# Data Collection and Preprocessing Phase

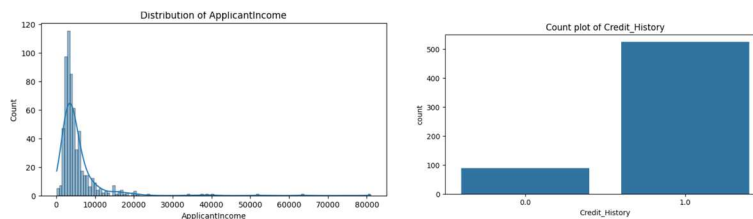| | |
|---|---|
| Date | 20 June 2025 |
| Team ID | SWTID1749791625 |
| Project Title | Smart Lender- Applicant Credibility Prediction for Loan Approval |
| Maximum Marks | 6 Marks |

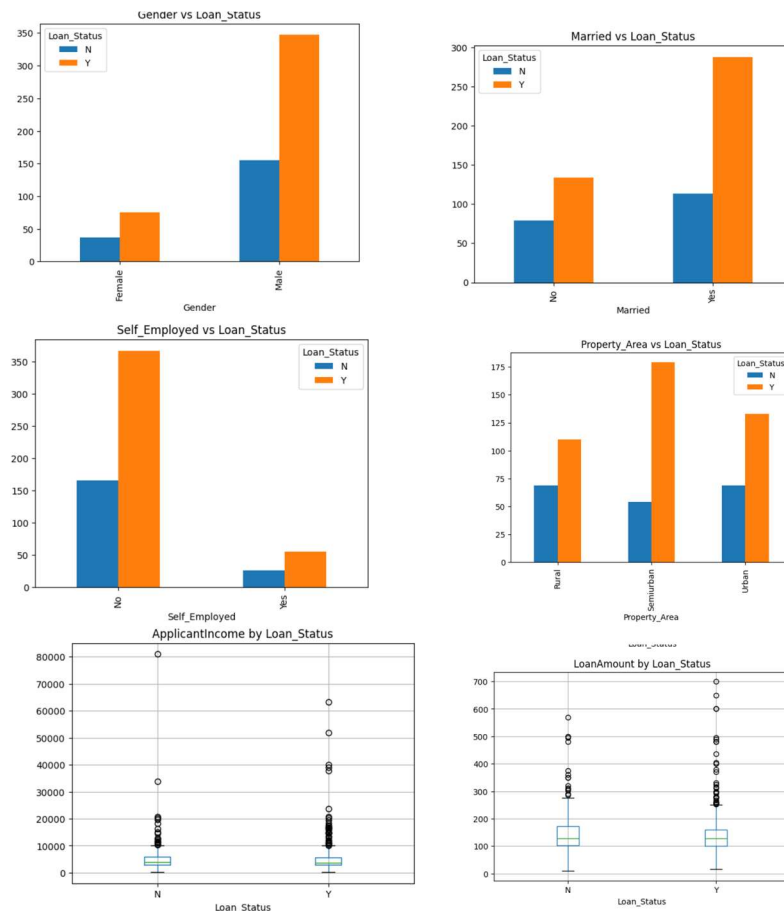**Data Exploration and Preprocessing Report**

Dataset variables will be statistically analyzed to identify patterns and outliers. Python will be used for preprocessing tasks such as normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring high-quality input for analysis and modeling, and providing a strong foundation for reliable insights and predictions.

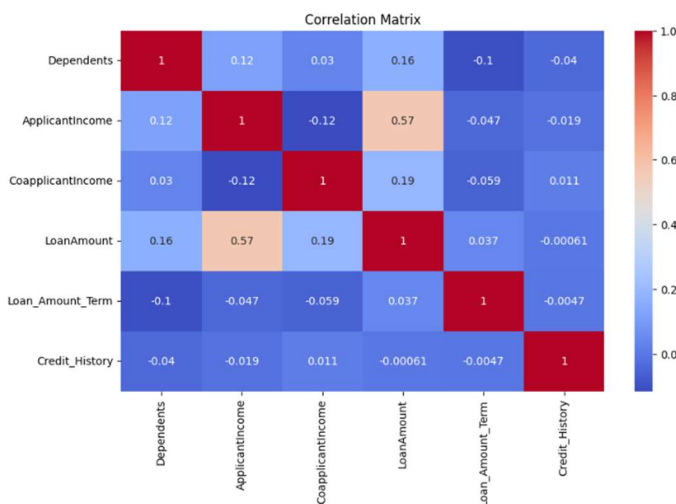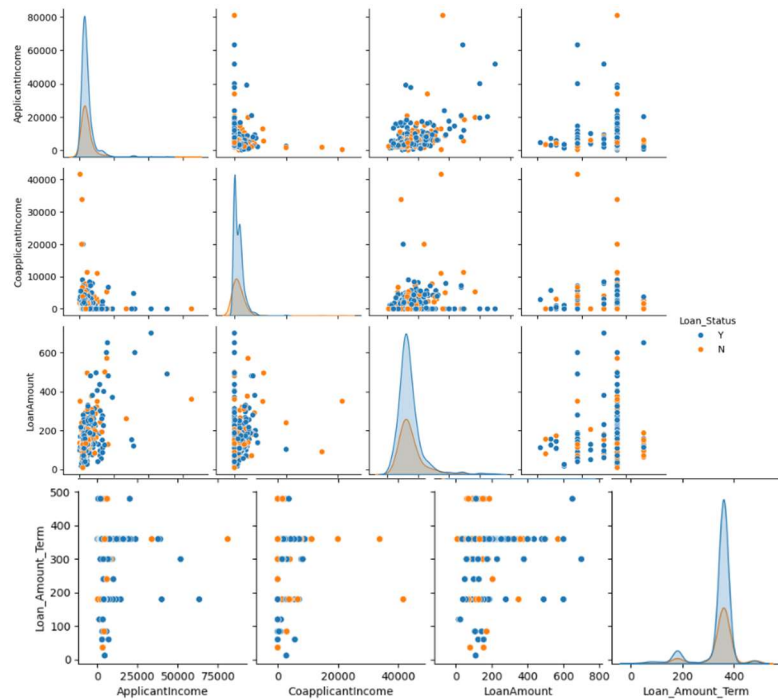| Section | Description |
|---|---|
| Data Overview | **Dimension:**<br>614rows×13columns.<br>**Descriptive Statistics:**<br><pre>       Dependents  ApplicantIncome  CoapplicantIncome  LoanAmount  \<br>count  614.000000       614.000000         614.000000  614.000000<br>mean     0.744300      5403.459283        1621.245798  145.752443<br>std      1.009623      6109.041673        2926.248369   84.107233<br>min      0.000000       150.000000           0.000000    9.000000<br>25%      0.000000      2877.500000           0.000000  100.250000<br>50%      0.000000      3812.500000        1188.500000  128.000000<br>75%      1.000000      5795.000000        2297.250000  164.750000<br>max      3.000000     81000.000000       41667.000000  700.000000<br><br>       Loan_Amount_Term  Credit_History<br>count        614.000000      614.000000<br>mean         342.410423        0.855049<br>std           64.428629        0.352339<br>min           12.000000        0.000000<br>25%          360.000000        1.000000<br>50%          360.000000        1.000000<br>75%          360.000000        1.000000<br>max          480.000000        1.000000<br>       Gender Married Education Self_Employed Property_Area Loan_Status<br>count     614     614       614           614           614         614<br>unique      2       2         2             2             3           2<br>top      Male     Yes  Graduate            No     Semiurban           Y<br>freq      502     401       480           532           233         422</pre> |

| | |
|---|---|
| Univariate Analysis |  |
| Bivariate Analysis |  |
| Multivariate Analysis |  |

| | |
|---|---|
| |  |
| Outliers and Anomalies | - |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data |  |
| Handling Missing Data |  |

| Data Transformation | Handling Categorical Dataset<br><br>```python<br>[ ] df['Loan_Status'] = df['Loan_Status'].map({'Y': 1, 'N': 0})<br>    df['Dependents'] = df['Dependents'].replace('3+', 3).astype(int)<br><br>    df = pd.get_dummies(df, columns=['Gender', 'Married', 'Education', 'Self_Employed', 'Property_Area'], drop_first=True)<br>```<br><br>Balancing And Scaling<br><br>```python<br>[ ] X = df.drop('Loan_Status', axis=1)<br>    y = df['Loan_Status']<br><br>    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)<br><br>    print("Train Target Distribution:")<br>    print(y_train.value_counts())<br><br>    print("Test Target Distribution:")<br>    print(y_test.value_counts())<br>    scaler = StandardScaler()<br>    X_train_scaled = scaler.fit_transform(X_train)<br>    X_test_scaled = scaler.transform(X_test)<br><br>    smote = SMOTE(random_state=42)<br>    X_train_smote, y_train_smote = smote.fit_resample(X_train_scaled, y_train)<br><br>    X_train_final, X_val, y_train_final, y_val = train_test_split(<br>        X_train_smote, y_train_smote, test_size=0.2, random_state=42<br>    )<br>``` |
|---|---|
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | ```python<br>import pickle<br>#Save model and scaler for Flask<br>with open('final_model.pkl', 'wb') as f:<br>    pickle.dump(voting_clf, f)<br><br>with open('scaler.pkl', 'wb') as f:<br>    pickle.dump(scaler, f)<br>``` |