

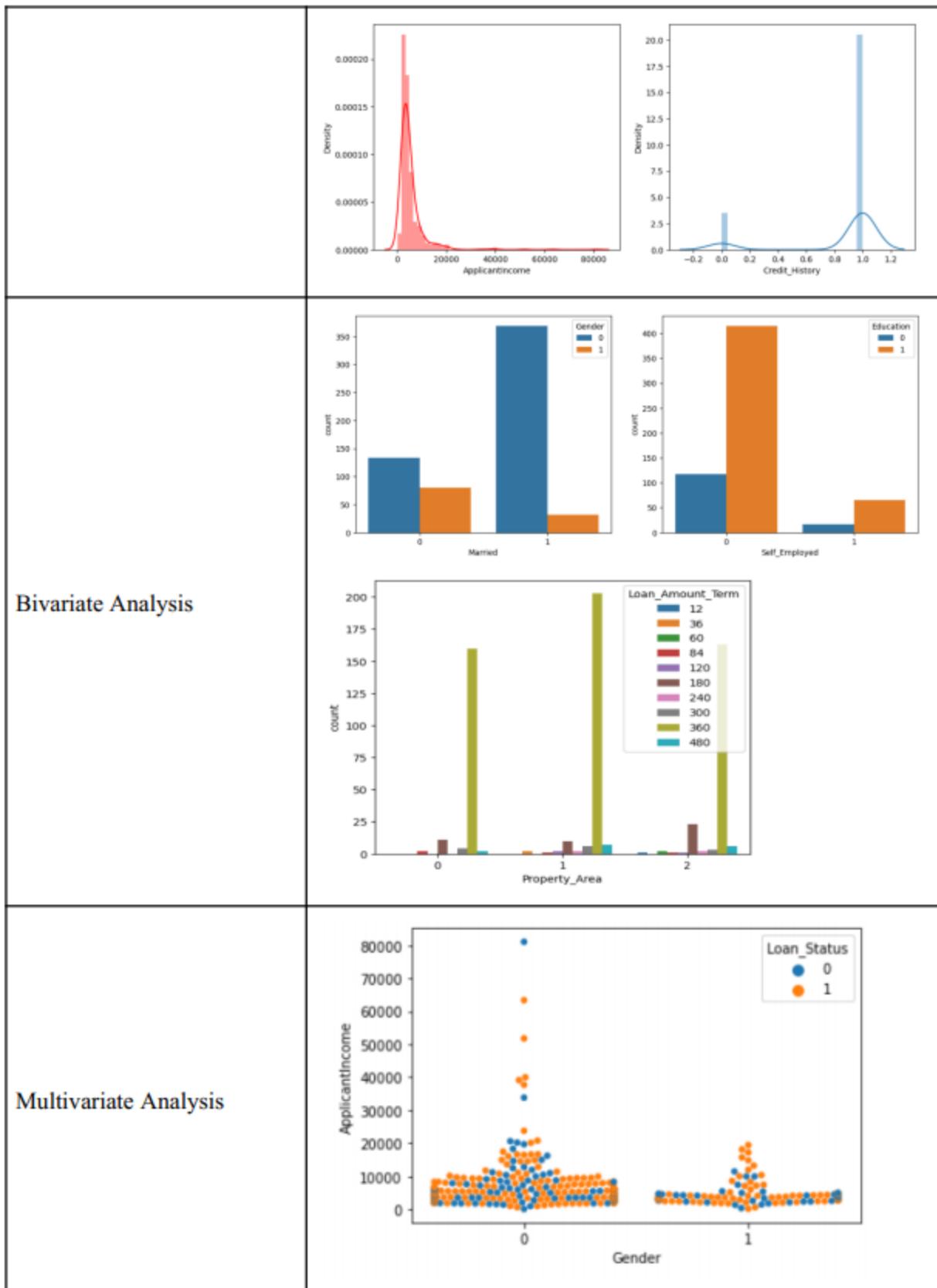
Data Collection and Preprocessing Phase

Date	20 February 2026
Team ID	LTVIP2026TMIDS46296
Project Title	Advancing nutrition science through geminial
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	<p><u>Dimension:</u> 614 rows × 13 columns</p> <p><u>Descriptive statistics:</u></p> <table border="1"> <thead> <tr> <th></th> <th>ApplicantIncome</th> <th>CosapplicantIncome</th> <th>LoanAmount</th> <th>Loan_Amount_Term</th> <th>Credit_History</th> </tr> </thead> <tbody> <tr> <td>count</td> <td>614.000000</td> <td>614.000000</td> <td>592.000000</td> <td>600.00000</td> <td>564.000000</td> </tr> <tr> <td>mean</td> <td>5403.459283</td> <td>1621.245798</td> <td>146.412162</td> <td>342.00000</td> <td>0.842199</td> </tr> <tr> <td>std</td> <td>6109.041673</td> <td>2926.248369</td> <td>85.587325</td> <td>65.12041</td> <td>0.364878</td> </tr> <tr> <td>min</td> <td>150.000000</td> <td>0.000000</td> <td>9.000000</td> <td>12.00000</td> <td>0.000000</td> </tr> <tr> <td>25%</td> <td>2877.500000</td> <td>0.000000</td> <td>100.000000</td> <td>360.00000</td> <td>1.000000</td> </tr> <tr> <td>50%</td> <td>3812.500000</td> <td>1188.500000</td> <td>128.000000</td> <td>360.00000</td> <td>1.000000</td> </tr> <tr> <td>75%</td> <td>5795.000000</td> <td>2297.250000</td> <td>168.000000</td> <td>360.00000</td> <td>1.000000</td> </tr> <tr> <td>max</td> <td>81000.000000</td> <td>41667.000000</td> <td>700.000000</td> <td>480.00000</td> <td>1.000000</td> </tr> </tbody> </table>		ApplicantIncome	CosapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	count	614.000000	614.000000	592.000000	600.00000	564.000000	mean	5403.459283	1621.245798	146.412162	342.00000	0.842199	std	6109.041673	2926.248369	85.587325	65.12041	0.364878	min	150.000000	0.000000	9.000000	12.00000	0.000000	25%	2877.500000	0.000000	100.000000	360.00000	1.000000	50%	3812.500000	1188.500000	128.000000	360.00000	1.000000	75%	5795.000000	2297.250000	168.000000	360.00000	1.000000	max	81000.000000	41667.000000	700.000000	480.00000	1.000000
	ApplicantIncome	CosapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History																																																		
count	614.000000	614.000000	592.000000	600.00000	564.000000																																																		
mean	5403.459283	1621.245798	146.412162	342.00000	0.842199																																																		
std	6109.041673	2926.248369	85.587325	65.12041	0.364878																																																		
min	150.000000	0.000000	9.000000	12.00000	0.000000																																																		
25%	2877.500000	0.000000	100.000000	360.00000	1.000000																																																		
50%	3812.500000	1188.500000	128.000000	360.00000	1.000000																																																		
75%	5795.000000	2297.250000	168.000000	360.00000	1.000000																																																		
max	81000.000000	41667.000000	700.000000	480.00000	1.000000																																																		
Univariate Analysis																																																							



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 border="1"> <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>CoaapplicantIncome</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	CoaapplicantIncome	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	CoaapplicantIncome																																															
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('+','') <ipython-input-71-6ac39c248773>: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 scaler 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.																																																						
Save Processed Data	-																																																						