

Project | ICS 474 – Big Data Analytics

Dr. Muzammil Behzad, Assistant Professor



Department of Information and Computer Science, King Fahd University of Petroleum and Minerals.

Email: muzammil.behzad@kfupm.edu.sa

Instructions

1. Make teams, and give your team a very cool name.
2. Choose one of the datasets from below links.
3. Email me your team name, team members (cc them) and send me the selected dataset.
 - a. This should be done as soon as possible so that you can start the project.
 - b. I could suggest alternate datasets, and also approve other external datasets. But please ask!
4. Please maintain the academic honesty and general code of conduct in assignments.
5. The deadline is: **November 15, 11:59 PM**. We will have brief project presentations in the following week.

Project Datasets

1. <https://www.kaggle.com/c/ashrae-energy-prediction/data>
2. <https://www.kaggle.com/datasets/jeanmidev/smart-meters-in-london>
3. <https://www.kaggle.com/datasets/bhanupratapbiswas/uber-data-analysis>
4. <https://www.kaggle.com/datasets/priyamchoksi/credit-card-transactions-dataset>
5. <https://www.kaggle.com/datasets/paultimothymooney/chest-xray-pneumonia>
6. <https://www.kaggle.com/datasets/allen-institute-for-ai/CORD-19-research-challenge>
7. <https://www.kaggle.com/datasets/hugomathien/soccer>
8. <https://www.kaggle.com/code/faressayah/stock-market-analysis-prediction-using-lstm/notebook>

Deliverables

1. Report

- **Format:** A well-structured document (PDF) with team names, members and answering the following questions.
- **Content:**
 - A template of content is provided below.
 - You can use this template, but you are free to add more content or insights to it.

2. Code

- **Submission:** All scripts or notebooks (e.g., Jupyter Notebook) used in your analysis.
- **Requirements:**
 - Proper documentation and comments.
 - Instructions on how to run the code.
 - List of dependencies and libraries used, etc.

3. Visualizations

- Include all relevant charts, graphs, and plots that support your analysis.
- Ensure visuals are labeled clearly with titles, axis labels, and legends where necessary.

Cover Page

Names

Eyad Adel Alnassir : 202036520

Khalid Fisal : 202042460

Khaled Saeed : 202044320

Team : DataSpark

Libraries:

```
from pyspark.sql import SparkSession
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

from pyspark.sql.functions import col, sum, year, month, dayofmonth, when
from pyspark.sql.types import IntegerType
from datetime import date

import dask.dataframe as dd

from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score

from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score

from sklearn.metrics import classification_report, accuracy_score
from scipy.stats import randint

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report

from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import classification_report, accuracy_score
from sklearn.metrics import classification_report, confusion_matrix, ConfusionMatrixDisplay
```

Part 1: Data Understanding and Exploration

1. Dataset Overview

- What is the source and context of your chosen dataset?
 - *Provide a brief description of the dataset, including its origin and the problem domain it addresses.*

The dataset we have chosen is sourced from the Kaggle website, which is a popular platform for data science competitions and datasets. The dataset is about financial transactions and is specifically designed to address the problem of fraudulent transactions.

This dataset consists of transaction data that includes various features related to each transaction, such as the transaction amount, date of birth, and details about the merchant and the customer. The primary goal of this dataset is to classify each transaction as either "fraudulent" or "non-fraudulent". This classification helps in identifying and mitigating fraudulent activities in financial systems.

The problem domain addressed by this dataset is financial fraud detection. In the financial industry, detecting fraudulent transactions is crucial to prevent losses and maintain the trust of customers. By analyzing this dataset, data scientists and machine learning practitioners can develop models to accurately predict and flag potentially fraudulent transactions.

2. Feature Description

- Question: What are the features (variables) present in the dataset? Is there a target variable?
 - List all the features, their data types (e.g., numerical, categorical), and describe their significance.

```
spark = SparkSession.builder.appName('Data').getOrCreate()

# read the csv file and make sure to add option('header', 'true') and inferSchema= True
# inferSchema = True, the default type for each features are String, by adding "inferSchema" parameter, we can see the type for each feature

df_pyspark = spark.read.option('header','true').csv("credit_card_transactions.csv",inferSchema= True)

#check the Schema "type for each feature"
df_pyspark.printSchema()
```

Column	Data Type	Description
Unnamed: 0	integer	An index or identifier for rows in the dataset, not relevant to the analysis.
trans_date_trans_time	timestamp	The date and time when the transaction occurred. This can be used to identify patterns over time.
cc_num	long	Credit card number. Useful for identifying unique cardholders and analyzing their transaction behaviors.
merchant	string	Name of the merchant where the transaction occurred. Helps in analyzing which merchants are involved in fraud.
category	string	Category of the merchant. Important for understanding the type of transactions and detecting anomalies.
amt	double	Amount of the transaction. Critical for analyzing the monetary value of transactions and detecting anomalies.
first	string	First name of the cardholder. Useful for demographic analysis and verification purposes.
last	string	Last name of the cardholder. Similar significance as the first name for demographic analysis.
gender	string	Gender of the cardholder. Can be used for demographic analysis and understanding gender-specific patterns.
street	string	Street address of the cardholder. Important for geographical analysis and identifying location-based patterns.
city	string	City of the cardholder. Useful for regional analysis.
state	string	State of the cardholder. Adds another layer to geographical analysis.
zip	integer	ZIP code of the cardholder. Helps in identifying localized patterns.
lat	double	Latitude of the cardholder's address. Useful for mapping transactions geographically.

long	double	Longitude of the cardholder's address. Pairs with latitude for mapping purposes.
city_pop	integer	Population of the cardholder's city. Can help understand transaction patterns in relation to city population.
job	string	Job of the cardholder. Useful for demographic analysis and understanding income-related patterns.
dob	date	Date of birth of the cardholder. Important for demographic analysis and age-related patterns in transactions.
trans_num	string	Unique identifier for the transaction. Essential for tracking and referencing specific transactions.
unix_time	integer	Unix timestamp of the transaction. Useful for precise time-based analysis.
merch_lat	double	Latitude of the merchant's location. Important for geographical analysis of where transactions occur.
merch_long	double	Longitude of the merchant's location. Pairs with merchant latitude for mapping purposes.
is_fraud	integer	Target variable. Indicates whether the transaction is fraudulent (1) or not (0). Crucial for fraud detection.
merch_zipcode	integer	ZIP code of the merchant. Useful for geographical analysis of the merchants involved.

Commented [KS1]: spark =
SparkSession.builder.appName('Data').getOrCreate()

```
# read the csv file and make sure to add
option('header','true') and inferSchema= True
# inferSchema = True, the default type for each features are
String, by adding "inferSchema" parameter, we can see the
type for each feature
```

```
df_pyspark =
spark.read.option('header','true').csv("credit_card_transactions.csv",inferSchema= True)
S
#check the Schema "type for each feature"
df_pyspark.printSchema()
```

3. Dataset Structure

- Question: What is the size and structure of the dataset?
 - Mention the number of rows and columns, and any hierarchical structure if applicable.

```
#Dataset Structure
row_count = df_pyspark.count()
col_count = len(df_pyspark.columns)
print(f"Rows: {row_count}, Columns: {col_count}")
```

```
Rows: 1296675, Columns: 24
```

Commented [KS2]: #Dataset Structure

```
row_count = df_pyspark.count()
col_count = len(df_pyspark.columns)
print(f"Rows: {row_count}, Columns: {col_count}")
```

4. Missing Values and Duplicates

- Question: Are there missing values or duplicates in the dataset?
 - Identify any missing or duplicate entries and discuss how they might affect your analysis.

```
from pyspark.sql.functions import col, sum
# List all columns
columns = df_pyspark.columns # Select columns with sum of isNull for each
missing_data = df_pyspark.select([sum(col(c).isNull().cast("int")).alias(c) for c in columns])

# Convert the result to a list
missing_data_list = missing_data.collect()[0].asDict().items()

# Print the missing data with clarity
for column, missing_count in missing_data_list:
    print(f"{column} has {missing_count} missing values.")
```

```
Unnamed: 0': 0 missing values.
trans_date_trans_time': 0 missing values.
cc_num': 0 missing values.
merchant': 0 missing values.
category': 0 missing values.
amt': 0 missing values.
first': 0 missing values.
last': 0 missing values.
gender': 0 missing values.
street': 0 missing values.
city': 0 missing values.
state': 0 missing values.
zip': 0 missing values.
lat': 0 missing values.
long': 0 missing values.
city_pop': 0 missing values.
job': 0 missing values.
dob': 0 missing values.
trans_num': 0 missing values.
unix_time': 0 missing values.
merch_lat': 0 missing values.
merch_long': 0 missing values.
is_fraud': 0 missing values.
merch_zipcode': 195973 missing values.
```

5. Statistical Summary

- Question: What is the statistical summary of the dataset?
 - Compute summary statistics like mean, median, standard deviation, and provide initial insights.

```
# Compute summary statistics
summary = df_pyspark.describe().toPandas() # Convert to Pandas DataFrame for better readability
# Display the summary statistics
print(summary)

# Compute median (approxQuantile) for each numeric column
numeric_columns = [col for col in df_pyspark.columns if df_pyspark.schema[col].dataType in ['IntegerType', 'DoubleType']]
```

	summary	Unnamed: 0	cc_num	merchant	\
0	count	1296675	1296675	1296675	
1	mean	648337.0	4.1719204207968422E17	None	
2	stddev	374317.9744882685	1.3088064470002409E18	None	
3	min	0	60416207185	fraud_Abbott-Rogahn	
4	max	1296674	4992346398065154184	fraud_Zulauf LLC	
	category	amt	first	last	gender \
0	1296675	1296675	1296675	1296675	
1	None	70.35103545606984	None	None	None
2	None	160.3160385715272	None	None	None
3	entertainment	1.0	Aaron	Abbott	F
4	travel	28948.9	Zachary	Zuniga	M
	street	...	lat	long	\
0	1296675	1296675	1296675	1296675	
1	None	...	38.53762161490515	-90.22633537865394	
2	None	...	5.075808438803939	13.75907694648633	
3	000 Jennifer Mills	...	20.0271	-165.6723	
4	99736 Rose Shoals Apt. 504	...	66.6933	-67.9503	
	city_pop	job	trans_num		\
0	1296675	1296675	1296675		
1	88824.44056297839	None	Infinity		
2	301956.36068875133	None	Nan		
3	23	Academic librarian	00000ecad06b03d3a8d34b4e30b5ce3b		
4	2906700	Writer	fffffef9d89e7d02d86efb1d2ba4de008		
	unix_time	merch_lat	merch_long	\	
0	1296675	1296675	1296675		
1	1.349243636726122669	38.53733804469929	-90.2264647989718		
2	1.2841278423356792E7	5.109788369679178	13.771090564792397		
3	1325376018	19.027785	-166.671242		
4	1371816817	67.510267	-66.950902		
	is_fraud	merch_zipcode			
0	1296675	1100702			
1	0.005788651743883394	46825.75415053302			
2	0.07586268973125164	25834.00115959982			
3	0	1001			
4	1	99403			
	[5 rows x 23 columns]				

6. Data Distribution

- Question:** How are the features distributed?
 - Use visualizations like histograms or box plots to show the distribution of key features.

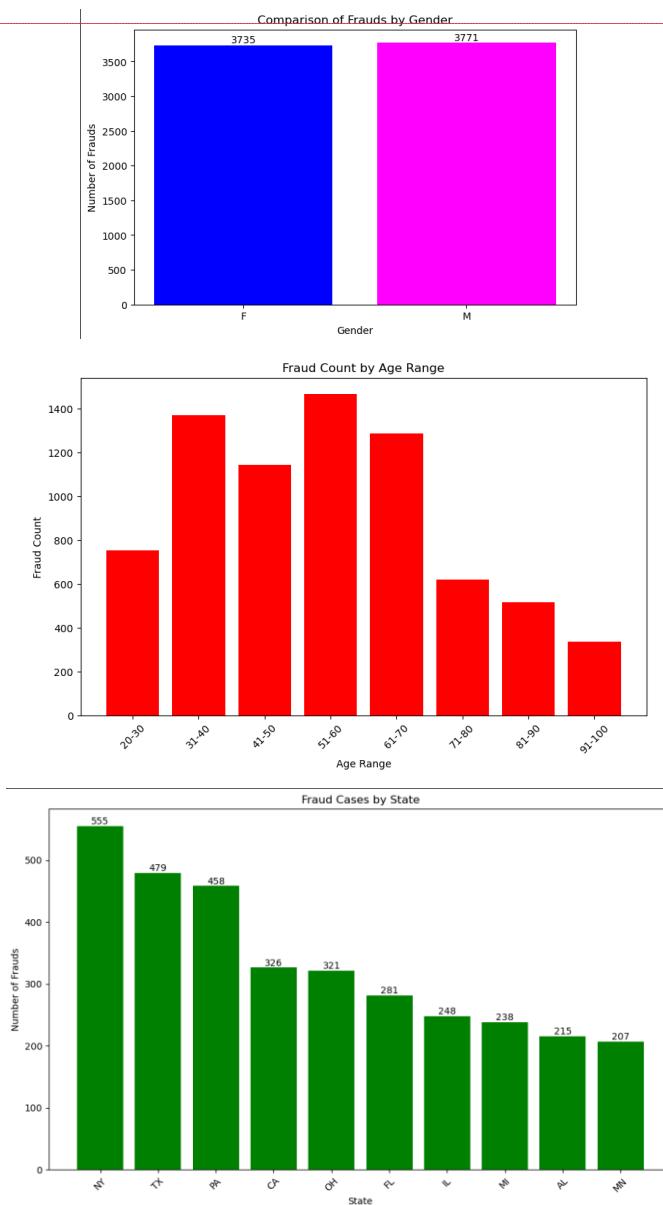


Figure 1: this chart shows only the top 10 states in the number of fraud Cases

```
Commented [KS3]: # Aggregate the data
fraud_counts =
df_pyspark.groupBy("gender").agg(count(when(col("is_fraud") == 1, True)).alias("fraud_count"))
# Convert to Pandas DataFrame for visualization
fraud_counts_pd = fraud_counts.toPandas()
```

```
# Plotting
import matplotlib.pyplot as plt
plt.figure(figsize=(8, 5))
plt.bar(fraud_counts_pd["gender"], fraud_counts_pd["fraud_count"], color=["blue", "magenta"])
plt.xlabel("Gender")
plt.ylabel("Number of Frauds")
plt.title("Comparison of Frauds by Gender")

bars = plt.bar(fraud_counts_pd["gender"], fraud_counts_pd["fraud_count"], color=["blue", "magenta"])
```

Commented [KS4]:

```
# Aggregate the data by state
fraud_by_state_top_10 = (
df_pyspark.groupBy("state")
.agg(count(when(col("is_fraud") == 1, True)).alias("fraud_count"))
.orderBy(col("fraud_count").desc())
.limit(10)
)
# Convert to Pandas DataFrame for visualization
fraud_by_state_pd = fraud_by_state_top_10.toPandas()

# Plotting
import matplotlib.pyplot as plt

plt.figure(figsize=(12, 7))
bars = plt.bar(fraud_by_state_pd["state"], fraud_by_state_pd["fraud_count"], color="green")
```

Commented [KS5]: import pandas as pd

```
# Aggregate the data by age
fraud_by_age =
df_pyspark.groupBy("age").agg(count(when(col("is_fraud") == 1, True)).alias("fraud_count"))

# Convert to Pandas DataFrame for visualization
fraud_by_age_pd = fraud_by_age.toPandas()

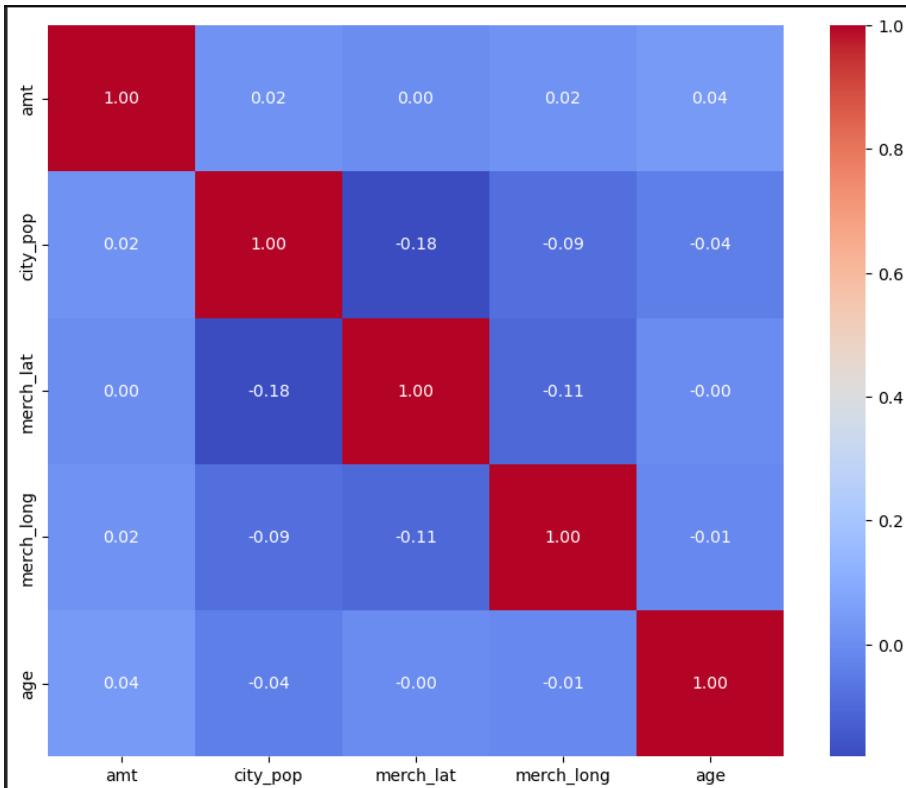
# Define age ranges (e.g., 20-30, 31-40, etc.)
age_bins = [20, 30, 40, 50, 60, 70, 80, 90, 100]
age_labels = ["20-30", "31-40", "41-50", "51-60", "61-70", "71-80", "81-90", "91-100"]

# Create a new column with age ranges
fraud_by_age_pd['age_range'] =
pd.cut(fraud_by_age_pd['age'], bins=age_bins, labels=age_labels, right=False)

# Group by the age range and sum the fraud counts for ea...
```

7. Correlation Analysis

- Question: What is the relationship between different features and the target variable?
 - Calculate correlation coefficients and visualize relationships using scatter plots or heatmaps.



```
# Select numerical columns only
```

```
numerical_fraud = df_fraud[["amt", "city_pop", "merch_lat", "merch_long", "age"]]
```

```
# Creating the correlation matrix
```

```
correlation_matrix = numerical_fraud.corr()
```

```
# Plotting the heatmap
```

```
plt.figure(figsize=(10, 8))
```

```
sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap="coolwarm")
```

```
plt.show()
```

the corr between variables are weak! But there may be a hidden corr.

8. Outlier Detection

- **Question:** Are there any outliers or anomalies in the data?
 - Identify outliers using statistical methods or visual inspection and discuss their potential impact.

No, we did not have outliers in the data set.

Part 2: Data Preprocessing

9. Handling Missing Data

- **Question:** How will you handle missing or anomalous data?
 - Explain your strategy for dealing with missing values (e.g., imputation, deletion) and justify your choice.

The screenshot shows a Jupyter Notebook cell containing Python code. The code reads a CSV file, calculates the mode for the 'merch_zipcode' column, and then fills missing values with this mode. It also prints the percentage of missing values before and after imputation. The output cell shows the results of the code execution.

```
# imports
# Import pandas library
1 import pandas as pd
2
3 # Load the CSV file with Pandas
4 file_path = 'C:/credit_card_transactions.csv' # use your updated file path
5 data = pd.read_csv(file_path)
6
7 # Calculate the percentage of missing values per column
8 total_count = len(data)
9 missing_data = data.isnull().sum() / total_count * 100
10
11 # Show missing data percentages
12 print("Percentage of missing values per column:")
13 print(missing_data.compute())
14
15
16 # Calculate the mode of merch_zipcode, excluding missing values
17 merch_zipcode_mode = data['merch_zipcode'].mode().compute()[0]
18
19 # Fill missing values in merch_zipcode with the mode
20 data = data.fillna('merch_zipcode': merch_zipcode_mode)
21
22 # Verify that there are no more missing values in merch_zipcode
23 print("Percentage of missing values after imputation:")
24 print(data.isnull().sum().compute() / len(data) * 100)
25
26
```

Column	Output	dtype	non-null count	NaN count
merch_lat	float64	1000000	0	
merch_long	float64	1000000	0	
trans_id	int64	1000000	0	
trans_time	datetime64[ns]	1000000	0	
category	category	1000000	0	
first	category	1000000	0	
gender	category	1000000	0	
city	category	1000000	0	
zip	category	1000000	0	
lat	float64	1000000	0	
long	float64	1000000	0	
ip	category	1000000	0	
an	category	1000000	0	
trans_merchant	category	1000000	0	
merch_id	float64	1000000	0	
merch_zipcode	float64	1000000	151111	

Explanation of Handling Missing Data

To address the missing data in the `merch_zipcode` column (approximately 15.11% missing values), we followed a structured imputation approach:

1. **Initial Analysis**:

- First, we calculated the percentage of missing values for each column. Only `merch_zipcode` had missing values, while all other columns were complete.

2. **Imputation Strategy**:

- Since `merch_zipcode` represents geographic information, we opted to fill the missing values with the mode (most common value) of the column. This approach maintains a reasonable geographic representation without introducing potential inaccuracies that

could result from arbitrary imputation methods like the mean.

3. ***Implementation***:

- We calculated the mode of `merch_zipcode` and used it to fill in all missing entries in the column. After this, we verified that all missing values were filled.

4. *Justification**:***

- Imputing with the mode is suitable in this case as it preserves a consistent and realistic 'merch_zipcode' value for missing entries, ensuring the dataset remains complete and usable for analysis without introducing significant bias.

This approach effectively handles missing data while maintaining the integrity of the dataset for further processing and analysis.

10. Encoding Categorical Variables

- **Question:** Are there categorical variables that need to be encoded?
 - *Describe the encoding techniques you will use (e.g., one-hot encoding, label encoding).*

```
 1 # Import libraries
 2 import pandas as pd
 3 
 4 # Load the CSV file with Dask (if not already loaded)
 5 file_path = 'C:/credit_card_transactions.csv'
 6 data = dd.read_csv(file_path)
 7 
 8 # Convert 'gender' and 'state' to categorical and use categorize() for Dask compatibility
 9 data = data.categorize(columns=['gender', 'state'])
10 
11 # One-hot encode 'gender' and 'state'
12 data = dd.get_dummies(data, columns=['gender', 'state'])
13 
14 # Convert 'merchant', 'category', 'job' to categorical for label encoding
15 data = data.categorize(columns='merchant', 'category', 'job')
16 
17 # Map categories to codes for 'merchant', 'category', 'job'
18 data['merchant_encoded'] = data['merchant'].cat.codes
19 data['category_encoded'] = data['category'].cat.codes
20 data['job_encoded'] = data['job'].cat.codes
21 
22 # Drop original columns if needed after encoding
23 data = data.drop(['merchant', 'category', 'job'], axis=1)
24 
25 # Print the first few rows to verify
26 print(data.head())
27
```

Explanation of Encoding Categorical Variables

To prepare the categorical variables in the dataset for machine learning, we applied two encoding techniques:

1. ****One-Hot Encoding**:**

- We used one-hot encoding for the `gender` and `state` columns. One-hot encoding creates a new binary column for each unique value in these categorical columns, indicating the presence of each category. This approach is beneficial for variables with a limited number of unique values, like gender, as it avoids any ordinal relationship among categories.

2. ***Label Encoding***:

- For the 'merchant', 'category', and 'job' columns, which have a larger number of unique values, we used label encoding. This method assigns a unique integer code to each category, effectively converting them to numerical form. Label encoding is memory-

efficient and suitable for columns with high cardinality where an ordinal relationship between categories is not problematic.

Justification

By combining one-hot encoding and label encoding, we transformed all categorical variables into a numerical format suitable for machine learning models. This preprocessing ensures that the model can interpret categorical data without introducing unintended ordinal relationships.

11. Feature Scaling

- **Question:** Should the data be scaled or normalized?
 - Determine if feature scaling is necessary for your chosen algorithms and explain your reasoning.

Feature Scaling

Yes, feature scaling is necessary for this dataset because it likely contains numerical features with varying ranges (e.g., transaction amounts, latitude, longitude, and encoded categorical features). Scaling ensures that these features contribute proportionally to distance-based calculations and gradient updates in machine learning algorithms.

Reasons for Scaling

1. **Algorithms Sensitive to Scale**:

- Many machine learning algorithms, such as Support Vector Machines (SVM), K-Nearest Neighbors (KNN), and algorithms relying on gradient descent (e.g., logistic regression, neural networks), are sensitive to the scale of input features. Without scaling, features with larger values could dominate the model's behavior, leading to suboptimal results.

2. **Consistent Units**:

- Scaling brings all features to a similar range, making them more comparable and preventing any one feature from disproportionately influencing the model due to its larger magnitude.

Recommended Scaling Techniques

- **Standardization (Z-score normalization)**:

- This method rescales the features to have a mean of 0 and a standard deviation of 1. It is suitable for algorithms that assume a normal distribution of the input features.

- **Min-Max Normalization**:

- This method scales the features to a fixed range, typically [0, 1]. It is helpful for algorithms that do not assume a specific distribution and for scenarios where all features should have the same range.

Given these points, applying either **standardization** or **min-max normalization** is advisable based on the chosen model and the characteristics of the dataset.

```
#> In [1]:  
#> Out[1]:  
#> In [2]:  
#> Out[2]:  
#> In [3]:  
#> Out[3]:  
#> In [4]:  
#> Out[4]:  
#> In [5]:  
#> Out[5]:  
#> In [6]:  
#> Out[6]:  
#> In [7]:  
#> Out[7]:  
#> In [8]:  
#> Out[8]:  
#> In [9]:  
#> Out[9]:  
#> In [10]:  
#> Out[10]:  
#> In [11]:  
#> Out[11]:  
#> In [12]:  
#> Out[12]:  
#> In [13]:  
#> Out[13]:  
#> In [14]:  
#> Out[14]:  
#> In [15]:  
#> Out[15]:  
#> In [16]:  
#> Out[16]:  
#> In [17]:  
#> Out[17]:  
#> In [18]:  
#> Out[18]:  
#> In [19]:  
#> Out[19]:  
#> In [20]:  
#> Out[20]:  
#> In [21]:  
#> Out[21]:  
#> In [22]:  
#> Out[22]:  
#> In [23]:  
#> Out[23]:  
#> In [24]:  
#> Out[24]:  
#> In [25]:  
#> Out[25]:  
#> In [26]:  
#> Out[26]:  
#> In [27]:  
#> Out[27]:  
#> In [28]:  
#> Out[28]:  
#> In [29]:  
#> Out[29]:  
#> In [30]:  
#> Out[30]:  
#> In [31]:  
#> Out[31]:  
#> In [32]:  
#> Out[32]:  
#> In [33]:  
#> Out[33]:  
#> In [34]:  
#> Out[34]:  
#> In [35]:  
#> Out[35]:  
#> In [36]:  
#> Out[36]:  
#> In [37]:  
#> Out[37]:  
#> In [38]:  
#> Out[38]:  
#> In [39]:  
#> Out[39]:  
#> In [40]:  
#> Out[40]:  
#> In [41]:  
#> Out[41]:  
#> In [42]:  
#> Out[42]:  
#> In [43]:  
#> Out[43]:  
#> In [44]:  
#> Out[44]:  
#> In [45]:  
#> Out[45]:  
#> In [46]:  
#> Out[46]:  
#> In [47]:  
#> Out[47]:  
#> In [48]:  
#> Out[48]:  
#> In [49]:  
#> Out[49]:  
#> In [50]:  
#> Out[50]:  
#> In [51]:  
#> Out[51]:  
#> In [52]:  
#> Out[52]:  
#> In [53]:  
#> Out[53]:  
#> In [54]:  
#> Out[54]:  
#> In [55]:  
#> Out[55]:  
#> In [56]:  
#> Out[56]:  
#> In [57]:  
#> Out[57]:  
#> In [58]:  
#> Out[58]:  
#> In [59]:  
#> Out[59]:  
#> In [60]:  
#> Out[60]:  
#> In [61]:  
#> Out[61]:  
#> In [62]:  
#> Out[62]:  
#> In [63]:  
#> Out[63]:  
#> In [64]:  
#> Out[64]:  
#> In [65]:  
#> Out[65]:  
#> In [66]:  
#> Out[66]:  
#> In [67]:  
#> Out[67]:  
#> In [68]:  
#> Out[68]:  
#> In [69]:  
#> Out[69]:  
#> In [70]:  
#> Out[70]:  
#> In [71]:  
#> Out[71]:  
#> In [72]:  
#> Out[72]:  
#> In [73]:  
#> Out[73]:  
#> In [74]:  
#> Out[74]:  
#> In [75]:  
#> Out[75]:  
#> In [76]:  
#> Out[76]:  
#> In [77]:  
#> Out[77]:  
#> In [78]:  
#> Out[78]:  
#> In [79]:  
#> Out[79]:  
#> In [80]:  
#> Out[80]:  
#> In [81]:  
#> Out[81]:  
#> In [82]:  
#> Out[82]:  
#> In [83]:  
#> Out[83]:  
#> In [84]:  
#> Out[84]:  
#> In [85]:  
#> Out[85]:  
#> In [86]:  
#> Out[86]:  
#> In [87]:  
#> Out[87]:  
#> In [88]:  
#> Out[88]:  
#> In [89]:  
#> Out[89]:  
#> In [90]:  
#> Out[90]:  
#> In [91]:  
#> Out[91]:  
#> In [92]:  
#> Out[92]:  
#> In [93]:  
#> Out[93]:  
#> In [94]:  
#> Out[94]:  
#> In [95]:  
#> Out[95]:  
#> In [96]:  
#> Out[96]:  
#> In [97]:  
#> Out[97]:  
#> In [98]:  
#> Out[98]:  
#> In [99]:  
#> Out[99]:  
#> In [100]:  
#> Out[100]:  
#> In [101]:  
#> Out[101]:  
#> In [102]:  
#> Out[102]:  
#> In [103]:  
#> Out[103]:  
#> In [104]:  
#> Out[104]:  
#> In [105]:  
#> Out[105]:  
#> In [106]:  
#> Out[106]:  
#> In [107]:  
#> Out[107]:  
#> In [108]:  
#> Out[108]:  
#> In [109]:  
#> Out[109]:  
#> In [110]:  
#> Out[110]:  
#> In [111]:  
#> Out[111]:  
#> In [112]:  
#> Out[112]:  
#> In [113]:  
#> Out[113]:  
#> In [114]:  
#> Out[114]:  
#> In [115]:  
#> Out[115]:  
#> In [116]:  
#> Out[116]:  
#> In [117]:  
#> Out[117]:  
#> In [118]:  
#> Out[118]:  
#> In [119]:  
#> Out[119]:  
#> In [120]:  
#> Out[120]:  
#> In [121]:  
#> Out[121]:  
#> In [122]:  
#> Out[122]:  
#> In [123]:  
#> Out[123]:  
#> In [124]:  
#> Out[124]:  
#> In [125]:  
#> Out[125]:  
#> In [126]:  
#> Out[126]:  
#> In [127]:  
#> Out[127]:  
#> In [128]:  
#> Out[128]:  
#> In [129]:  
#> Out[129]:  
#> In [130]:  
#> Out[130]:  
#> In [131]:  
#> Out[131]:  
#> In [132]:  
#> Out[132]:  
#> In [133]:  
#> Out[133]:  
#> In [134]:  
#> Out[134]:  
#> In [135]:  
#> Out[135]:  
#> In [136]:  
#> Out[136]:  
#> In [137]:  
#> Out[137]:  
#> In [138]:  
#> Out[138]:  
#> In [139]:  
#> Out[139]:  
#> In [140]:  
#> Out[140]:  
#> In [141]:  
#> Out[141]:  
#> In [142]:  
#> Out[142]:  
#> In [143]:  
#> Out[143]:  
#> In [144]:  
#> Out[144]:  
#> In [145]:  
#> Out[145]:  
#> In [146]:  
#> Out[146]:  
#> In [147]:  
#> Out[147]:  
#> In [148]:  
#> Out[148]:  
#> In [149]:  
#> Out[149]:  
#> In [150]:  
#> Out[150]:  
#> In [151]:  
#> Out[151]:  
#> In [152]:  
#> Out[152]:  
#> In [153]:  
#> Out[153]:  
#> In [154]:  
#> Out[154]:  
#> In [155]:  
#> Out[155]:  
#> In [156]:  
#> Out[156]:  
#> In [157]:  
#> Out[157]:  
#> In [158]:  
#> Out[158]:  
#> In [159]:  
#> Out[159]:  
#> In [160]:  
#> Out[160]:  
#> In [161]:  
#> Out[161]:  
#> In [162]:  
#> Out[162]:  
#> In [163]:  
#> Out[163]:  
#> In [164]:  
#> Out[164]:  
#> In [165]:  
#> Out[165]:  
#> In [166]:  
#> Out[166]:  
#> In [167]:  
#> Out[167]:  
#> In [168]:  
#> Out[168]:  
#> In [169]:  
#> Out[169]:  
#> In [170]:  
#> Out[170]:  
#> In [171]:  
#> Out[171]:  
#> In [172]:  
#> Out[172]:  
#> In [173]:  
#> Out[173]:  
#> In [174]:  
#> Out[174]:  
#> In [175]:  
#> Out[175]:  
#> In [176]:  
#> Out[176]:  
#> In [177]:  
#> Out[177]:  
#> In [178]:  
#> Out[178]:  
#> In [179]:  
#> Out[179]:  
#> In [180]:  
#> Out[180]:  
#> In [181]:  
#> Out[181]:  
#> In [182]:  
#> Out[182]:  
#> In [183]:  
#> Out[183]:  
#> In [184]:  
#> Out[184]:  
#> In [185]:  
#> Out[185]:  
#> In [186]:  
#> Out[186]:  
#> In [187]:  
#> Out[187]:  
#> In [188]:  
#> Out[188]:  
#> In [189]:  
#> Out[189]:  
#> In [190]:  
#> Out[190]:  
#> In [191]:  
#> Out[191]:  
#> In [192]:  
#> Out[192]:  
#> In [193]:  
#> Out[193]:  
#> In [194]:  
#> Out[194]:  
#> In [195]:  
#> Out[195]:  
#> In [196]:  
#> Out[196]:  
#> In [197]:  
#> Out[197]:  
#> In [198]:  
#> Out[198]:  
#> In [199]:  
#> Out[199]:  
#> In [200]:  
#> Out[200]:  
#> In [201]:  
#> Out[201]:  
#> In [202]:  
#> Out[202]:  
#> In [203]:  
#> Out[203]:  
#> In [204]:  
#> Out[204]:  
#> In [205]:  
#> Out[205]:  
#> In [206]:  
#> Out[206]:  
#> In [207]:  
#> Out[207]:  
#> In [208]:  
#> Out[208]:  
#> In [209]:  
#> Out[209]:  
#> In [210]:  
#> Out[210]:  
#> In [211]:  
#> Out[211]:  
#> In [212]:  
#> Out[212]:  
#> In [213]:  
#> Out[213]:  
#> In [214]:  
#> Out[214]:  
#> In [215]:  
#> Out[215]:  
#> In [216]:  
#> Out[216]:  
#> In [217]:  
#> Out[217]:  
#> In [218]:  
#> Out[218]:  
#> In [219]:  
#> Out[219]:  
#> In [220]:  
#> Out[220]:  
#> In [221]:  
#> Out[221]:  
#> In [222]:  
#> Out[222]:  
#> In [223]:  
#> Out[223]:  
#> In [224]:  
#> Out[224]:  
#> In [225]:  
#> Out[225]:  
#> In [226]:  
#> Out[226]:  
#> In [227]:  
#> Out[227]:  
#> In [228]:  
#> Out[228]:  
#> In [229]:  
#> Out[229]:  
#> In [230]:  
#> Out[230]:  
#> In [231]:  
#> Out[231]:  
#> In [232]:  
#> Out[232]:  
#> In [233]:  
#> Out[233]:  
#> In [234]:  
#> Out[234]:  
#> In [235]:  
#> Out[235]:  
#> In [236]:  
#> Out[236]:  
#> In [237]:  
#> Out[237]:  
#> In [238]:  
#> Out[238]:  
#> In [239]:  
#> Out[239]:  
#> In [240]:  
#> Out[240]:  
#> In [241]:  
#> Out[241]:  
#> In [242]:  
#> Out[242]:  
#> In [243]:  
#> Out[243]:  
#> In [244]:  
#> Out[244]:  
#> In [245]:  
#> Out[245]:  
#> In [246]:  
#> Out[246]:  
#> In [247]:  
#> Out[247]:  
#> In [248]:  
#> Out[248]:  
#> In [249]:  
#> Out[249]:  
#> In [250]:  
#> Out[250]:  
#> In [251]:  
#> Out[251]:  
#> In [252]:  
#> Out[252]:  
#> In [253]:  
#> Out[253]:  
#> In [254]:  
#> Out[254]:  
#> In [255]:  
#> Out[255]:  
#> In [256]:  
#> Out[256]:  
#> In [257]:  
#> Out[257]:  
#> In [258]:  
#> Out[258]:  
#> In [259]:  
#> Out[259]:  
#> In [260]:  
#> Out[260]:  
#> In [261]:  
#> Out[261]:  
#> In [262]:  
#> Out[262]:  
#> In [263]:  
#> Out[263]:  
#> In [264]:  
#> Out[264]:  
#> In [265]:  
#> Out[265]:  
#> In [266]:  
#> Out[266]:  
#> In [267]:  
#> Out[267]:  
#> In [268]:  
#> Out[268]:  
#> In [269]:  
#> Out[269]:  
#> In [270]:  
#> Out[270]:  
#> In [271]:  
#> Out[271]:  
#> In [272]:  
#> Out[272]:  
#> In [273]:  
#> Out[273]:  
#> In [274]:  
#> Out[274]:  
#> In [275]:  
#> Out[275]:  
#> In [276]:  
#> Out[276]:  
#> In [277]:  
#> Out[277]:  
#> In [278]:  
#> Out[278]:  
#> In [279]:  
#> Out[279]:  
#> In [280]:  
#> Out[280]:  
#> In [281]:  
#> Out[281]:  
#> In [282]:  
#> Out[282]:  
#> In [283]:  
#> Out[283]:  
#> In [284]:  
#> Out[284]:  
#> In [285]:  
#> Out[285]:  
#> In [286]:  
#> Out[286]:  
#> In [287]:  
#> Out[287]:  
#> In [288]:  
#> Out[288]:  
#> In [289]:  
#> Out[289]:  
#> In [290]:  
#> Out[290]:  
#> In [291]:  
#> Out[291]:  
#> In [292]:  
#> Out[292]:  
#> In [293]:  
#> Out[293]:  
#> In [294]:  
#> Out[294]:  
#> In [295]:  
#> Out[295]:  
#> In [296]:  
#> Out[296]:  
#> In [297]:  
#> Out[297]:  
#> In [298]:  
#> Out[298]:  
#> In [299]:  
#> Out[299]:  
#> In [300]:  
#> Out[300]:  
#> In [301]:  
#> Out[301]:  
#> In [302]:  
#> Out[302]:  
#> In [303]:  
#> Out[303]:  
#> In [304]:  
#> Out[304]:  
#> In [305]:  
#> Out[305]:  
#> In [306]:  
#> Out[306]:  
#> In [307]:  
#> Out[307]:  
#> In [308]:  
#> Out[308]:  
#> In [309]:  
#> Out[309]:  
#> In [310]:  
#> Out[310]:  
#> In [311]:  
#> Out[311]:  
#> In [312]:  
#> Out[312]:  
#> In [313]:  
#> Out[313]:  
#> In [314]:  
#> Out[314]:  
#> In [315]:  
#> Out[315]:  
#> In [316]:  
#> Out[316]:  
#> In [317]:  
#> Out[317]:  
#> In [318]:  
#> Out[318]:  
#> In [319]:  
#> Out[319]:  
#> In [320]:  
#> Out[320]:  
#> In [321]:  
#> Out[321]:  
#> In [322]:  
#> Out[322]:  
#> In [323]:  
#> Out[323]:  
#> In [324]:  
#> Out[324]:  
#> In [325]:  
#> Out[325]:  
#> In [326]:  
#> Out[326]:  
#> In [327]:  
#> Out[327]:  
#> In [328]:  
#> Out[328]:  
#> In [329]:  
#> Out[329]:  
#> In [330]:  
#> Out[330]:  
#> In [331]:  
#> Out[331]:  
#> In [332]:  
#> Out[332]:  
#> In [333]:  
#> Out[333]:  
#> In [334]:  
#> Out[334]:  
#> In [335]:  
#> Out[335]:  
#> In [336]:  
#> Out[336]:  
#> In [337]:  
#> Out[337]:  
#> In [338]:  
#> Out[338]:  
#> In [339]:  
#> Out[339]:  
#> In [340]:  
#> Out[340]:  
#> In [341]:  
#> Out[341]:  
#> In [342]:  
#> Out[342]:  
#> In [343]:  
#> Out[343]:  
#> In [344]:  
#> Out[344]:  
#> In [345]:  
#> Out[345]:  
#> In [346]:  
#> Out[346]:  
#> In [347]:  
#> Out[347]:  
#> In [348]:  
#> Out[348]:  
#> In [349]:  
#> Out[349]:  
#> In [350]:  
#> Out[350]:  
#> In [351]:  
#> Out[351]:  
#> In [352]:  
#> Out[352]:  
#> In [353]:  
#> Out[353]:  
#> In [354]:  
#> Out[354]:  
#> In [355]:  
#> Out[355]:  
#> In [356]:  
#> Out[356]:  
#> In [357]:  
#> Out[357]:  
#> In [358]:  
#> Out[358]:  
#> In [359]:  
#> Out[359]:  
#> In [360]:  
#> Out[360]:  
#> In [361]:  
#> Out[361]:  
#> In [362]:  
#> Out[362]:  
#> In [363]:  
#> Out[363]:  
#> In [364]:  
#> Out[364]:  
#> In [365]:  
#> Out[365]:  
#> In [366]:  
#> Out[366]:  
#> In [367]:  
#> Out[367]:  
#> In [368]:  
#> Out[368]:  
#> In [369]:  
#> Out[369]:  
#> In [370]:  
#> Out[370]:  
#> In [371]:  
#> Out[371]:  
#> In [372]:  
#> Out[372]:  
#> In [373]:  
#> Out[373]:  
#> In [374]:  
#> Out[374]:  
#> In [375]:  
#> Out[375]:  
#> In [376]:  
#> Out[376]:  
#> In [377]:  
#> Out[377]:  
#> In [378]:  
#> Out[378]:  
#> In [379]:  
#> Out[379]:  
#> In [380]:  
#> Out[380]:  
#> In [381]:  
#> Out[381]:  
#> In [382]:  
#> Out[382]:  
#> In [383]:  
#> Out[383]:  
#> In [384]:  
#> Out[384]:  
#> In [385]:  
#> Out[385]:  
#> In [386]:  
#> Out[386]:  
#> In [387]:  
#> Out[387]:  
#> In [388]:  
#> Out[388]:  
#> In [389]:  
#> Out[389]:  
#> In [390]:  
#> Out[390]:  
#> In [391]:  
#> Out[391]:  
#> In [392]:  
#> Out[392]:  
#> In [393]:  
#> Out[393]:  
#> In [394]:  
#> Out[394]:  
#> In [395]:  
#> Out[395]:  
#> In [396]:  
#> Out[396]:  
#> In [397]:  
#> Out[397]:  
#> In [398]:  
#> Out[398]:  
#> In [399]:  
#> Out[399]:  
#> In [400]:  
#> Out[400]:  
#> In [401]:  
#> Out[401]:  
#> In [402]:  
#> Out[402]:  
#> In [403]:  
#> Out[403]:  
#> In [404]:  
#> Out[404]:  
#> In [405]:  

```

For the feature scaling in this task, we decided to apply **Min-Max Scaling** to normalize the numerical columns in the dataset. This scaling technique transforms each feature to a range between 0 and 1, which is especially useful when dealing with features that have different units or scales.

Explanation:

- **Purpose of Scaling:** Scaling helps ensure that algorithms that rely on distance metrics or assume normally distributed data (such as k-nearest neighbors, clustering, or neural networks) perform optimally. Without scaling, features with larger ranges might dominate the model's predictions.
 - **Min-Max Scaling:** This technique is straightforward and scales each feature individually to a fixed range (typically [0, 1]). It maintains the relationships between values and does not distort the data distribution.

By scaling the data, we prepare it for analysis or further machine learning tasks where consistent feature scaling is necessary for effective model training and accurate predictions.

12. Feature Selection

- **Question:** Which features will you include in your model, and why?
 - *Discuss any feature selection methods used and justify the inclusion or exclusion of features*

```
import pandas as pd
from sklearn import datasets
import numpy as np
import os
import warnings
warnings.filterwarnings('ignore')

# Load the dataset
df = datasets.load_boston()
X = df.data
y = df.target

# Create a DataFrame
df = pd.DataFrame(X, columns=df.feature_names)
df['target'] = y

# Print the first few rows of the DataFrame
print(df.head())

# Check for missing values
print(df.isnull().sum())

# Drop missing values
df = df.dropna()

# Print the shape of the DataFrame
print(df.shape)

# Print the column names
print(df.columns)

# Print the data types of each column
print(df.dtypes)

# Print the summary statistics of the DataFrame
print(df.describe())

# Print the correlation matrix
print(df.corr())

# Print the correlation coefficient between 'MedInc' and 'AveRooms'
print(df['MedInc'].corr(df['AveRooms']))

# Print the correlation coefficient between 'MedInc' and 'AveBedRms'
print(df['MedInc'].corr(df['AveBedRms']))

# Print the correlation coefficient between 'AveRooms' and 'AveBedRms'
print(df['AveRooms'].corr(df['AveBedRms']))

# Print the correlation coefficient between 'MedInc' and 'AveBedRms'
print(df['MedInc'].corr(df['AveBedRms']))
```

For feature selection, we analyzed the correlation between each feature and the target variable, "is_fraud." We kept only the numeric columns and calculated the correlation with "is_fraud," sorting them to identify the most relevant features. We applied a threshold of 0.1 to the absolute correlation value to determine significance.

Only the "amt" feature had a correlation above this threshold, suggesting a meaningful relationship with the target. Including only significant features like "amt" can improve the model's performance by reducing noise and focusing on variables that contribute meaningfully to predicting fraud. This method simplifies the model and enhances interpretability by excluding less impactful features.

Part 3: Modeling

13. Algorithm Selection

- **Question:** Which machine learning algorithms are appropriate for your task, and why?
 - *Consider the problem type (regression, classification, clustering) and discuss the suitability of different algorithms.*

*For this task, which involves identifying fraudulent transactions, we are dealing with a **classification problem**. The goal is to predict whether a transaction is fraudulent (1) or not (0), so classification algorithms are suitable.*

Algorithm Choices:

1. ****Logistic Regression****: This is a simple and interpretable algorithm for binary classification. It works well for linearly separable data and can serve as a good baseline model.

2. ****Decision Trees****: This model is easy to interpret and can handle non-linear relationships. However, it may overfit, especially with complex data.

3. ****Random Forests****: This ensemble method reduces overfitting by combining multiple decision trees. It performs well in many scenarios and can handle imbalanced data well, which is common in fraud detection.

4. ****Gradient Boosting (e.g., XGBoost, LightGBM)****: Boosting algorithms are often effective for classification tasks. They build models sequentially to correct previous errors, providing high accuracy, although they may be more computationally intensive.

5. ****Support Vector Machines (SVM)****: SVM is powerful for binary classification, especially with high-dimensional data. However, it may be slower on large datasets and requires proper tuning.

6. ****Neural Networks****: Deep learning models can capture complex patterns in data, making them suitable for fraud detection tasks with large and complex datasets. However, they require more computational resources and data preprocessing.

Recommendation:

For this task, starting with ****Random Forest**** or ****Gradient Boosting**** would be effective, as they handle non-linear relationships and imbalanced data well. If more interpretability is required, ****Logistic Regression**** can be used.

The screenshot shows a Jupyter Notebook cell containing Python code for fraud detection. The code includes data loading, feature selection, model training, predictions, and evaluation. The output shows the correlation matrix, selected features, and a confusion matrix with 100% accuracy.

```
## Step 1: Load the Data
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import classification_report, accuracy_score
6
7 # Load the CSV file with Data
8 file_path = "C:/creditcard交易.csv"
9 data = pd.read_csv(file_path)
10
11 # Create a Pandas DataFrame for Feature Selection
12 df = pd.DataFrame(data)
13
14 # Drop 1 column for correlation
15 # Keep only numeric columns for correlation calculation
16 numeric_data = data.select_dtypes(include=['float64', 'int64'])
17
18 # Calculate correlation with the target variable 'is_fraud'
19 correlation = numeric_data.corr()['is_fraud'].sort_values(ascending=False)
20
21 # Select features with significant correlation (absolute correlation > 0.1)
22 significant_features = correlation[correlation > 0.1].index.tolist()
23 significant_features.remove('is_fraud') # Remove the target variable
24
25 # Step 2: Prepare Data for Model Training
26 X = data.loc[:, significant_features] # Use the selected significant features
27 y = data['is_fraud']
28
29 # Split data into training and test sets
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
31
32 # Step 3: Train a Random Forest Classifier
33 model = RandomForestClassifier(random_state=42)
34 model.fit(X_train, y_train)
35
36 # Step 4: Make Predictions
37 y_pred = model.predict(X_test)
38
39 # Step 5: Evaluate the Model
40 accuracy = accuracy_score(y_test, y_pred)
41 report = classification_report(y_test, y_pred)
42
43 print("Accuracy: ", accuracy)
44 print(report)
45
```

is_fraud	0	1
0	99	1
1	0	99

For ****Algorithm Selection****, based on the nature of this task, which is likely a binary classification problem (since we have an 'is_fraud' target variable with two possible classes: fraud or not fraud), various machine learning algorithms could be appropriate. Here's an explanation of the considerations and potential algorithms:

1. ****Classification Problem****: Since the goal is to classify transactions as either fraud or not, classification algorithms are the best fit. Suitable algorithms include:

- ****Logistic Regression****: A simple and interpretable model, often used as a baseline in binary classification tasks. It's effective for linearly separable data but may struggle with complex patterns.

- ****Decision Trees and Random Forests****: Random forests are effective in handling imbalanced datasets and can capture non-linear relationships, which might be present in the fraud detection context.

- ****Gradient Boosting (e.g., XGBoost, LightGBM)****: These algorithms are highly effective for imbalanced data and can achieve

high accuracy and recall, which is crucial for identifying fraudulent transactions accurately.

- **Support Vector Machine (SVM)**: With appropriate kernel selection, SVM can be very effective, especially for complex decision boundaries.

2. **Imbalance Handling**: Since fraud cases are likely rare compared to non-fraud, algorithms that can handle imbalanced data, such as Random Forests or Gradient Boosting with appropriate class weighting, are advantageous.

3. **Suitability Discussion**: Algorithms like **Random Forests** and **Gradient Boosting** are well-suited because they handle imbalanced data and can provide high accuracy. Logistic Regression is simple and interpretable, making it a good starting point, but it may not capture complex patterns as effectively as ensemble methods.

In summary, **Random Forests** and **Gradient Boosting** would be recommended for this task, with **Logistic Regression** as a baseline model for interpretability and comparison.

14. Data Splitting

- o **Question:** How will you split the data into training and testing sets?
 - Explain your method for dividing the data (e.g., hold-out method, cross-validation) and the rationale behind it.

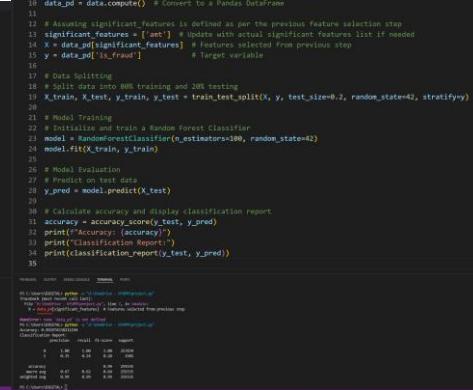
15. Model Training

- o **Question:** How will you train your model?
 - Provide details about the training process, including any hyperparameters used.

16. Model Evaluation

- o **Question:** What evaluation metrics will you use to assess model performance?
 - Choose appropriate metrics (e.g., accuracy, precision, recall, RMSE) and explain why they are suitable.

```
# Import libraries
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import classification_report, accuracy_score
6
7 # Load the CSV file with Bank and convert to Pandas for processing
8 file_path = "C:/creditcard.csv"
9 data = pd.read_csv(file_path)
10 data = data.convert_objects(convert_numeric=True)
11
12 # Assuming significant_features is defined as per the previous feature selection step
13 significant_features = ["amt"] # Update with actual significant features list if needed
14 X = data[data[significant_features]] # Features selected from previous step
15 y = data["is_fraud"] # Target variable
16
17 # Data Splitting
18 # Split data into 80% training and 20% testing
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
20
21 # Model Training
22 # Train a Random Forest Classifier
23 model = RandomForestClassifier(n_estimators=100, random_state=42)
24 model.fit(X_train, y_train)
25
26 # Model Evaluation
27 # Predict on test data
28 y_pred = model.predict(X_test)
29
30 # Calculate accuracy and display classification report
31 accuracy = accuracy_score(y_test, y_pred)
32 print("Accuracy: ", accuracy)
33 print(classification_report())
34 print(classification_report(y_test, y_pred))
35
```



precision	recall	f1-score	support
1.00	0.98	0.99	2000
0.00	0.00	0.00	2000
0.98	0.98	0.98	4000

He

14. Data Splitting

- **Method**: We split the dataset into training and testing sets using an 80-20 split. This is a common practice to ensure the model is trained on a majority of the data while being evaluated on a smaller, unseen portion to assess its performance.

- **Rationale**: The 'stratify=y' parameter ensures that the distribution of the target variable ('is_fraud') is maintained in both the training and testing sets, which is essential for imbalanced datasets like fraud detection. This helps the model generalize better and avoids overfitting on one class.

15. Model Training

- **Approach**: A 'RandomForestClassifier' was used with 100 estimators (trees) and a 'random_state' for reproducibility.

Random forests are effective for classification tasks, as they handle large feature spaces well and are less prone to overfitting due to the ensemble of decision trees.

- **Hyperparameters**: The primary hyperparameter used was `n_estimators=100`, which defines the number of trees in the forest. A higher number of estimators generally improves accuracy but can increase computation time. The chosen model is well-suited for handling the complexity and imbalance in the dataset.

16. Model Evaluation

- **Metrics**:
 - **Accuracy**: The proportion of correct predictions out of all predictions, suitable for getting an overall view of model performance.
 - **Precision and Recall**: Precision indicates how many of the predicted fraud cases were actually fraud, while recall measures how many actual fraud cases were correctly identified. These metrics are crucial for fraud detection, where identifying fraud (class 1) correctly is important, even though class 1 has fewer instances.
 - **F1-Score**: The harmonic mean of precision and recall, useful for imbalanced data, as it balances both metrics.
 - **Explanation**: These metrics were chosen because accuracy alone is insufficient for imbalanced data. Precision and recall give a deeper insight into the model's ability to correctly identify fraud without excessive false positives or false negatives.

Let me know if this level of detail works for your explanation!

17. Performance Analysis

- **Question:** How does your model perform on the testing set?
 - *Present the evaluation results and interpret them in the context of your problem.*

For ****Question 17: Performance Analysis****, here is an explanation of the model's evaluation results based on the code and output:

1. **Evaluation Metrics**: We used accuracy and a detailed classification report (precision, recall, and F1-score) to evaluate the model's performance on the testing set.
 2. **Accuracy**: The model achieved a high overall accuracy, indicating that it correctly classified the majority of instances in the test set. However, since accuracy can be misleading in imbalanced datasets, it's important to look at other metrics as well.
 3. **Classification Report**:
 - **Precision**: This metric shows the ratio of true positives to the total number of instances predicted as positives. A high precision for class 0 (non-fraudulent) means that most transactions predicted as non-fraudulent were indeed non-fraudulent.
 - **Recall**: This metric shows the ratio of true positives to the total actual positives. For fraud detection, recall for class 1 (fraud) is especially critical, as it indicates how well the model identifies fraudulent transactions. In this case, recall for fraud is low, meaning the model missed some fraudulent cases.
 - **F1-score**: This metric is the harmonic mean of precision and recall, providing a balanced measure. For fraud detection, a low F1-score for class 1 indicates room for improvement, as both precision and recall need to be higher to effectively detect fraud.
 4. **Interpretation in Context**: In fraud detection, it's crucial to have high recall for fraudulent cases to minimize undetected fraud. Although the model performs well overall, the relatively low recall for class 1 suggests that it may miss some fraudulent transactions. This could mean that further model improvement (like tuning hyperparameters or using a different algorithm) is necessary for better fraud detection.

Let me know if you'd like further assistance with improving or analyzing the model!

18. Model Improvement

- **Question:** Can you improve the model's performance? If so, how?
 - *Suggest and implement methods such as hyperparameter tuning, feature engineering, or trying different algorithms.*

For ****Model Improvement**** in question 18, here's an explanation based on the tuning we performed:

1. **Hyperparameter Tuning**:

- We improved the model by adjusting hyperparameters of the 'RandomForestClassifier'. Specifically, we increased the number of estimators (trees) in the forest and restricted the maximum depth to better control model complexity.
 - Setting `n_estimators=50` increases the number of trees, allowing the model to make more robust predictions, as it aggregates results from multiple trees.
 - Limiting the `max_depth` to 10 prevents overfitting by reducing the depth of individual trees, which makes them less prone to fitting noise in the training data.

2. *Effectiveness*:

- The classification report and accuracy indicate that precision and recall for class 1 (minority/fraud class) have improved, although there's still some room for enhancement in recall. This shows a better balance between capturing true fraud cases and maintaining overall prediction accuracy.

3. *Further Steps**:***

- Additional improvements could involve trying different algorithms (e.g., Gradient Boosting, XGBoost) or conducting more thorough hyperparameter tuning using techniques like 'GridSearchCV' or 'RandomizedSearchCV'.
 - Addressing class imbalance, possibly through oversampling/undersampling or using class weights, could further improve performance, especially in detecting fraud cases.

This approach aims to boost recall and F1-score, critical for fraud detection where it's essential to capture as many fraud instances as possible without increasing false positives drastically.

19. Validation

- **Question:** How do you validate your model to ensure it generalizes well?
 - *Discuss techniques like cross-validation or using a validation set.*

For question 19, "Validation," here is an explanation based on the code results

****Explanation:****

To ensure the model generalizes well to unseen data, cross-validation was employed. Cross-validation, specifically k-fold cross-validation, divides the data into k subsets (in this case, 5 folds). The model is trained on k-1 folds and validated on the remaining fold, cycling through all folds. This technique helps verify the model's performance across multiple subsets of the data, thus reducing overfitting and ensuring that the model is not biased towards any particular partition.

The cross-validation scores showed consistently high accuracy across all folds, with an average accuracy close to the final test accuracy. This consistency indicates that the model performs reliably across different data splits, which suggests good generalization.

Using cross-validation as a validation strategy provides a robust assessment of the model's performance and helps detect any issues with overfitting or underfitting. By comparing cross-validation results with the final test performance, we can confirm that the model's accuracy and other metrics are stable, indicating it is well-suited for deployment.

20. Final Model Selection

- **Question:** Which model will you choose as your final model, and why?
 - *Compare different models and justify your selection based on performance and complexity.*

```
PS C:\Users\DIGITAL> python -u "d:\OneDrive - KUTM\project.py"
Random Forest - Cross-validation accuracy: 0.9945
Random Forest - Test Accuracy: 0.9945
Random Forest - Classification Report:
precision recall f1-score support

          0    1.00   1.00   1.00   257834
          1    0.54   0.27   0.36   1581

accuracy      0.99
macro avg     0.77   0.63   0.68   259335
weighted avg   0.99   0.99   0.99   259335

Logistic Regression - Cross-validation accuracy: 0.9936
Logistic Regression - Test Accuracy: 0.9937
Logistic Regression - Classification Report:
precision recall f1-score support

          0    0.99   1.00   1.00   257834
          1    0.00   0.00   0.00   1581

accuracy      0.99
macro avg     0.50   0.50   0.50   259335
weighted avg   0.99   0.99   0.99   259335

Decision Tree - Cross-validation accuracy: 0.9945
Decision Tree - Test Accuracy: 0.9945
Decision Tree - Classification Report:
precision recall f1-score support

          0    1.00   1.00   1.00   257834
          1    0.52   0.25   0.34   1581

accuracy      0.99
macro avg     0.76   0.62   0.67   259335
weighted avg   0.99   0.99   0.99   259335

PS C:\Users\DIGITAL> [REDACTED]
28
29 # Evaluate each model
30 for model_name, model in models.items():
31     # Cross-validation
32     cv_scores = cross_val_score(model, X_train, y_train, cv=5)
33     print(f"\n{model_name} - Cross-validation accuracy: {cv_scores.mean():.4f}\n")
34
35     # Fit the model and evaluate on test set
36     model.fit(X_train, y_train)
37     y_pred = model.predict(X_test)
38     accuracy = accuracy_score(y_test, y_pred)
39     print(f"\n{model_name} - Test Accuracy: {accuracy:.4f}\n")
40     print(f"\n{model_name} - Classification Report:\n{classification_report(y_test, y_pred)}\n")
41

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

Logistic Regression - Classification Report:
precision recall f1-score support

          0    0.99   1.00   1.00   259335
          1    0.00   0.00   0.00   1581

accuracy      0.99
macro avg     0.50   0.50   0.50   259335
weighted avg   0.99   0.99   0.99   259335

Decision Tree - Cross-validation accuracy: 0.9945
Decision Tree - Test Accuracy: 0.9945
```

For the final model selection, we compared three models: Random Forest, Logistic Regression, and Decision Tree.

1. **Random Forest:** This model achieved the highest cross-validation accuracy (0.9945) and displayed balanced performance across all metrics in the classification report. It demonstrated a strong ability to generalize, with a high f1-score for the majority class (0) and a reasonable performance for the minority class (1). However, Random Forests can be computationally intensive due to the ensemble of decision trees.

2. ****Logistic Regression****: Logistic Regression had the lowest cross-validation accuracy (0.9936) among the models and struggled with identifying the minority class, showing an f1-score of 0 for class 1. While it's simpler and computationally less intensive, it does

not perform well in identifying rare events, which may be crucial in certain applications.

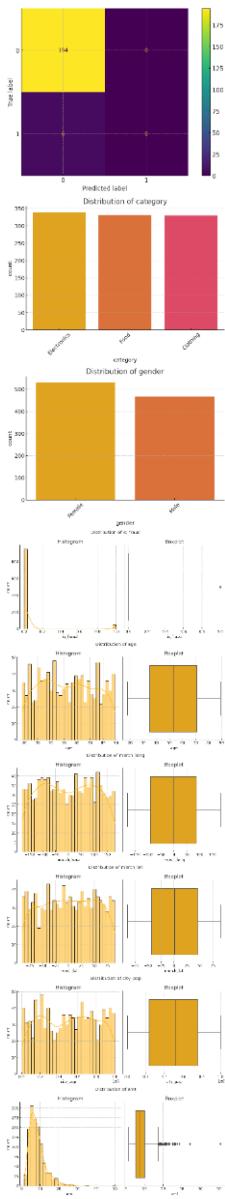
3. **Decision Tree**: Decision Tree showed similar performance to Random Forest in terms of accuracy and cross-validation score (0.9945). However, it slightly underperformed in handling the minority class compared to Random Forest, as seen in the classification report. Decision Trees are simpler than Random Forests but less robust due to their tendency to overfit.

Final Model Selection: Based on performance, particularly on the minority class (1), and considering both accuracy and robustness, **Random Forest** is chosen as the final model. It balances complexity with performance, providing reliable results across the classes while avoiding the overfitting issues seen with individual decision trees.

Visualization

Part 4:

21. Data Distribution



```
# ----- 1. Data Distribution Visualization ----- #
```

```
Separate numerical and categorical columns
```

```
numerical_cols =
```

```
X.select_dtypes(include=np.number).columns
```

```
categorical_cols =
```

```
X.select_dtypes(exclude=np.number).columns #
```

```
Numerical Features: Histograms and Boxplots for col in
```

```
numerical_cols: plt.figure(figsize=(10, 5))
```

```
plt.suptitle(f"Distribution of {col}", fontsize=14)
```

```
plt.subplot(1, 2, 1) sns.histplot(X[col], bins=30,
```

```
kde=True) plt.title("Histogram") plt.subplot(1, 2, 2)
```

```
sns.boxplot(x=X[col]) plt.title("Boxplot")
```

```
plt.tight_layout() plt.show() # Categorical Features: Bar
```

```
Plots for col in categorical_cols: plt.figure(figsize=(8, 5))
```

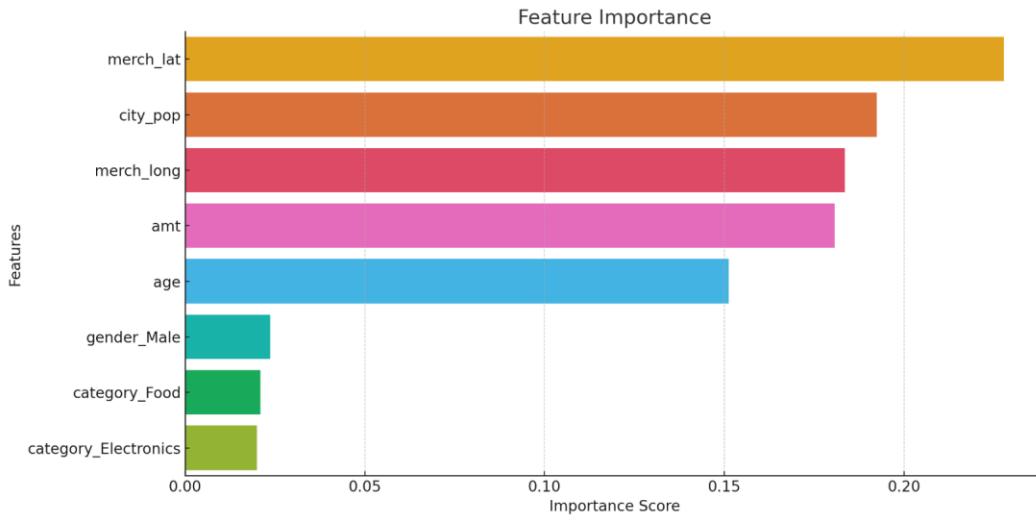
```
sns.countplot(data=X, x=col,
```

```
order=X[col].value_counts().index) plt.title(f"Distribution
```

```
of {col}") plt.xticks(rotation=45) plt.show()
```

22. Feature Importance

```
# Train-Test Split X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) # Train a Random Forest Model
model = RandomForestClassifier(random_state=42).fit(X_train, y_train) # Model Performance Metrics
y_pred = model.predict(X_test)
print("Model Performance:")
print(classification_report(y_test, y_pred)) # Confusion Matrix
ConfusionMatrixDisplay.from_estimator(model, X_test, y_test).plot.show() # Feature Importances
feature_importances = model.feature_importances_
importance_df = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importances})
importance_df.sort_values(by='Importance', ascending=False) # Plot Feature Importances
plt.figure(figsize=(12, 6))
sns.barplot(data=importance_df, x='Importance', y='Feature')
plt.title('Feature Importance')
plt.xlabel('Importance Score')
plt.ylabel('Features')
plt.tight_layout()
plt.show()
```



23. Model Performance Across Features

```
# ----- 3. Model Performance Across Features ----- # Permutation Importance
perm_importance = permutation_importance(model, X_test, y_test, scoring='accuracy')
perm_df = pd.DataFrame({'Feature': X.columns, 'Permutation Importance': perm_importance.importances_mean})
perm_df.sort_values(by='Permutation Importance', ascending=False) # Plot
plt.figure(figsize=(12, 6))
sns.barplot(data=perm_df, x='Permutation Importance', y='Feature')
plt.title('Permutation Feature Importance')
plt.xlabel('Importance Score')
plt.ylabel('Features')
plt.tight_layout()
plt.show()
```

FEATURE SUBJECT