# Group 12 Project Documentation
# Sales Analytics

## Problem Statement 1: Advanced Data Cleaning

### Objective
To ensure data quality by identifying and handling outliers, addressing missing values, and correcting data type inconsistencies.

### Steps and Justifications:
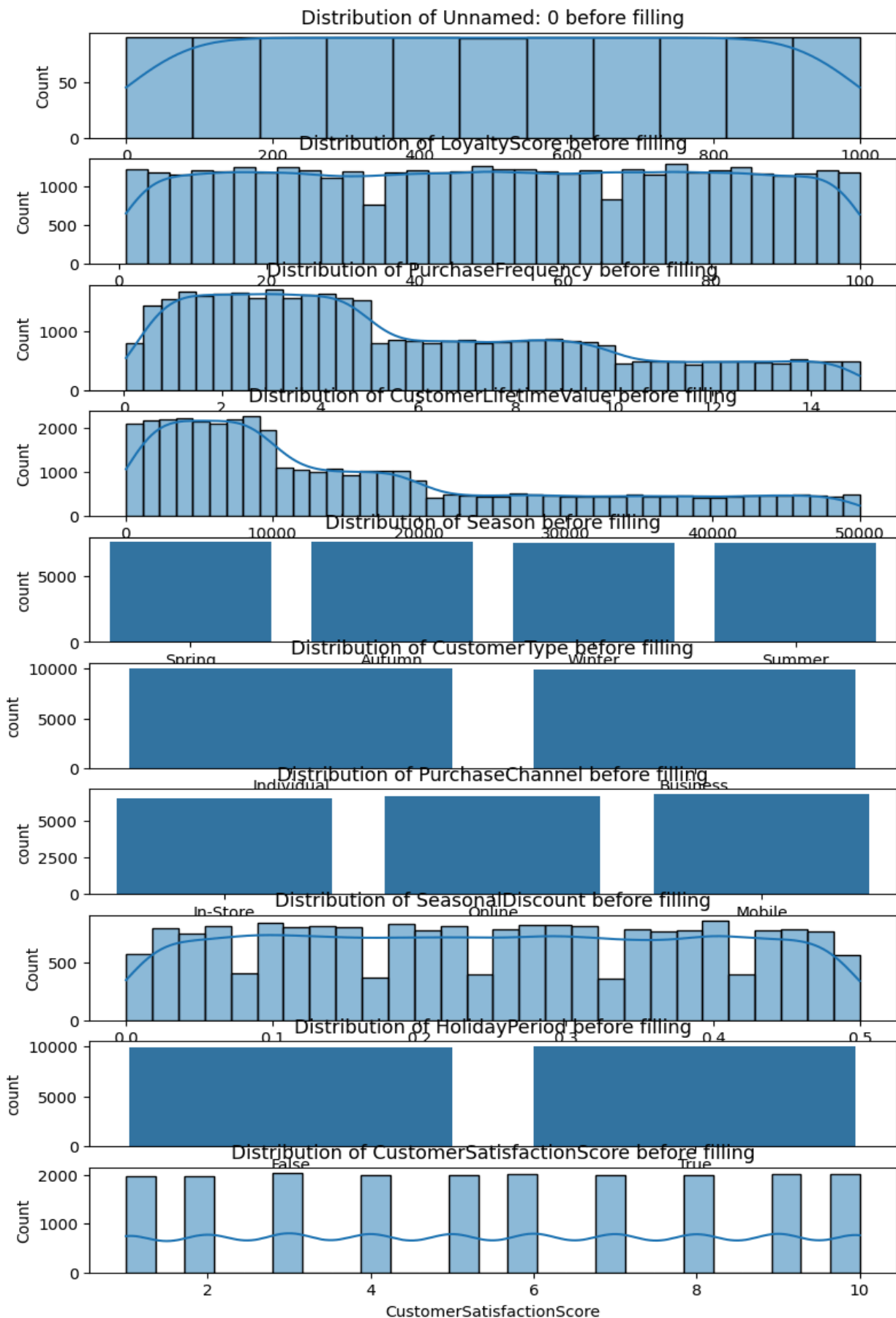
1. **Loading the Dataset**
   - Loaded the dataset using Pandas to facilitate data manipulation and cleaning.

2. **Identifying Outliers**
   - Used the Interquartile Range (IQR) method to detect outliers. This method is chosen for its robustness in identifying extreme values without being influenced by them.
   - Calculated the first quartile (Q1) and third quartile (Q3) and determined the IQR as Q3 - Q1.
   - Outliers are defined as data points that fall below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR.

**Output:**

```
Missing values before imputation:
 Unnamed: 0                      40000
CustomerID                          0
Age                                 0
Gender                              0
Location                            0
ProductCategory                     0
PurchaseDate                        0
PurchaseAmount                      0
PaymentMethod                       0
Quantity                            0
DiscountPercentage                  0
IsReturned                          0
Rating                              0
IsPromotion                         0
CustomerSegment                     0
ShippingDuration                    0
Region                              0
LoyaltyScore                     1000
PurchaseFrequency                1000
CustomerLifetimeValue            1000
Season                          11000
CustomerType                    21000
PurchaseChannel                 21000
SeasonalDiscount                21000
HolidayPeriod                   21000
CustomerSatisfactionScore       21000
dtype: int64
```
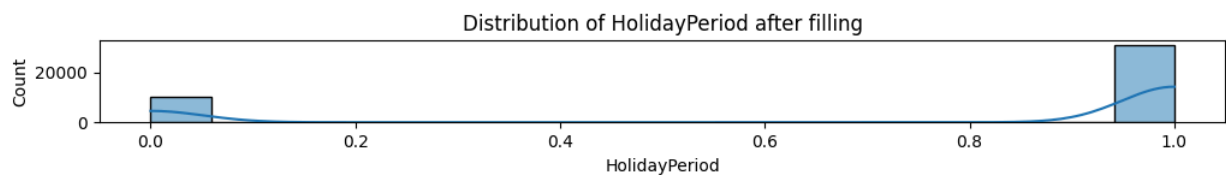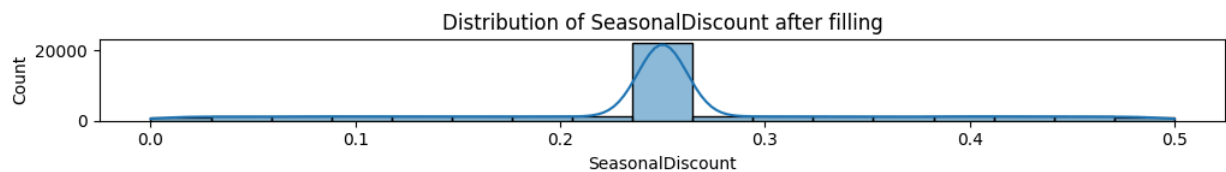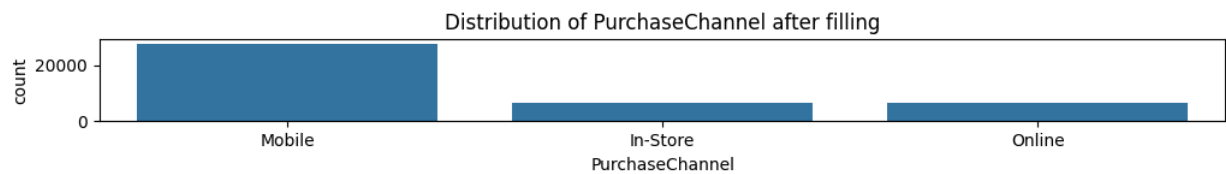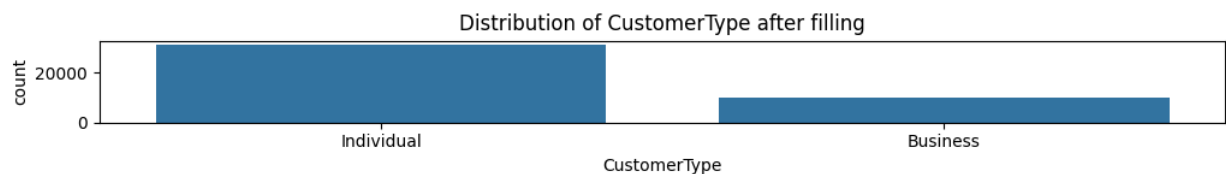
Distribution of Unnamed: 0 before filling

Distribution of LoyaltyScore before filling

Distribution of PurchaseFrequency before filling

Distribution of CustomerLifetimeValue before filling

Distribution of Season before filling

Distribution of CustomerType before filling

Distribution of PurchaseChannel before filling

Distribution of SeasonalDiscount before filling

Distribution of HolidayPeriod before filling

Distribution of CustomerSatisfactionScore before filling

3. **Handling Missing Values**
   - Applied different imputation techniques based on the nature of the data:
      - **Mean Imputation**: Used for numerical data where the mean is appropriate.
      - **Median Imputation**: Chosen for skewed numerical data to avoid mean distortion.
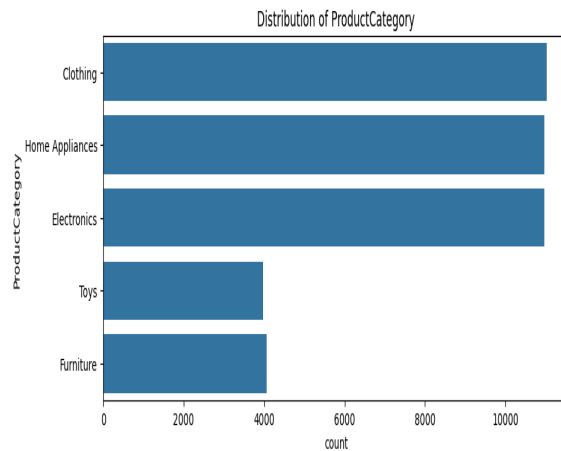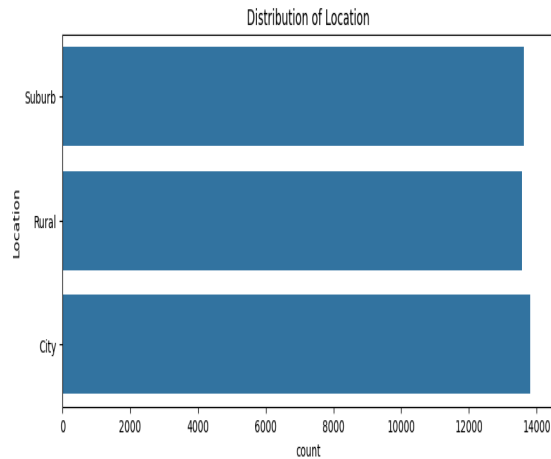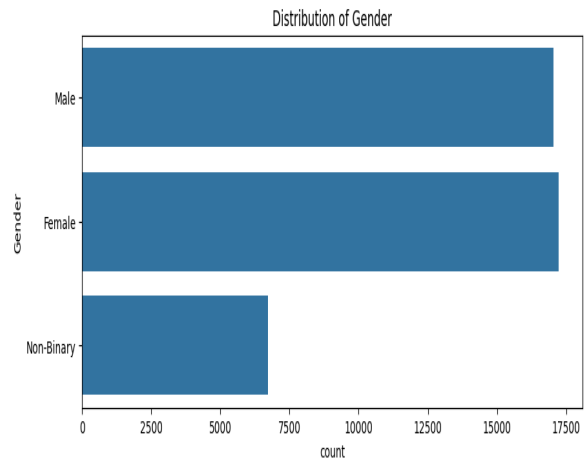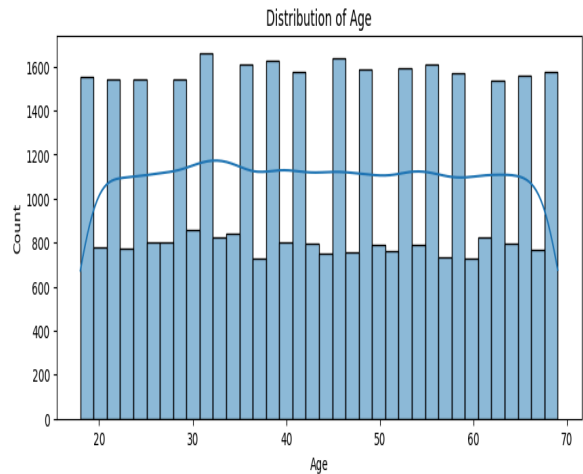      - **Mode Imputation**: Used for categorical data to fill in the most frequent value.

**Output:**
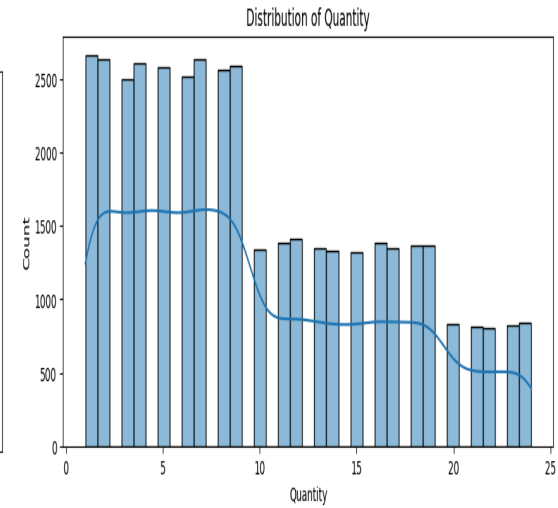
```
Missing values after imputation:
 Unnamed: 0                        0
CustomerID                         0
Age                                0
Gender                             0
Location                           0
ProductCategory                    0
PurchaseDate                       0
PurchaseAmount                     0
PaymentMethod                      0
Quantity                           0
DiscountPercentage                 0
IsReturned                         0
Rating                             0
IsPromotion                        0
CustomerSegment                    0
ShippingDuration                   0
Region                             0
LoyaltyScore                       0
PurchaseFrequency                  0
CustomerLifetimeValue              0
Season                             0
CustomerType                       0
PurchaseChannel                    0
SeasonalDiscount                   0
HolidayPeriod                      0
CustomerSatisfactionScore          0
dtype: int64
```
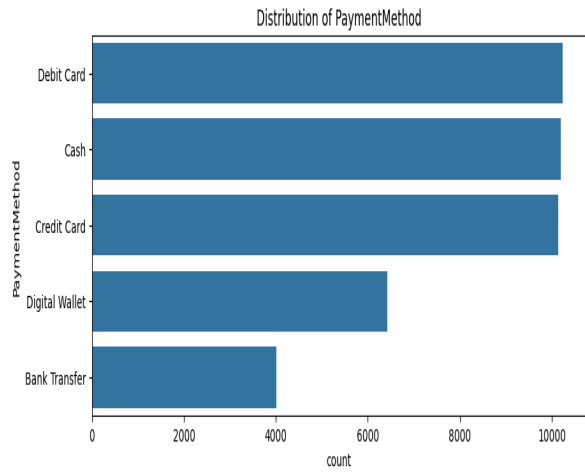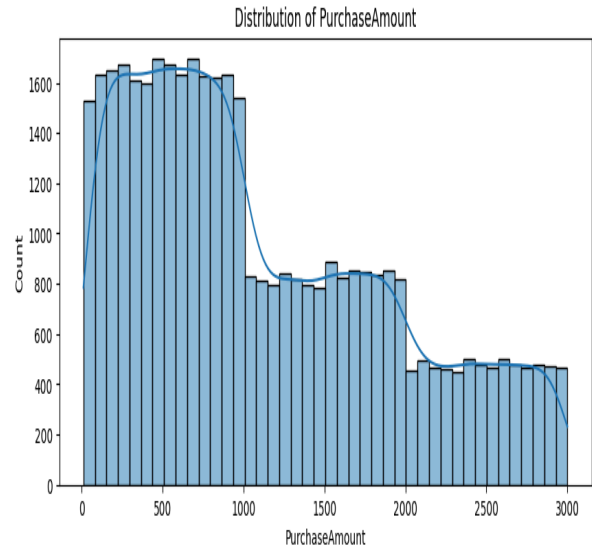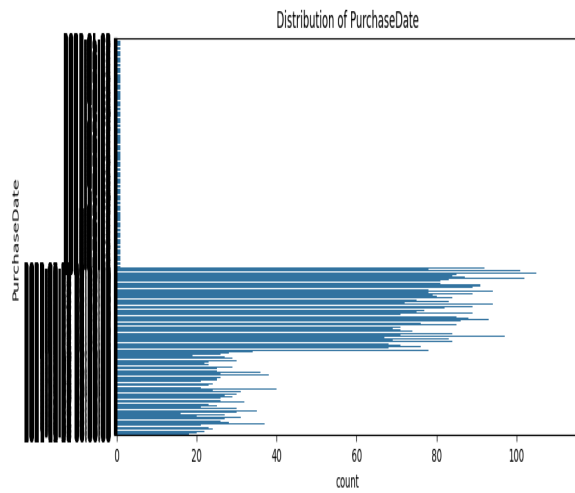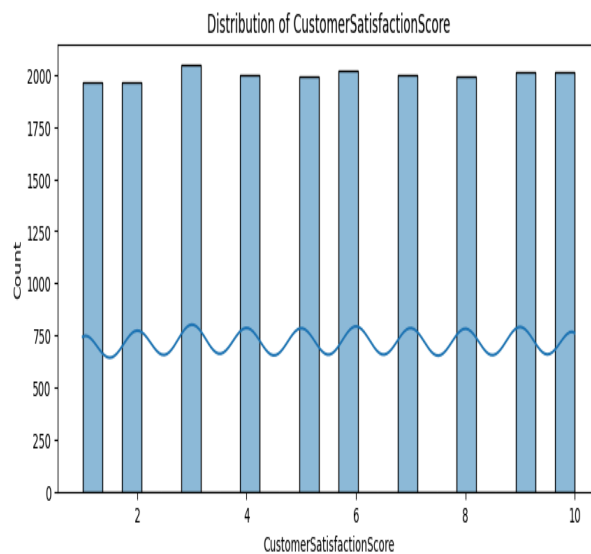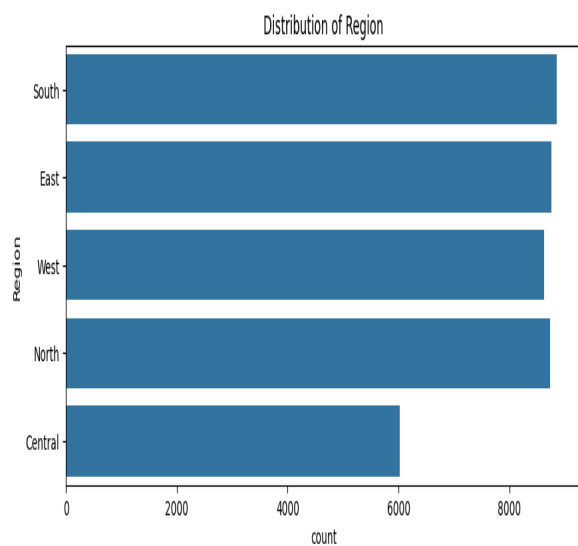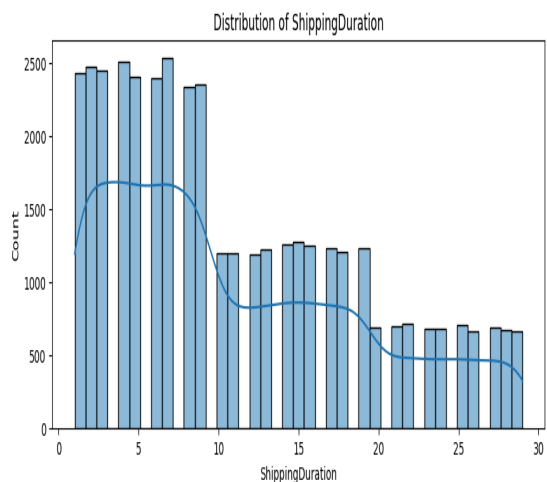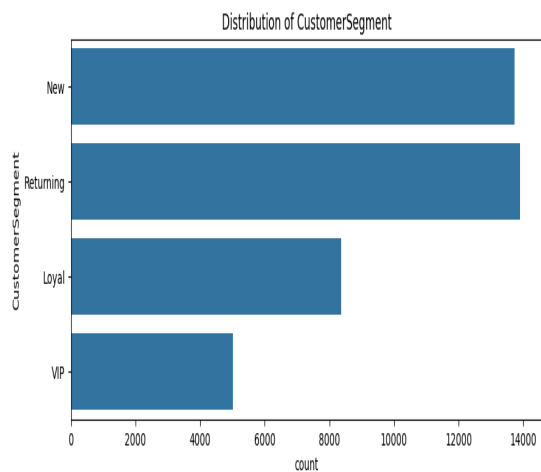
## Distribution of Unnamed: 0 after filling

## Distribution of LoyaltyScore after filling

## Distribution of PurchaseFrequency after filling

## Distribution of CustomerLifetimeValue after filling

## Distribution of Season after filling

## Distribution of CustomerType after filling

## Distribution of PurchaseChannel after filling

## Distribution of SeasonalDiscount after filling

## Distribution of HolidayPeriod after filling

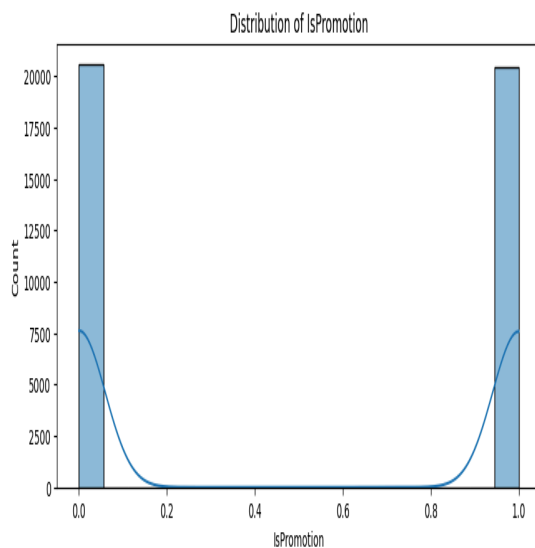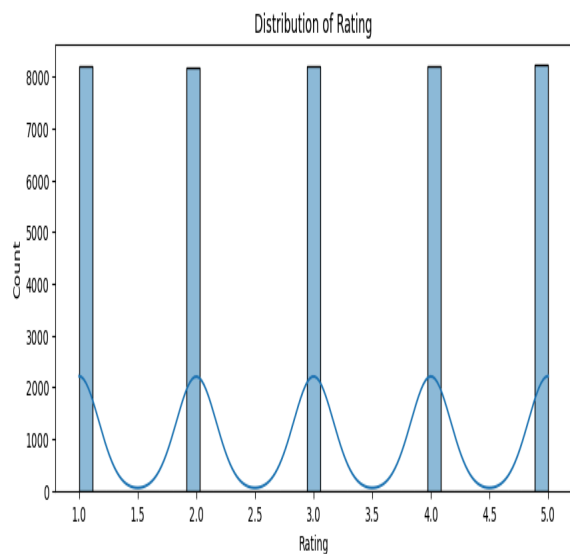## Distribution of CustomerSatisfactionScore after filling

## 4. **Correcting Data Type Inconsistencies**
- Ensured that each column had a consistent data type:
  - Converted date columns to datetime objects.
  - Changed numerical columns stored as strings to appropriate numerical types.
- Verified and corrected any misclassified data types.

## Distribution of PurchaseDate

## Distribution of PurchaseAmount

## Distribution of PaymentMethod

## Distribution of Quantity

## Distribution of DiscountPercentage

## Distribution of IsReturned

Distribution of Rating

Distribution of IsPromotion

Distribution of CustomerSegment

Distribution of ShippingDuration

Distribution of Region

Distribution of CustomerSatisfactionScore

# Problem Statement 2: Data Augmentation

## Objective
To enhance the dataset by generating additional samples while maintaining the statistical properties of the original data.

## Steps and Justifications

### 1. Analyzing Data Distribution
   - Conducted an in-depth analysis of the existing data's distribution to understand its characteristics.
   - Identified key statistical properties such as mean, variance, skewness, and kurtosis.

### 2. Data Augmentation Techniques
   - Applied bootstrapping to create additional samples. Bootstrapping is a resampling technique that generates new data points by sampling with replacement from the existing data.
   - Ensured that the augmented data followed the original dataset's statistical distribution.
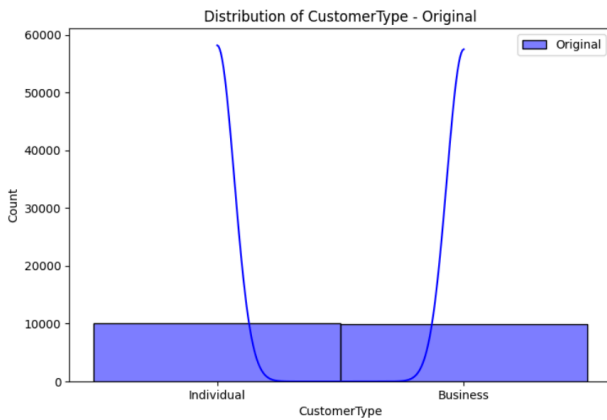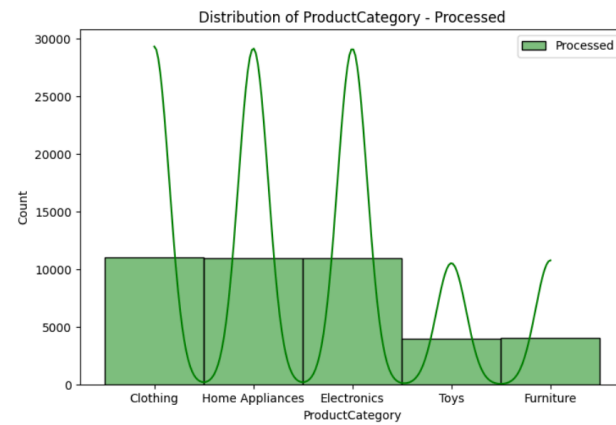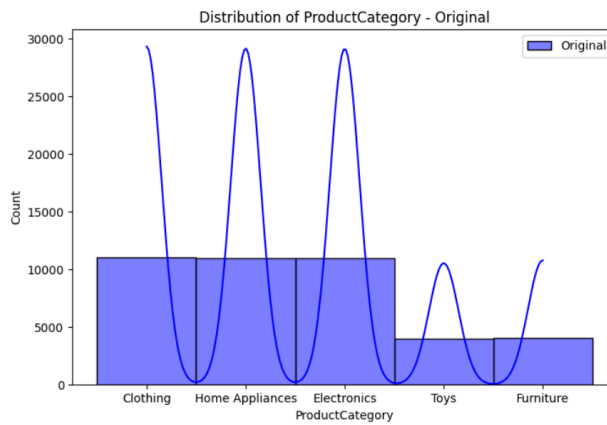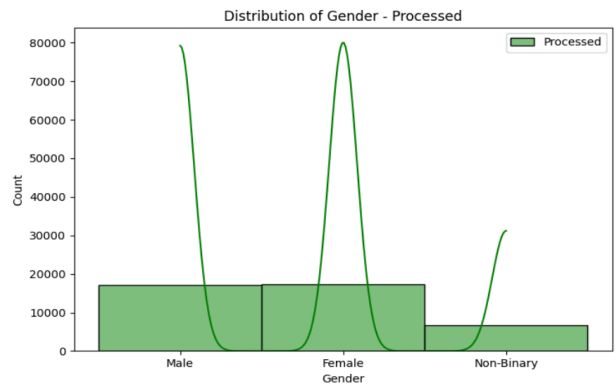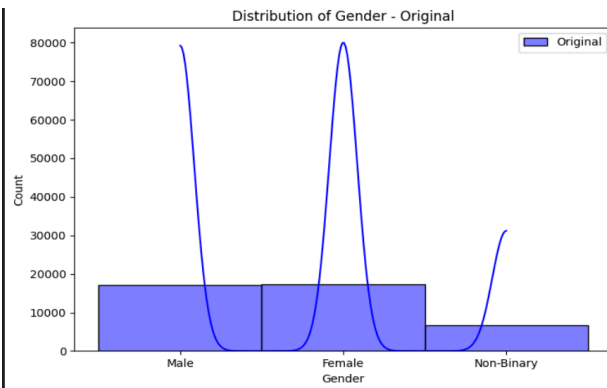
### 3. Integrating Augmented Data
   - Merged the newly generated samples with the original dataset to create an expanded dataset.
   - Maintained the integrity of the original data while ensuring the augmented data enhanced the dataset.

### 4. Validation
   - Performed rigorous validation to ensure the augmented data met quality standards.
   - Compared statistical properties of the augmented dataset with the original to confirm consistency.

Distribution of Gender - Original

Distribution of Gender - Processed

Distribution of ProductCategory - Original

Distribution of ProductCategory - Processed

Distribution of CustomerType - Original

Distribution of CustomerType - Processed

Distribution of CustomerSatisfactionScore - Original

Distribution of CustomerSatisfactionScore - Processed

# Problem Statement 3: Real-time Data Ingestion

## Objective

To set up a real-time data ingestion pipeline using Apache Kafka and ensure optimized data flow into SQL databases.

## Steps and Justifications

1. **Setting Up Apache Kafka**
   - Configured an Apache Kafka environment to manage real-time data streams.
   - Established Kafka brokers, topics, and partitions to facilitate efficient data flow.

2. **Creating Kafka Producers**
   - Developed Kafka producers to simulate real-time data streams.
   - Configured producers to send data to the appropriate Kafka topics.

3. **Developing Kafka Consumers**
   - Used Python to create Kafka consumers that ingest data from Kafka topics into SQL databases.
   - Ensured consumers were optimized for high throughput and low latency to handle real-time data efficiently.

4. **Optimizing Data Ingestion**
   - Implemented strategies to minimize latency and maximize throughput.
   - Used batching and compression techniques to enhance performance.

# Problem Statement 4: Storage Optimization

## Objective
To evaluate and optimize storage formats for better efficiency and performance.

## Steps and Justifications

1. **Evaluating Columnar Storage Formats**
   - Assessed columnar storage formats such as Parquet and ORC for their storage efficiency and performance.
   - Compared these formats with traditional row-based storage.

2. **Converting Dataset**
   - Converted the dataset to Parquet and ORC formats.
   - Evaluated the storage space required and the query performance for each format.

3. **Comparison and Analysis**
   - Conducted a detailed comparison of storage efficiency and query performance between columnar and row-based storage.
   - Analyzed metrics such as storage size, read/write speeds, and query response times.

This documentation provides a clear and comprehensive overview of each problem statement, the methods used, and the justifications for these methods, ensuring a thorough understanding of the tasks and their execution.

# Agile Model - Jira (Progress Dashboard)