# Data Collection and Preprocessing Phase
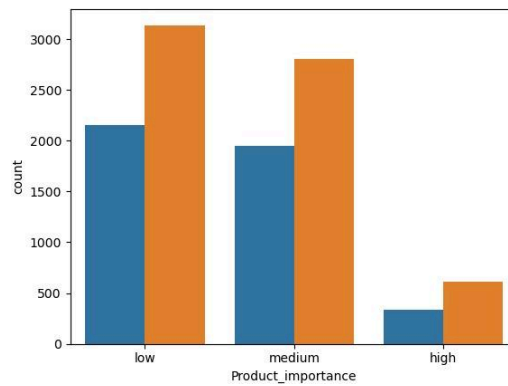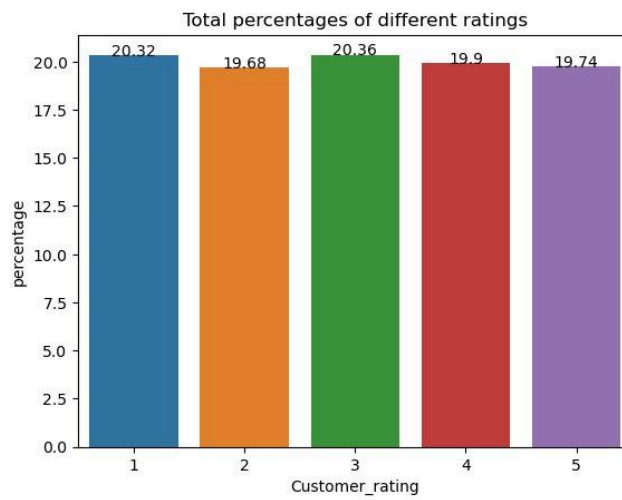
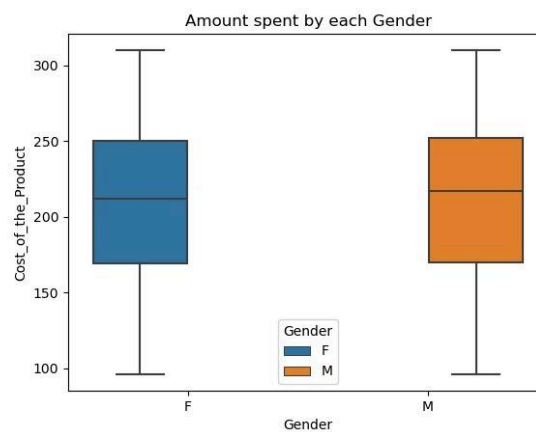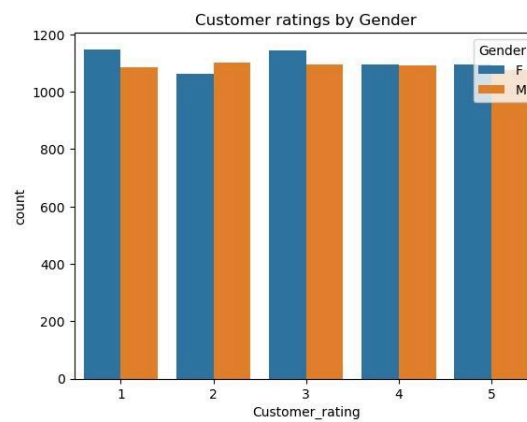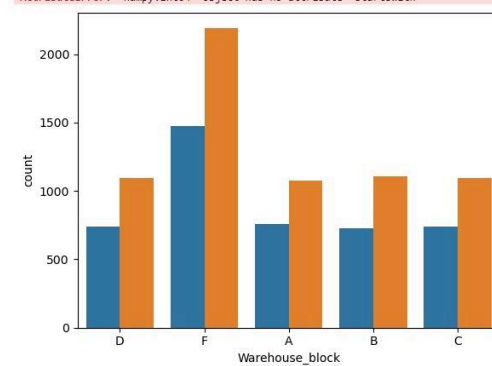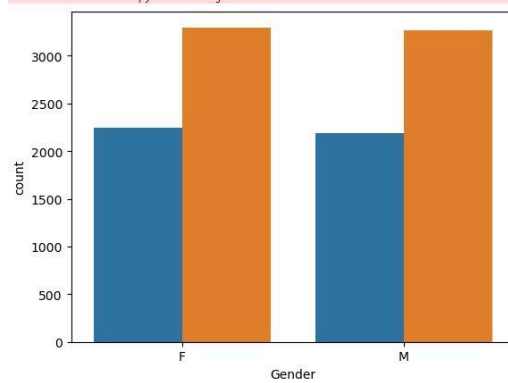| Date | 11 JULY 2024 |
|---|---|
| Team ID | SWTID1720116037 |
| Project Title | Ecommerce Shipping Prediction Using Machine Learning |
| Maximum Marks | 6 Marks |

## 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.
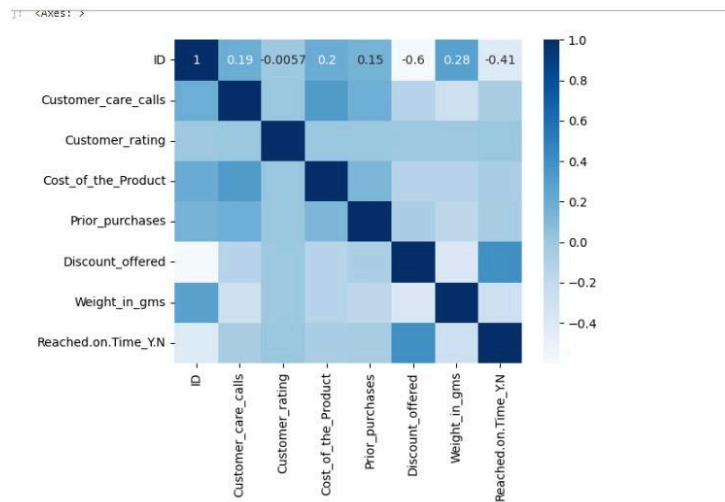
| Section | Description |
|---|---|
| Data Overview |  |
| Univariate Analysis |  |

**Total number of delayed deliveries vs Warehouse block**



**Total percentages of different ratings**



Bivariate Analysis

Customer ratings by Gender



Amount spent by each Gender
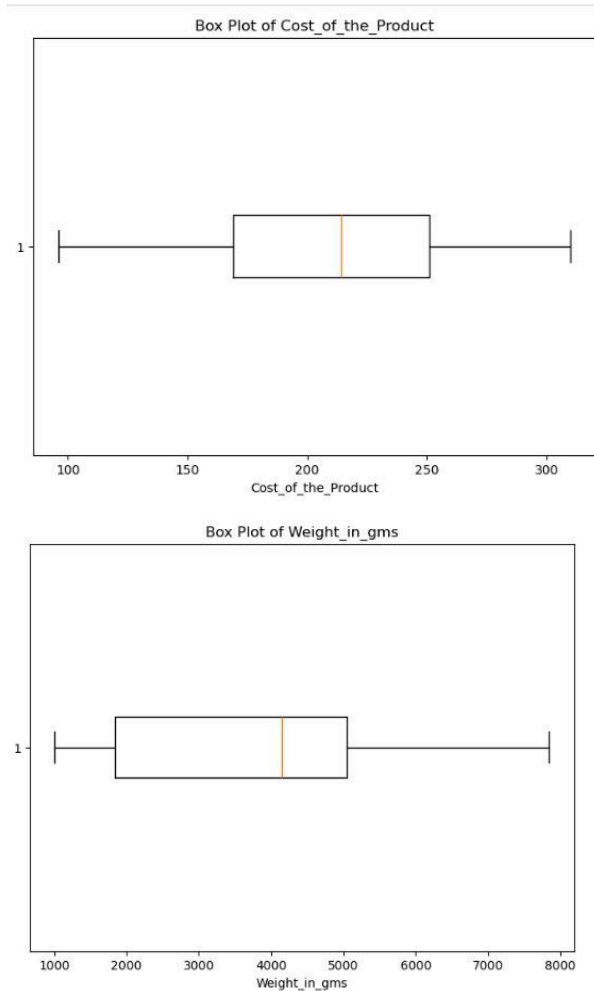
| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies |  |

**Data Preprocessing Code Screenshots**

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