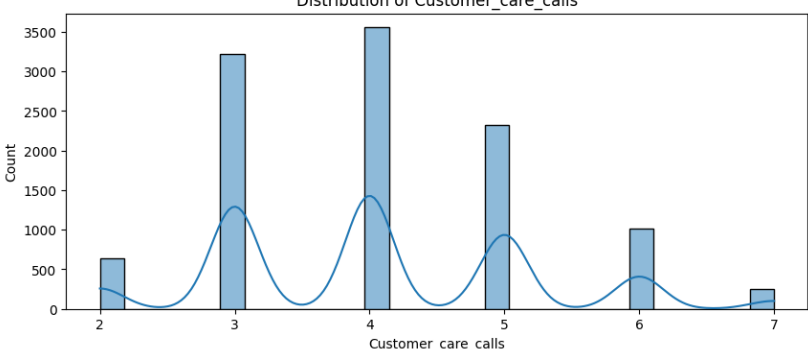


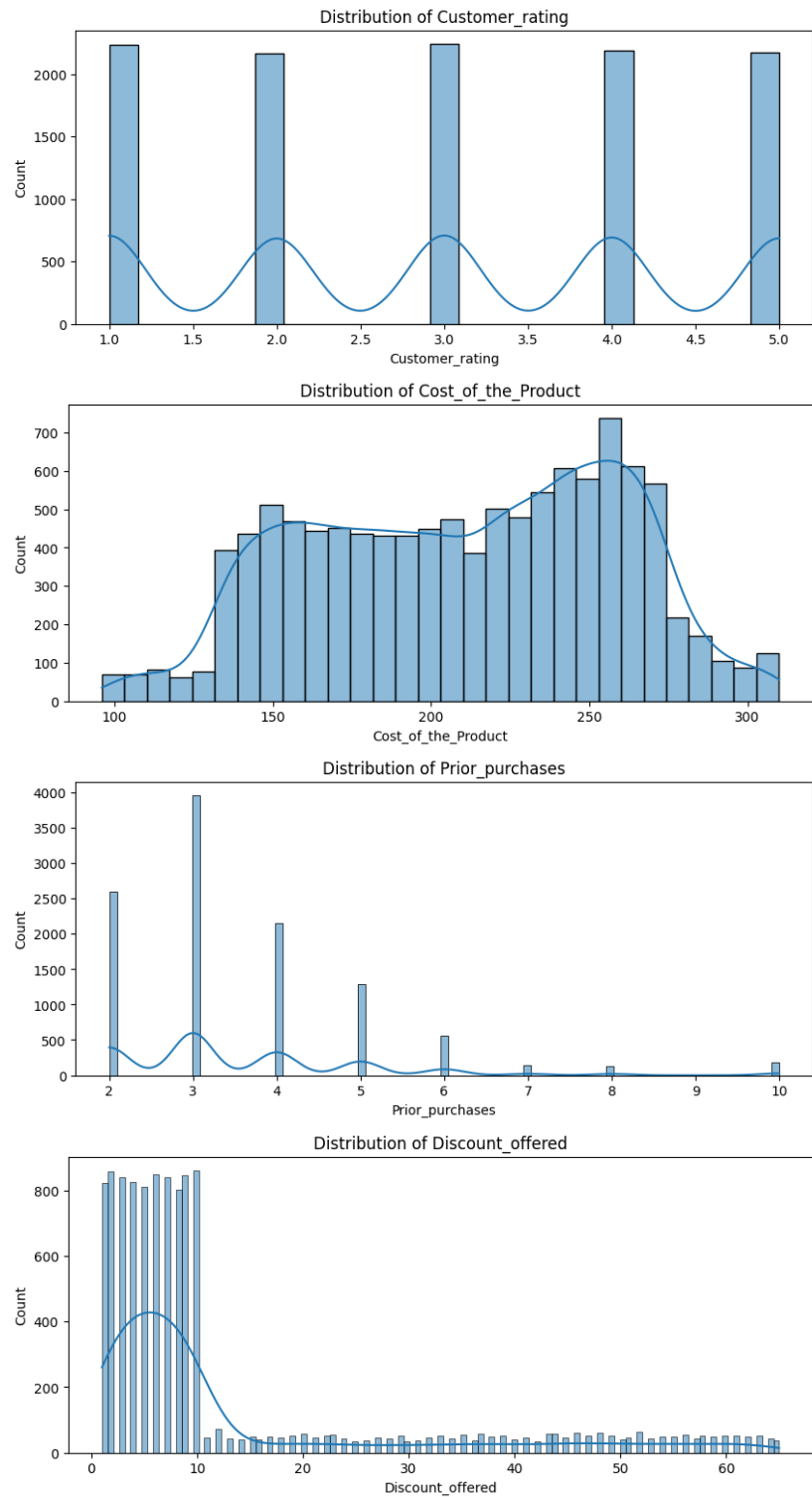
Data Collection and Preprocessing Phase

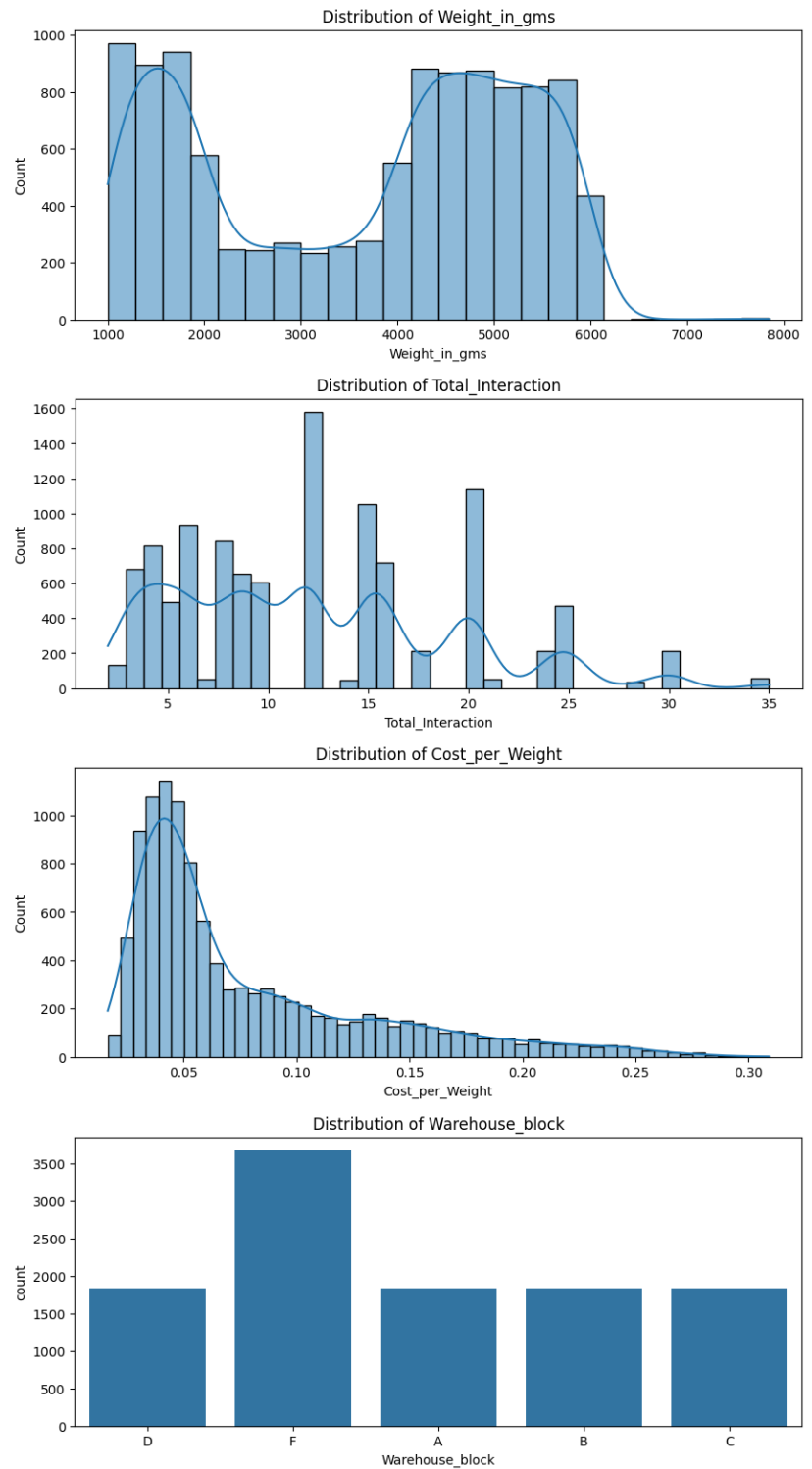
Date	20 July 2024
Team ID	SWTID1720110595
Project Title	Ecommerce Shipping Prediction Using Machine Learning
Maximum Marks	6 Marks

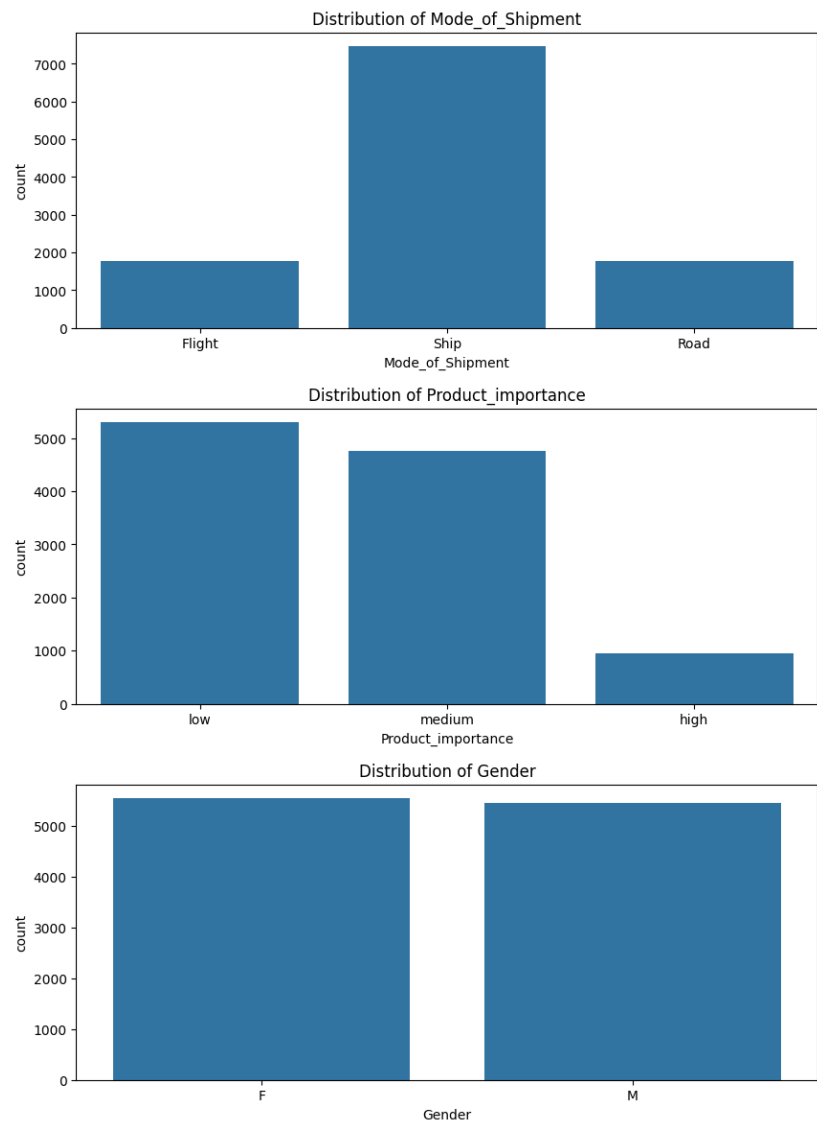
Data Exploration and Preprocessing Template

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

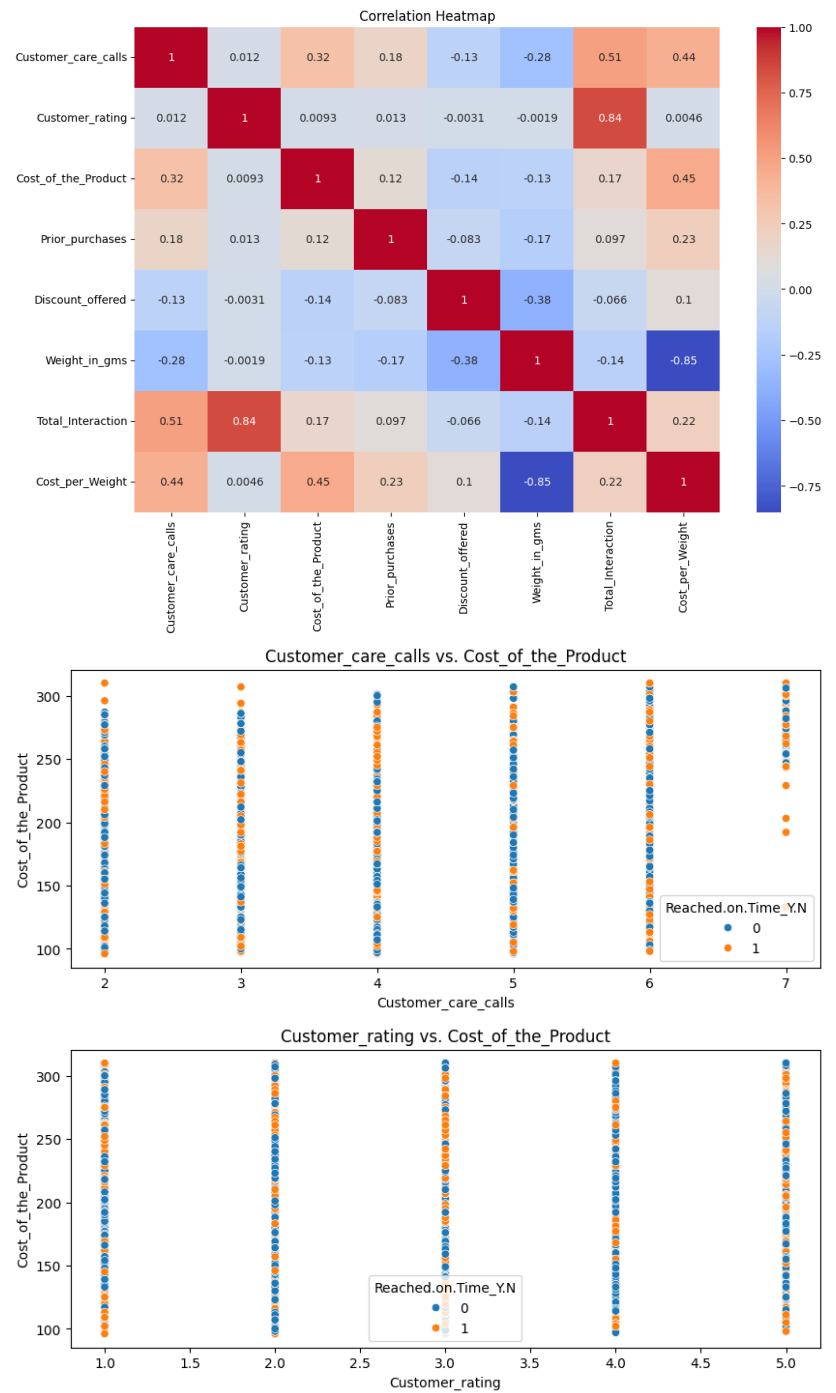
Section	Description
Data Overview	<pre> ... count ID Customer_care_calls Customer_rating Cost_of_the_Product \ mean 5500.00000 4.054459 2.990545 210.196836 std 3175.28214 1.141490 1.413603 48.063272 min 1.00000 2.000000 1.000000 96.000000 25% 2750.50000 3.000000 2.000000 169.000000 50% 5500.00000 4.000000 3.000000 214.000000 75% 8249.50000 5.000000 4.000000 251.000000 max 10999.00000 7.000000 5.000000 310.000000 Prior_purchases Discount_offered Weight_in_gms Reached.on.Time_Y.N count 10999.000000 10999.000000 10999.000000 10999.000000 mean 3.567597 13.373216 3634.016729 0.596691 std 1.522860 16.205527 1635.377251 0.490584 min 2.000000 1.000000 1001.000000 0.000000 25% 3.000000 4.000000 1839.500000 0.000000 50% 3.000000 7.000000 4149.000000 1.000000 75% 4.000000 10.000000 5050.000000 1.000000 max 10.000000 65.000000 7846.000000 1.000000 </pre>
Univariate Analysis	<p>Distribution of Customer_care_calls</p> 

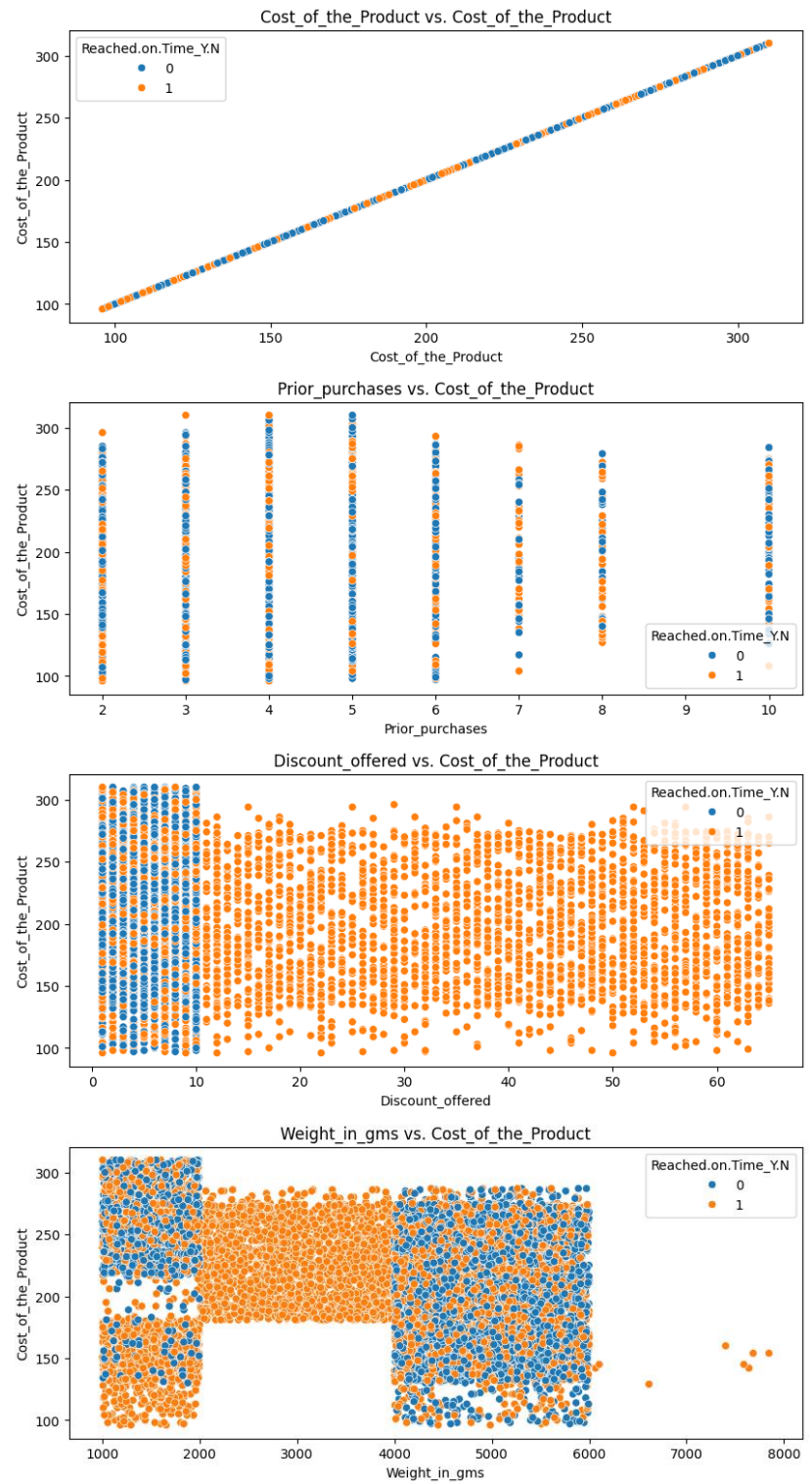


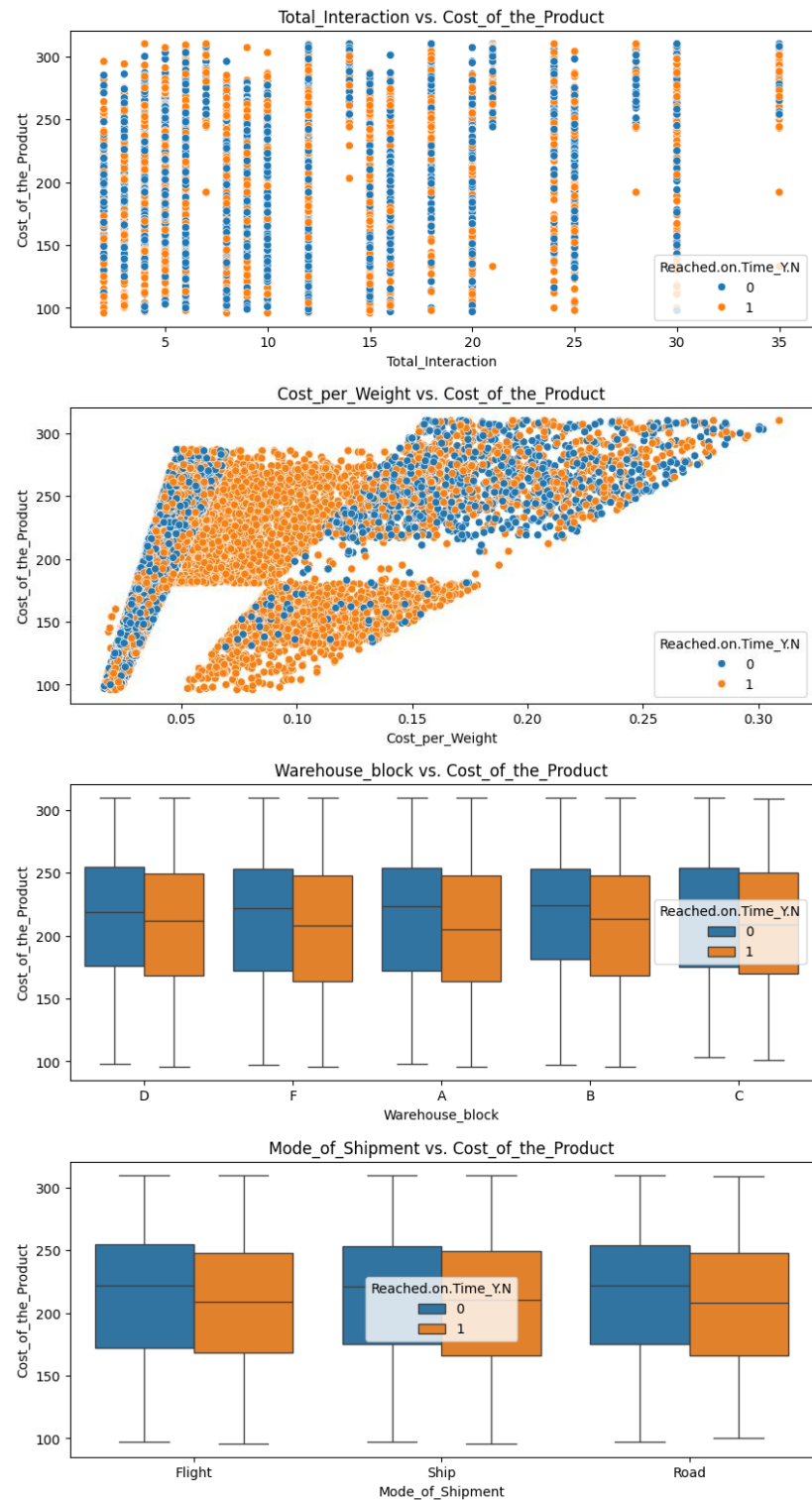


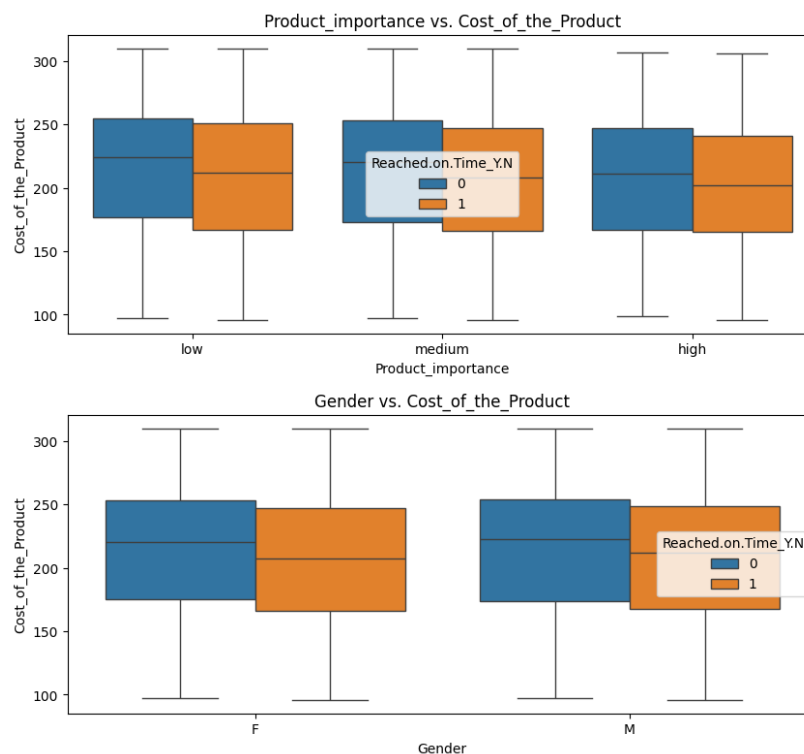


Bivariate Analysis

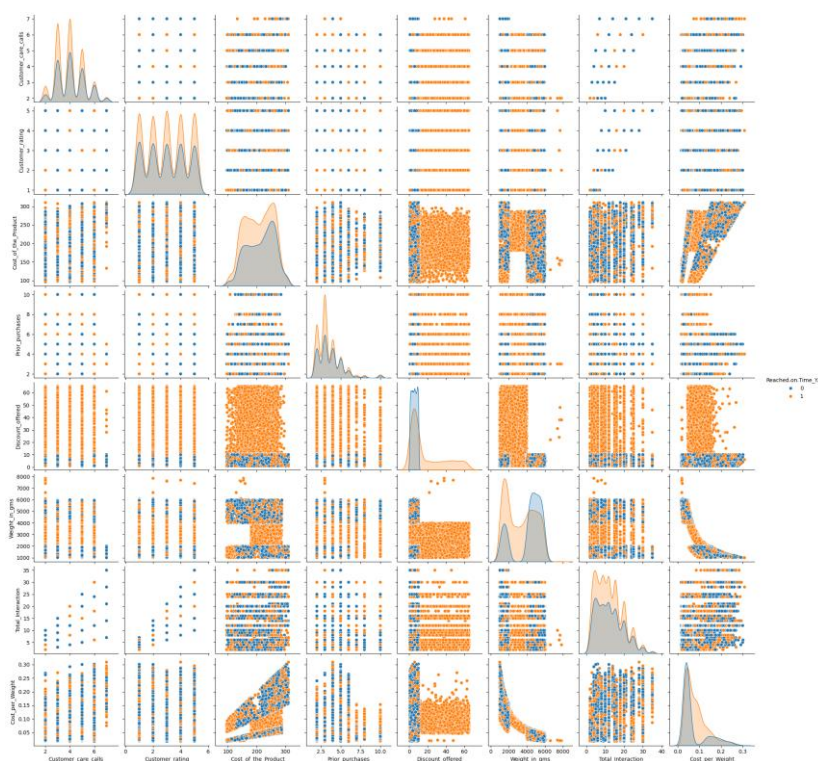


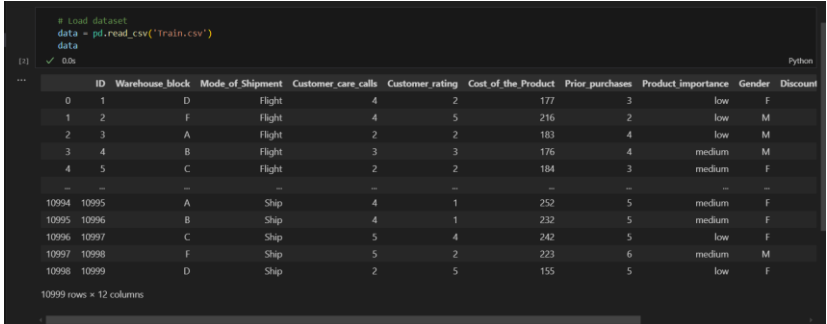






Multivariate Analysis



Outliers and Anomalies	<pre># Outliers and Anomalies z_scores = np.abs(stats.zscore(data[numerical_features])) outliers = np.where(z_scores > 3) print("Outliers detected:", outliers) for feature in numerical_features: plt.figure(figsize=(10, 4)) sns.boxplot(data[feature]) plt.title(f'Boxplot of {feature}') plt.show() data_cleaned = data[(z_scores < 3).all(axis=1)]</pre>
Data Preprocessing Code Screenshots	
Loading Data	 <p>The screenshot shows a Jupyter Notebook cell where a dataset is loaded using <code>pd.read_csv('train.csv')</code>. Below the code, a preview of the data is shown as a table with 12 columns: ID, Warehouse block, Mode of Shipment, Customer care calls, Customer rating, Cost of the Product, Prior purchases, Product importance, Gender, and Discount. The table displays rows 0 to 4, followed by an ellipsis, and then rows 10994 to 10999. The bottom of the preview indicates '10999 rows x 12 columns'.</p>
Handling Missing Data	<pre># Data preprocessing data = data.fillna(method='ffill')</pre>
Data Transformation	<pre># Define numerical and categorical features numerical_features = ['Customer_care_calls', 'Customer_rating', 'Cost_of_the_Product', 'Prior_purchases', 'Discount_offered', 'Weight_in_gms', 'Total_in'] categorical_features = ['Warehouse_block', 'Mode_of_Shipment', 'Product_importance', 'Gender'] numerical_transformer = StandardScaler() categorical_transformer = OneHotEncoder(handle_unknown='ignore', drop='first') preprocessor = ColumnTransformer(transformers=[('num', numerical_transformer, numerical_features), ('cat', categorical_transformer, categorical_features)]) # Preprocess data before applying SMOTE X_preprocessed = preprocessor.fit_transform(x)</pre>
Feature Engineering	<pre># Feature engineering data['Total_Interaction'] = data['Customer_care_calls'] * data['Customer_rating'] data['Cost_per_Weight'] = data['Cost_of_the_Product'] / data['Weight_in_gms']</pre>

Save Processed Data

```
# Serialize the preprocessor
with open('preprocessor.pkl', 'wb') as f:
    pickle.dump(preprocessor, f)

# Serialize the trained stacking pipeline
with open('model.pkl', 'wb') as f:
    pickle.dump(stacking_pipeline, f)
```