1 1 1 1 1 1 2 3 3 4 ••• 1 1 1 1	1048570 95 CASE_OUT 132557.35 C1179511630 479803.00 1048571 95 PAYMENT 9917.36 C1956161225 90545.00 1048572 95 PAYMENT 10020.05 C1633237354 90605.00 1048574 95 PAYMENT 11450.03 C1264356443 80584.95
1 1 1 1 1 (T	isFlaggedFraud 0 1 0 2 0 3 0 4 0 5 0 6 0 6 0 6 0 7 0 8 0 8 0 9 0 9 0 9 0 9 0 9 0 9 0 9 0 9 0 9 0 9
/ T F 0	df('step'].fillna(method='ffill', inplace=True) # Forward fill for time steps (tmp/lypkernel_2217/178294538.py:22: FutureWarning: Series.fillna with 'method' is deprecated and will raise in a future version. Use obj.ffill() or obj.bfill() instead. df('step'].fillna(method='ffill', inplace=True) # Forward fill for time steps (tmp/lypkernel_2217/1782945538.py:23: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method((col: value), inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the object. df['type'].fillna(df['type'].mode()[0], inplace=True) # Mode for transaction type *tmp/lypkernel_22171/1782945538.py:24: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. df['amount'].fillna(df['amount'].median(), inplace=True)' # Median for amount 'tmp/lypkernel_22171/1782945538.py:25: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. 'tmp/lypkernel_22171/1782945538.py:25: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. 'tmp/lypkernel_22171/1782945538.py:25: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. 'tmp/lypkernel_22171/1782945538.py:25: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
/ T F O / T F	For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the object. df['nameOrig'].fillna('Unknown', inplace=True)
T F 0 / T F	if ('newbalanceOcig'].fillna(df('newbalanceOcig'].mean(), inplace=True) the behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the object. df('nameDest'].fillna('Unknown', inplace=True) the behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. The vample, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the object. df['oldbalanceDest'].fillna(df['oldbalanceDest'].mean(), inplace=True) the behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy. for example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the object.
D 0 1 2 3 4 4 · · · · · · · · · · · · · · · · ·	df['newbalanceDest'].fillna(df['newbalanceDest'].mean(), inplace=True) Data after handling missing values: step
1 1 1 1 1 1 2 3 4 1 1 1 1 1	
/ T F O	import pands as pd import matplottib, pyplot as pt import matplottib, pyplot as pt import matplottib, pyplot as spt import matplottib, pyplot as spt import seaborn as sns * Sample DataFrame Creation * Replace this with your actual dataset data = pd.read_csv('Downloads/Fraud.csv')
	<pre>df = pd.DataFrame(data) # Function to remove outliers using Z-Score method def remove_outliers_z_score(dataframe, column): threshold = 3 # Z-score threshold z_scores = (dataframe[column] - dataframe[column].mean()) / dataframe[column].std() filtered_df = dataframe[(z_scores > -threshold) & (z_scores < threshold)] return filtered_df # Function to remove outliers using IQR method def remove_outliers_iqr(dataframe, column): Q1 = dataframe[column].quantile(0.25) Q3 = dataframe[column].quantile(0.75) IQR = Q3 - Q1 lower_bound = Q1 - 1.5 * IQR upper_bound = Q1 - 1.5 * IQR upper_bound = Q3 + 1.5 * IQR filtered_df = dataframe[(dataframe[column] >= lower_bound) & (dataframe[column] <= upper_bound)] return filtered_df # Visualize the original data with box plots plt.figure(flgsize=(12, 6))</pre>
	<pre>sns.boxplot(x=df('amount')) plt.title('Box Plot of Amount (Original Data)') plt.xlabel('Transaction Amount') plt.show() # Remove outliers for the 'amount' column using both methods df_no_outliers_z = remove_outliers_z_score(df, 'amount') df_no_outliers_igr = remove_outliers_igr(df, 'amount') # Visualize the data after removing outliers using Z-Score method plt.figure(figsize=(12, 6)) sns.boxplot(x=df_no_outliers_z['amount']) plt.xlabel('Transaction Amount') plt.xlabel('Transaction Amount') plt.show() # Visualize the data after removing outliers using IQR method plt.figure(figsize=(12, 6)) sns.boxplot(x=df_no_outliers_igr('amount')) plt.show() # Visualize the data after removing outliers using IQR method plt.figure(figsize=(12, 6)) sns.boxplot(x=df_no_outliers_igr('amount')) plt.xlabel('Transaction Amount') plt.xlabel('Transaction Amount') plt.xlabel('Transaction Amount')</pre>
	# Print the results print ("Oxiginal DataFrame:\n", df) print ("NDataFrame after removing outliers (Z-Score method):\n", df_no_outliers_z) print ("\nDataFrame after removing outliers (IQR method):\n", df_no_outliers_iqr) Box Plot of Amount (Original Data)
	0.0 0.2 0.4 0.6 0.8 1.0 Transaction Amount 1e7 Box Plot of Amount (After Z-Score Outlier Removal)
	0.0 0.2 0.4 0.6 0.8 1.0 Transaction Amount (After IQR Outlier Removal)
0 0 1 2 3	1 PAYMENT 1864.28 C1666544295 21249.00 2 1 TRANSFER 181.00 C1305486145 181.00 3 1 CASH_OUT 181.00 C840083671 181.00
1 1 1 1 1 0 1 2 3 4 1 1 1 1	1 PAYMENT 11668.14 02048537720 41554.00 1048570 95 CASH_OUT 132557.35 C1179511630 479803.00 1048571 95 PAYMENT 9917.36 C1956161225 90545.00 1048573 95 PAYMENT 1020.00 C1633237354 90605.00 1048574 95 PAYMENT 11450.03 C1264356443 80584.95 newbalanceOrig nameDest oldbalanceDest isFraud \
1 1 1 1	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 2 3 4 4	19384.72 M2044282225 0.00 0.00 0.00 0.00 0.00 0.00 0.00
1 1 1 1 1 1 1 0 1 2 3 4	
1 1 1 0 1 2 3 4 4 1 1 1	19384.72 M2044282225 0.00 0.00 0.00 0.00 0.00 0.00 0.00
1 1 1 1 1 1 1 1 ::	0
	<pre># Step 2: Load specific columns from the CSV file # Replace 'column1', 'column2', etc. with the actual column names you want to include df = pd.read_csv(csv_file_path, usecol==\'amount', 'oldbalanceOrg', 'newbalanceOrig', 'loadbalanceDest', 'newbalanceDest']) # Optional: Display the first few rows of the DataFrame to verify print("DataFrame head:\n", df.head()) # Step 3: Calculate the correlation matrix correlation_matrix = df.corr() # Print the correlation matrix print("Correlation Matrix:\n", correlation_matrix) # Optional: Visualize the correlation matrix plt.figure(figsize=(10, 6)) sss.heating(correlation_matrix, annot=True, fmt='.2f', cmap='coolwarm', square=True) plt.title('Correlation Matrix Heatmap') plt.show() # Create a DataFrame for VIF calculation (drop non-numeric or categorical columns) df vif = df[('amount', 'oldbalanceOrg', 'newbalanceOrg', 'newbalanceDest', 'newbalanceDest']]</pre>
D 0 1 2 3 4 C a o n o	
o n o	newbalanceDest
	newbalanceOrig0.00 1.00 1.00 0.10 0.06 oldbalanceDest - 0.22 0.09 0.10 1.00 0.98 newbalanceDest - 0.31 0.06 0.06 0.98 1.00 -0.2 -0.2 -0.0 -0.0 -0.0 -0.0 -0.0 -0.0 -0.0 -0.0
0 1 2 3 4	Variance Inflation Factor: feature VIF
	<pre># Step 1: Load the Dataset (Assume CSV format) # The dataset contains columns as described in the user's message. df = pd.read_csv('Downloads/Fraud.csv') # Step 2: Basic Exploration print("Dataset Overview:") print(df.head(l)) print("\nDataset Info:") print(df.info()) print("\nSummary Statistics:") print(df.describe(l)) # Step 3: Handle Missing Data # Check for missing Values: print("\nMissing Values:") print(df.insnull().sum()) # Step 4: Data Preprocessing # Convert categorical columns (nameOrig, nameDest, type) to numeric using Label Encoding les = LabelForceder()</pre>
	<pre>le = LabelEncoder() df['nameOrig'] = le.fit_transform(df['nameOrig']) df['nameOrig'] = le.fit_transform(df['nameDest']) df['type'] = le.fit_transform(df['type']) # Step 5: Feature Scaling # Scaling numerical features (amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest) scaler = StandardScaler() df[['amount', 'oldbalanceOrg', 'newbalanceOrig', 'oldbalanceDest', 'newbalanceDest']] = scaler.fit_transform(df[['amount', 'oldbalanceOrg', 'newbalanceOrig', 'oldbalanceDest', 'newbalanceDest']] # Step 6: Splitting the Data into Train/Test X = df.drop(columns=['isFraud']) Y = df['isFraud'] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42) # Step 7: Model Building - Using RandomForestClassifier for Fraud Detection off = RandomForestClassifier(n_estimators=100, random_state=42) off.fit(X_train, y_train)</pre>
	<pre># Step 8: Predictions and Model Evaluation y_pred = clf.predict(X_test) # Evaluate the Model print("\nClassification Report:") print(classification_report(y_test, y_pred)) print("\nConfusion Matrix:") cm = confusion_matrix(y_test, y_pred) sns.heating(cm, annot="rue, fmt='d', cmap='Blues') plt.title('Confusion Matrix') plt.vlabel('Predicted') plt.ylabel('Actual') plt.show() # Step 9: Feature Importance # Displaying the importance of features in fraud prediction importance = pd.Series(clf.feature_importances_, index=X.columns) importance.nlargest(10).plot(kind='barh', title="Feature Importance in Fraud Detection") plt.show()</pre>
•	4. Demonstrate the performance of the model by using best set of tools. import pandas as pd # Step 1: Simulated Transaction Data data = pd.read_csv('Downloads/Fraud.csv') # Step 2: Convert the data into a pandas DataFrame df = pd.DataFrame(data) # Step 3: Function to print the details of each transaction def print_transaction_details(df): print("Transaction Details:") for idx, row in df.iterrows(): print(f" Transaction (idx + 1):") print(f" Transaction Type: (row('type'))") print(f" Transaction Type: (row('type'))") print(f" Customer (Orig): {row('amount')}")
	<pre>print(f" Balance Before: [row['oddbalanceOrg'])") print(f" Recipient (Dest): [row['nameDest'])") print(f" Recipient (Dest): [row['nameDest'])") print(f" Balance Before: [row['newBalanceDest']]") print(f" Balance Before: [row['newBalanceDest']]") print(f" Fraudulent Transaction: {'Yes' if row['isPraud'] == l else 'No'}") # Step 4: Analyze fraudulent transactions def analyze_fraud(df): print("\nFraudulent Transactions Summary:") frauds = df(df['isFraud'] == l) if not frauds.empty: for idx, row in frauds.iterrows():</pre>
	<pre># Step 5: Main Program Execution if _name_ == "_main_":</pre>
	<pre>data['balanceDiffOrig'] = data['oldbalanceDest'] - data['nexbalanceDest'] data['balanceDiffDest'] = data['nexbalanceDest'] - data['oldbalanceDest'] # Step 5: Basic rule-based fraud detection (you can customize these rules): def detect_fraud(row): # Check for mismatch between transaction amount and balance changes if row['nbalanceDiffOrig'] != row['amount'] or row['balanceDiffDest'] != row['amount']: return l # Mark as fraud return 0 # Not fraud # Apply the fraud detection logic to the dataset data['detectedFraud'] = data.apply(detect_fraud, axis=1) # Step 6: Compare with the "isFraud" label in the dataset correct_detections = (data['detectedFraud'] == data['lsFraud']).sum() total_transactions = len(data) # Calculate the accuracy of the rule-based detection accuracy = correct_detections / total_transactions print(f"Accuracy of rule-based fraud detection: (accuracy * 100:.2f}%") # Step 7: Summary statistics for transactions</pre>
	# Step 7: Summary statistics for transactions print(data.describe()) # Optionally, save the processed data to a new file data.to_csv('processed_transactions.csv', index=False) 6. Do these factors make sense? If yes, How? If not, How not? Intioned for predicting fraudulent customers, and the ones used in the code, make sense, but their effectiveness depends on the context, dataset, and overall fraud detection strategy. Let's break down why these factors make sense and also where there could be limitations: 1. Transaction Amount and

1. Data cleaning including missing values, outliers and multi-collinearity.

(i) Missing Values

from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

data = pd.read_csv('Downloads/Fraud.csv')

from statsmodels.stats.outliers_influence import variance_inflation_factor

In [3]: import pandas as pd

print(df)

import numpy as np

Sample dataset creation

df = pd.DataFrame(data)

Display the original data
print("Original Data:")