is_fraud'].value_counts()

# Converting fraud labels to readable text
fraud_labels = {0: "Non-Fraud", 1: "Fraud"}

# Plotting the graph
plt.figure(figsize=(6, 4))
sns.barplot(x=fraud_counts.index.map(fraud_labels), y=fraud_counts.values,
            palette=["blue", "red"], hue=fraud_counts.index, legend=False)

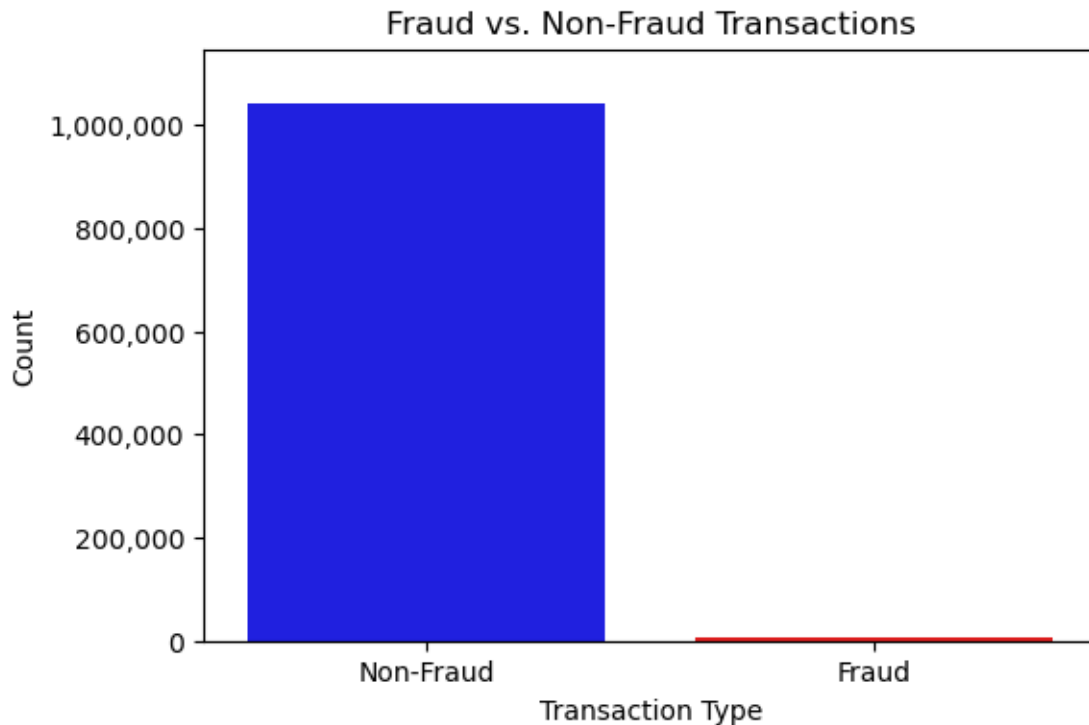
# Setting y-axis limits to make the fraud bar visible
plt.ylim(0, max(fraud_counts.values) * 1.1) # Adjusting the upper limit to
            ↪ make fraud visible

# Formatting y-axis labels with commas for better readability
def comma_format(x, pos):
    return f'{int(x):,}'

plt.gca().yaxis.set_major_formatter(FuncFormatter(comma_format))

# Labeling and titling the graph
plt.xlabel("Transaction Type")
plt.ylabel("Count")
plt.title("Fraud vs. Non-Fraud Transactions")
plt.show()

```



Explanation: “The Fraud vs. Non-Fraud Transactions” bar chart shows the number of Non-Fraud (highlighted in blue) and Fraud (highlighted in red) transactions in the dataset, with count displayed on the left y-axis. As seen in the chart, the majority of transactions (over 1 million) are Non-Fraud, with a significantly smaller number being Fraud. This indicates that fraud is a relatively rare occurrence in the dataset.

2. Transaction Amounts by Fraud vs Non-Fraud (Box Plot):

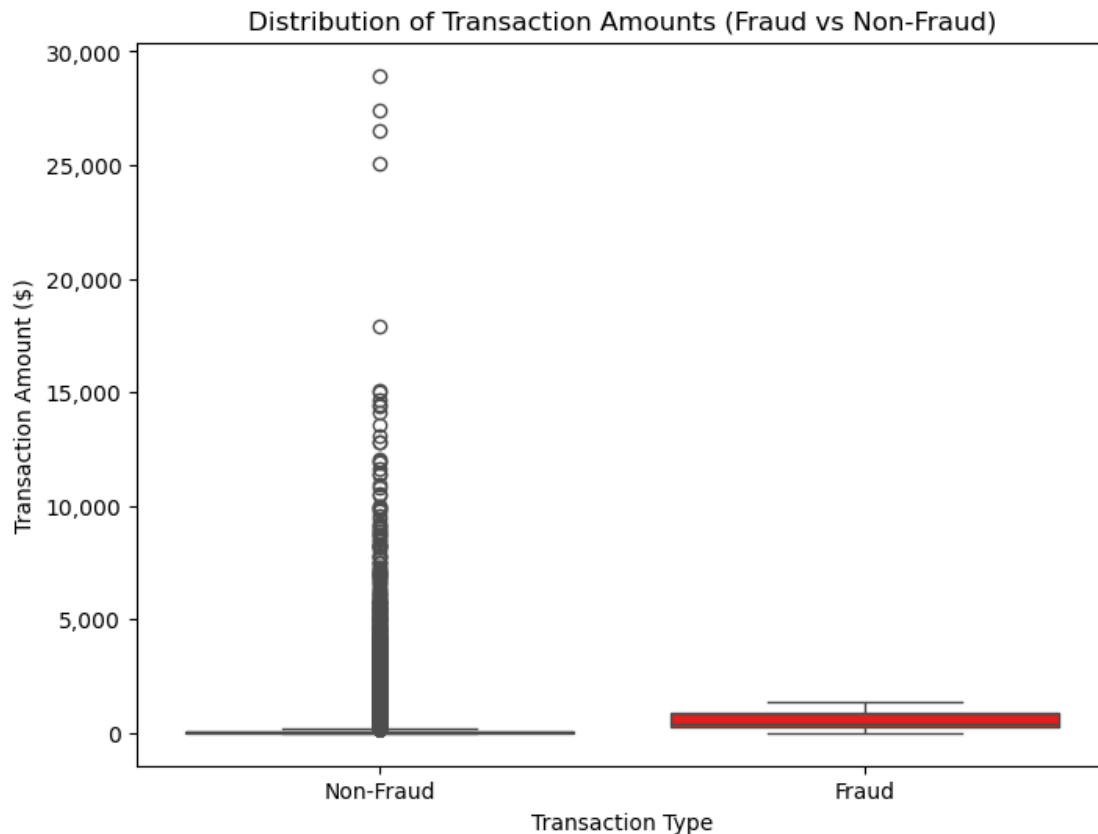
```
[152]: # Plotting Transaction Amounts by Fraud vs Non-Fraud
plt.figure(figsize=(8, 6))
sns.boxplot(x="is_fraud", y="amt", data=df, hue="is_fraud", palette=["blue", "red"], legend=False)

# Converting 0 and 1 to "Non-Fraud" and "Fraud"
plt.xticks(ticks=[0, 1], labels=["Non-Fraud", "Fraud"])

plt.xlabel("Transaction Type")
plt.ylabel("Transaction Amount ($)")
plt.title("Distribution of Transaction Amounts (Fraud vs Non-Fraud)")

# Format the y-axis with commas for readability
plt.gca().yaxis.set_major_formatter(plt.FuncFormatter(lambda x, loc: f'{x:,}'))
```

```
plt.show()
```



Explanation: The “Distribution of Transaction Amounts (Fraud vs Non-Fraud)” graph compares the transaction amounts of Non-Fraud and Fraud transactions in the dataset. The graph shows the transaction amounts on the y-axis and the transaction type (Non-Fraud or Fraud) on the x-axis. As shown, Non-Fraud transactions vary widely, ranging from 1 to 28,948.90 dollars. In contrast, Fraud transactions are generally smaller, ranging from 1.18 to 1,371.81 dollars. This suggests that fraudsters tend to make smaller credit card transactions, possibly hoping that small amounts will go unnoticed by banks.

Category-wise Distribution of Transactions

```
[153]: import matplotlib.pyplot as plt
import seaborn as sns

# Plotting Category-wise Fraud vs Non-Fraud Distribution
plt.figure(figsize=(10, 6))
sns.countplot(x="category", hue="is_fraud", data=df, palette=["blue", "red"],
             legend=False)

# Labeling and rotating the x-axis ticks
```

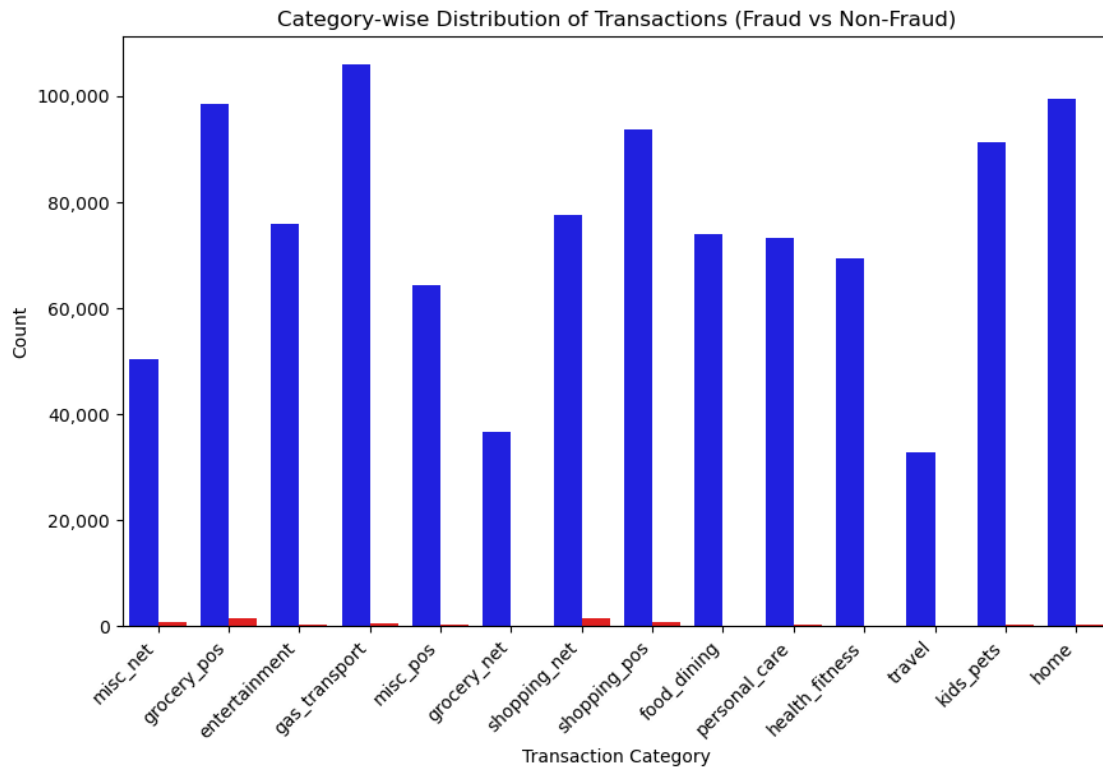
```

plt.xlabel("Transaction Category")
plt.ylabel("Count")
plt.title("Category-wise Distribution of Transactions (Fraud vs Non-Fraud)")
plt.xticks(rotation=45, ha='right') # Adjusting the horizontal alignment for
    ↪ better readability

# Formatting the y-axis with commas for readability
plt.gca().yaxis.set_major_formatter(plt.FuncFormatter(lambda x, loc: f'{int(x):
    ↪ ,}')))

plt.show()

```



Explanation: The “Category-wise Distribution of Transactions (Fraud vs Non-Fraud)” graph shows transactions by category for both Non-Fraud and Fraud. It separates transactions by each category and indicates whether fraud occurred online or at the point of sale. For example, the red bar for misc_net shows online fraud transactions, while the red bar for misc_pos shows fraudulent transactions that occurred at the point of sale.

In line with the common belief that fraudulent credit card transactions mostly happen online, the graph partly validates this. We can see that for misc_net, fraudulent transactions (highlighted in red) are slightly higher than for misc_pos. The same trend is observed for shopping_net vs. shopping_pos, with the red bar for shopping_net being marginally higher than for shopping_pos.

However, the data for grocery_net contradicts this assumption. No fraud was detected online

(grocery_net), while all fraudulent transactions occurred at the point of sale (grocery_pos). This is counterintuitive to the common belief that fraud mostly occurs online.

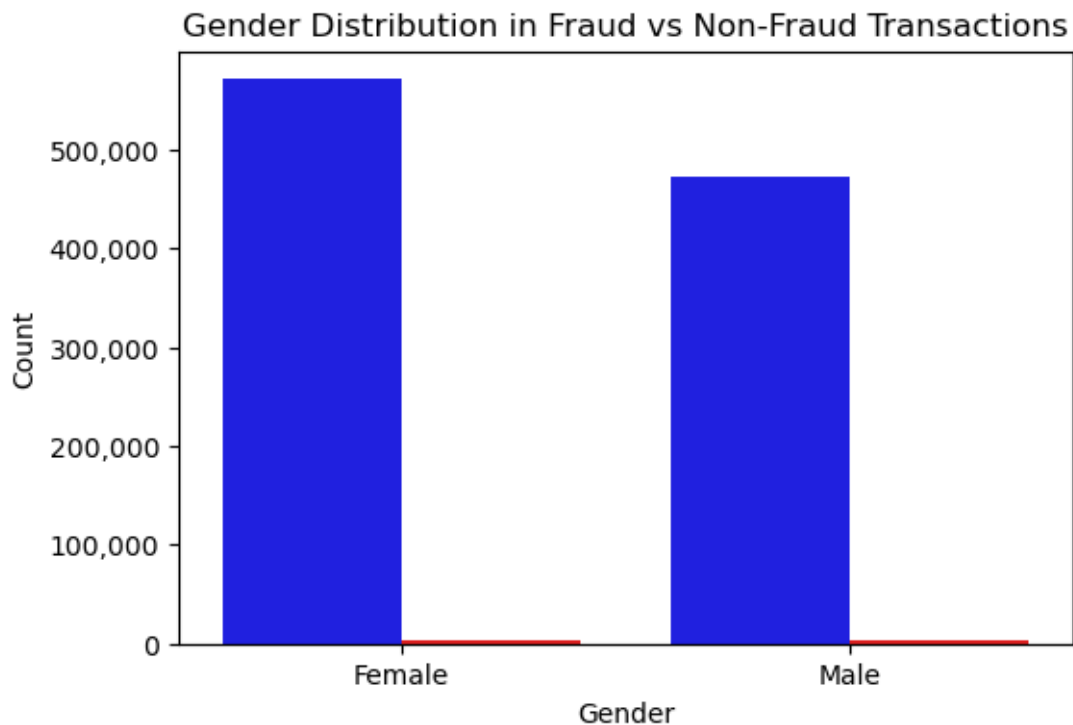
Gender Distribution in Fraud vs Non-Fraud Transactions

```
[154]: # Plot: Gender Distribution in Fraud vs Non-Fraud Transactions
plt.figure(figsize=(6, 4))
sns.countplot(x="gender", hue="is_fraud", data=df, palette=["blue", "red"],
             legend=False)

# Replace 'F' with 'Female' and 'M' with 'Male'
gender_labels = {'F': 'Female', 'M': 'Male'}
plt.xticks(ticks=[0, 1], labels=[gender_labels.get(x, x) for x in df['gender'].
             unique()])

# Format the y-axis with commas for readability
plt.gca().yaxis.set_major_formatter(plt.FuncFormatter(lambda x, loc: f'{int(x):
             },}'))

# Labels and title
plt.xlabel("Gender")
plt.ylabel("Count")
plt.title("Gender Distribution in Fraud vs Non-Fraud Transactions")
plt.show()
```



Explanation: The “Gender Distribution in Fraud vs Non-Fraud Transactions” graph shows the number of fraudulent transactions made by females and males. On the y-axis, the graph counts the number of transactions made by each gender, while the x-axis shows the gender. An interesting insight from the dataset is that the number of fraudulent transactions is nearly the same for both genders. However, it is evident that a higher number of Non-Fraud transactions were made by females. Given that, one might expect more fraudulent transactions from females due to the higher overall transaction count. Despite this, the graph suggests that males are more prone to commit credit card fraud.

Overview/Conclusion:

By conducting graphical analysis, I’ve gained several insights from the dataset. For instance, I initially believed that there would be a much higher number of fraudulent credit card transactions, but the data revealed that fraud is actually quite minimal, amounting to only a few thousand dollars. Additionally, the analysis showed that males tend to commit more credit card fraud.

Furthermore, the data partially supports the common belief that more fraudulent transactions occur online than at the point of sale. From the visualizations, we can see that the highest credit card fraud occurred in the “grocery_pos” and “shopping_net” categories. This suggests that in the “grocery_pos” category, fraudsters went to physical stores and used stolen cards to make purchases, while in the “shopping_net” category, fraudsters made purchases online.

Finally, the transaction amount graph highlighted that fraudulent transactions tend to involve smaller amounts, with the highest being \$1,371.81. This could suggest that fraudsters prefer to purchase cheaper items, possibly hoping that banks won’t notice the fraud because of the smaller transaction sizes.

2 DSC 550 Term Project - Milestone 2

2.0.1 Data Preparation for Modeling

In this section I will include all the necessary steps to completely prepare the credit card transactions dataset for model building.

This will include data transformation, feature engineering, removing duplicates, handling missing values, and addressing class imbalance.

2.0.2 Importing Necessary Libraries

In this step, I import the necessary libraries and load the dataset for data manipulation and analysis.

```
[155]: import pandas as pd
import numpy as np

# Loading the dataset
file_path = r"C:\Users\maxim\OneDrive\Desktop\BU\DSC 550\Term_1\
↳Project\credit_card_transactions.csv"
df = pd.read_csv(file_path)
```

2.0.3 Dropping Unnecessary Features

In this step, I am dropping columns that are not useful for modeling, such as ‘unix_time’, or any redundant identifiers.

```
[156]: # Dropping unnecessary columns/features not useful for modeling (e.g.,  
        ↪ 'unix_time' or redundant identifiers)  
df.drop(['unix_time'], axis=1, inplace=True)  
  
# Displaying first two rows to confirm the unnecessary columns/features dropped  
df.head(2)
```

```
[156]: Unnamed: 0 trans_date_trans_time      cc_num \  
0          0          1/1/2019 0:00  2.703190e+15  
1          1          1/1/2019 0:00  6.304230e+11  
  
        merchant      category      amt      first      last \  
0      fraud_Rippin, Kub and Mann      misc_net      4.97      Jennifer      Banks  
1  fraud_Heller, Gutmann and Zieme      grocery_pos      107.23      Stephanie      Gill  
  
        gender      street      ...      lat      long      city_pop \  
0          F          561 Perry Cove      ...      36.0788      -81.1781          3495  
1          F  43039 Riley Greens Suite 393      ...      48.8878      -118.2105          149  
  
        job      dob \  
0      Psychologist, counselling      3/9/1988  
1  Special educational needs teacher      6/21/1978  
  
        trans_num      merch_lat      merch_long      is_fraud \  
0  0b242abb623afc578575680df30655b9      36.011293      -82.048315          0  
1  1f76529f8574734946361c461b024d99      49.159047      -118.186462          0  
  
        merch_zipcode  
0          28705.0  
1          NaN  
  
[2 rows x 23 columns]
```

2.0.4 Handling Missing Values

This step handles missing values in the dataset. I replace “?” with NaN and drop rows with missing data.

Alternatively, imputation can be used for missing data, but here I drop them.

```
[157]: # Checking and handling missing values appropriately  
df.replace("?", np.nan, inplace=True)  
df = df.dropna()
```



```
# Displaying first two rows to confirm the missing values handled appropriately
df.head(2)
```

```
[157]: Unnamed: 0 trans_date_trans_time cc_num merchant \
0 0 1/1/2019 0:00 2.703190e+15 fraud_Rippin, Kub and Mann
2 2 1/1/2019 0:00 3.885950e+13 fraud_Lind-Buckridge

category amt first last gender street \
0 misc_net 4.97 Jennifer Banks F 561 Perry Cove
2 entertainment 220.11 Edward Sanchez M 594 White Dale Suite 530

... lat long city_pop job dob \
0 ... 36.0788 -81.1781 3495 Psychologist, counselling 3/9/1988
2 ... 42.1808 -112.2620 4154 Nature conservation officer 1/19/1962

trans_num merch_lat merch_long is_fraud \
0 0b242abb623afc578575680df30655b9 36.011293 -82.048315 0
2 a1a22d70485983eac12b5b88dad1cf95 43.150704 -112.154481 0

merch_zipcode
0 28705.0
2 83236.0

[2 rows x 23 columns]
```

2.0.5 Data Type Conversion

I ensure that the data types of features are correct for modeling, specifically converting the 'amt' column to a float.

```
[158]: # Ensuring correct data types for modeling
df['amt'] = df['amt'].astype(float)

# Displaying the first two rows to confirm the 'amt' column converted to a float
df.head(2)
```

```
[158]: Unnamed: 0 trans_date_trans_time cc_num merchant \
0 0 1/1/2019 0:00 2.703190e+15 fraud_Rippin, Kub and Mann
2 2 1/1/2019 0:00 3.885950e+13 fraud_Lind-Buckridge

category amt first last gender street \
0 misc_net 4.97 Jennifer Banks F 561 Perry Cove
2 entertainment 220.11 Edward Sanchez M 594 White Dale Suite 530

... lat long city_pop job dob \
0 ... 36.0788 -81.1781 3495 Psychologist, counselling 3/9/1988
2 ... 42.1808 -112.2620 4154 Nature conservation officer 1/19/1962
```

```

            trans_num  merch_lat  merch_long  is_fraud  \
0  0b242abb623afc578575680df30655b9  36.011293  -82.048315      0
2  a1a22d70485983eac12b5b88dad1cf95  43.150704  -112.154481      0

    merch_zipcode
0          28705.0
2          83236.0

[2 rows x 23 columns]

```

2.0.6 Feature Engineering

In this step, I create new features from existing ones. Here, I create a ‘merchant_category’ feature by combining ‘merchant’ and ‘category’.

```

[159]: # Creating a new feature from 'merchant' or 'category'
df['merchant_category'] = df['merchant'] + '_' + df['category']

# Displaying output anfter creating new feature
df.head(2)

```

```

[159]: Unnamed: 0  trans_date_trans_time      cc_num      merchant  \
0          0      1/1/2019 0:00  2.703190e+15  fraud_Rippin, Kub and Mann
2          2      1/1/2019 0:00  3.885950e+13      fraud_Lind-Buckridge

            category    amt    first    last gender      street  \
0      misc_net    4.97  Jennifer    Banks      F      561 Perry Cove
2  entertainment  220.11    Edward  Sanchez      M  594 White Dale Suite 530

...      long city_pop      job      dob  \
0 ... -81.1781    3495  Psychologist, counselling  3/9/1988
2 ... -112.2620    4154  Nature conservation officer  1/19/1962

            trans_num  merch_lat  merch_long  is_fraud  \
0  0b242abb623afc578575680df30655b9  36.011293  -82.048315      0
2  a1a22d70485983eac12b5b88dad1cf95  43.150704  -112.154481      0

    merch_zipcode      merchant_category
0      28705.0  fraud_Rippin, Kub and Mann_misc_net
2      83236.0  fraud_Lind-Buckridge_entertainment

[2 rows x 24 columns]

```

2.0.7 Encoding Categorical Variables

I convert categorical columns into dummy variables, which is necessary for many machine learning algorithms.

```
[160]: # Converting categorical columns to dummy variables
df = pd.get_dummies(df, columns=['gender', 'category'], drop_first=True)

# Displaying the first two rows to confirm the conversion
df.head(2)
```

```
[160]: Unnamed: 0  trans_date_trans_time      cc_num      merchant \
0          0      1/1/2019 0:00  2.703190e+15  fraud_Rippin, Kub and Mann
2          2      1/1/2019 0:00  3.885950e+13      fraud_Lind-Buckridge

      amt  first  last      street      city state \
0   4.97  Jennifer  Banks      561 Perry Cove  Moravian Falls  NC
2  220.11  Edward  Sanchez  594 White Dale Suite 530      Malad City  ID

...  category_grocery_pos  category_health_fitness  category_home \
0  ...                  False                  False          False
2  ...                  False                  False          False

      category_kids_pets  category_misc_net  category_misc_pos \
0                  False                  True          False
2                  False                  False          False

      category_personal_care  category_shopping_net  category_shopping_pos \
0                  False                  False          False
2                  False                  False          False

      category_travel
0                  False
2                  False

[2 rows x 36 columns]
```

2.0.8 Displaying the Cleaned and Transformed Dataset

Here I display the first few rows of the cleaned and transformed dataset to confirm all the changes.

```
[161]: # Printing the cleaned and transformed dataset
df.head()
```

```
[161]: Unnamed: 0  trans_date_trans_time      cc_num \
0          0      1/1/2019 0:00  2.703190e+15
2          2      1/1/2019 0:00  3.885950e+13
4          4      1/1/2019 0:03  3.755340e+14
5          5      1/1/2019 0:04  4.767270e+15
7          7      1/1/2019 0:05  6.011360e+15

      merchant      amt  first  last \
0  fraud_Rippin, Kub and Mann  4.97  Jennifer  Banks
```

2	fraud_Lind-Buckridge	220.11	Edward	Sanchez
4	fraud_Keeling-Crist	41.96	Tyler	Garcia
5	fraud_Stroman, Hudson and Erdman	94.63	Jennifer	Conner
7	fraud_Corwin-Collins	71.65	Steven	Williams

	street	city	state	...	category_grocery_pos \
0	561 Perry Cove	Moravian Falls	NC	...	False
2	594 White Dale Suite 530	Malad City	ID	...	False
4	408 Bradley Rest	Doe Hill	VA	...	False
5	4655 David Island	Dublin	PA	...	False
7	231 Flores Pass Suite 720	Edinburg	VA	...	False

	category_health_fitness	category_home	category_kids_pets \
0	False	False	False
2	False	False	False
4	False	False	False
5	False	False	False
7	False	False	False

	category_misc_net	category_misc_pos	category_personal_care \
0	True	False	False
2	False	False	False
4	False	True	False
5	False	False	False
7	False	False	False

	category_shopping_net	category_shopping_pos	category_travel
0	False	False	False
2	False	False	False
4	False	False	False
5	False	False	False
7	False	False	False

[5 rows x 36 columns]