

Data Collection and Preprocessing Phase

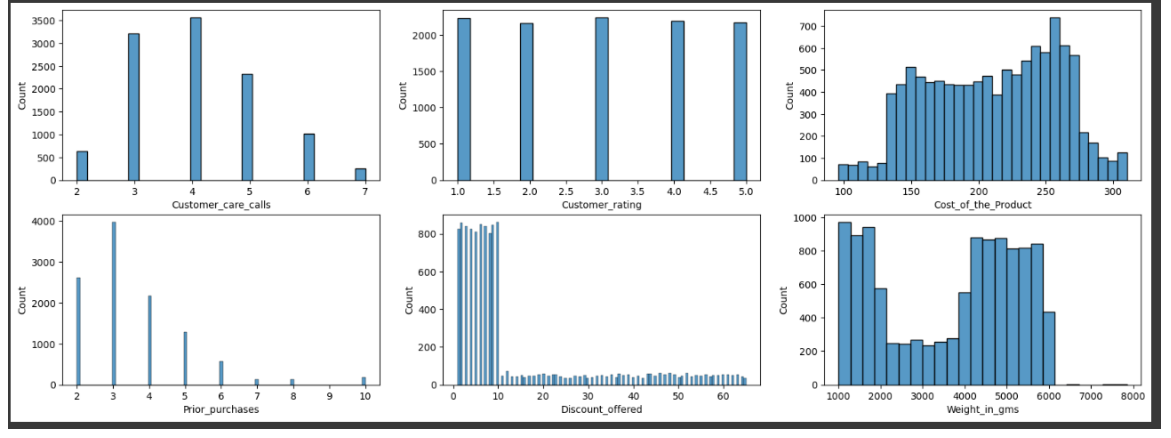
Date	5 July 2024
Team ID	SWTID1720082658
Project Title	Ecommerce Shipping Prediction Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

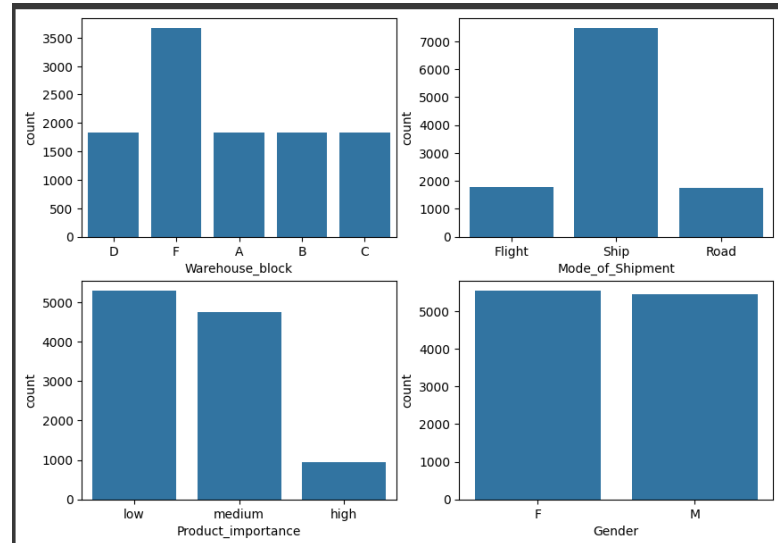
Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong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

Section	Description																																																															
Data Overview	Dimensions : 10999 rows x 12 columns																																																															
	<table><tr><th></th><th>ID</th><th>Warehouse_block</th><th>Mode_of_Shipment</th><th>Customer_care_calls</th><th>Customer_rating</th><th>Cost_of_the_Product</th></tr><tr><td>count</td><td>10999.00000</td><td>10999.00000</td><td>10999.00000</td><td>10999.00000</td><td>10999.00000</td><td>10999.00000</td></tr><tr><td>mean</td><td>5500.00000</td><td>2.333394</td><td>1.516865</td><td>4.054459</td><td>2.990545</td><td>210.196836</td></tr><tr><td>std</td><td>3175.28214</td><td>1.490726</td><td>0.756894</td><td>1.141490</td><td>1.413603</td><td>48.063272</td></tr><tr><td>min</td><td>1.00000</td><td>0.00000</td><td>0.00000</td><td>2.00000</td><td>1.00000</td><td>96.00000</td></tr><tr><td>25%</td><td>2750.50000</td><td>1.00000</td><td>1.00000</td><td>3.00000</td><td>2.00000</td><td>169.00000</td></tr><tr><td>50%</td><td>5500.00000</td><td>3.00000</td><td>2.00000</td><td>4.00000</td><td>3.00000</td><td>214.00000</td></tr><tr><td>75%</td><td>8249.50000</td><td>4.00000</td><td>2.00000</td><td>5.00000</td><td>4.00000</td><td>251.00000</td></tr><tr><td>max</td><td>10999.00000</td><td>4.00000</td><td>2.00000</td><td>7.00000</td><td>5.00000</td><td>310.00000</td></tr></table>		ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product	count	10999.00000	10999.00000	10999.00000	10999.00000	10999.00000	10999.00000	mean	5500.00000	2.333394	1.516865	4.054459	2.990545	210.196836	std	3175.28214	1.490726	0.756894	1.141490	1.413603	48.063272	min	1.00000	0.00000	0.00000	2.00000	1.00000	96.00000	25%	2750.50000	1.00000	1.00000	3.00000	2.00000	169.00000	50%	5500.00000	3.00000	2.00000	4.00000	3.00000	214.00000	75%	8249.50000	4.00000	2.00000	5.00000	4.00000	251.00000	max	10999.00000	4.00000	2.00000	7.00000	5.00000	310.00000
		ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product																																																									
	count	10999.00000	10999.00000	10999.00000	10999.00000	10999.00000	10999.00000																																																									
	mean	5500.00000	2.333394	1.516865	4.054459	2.990545	210.196836																																																									
	std	3175.28214	1.490726	0.756894	1.141490	1.413603	48.063272																																																									
	min	1.00000	0.00000	0.00000	2.00000	1.00000	96.00000																																																									
	25%	2750.50000	1.00000	1.00000	3.00000	2.00000	169.00000																																																									
	50%	5500.00000	3.00000	2.00000	4.00000	3.00000	214.00000																																																									
	75%	8249.50000	4.00000	2.00000	5.00000	4.00000	251.00000																																																									
max	10999.00000	4.00000	2.00000	7.00000	5.00000	310.00000																																																										

Univariate Analysis



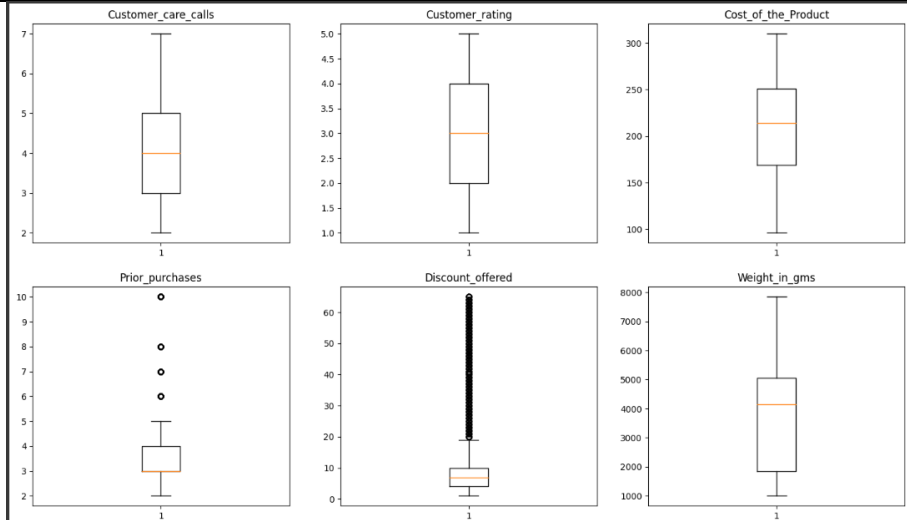
Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data

```
[5] data=pd.read_csv("Train.csv")
data.head()
```

	ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Product_importance	Gender	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N
0	1	D	Flight	4	2	177	3	low	F	44	1233	1
1	2	F	Flight	4	5	216	2	low	M	59	3088	1
2	3	A	Flight	2	2	183	4	low	M	48	3374	1
3	4	B	Flight	3	3	176	4	medium	M	10	1177	1
4	5	C	Flight	2	2	184	3	medium	F	46	2484	1

Handling Missing Data

```
[7] data.isnull().sum()
```

	sum
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
dtype: int64	

Data Encoding	<pre>[10] from sklearn.preprocessing import LabelEncoder le = LabelEncoder() data.Warehouse_block = le.fit_transform(data.Warehouse_block) data.Mode_of_Shipment = le.fit_transform(data.Mode_of_Shipment) data.Product_importance = le.fit_transform(data.Product_importance) data.Gender = le.fit_transform(data.Gender) data.head()</pre> <table><thead><tr><th></th><th>ID</th><th>Warehouse_block</th><th>Mode_of_Shipment</th><th>Customer_care_calls</th><th>Customer_rating</th><th>Cost_of_the_Product</th><th>Prior_purchases</th><th>Product_importance</th><th>Gender</th><th>Discount_offered</th><th>Weight_in_gms</th><th>Reached.on.Time_Y.N</th></tr></thead><tbody><tr><td>0</td><td>1</td><td>3</td><td>0</td><td>4</td><td>2</td><td>177</td><td>3</td><td>1</td><td>0</td><td>44</td><td>1233</td><td>1</td></tr><tr><td>1</td><td>2</td><td>4</td><td>0</td><td>4</td><td>5</td><td>216</td><td>2</td><td>1</td><td>1</td><td>59</td><td>3088</td><td>1</td></tr><tr><td>2</td><td>3</td><td>0</td><td>0</td><td>2</td><td>2</td><td>183</td><td>4</td><td>1</td><td>1</td><td>48</td><td>3374</td><td>1</td></tr><tr><td>3</td><td>4</td><td>1</td><td>0</td><td>3</td><td>3</td><td>176</td><td>4</td><td>2</td><td>1</td><td>10</td><td>1177</td><td>1</td></tr><tr><td>4</td><td>5</td><td>2</td><td>0</td><td>2</td><td>2</td><td>184</td><td>3</td><td>2</td><td>0</td><td>46</td><td>2484</td><td>1</td></tr></tbody></table>		ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Product_importance	Gender	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N	0	1	3	0	4	2	177	3	1	0	44	1233	1	1	2	4	0	4	5	216	2	1	1	59	3088	1	2	3	0	0	2	2	183	4	1	1	48	3374	1	3	4	1	0	3	3	176	4	2	1	10	1177	1	4	5	2	0	2	2	184	3	2	0	46	2484	1
	ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purchases	Product_importance	Gender	Discount_offered	Weight_in_gms	Reached.on.Time_Y.N																																																																			
0	1	3	0	4	2	177	3	1	0	44	1233	1																																																																			
1	2	4	0	4	5	216	2	1	1	59	3088	1																																																																			
2	3	0	0	2	2	183	4	1	1	48	3374	1																																																																			
3	4	1	0	3	3	176	4	2	1	10	1177	1																																																																			
4	5	2	0	2	2	184	3	2	0	46	2484	1																																																																			
Data Transformation	<pre>[18] from sklearn.preprocessing import StandardScaler scale=StandardScaler() xnorm_train = scale.fit_transform(x_train) xnorm_test = scale.fit_transform(x_test)</pre> <pre>[19] from sklearn.preprocessing import MinMaxScaler norm=MinMaxScaler() x=norm.fit_transform(x) x</pre>																																																																														
Feature Engineering	Code is in the final code submitted.																																																																														
Save Processed Data	-																																																																														