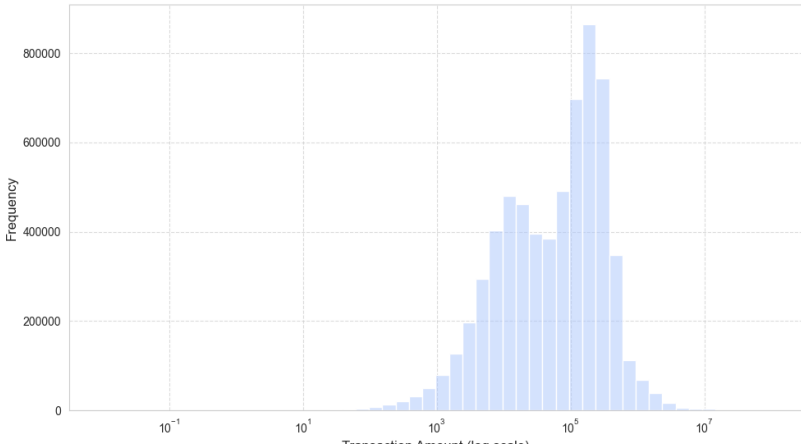


Data Collection and Preprocessing Phase

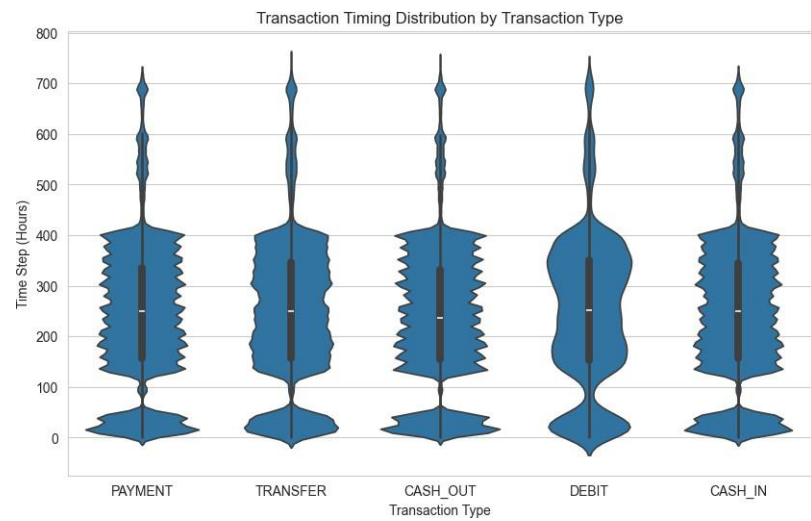
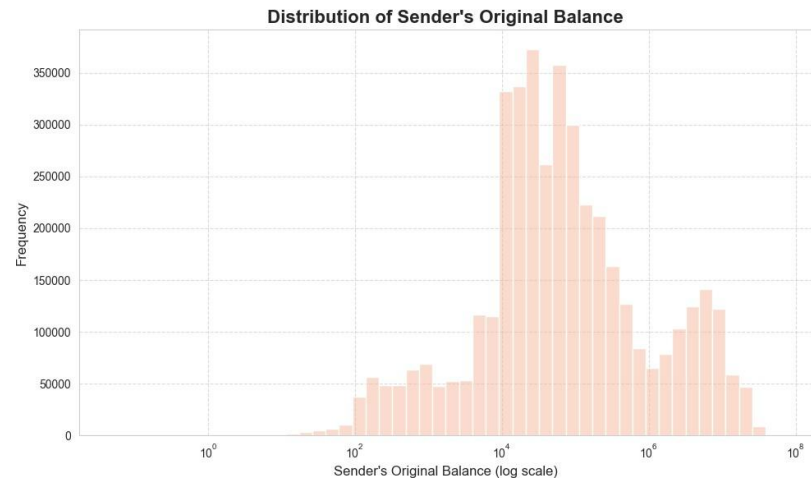
Date	19 Feb 2026
Team ID	LTVIP2026TMIDS80731
Project Title	Online Payment Fraud Detection using ML
Maximum Marks	6 Marks

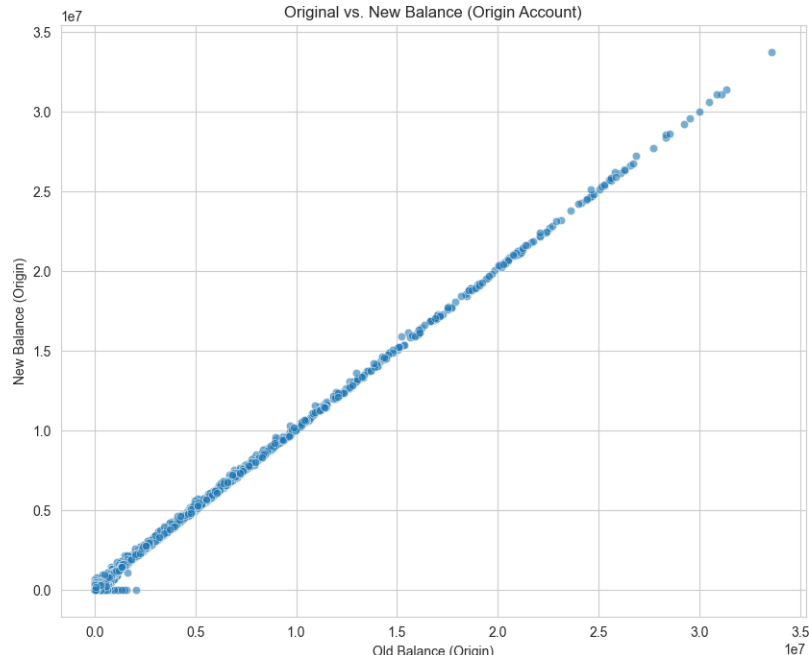
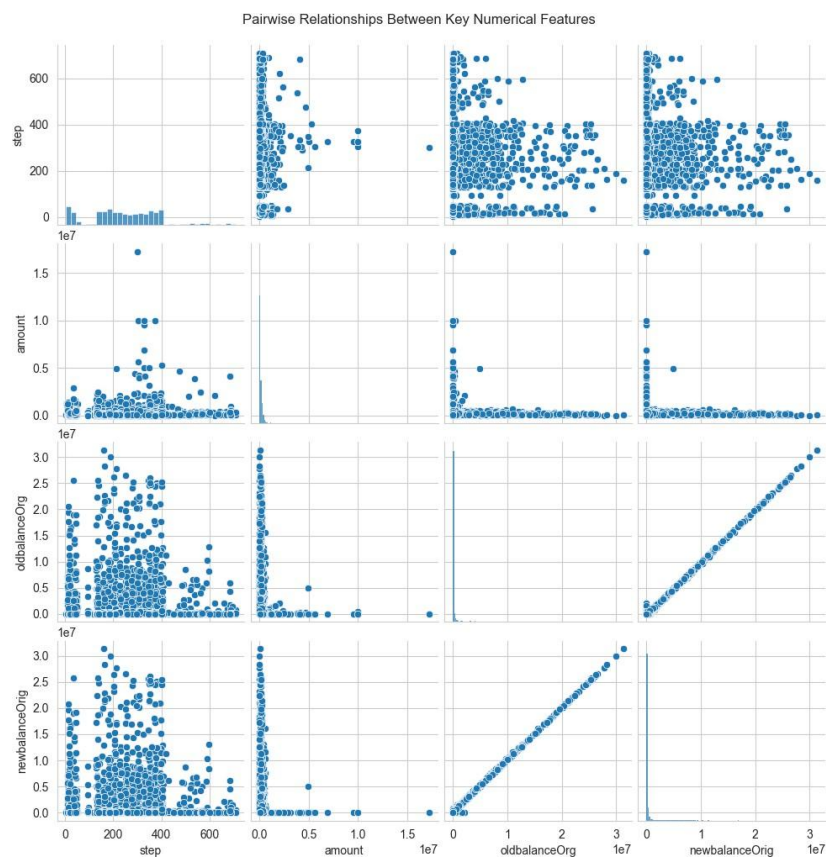
Data Exploration and Preprocessing

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																																																	
Data Overview	<table><tr><th></th><th>step</th><th>amount</th><th>oldbalanceOrig</th><th>newbalanceOrig</th><th>oldbalanceDest</th><th>newbalanceDest</th><th>isFraud</th><th>isFlaggedFraud</th></tr><tr><td>count</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td><td>6.362620e+06</td></tr><tr><td>mean</td><td>2.433972e+02</td><td>1.798619e+05</td><td>8.338831e+05</td><td>8.551137e+05</td><td>1.100702e+06</td><td>1.224996e+06</td><td>1.290820e-03</td><td>2.514687e-06</td></tr><tr><td>std</td><td>1.423320e+02</td><td>6.038582e+05</td><td>2.888243e+06</td><td>2.924049e+06</td><td>3.399180e+06</td><td>3.674129e+06</td><td>3.590480e-02</td><td>1.585775e-03</td></tr><tr><td>min</td><td>1.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>25%</td><td>1.560000e+02</td><td>1.338957e+04</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>50%</td><td>2.390000e+02</td><td>7.487194e+04</td><td>1.420800e+04</td><td>0.000000e+00</td><td>1.327057e+05</td><td>2.146614e+05</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>75%</td><td>3.350000e+02</td><td>2.087215e+05</td><td>1.073152e+05</td><td>1.442584e+05</td><td>9.430367e+05</td><td>1.111909e+06</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>max</td><td>7.430000e+02</td><td>9.244552e+07</td><td>5.958504e+07</td><td>4.958504e+07</td><td>3.560159e+08</td><td>3.561793e+08</td><td>1.000000e+00</td><td>1.000000e+00</td></tr></table>		step	amount	oldbalanceOrig	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud	count	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	mean	2.433972e+02	1.798619e+05	8.338831e+05	8.551137e+05	1.100702e+06	1.224996e+06	1.290820e-03	2.514687e-06	std	1.423320e+02	6.038582e+05	2.888243e+06	2.924049e+06	3.399180e+06	3.674129e+06	3.590480e-02	1.585775e-03	min	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	25%	1.560000e+02	1.338957e+04	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	50%	2.390000e+02	7.487194e+04	1.420800e+04	0.000000e+00	1.327057e+05	2.146614e+05	0.000000e+00	0.000000e+00	75%	3.350000e+02	2.087215e+05	1.073152e+05	1.442584e+05	9.430367e+05	1.111909e+06	0.000000e+00	0.000000e+00	max	7.430000e+02	9.244552e+07	5.958504e+07	4.958504e+07	3.560159e+08	3.561793e+08	1.000000e+00	1.000000e+00
	step	amount	oldbalanceOrig	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud																																																																										
count	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06	6.362620e+06																																																																										
mean	2.433972e+02	1.798619e+05	8.338831e+05	8.551137e+05	1.100702e+06	1.224996e+06	1.290820e-03	2.514687e-06																																																																										
std	1.423320e+02	6.038582e+05	2.888243e+06	2.924049e+06	3.399180e+06	3.674129e+06	3.590480e-02	1.585775e-03																																																																										
min	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00																																																																										
25%	1.560000e+02	1.338957e+04	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00																																																																										
50%	2.390000e+02	7.487194e+04	1.420800e+04	0.000000e+00	1.327057e+05	2.146614e+05	0.000000e+00	0.000000e+00																																																																										
75%	3.350000e+02	2.087215e+05	1.073152e+05	1.442584e+05	9.430367e+05	1.111909e+06	0.000000e+00	0.000000e+00																																																																										
max	7.430000e+02	9.244552e+07	5.958504e+07	4.958504e+07	3.560159e+08	3.561793e+08	1.000000e+00	1.000000e+00																																																																										
Univariate Analysis	<div><div>Distribution of Transaction Amounts</div></div>																																																																																	

Bivariate Analysis



	
Multivariate Analysis	
Outliers and Anomalies	<p>Due to the extremely wide distribution and heavy-tailed nature of financial transaction data in this dataset, traditional outlier</p>

detection methods are not effective. The data spans several orders of magnitude (e.g., transaction amounts range from \$0 to \$92+ million), making it difficult to distinguish between legitimate large transactions and true anomalies using standard statistical methods.

Data Preprocessing Code Screenshots

Loading Data

```
#Loading Data
data=pd.read_csv("data.csv")
data_og=data.copy()
data.head()
```

	step	type	amount	nameOrig	oldbalanceOrg	newbalanceOrig	nameDest	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud
0	1	PAYMENT	9839.64	C1231006815	170136.0	160296.36	M1979787155	0.0	0.0	0	0
1	1	PAYMENT	1864.28	C1666544295	21249.0	19384.72	M2044282225	0.0	0.0	0	0
2	1	TRANSFER	181.00	C1305486145	181.0	0.00	C553264065	0.0	0.0	1	0
3	1	CASH_OUT	181.00	C840083671	181.0	0.00	C38997010	21182.0	0.0	1	0
4	1	PAYMENT	11668.14	C2048537720	41554.0	29885.86	M1230701703	0.0	0.0	0	0

Handling Missing Data

```
#Handling Missing Data (There isnt any missing data)
data.isnull().sum()
```

```
step          0
type          0
amount        0
nameOrig      0
oldbalanceOrg 0
newbalanceOrig 0
nameDest      0
oldbalanceDest 0
newbalanceDest 0
isFraud       0
isFlaggedFraud 0
dtype: int64
```

Data Transformation

```
df_clean = data.drop(['nameOrig', 'nameDest', 'isFlaggedFraud'], axis=1)
df_clean.columns.tolist()
```

```
['step',
 'type',
 'amount',
 'oldbalanceOrg',
 'newbalanceOrig',
 'oldbalanceDest',
 'newbalanceDest',
 'isFraud']
```

	<pre>df_encoded = pd.get_dummies(df_clean, columns=['type'], prefix='type') df_encoded.head()</pre> <table><thead><tr><th>amount</th><th>oldbalanceOrig</th><th>newbalanceOrig</th><th>oldbalanceDest</th><th>newbalanceDest</th><th>isFraud</th><th>type_CASH_IN</th><th>type_CASH_OUT</th><th>type_DEBIT</th><th>type_PAYMENT</th><th>type_TRANSFER</th></tr></thead><tbody><tr><td>9839.64</td><td>170136.0</td><td>160296.36</td><td>0.0</td><td>0.0</td><td>0</td><td>False</td><td>False</td><td>False</td><td>True</td><td>False</td></tr><tr><td>1864.28</td><td>21249.0</td><td>19384.72</td><td>0.0</td><td>0.0</td><td>0</td><td>False</td><td>False</td><td>False</td><td>True</td><td>False</td></tr><tr><td>181.00</td><td>181.0</td><td>0.00</td><td>0.0</td><td>0.0</td><td>1</td><td>False</td><td>False</td><td>False</td><td>False</td><td>True</td></tr><tr><td>181.00</td><td>181.0</td><td>0.00</td><td>21182.0</td><td>0.0</td><td>1</td><td>False</td><td>True</td><td>False</td><td>False</td><td>False</td></tr><tr><td>11668.14</td><td>41554.0</td><td>29885.86</td><td>0.0</td><td>0.0</td><td>0</td><td>False</td><td>False</td><td>False</td><td>True</td><td>False</td></tr></tbody></table> <pre>X = df_encoded.drop('isFraud', axis=1) y = df_encoded['isFraud'] scaler = StandardScaler() X_scaled = scaler.fit_transform(X)</pre>	amount	oldbalanceOrig	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud	type_CASH_IN	type_CASH_OUT	type_DEBIT	type_PAYMENT	type_TRANSFER	9839.64	170136.0	160296.36	0.0	0.0	0	False	False	False	True	False	1864.28	21249.0	19384.72	0.0	0.0	0	False	False	False	True	False	181.00	181.0	0.00	0.0	0.0	1	False	False	False	False	True	181.00	181.0	0.00	21182.0	0.0	1	False	True	False	False	False	11668.14	41554.0	29885.86	0.0	0.0	0	False	False	False	True	False
amount	oldbalanceOrig	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud	type_CASH_IN	type_CASH_OUT	type_DEBIT	type_PAYMENT	type_TRANSFER																																																									
9839.64	170136.0	160296.36	0.0	0.0	0	False	False	False	True	False																																																									
1864.28	21249.0	19384.72	0.0	0.0	0	False	False	False	True	False																																																									
181.00	181.0	0.00	0.0	0.0	1	False	False	False	False	True																																																									
181.00	181.0	0.00	21182.0	0.0	1	False	True	False	False	False																																																									
11668.14	41554.0	29885.86	0.0	0.0	0	False	False	False	True	False																																																									
Feature Engineering	New features such as balance changes and ratios were considered but not used in the baseline due to time/resource constraints.																																																																		
Save Processed Data	<pre>X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns) X_scaled_df.describe()</pre>																																																																		