# Data Collection and Preprocessing Phase

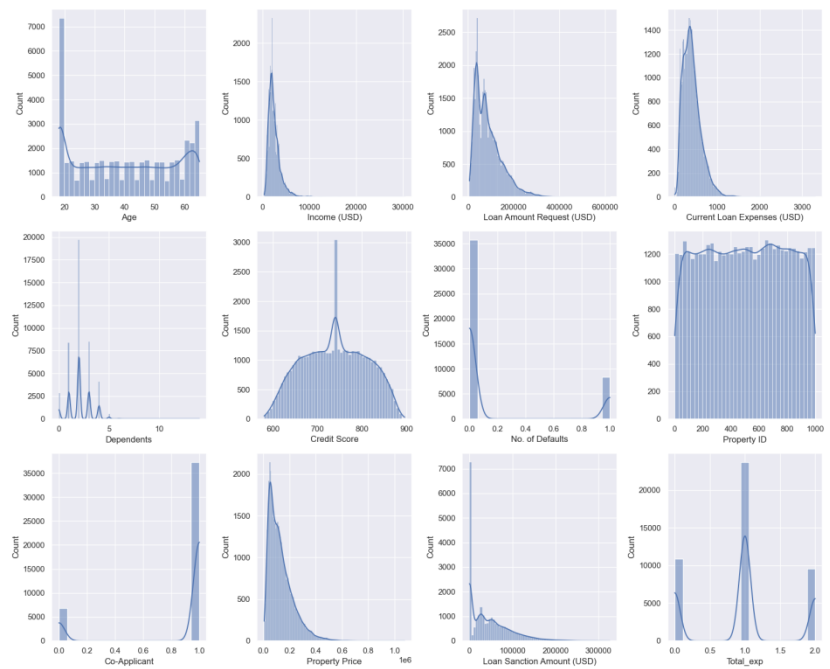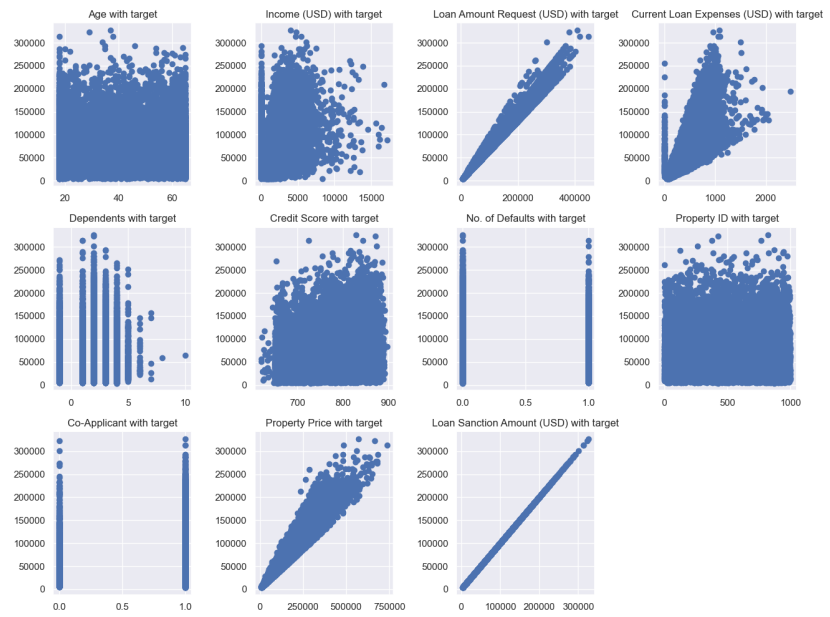| | |
|---|---|
| Date | 15 June 2024 |
| Team ID | 740185 |
| Project Title | Loan Sanction Amount Prediction Data With Ml |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**
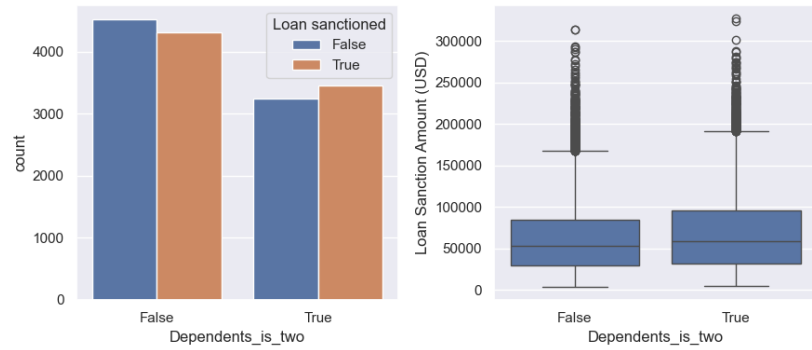
**Data Exploration and Preprocessing Template**

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

| Section | Description |
|---|---|
| Data Overview | Dimension:<br>8 rows x 131 columns<br>Descriptive stastistics:<br> |

| Univariate Analysis |  |

| | |
|---|---|
| Bivariate Analysis |  |
| | |
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

**Data Preprocessing Code Screenshots**

| Loading train Data |  |
| --- | --- |
| Loading test Data |  |
| Handling Missing Data In train and test |  |

| Data Transformation |  |
|---|---|
| Feature Engineering | - |
| Save Processed Data | - |