

Data Collection and Preprocessing Phase

Date	10 July 2024
Team ID	SWTID1720369851
Project Title	Ecommerce Shipping Prediction
Maximum Marks	6 Marks

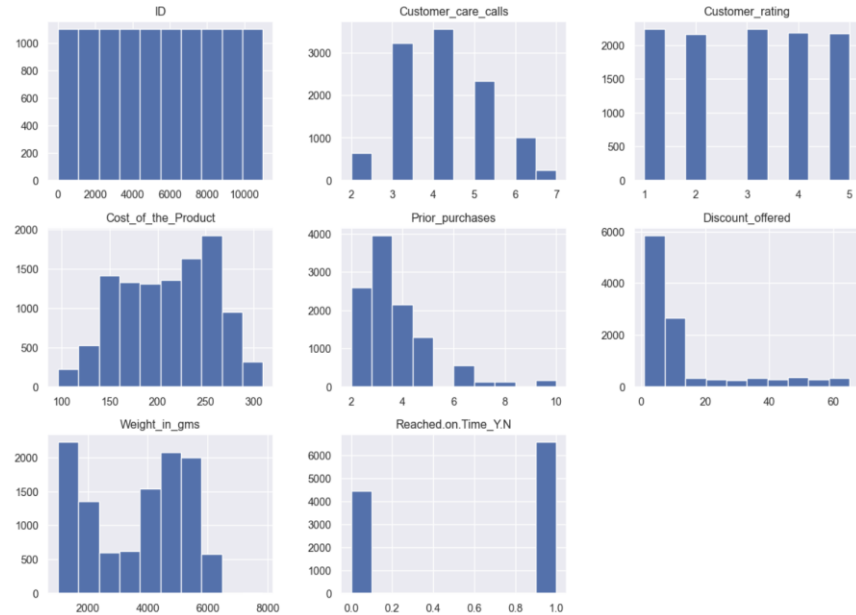
Data Exploration and Preprocessing Template

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	<pre>data.shape</pre>																																																																																	
	<pre>(10999, 12)</pre>																																																																																	
	<pre>data.describe()</pre>																																																																																	
	<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>5500.000000</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	5500.000000	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	5500.000000	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																																																																										
<pre>data.info()</pre>																																																																																		
<pre><class 'pandas.core.frame.DataFrame'> RangeIndex: 10999 entries, 0 to 10998 Data columns (total 12 columns): # Column Non-Null Count Dtype --- - 0 ID 10999 non-null int64 1 Warehouse_block 10999 non-null object 2 Mode_of_Shipment 10999 non-null object 3 Customer_care_calls 10999 non-null int64 4 Customer_rating 10999 non-null int64 5 Cost_of_the_Product 10999 non-null int64 6 Prior_purchases 10999 non-null int64 7 Product_importance 10999 non-null object 8 Gender 10999 non-null object 9 Discount_offered 10999 non-null int64 10 Weight_in_gms 10999 non-null int64 11 Reached.on.Time_Y.N 10999 non-null int64 dtypes: int64(8), object(4) memory usage: 1.0+ MB</pre>																																																																																		

Univariate Analysis

```
data.select_dtypes(include=[np.number]).hist(bins=10, figsize=(15,10))
plt.show()
```

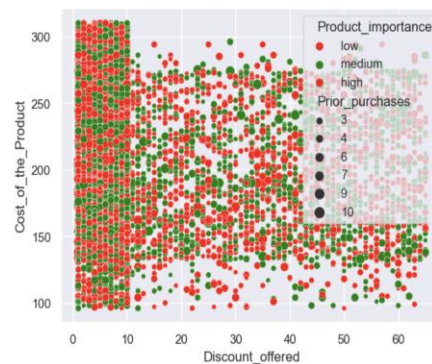


Bivariate Analysis

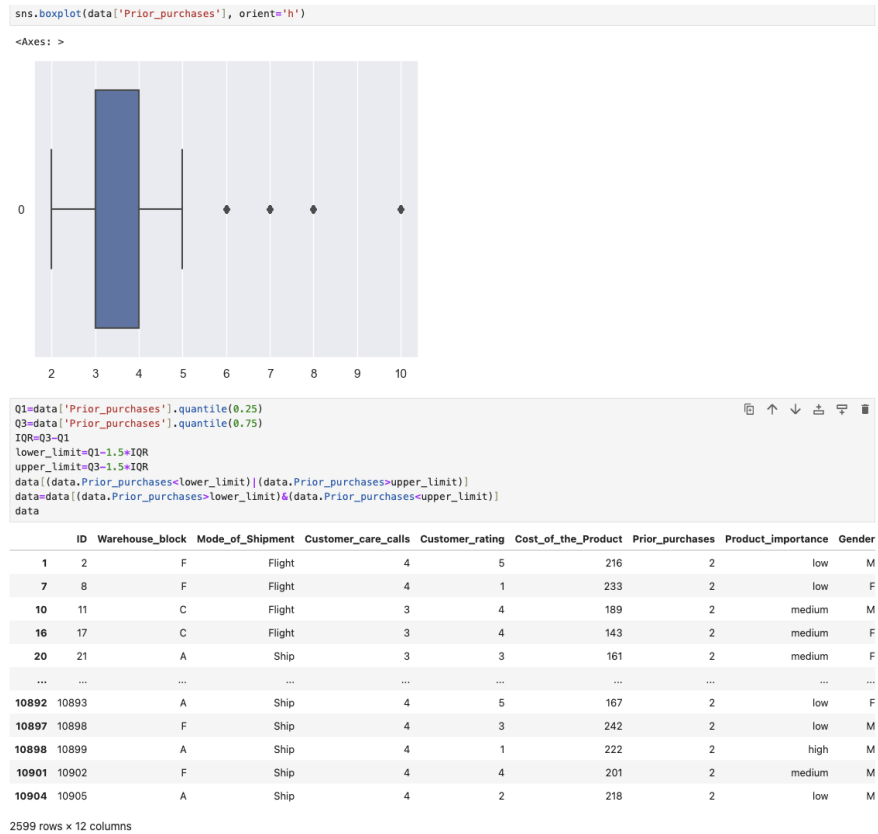
```
data.corr(numeric_only=True)
```

	ID	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N
ID	1.000000	0.188998	-0.005722	0.196791	0.145369	-0.598278	0.278312	-0.411822
Customer_care_calls	0.188998	1.000000	0.012209	0.323182	0.180771	-0.130750	-0.276615	-0.067126
Customer_rating	-0.005722	0.012209	1.000000	0.009270	0.013179	-0.003124	-0.001897	0.013119
Cost_of_the_Product	0.196791	0.323182	0.009270	1.000000	0.123676	-0.138312	-0.132604	-0.073587
Prior_purchases	0.145369	0.180771	0.013179	0.123676	1.000000	-0.082769	-0.168213	-0.055515
Discount_offered	-0.598278	-0.130750	-0.003124	-0.138312	-0.082769	1.000000	-0.376067	0.397108
Weight_in_gms	0.278312	-0.276615	-0.001897	-0.132604	-0.168213	-0.376067	1.000000	-0.268793
Reached.on.Time_Y.N	-0.411822	-0.067126	0.013119	-0.073587	-0.055515	0.397108	-0.268793	1.000000

```
# Scatterplot
sns.scatterplot(x='Discount_offered',y='Cost_of_the_Product',data=data, hue='Product_importance', size='Prior_purchases', palette=['red','green'])
<Axes: xlabel='Discount_offered', ylabel='Cost_of_the_Product'>
```



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data

```
# Load the data

data = pd.read_csv("Train.csv")
```

Handling Missing Data	<pre>data.isnull().sum()</pre> <table><tr><td>ID</td><td>0</td></tr><tr><td>Warehouse_block</td><td>0</td></tr><tr><td>Mode_of_Shipment</td><td>0</td></tr><tr><td>Customer_care_calls</td><td>0</td></tr><tr><td>Customer_rating</td><td>0</td></tr><tr><td>Cost_of_the_Product</td><td>0</td></tr><tr><td>Prior_purchases</td><td>0</td></tr><tr><td>Product_importance</td><td>0</td></tr><tr><td>Gender</td><td>0</td></tr><tr><td>Discount_offered</td><td>0</td></tr><tr><td>Weight_in_gms</td><td>0</td></tr><tr><td>Reached.on.Time_Y.N</td><td>0</td></tr></table> <p>dtype: int64</p> <pre>data.duplicated().sum()</pre> <p>0</p>	ID	0	Warehouse_block	0	Mode_of_Shipment	0	Customer_care_calls	0	Customer_rating	0	Cost_of_the_Product	0	Prior_purchases	0	Product_importance	0	Gender	0	Discount_offered	0	Weight_in_gms	0	Reached.on.Time_Y.N	0																								
ID	0																																																
Warehouse_block	0																																																
Mode_of_Shipment	0																																																
Customer_care_calls	0																																																
Customer_rating	0																																																
Cost_of_the_Product	0																																																
Prior_purchases	0																																																
Product_importance	0																																																
Gender	0																																																
Discount_offered	0																																																
Weight_in_gms	0																																																
Reached.on.Time_Y.N	0																																																
Data Transformation	<p>Encoding the categorical variables</p> <pre>le = LabelEncoder()</pre> <pre>def Label_Enc(col): Categorical_col[col] = le.fit_transform(Categorical_col[col])</pre> <pre>for i in ['Warehouse_block', 'Mode_of_Shipment', 'Product_importance', 'Gender']: Label_Enc(i)</pre> <pre>Categorical_col.head()</pre> <table><tr><th></th><th>Warehouse_block</th><th>Mode_of_Shipment</th><th>Product_importance</th><th>Gender</th></tr><tr><td>0</td><td>3</td><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>4</td><td>0</td><td>1</td><td>1</td></tr><tr><td>2</td><td>0</td><td>0</td><td>1</td><td>1</td></tr><tr><td>3</td><td>1</td><td>0</td><td>2</td><td>1</td></tr><tr><td>4</td><td>2</td><td>0</td><td>2</td><td>0</td></tr></table>		Warehouse_block	Mode_of_Shipment	Product_importance	Gender	0	3	0	1	0	1	4	0	1	1	2	0	0	1	1	3	1	0	2	1	4	2	0	2	0																		
	Warehouse_block	Mode_of_Shipment	Product_importance	Gender																																													
0	3	0	1	0																																													
1	4	0	1	1																																													
2	0	0	1	1																																													
3	1	0	2	1																																													
4	2	0	2	0																																													
Feature Engineering	<pre>Numerical_col.drop(columns = ["ID"],axis = 1,inplace = True) Numerical_col.head()</pre> <table><tr><th></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><tr><td>0</td><td>4</td><td>2</td><td>177</td><td>3</td><td>44</td><td>1233</td><td>1</td></tr><tr><td>1</td><td>4</td><td>5</td><td>216</td><td>2</td><td>59</td><td>3088</td><td>1</td></tr><tr><td>2</td><td>2</td><td>2</td><td>183</td><td>4</td><td>48</td><td>3374</td><td>1</td></tr><tr><td>3</td><td>3</td><td>3</td><td>176</td><td>4</td><td>10</td><td>1177</td><td>1</td></tr><tr><td>4</td><td>2</td><td>2</td><td>184</td><td>3</td><td>46</td><td>2484</td><td>1</td></tr></table>		Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N	0	4	2	177	3	44	1233	1	1	4	5	216	2	59	3088	1	2	2	2	183	4	48	3374	1	3	3	3	176	4	10	1177	1	4	2	2	184	3	46	2484	1
	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N																																										
0	4	2	177	3	44	1233	1																																										
1	4	5	216	2	59	3088	1																																										
2	2	2	183	4	48	3374	1																																										
3	3	3	176	4	10	1177	1																																										
4	2	2	184	3	46	2484	1																																										
Save Processed Data	<pre># Save to a CSV file data.to_csv('Train(new).csv', index=False)</pre>																																																