

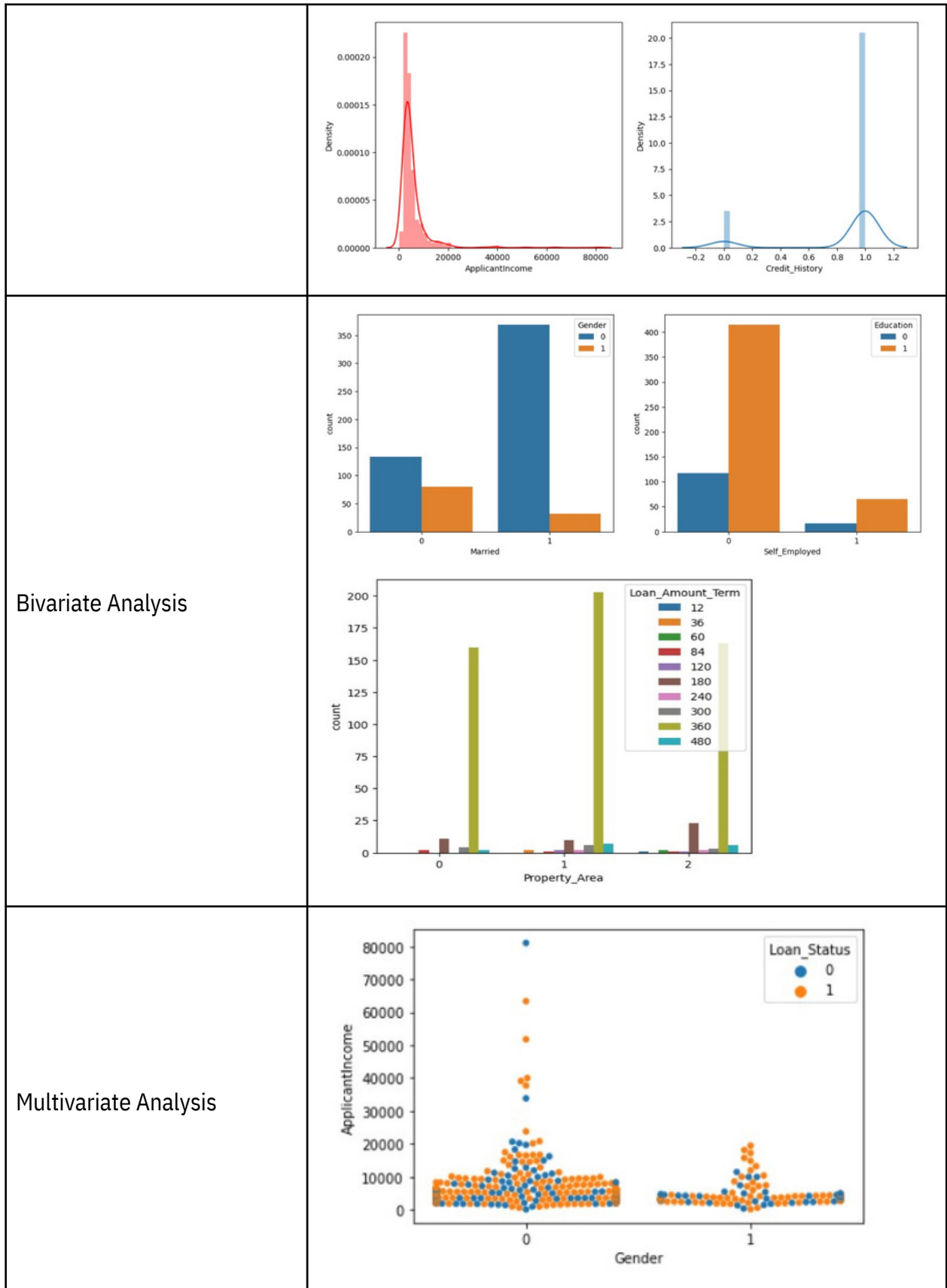
## Data Collection and Preprocessing Phase

Date	27 January 2025
Team ID	SWUID20240011509
Project Title	Restaurant Recommendation System
Maximum Marks	6 Marks

### Data Exploration and Preprocessing Report

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	<u>Dimension:</u> 614 rows × 13 columns
	<u>Descriptive statistics:</u>



Outliers and Anomalies	-																																																						
Data Preprocessing Code Screenshots																																																							
Loading Data	<pre>#importing the dataset which is in csv file data = pd.read_csv('/content/Dataset/loan_prediction.csv') data</pre> <table><thead><tr><th></th><th>Loan_ID</th><th>Gender</th><th>Married</th><th>Dependents</th><th>Education</th><th>Self_Employed</th><th>ApplicantIncome</th><th>CoapplicantIncome</th></tr></thead><tbody><tr><td>0</td><td>LP001002</td><td>Male</td><td>No</td><td>0</td><td>Graduate</td><td>No</td><td>5849</td><td>0.0</td></tr><tr><td>1</td><td>LP001003</td><td>Male</td><td>Yes</td><td>1</td><td>Graduate</td><td>No</td><td>4583</td><td>1508.0</td></tr><tr><td>2</td><td>LP001005</td><td>Male</td><td>Yes</td><td>0</td><td>Graduate</td><td>Yes</td><td>3000</td><td>0.0</td></tr><tr><td>3</td><td>LP001006</td><td>Male</td><td>Yes</td><td>0</td><td>Not Graduate</td><td>No</td><td>2583</td><td>2358.0</td></tr><tr><td>4</td><td>LP001008</td><td>Male</td><td>No</td><td>0</td><td>Graduate</td><td>No</td><td>6000</td><td>0.0</td></tr></tbody></table>		Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	0	LP001002	Male	No	0	Graduate	No	5849	0.0	1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	4	LP001008	Male	No	0	Graduate	No	6000	0.0
	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome																																															
0	LP001002	Male	No	0	Graduate	No	5849	0.0																																															
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0																																															
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0																																															
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0																																															
4	LP001008	Male	No	0	Graduate	No	6000	0.0																																															
Handling Missing Data	<pre>data['Gender'] = data['Gender'].fillna(data['Gender'].mode()[0])  data['Married'] = data['Married'].fillna(data['Married'].mode()[0])  #replacing + with space for filling the nan values data['Dependents']=data['Dependents'].str.replace('+','')  &lt;ipython-input-71-6ac39c248773&gt;:2: FutureWarning: The default value of regex will change from data['Dependents']=data['Dependents'].str.replace('+','')  data['Dependents'] = data['Dependents'].fillna(data['Dependents'].mode()[0])  data['Self_Employed'] = data['Self_Employed'].fillna(data['Self_Employed'].mode()[0])  data['LoanAmount'] = data['LoanAmount'].fillna(data['LoanAmount'].mode()[0])  data['Loan_Amount_Term'] = data['Loan_Amount_Term'].fillna(data['Loan_Amount_Term'].mode()[0])  data['Credit_History'] = data['Credit_History'].fillna(data['Credit_History'].mode()[0])</pre>																																																						
Data Transformation	<pre>data['Gender']=data['Gender'].map({'Female':1,'Male':0}) data['Property_Area']=data['Property_Area'].map({'Urban':2,'Semiurban': 1,'Rural':0}) data['Married']=data['Married'].map({'Yes':1,'No':0}) data['Education']=data['Education'].map({'Graduate':1,'Not Graduate':0}) data['Loan_Status']=data['Loan_Status'].map({'Y':1,'N':0})  # performing feature Scaling operation using standard scaller on X part of the dataset because # there different type of values in the columns sc=StandardScaler() x_bal=sc.fit_transform(x_bal)</pre>																																																						
Feature Engineering	Attached the codes in final submission.																																																						