

Data Collection and Preprocessing Phase

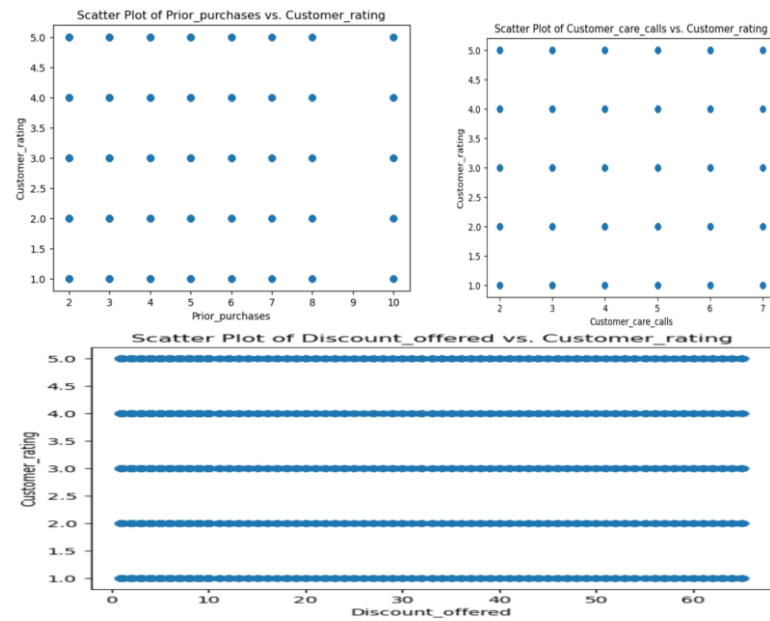
Date	11 JULY 2024
Team ID	SWTID1720115788
Project Title	Ecommerce Shipping Prediction Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

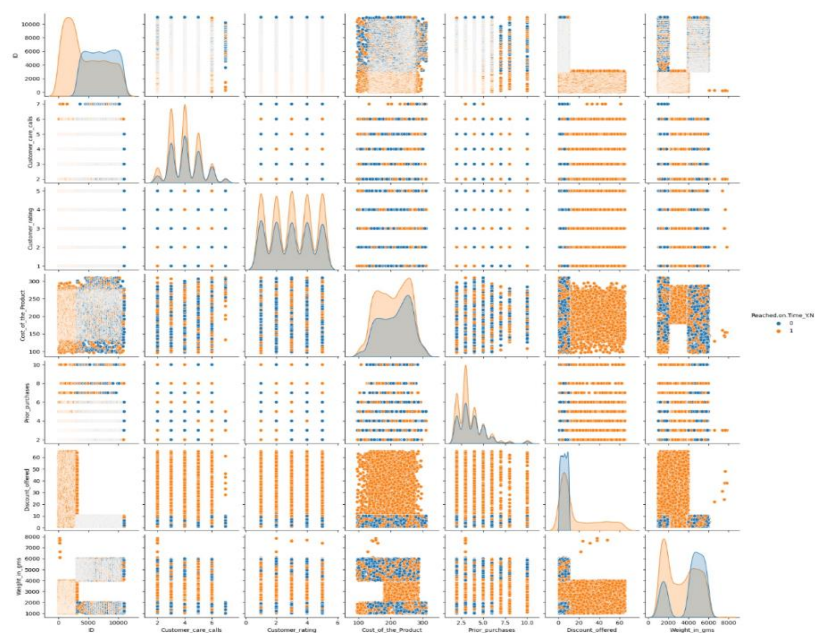
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																																																	
Data Overview	<table><thead><tr><th></th><th>ID</th><th>Customer_care_calls</th><th>Customer_rating</th><th>Cost_of_the_Product</th><th>Prior_purchases</th><th>Discount_offered</th><th>Weight_in_gms</th><th>Reached.on.Time_Y_N</th></tr></thead><tbody><tr><td>count</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td><td>10999.000000</td></tr><tr><td>mean</td><td>5500.00000</td><td>4.054459</td><td>2.990545</td><td>210.196836</td><td>3.567597</td><td>13.373216</td><td>3634.016729</td><td>0.596691</td></tr><tr><td>std</td><td>3175.28214</td><td>1.141490</td><td>1.413603</td><td>48.063272</td><td>1.522860</td><td>16.205527</td><td>1635.377251</td><td>0.490584</td></tr><tr><td>min</td><td>1.000000</td><td>2.000000</td><td>1.000000</td><td>96.000000</td><td>2.000000</td><td>1.000000</td><td>1001.000000</td><td>0.000000</td></tr><tr><td>25%</td><td>2750.500000</td><td>3.000000</td><td>2.000000</td><td>169.000000</td><td>3.000000</td><td>4.000000</td><td>1839.500000</td><td>0.000000</td></tr><tr><td>50%</td><td>5500.000000</td><td>4.000000</td><td>3.000000</td><td>214.000000</td><td>3.000000</td><td>7.000000</td><td>4149.000000</td><td>1.000000</td></tr><tr><td>75%</td><td>8249.500000</td><td>5.000000</td><td>4.000000</td><td>251.000000</td><td>4.000000</td><td>10.000000</td><td>5050.000000</td><td>1.000000</td></tr><tr><td>max</td><td>10999.000000</td><td>7.000000</td><td>5.000000</td><td>310.000000</td><td>10.000000</td><td>65.000000</td><td>7846.000000</td><td>1.000000</td></tr></tbody></table>		ID	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Discount_offered	Weight_in_gms	Reached.on.Time_Y_N	count	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	mean	5500.00000	4.054459	2.990545	210.196836	3.567597	13.373216	3634.016729	0.596691	std	3175.28214	1.141490	1.413603	48.063272	1.522860	16.205527	1635.377251	0.490584	min	1.000000	2.000000	1.000000	96.000000	2.000000	1.000000	1001.000000	0.000000	25%	2750.500000	3.000000	2.000000	169.000000	3.000000	4.000000	1839.500000	0.000000	50%	5500.000000	4.000000	3.000000	214.000000	3.000000	7.000000	4149.000000	1.000000	75%	8249.500000	5.000000	4.000000	251.000000	4.000000	10.000000	5050.000000	1.000000	max	10999.000000	7.000000	5.000000	310.000000	10.000000	65.000000	7846.000000	1.000000
	ID	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Discount_offered	Weight_in_gms	Reached.on.Time_Y_N																																																																										
count	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000																																																																										
mean	5500.00000	4.054459	2.990545	210.196836	3.567597	13.373216	3634.016729	0.596691																																																																										
std	3175.28214	1.141490	1.413603	48.063272	1.522860	16.205527	1635.377251	0.490584																																																																										
min	1.000000	2.000000	1.000000	96.000000	2.000000	1.000000	1001.000000	0.000000																																																																										
25%	2750.500000	3.000000	2.000000	169.000000	3.000000	4.000000	1839.500000	0.000000																																																																										
50%	5500.000000	4.000000	3.000000	214.000000	3.000000	7.000000	4149.000000	1.000000																																																																										
75%	8249.500000	5.000000	4.000000	251.000000	4.000000	10.000000	5050.000000	1.000000																																																																										
max	10999.000000	7.000000	5.000000	310.000000	10.000000	65.000000	7846.000000	1.000000																																																																										
Univariate Analysis	<div><div><div>count 18999.000000 mean 2.998545 std 1.413683 min 1.000000 25% 2.000000 50% 3.000000 75% 4.000000 max 5.000000 Name: Customer_rating, dtype: float64</div><div></div></div><div><div>count 18999 unique 3 top Ship freq 7462 Name: Mode_of_Shipment, dtype: object</div><div></div></div><div><div>count 10999 unique 2 top F freq 5545 Name: Gender, dtype: object</div><div></div></div></div>																																																																																	

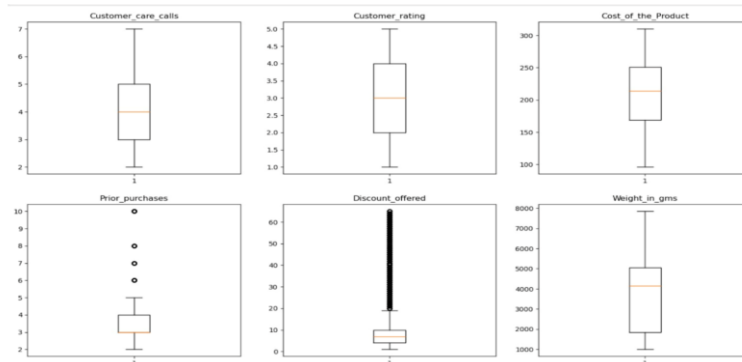
Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots	
Loading Data	<pre>import pandas as pd df = pd.read_csv('/Users/shanmukhanandudu/Downloads/train (3).csv') df</pre>
Handling Missing Data	<pre>df.isnull().sum()</pre>
Data Transformation	<pre># Encode categorical variables le = LabelEncoder() df['Warehouse_block'] = le.fit_transform(df['Warehouse_block']) df['Mode_of_Shipment'] = le.fit_transform(df['Mode_of_Shipment']) df['Product_importance'] = le.fit_transform(df['Product_importance']) df['Gender'] = le.fit_transform(df['Gender']) # Scale/normalize features scaler = StandardScaler() columns_to_scale = ['Customer_care_calls', 'Customer_rating', 'Cost_of_the_Product', 'Prior_purchases', 'Discount_offered', 'Weight_in_gms'] df[columns_to_scale] = scaler.fit_transform(df[columns_to_scale])</pre>
Feature Engineering	<pre>import pandas as pd # create a sample dataframe data = {'priority': ['low', 'medium', 'high', 'low', 'medium', 'high']} df = pd.DataFrame(data) # create a new column with the mapped values df['priority_code'] = df['priority'].map({'low': 0, 'medium': 1, 'high': 2}) print(df)</pre>
Save Processed Data	<pre>df.to_csv('my_dataset.csv', index=False)</pre>

