

---

# Supermarket Sales Analysis

## Research Project Report

---

*By*  
**Ibrahim Mahmood**  
ID No. 2021748



Ghulam Ishaq Khan Institute of Engineering Sciences and Technology

February 2024

# *Abstract*

This project navigates the supermarket sales spectrum using Python and data mining techniques. Key aspects include association rule mining, classification, regression, and outlier detection. Beyond conventional analyses, it explores temporal patterns in customer behavior, utilizing sequential pattern mining and tracking changes in item sets over time. The investigation extends to clustering algorithms for grouping customers based on purchase history. The project offers concise insights into customer behavior regarding the Supermarket Sales dataset.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Contents</b>	<b>ii</b>
<b>1 Exploratory Data Analysis (EDA)</b>	<b>1</b>
<b>2 Data Preprocessing</b>	<b>6</b>
2.1 Data Cleaning . . . . .	6
2.2 DateTime . . . . .	7
2.3 Encoding . . . . .	7
2.3.1 Label Encoding . . . . .	8
2.3.2 One Hot Encoding . . . . .	8
2.4 Standardization . . . . .	9
2.4.1 Why was Standardization required here? . . . . .	9
2.4.1.1 Visual Insepection . . . . .	9
2.4.1.2 Statistical Test (Shapiro-Wilk Test) . . . . .	10
2.4.2 Which Features require Standardization and why? . . . . .	10
2.4.3 Choosing between Normalization and Standardization . . . . .	11
2.4.4 Using StandardScaler . . . . .	11
2.4.5 How did this step contribute to effective data mining? . . . . .	11
2.5 Feature Selection . . . . .	12
2.5.1 Choosing K . . . . .	12
2.6 Splitting Data . . . . .	13
2.6.1 Importance of splitting . . . . .	13
2.6.2 Applying Train-Test Split . . . . .	14
<b>3 Visualization</b>	<b>15</b>
3.1 Variable Distributions and Patterns . . . . .	15
3.1.1 Numerical Variable Distributions: . . . . .	15
3.1.2 Categorical Variable Distributions . . . . .	15
3.1.3 Bivariate Relationships . . . . .	16
3.1.4 Time Trends . . . . .	16
3.2 Why is Visualization required here? . . . . .	16
3.2.1 Unveiling Hidden Patterns and Relationships . . . . .	16
3.2.2 Enhancing Understanding and Communication . . . . .	16

3.2.3	Guiding Further Analysis and Decision-Making:	17
3.3	Visualizations	17
<b>4</b>	<b>Association Rule Mining</b>	<b>22</b>
4.1	Data preparation	22
4.1.1	Creating a list of lists	22
4.1.2	Grouping products by count	22
4.1.3	Encoding the DataFrame using Transaction Encoder	23
4.2	Frequent Itemsets	23
4.2.1	Calculation	23
4.2.2	Choosing Minimum Support	23
4.2.3	The frequent itemset produced	24
4.3	Conducting Association Rule Mining	25
4.3.1	Applying association rule mining:	25
4.3.2	Setting the confidence	25
4.3.2.1	What is confidence?	25
4.3.2.2	Choosing Minimum Confidence	25
4.3.3	The obtained association rules	26
4.3.3.1	Lift	26
4.3.3.2	Leverage	26
4.3.3.3	The obtained association rules:	26
<b>5</b>	<b>Comparative Analysis with FP Growth Algorithm</b>	<b>27</b>
5.1	Applying FP Growth Algorithm:	27
5.2	Comparison of Apriori and FP Growth Algorithms:	28
5.2.1	Functionality:	28
5.2.2	Time Complexity:	28
5.2.3	Memory Usage:	29
5.2.4	Summarized Differences:	30
<b>6</b>	<b>Tracking Patterns and Customer Behavior Analysis</b>	<b>31</b>
6.1	Creating Time-Series Data:	31
6.2	Visualizations:	32
6.2.1	Seasonal Decompose Graphs:	33
6.2.2	Predictive Modeling using ARIMA:	34
6.2.2.1	How is it beneficial here?	34
6.2.2.2	Applying the model:	34
6.2.2.3	Prediction and analysis	35
6.2.3	Feature Engineering to calculate and visualize Rolling Mean:	35
6.2.3.1	The visualization:	36
6.2.3.2	How the graph works:	36
6.2.3.3	Analyzing the visualization	36
<b>7</b>	<b>Sequential Pattern Mining</b>	<b>37</b>
7.1	Relevance to this project	37
7.2	How is this different from using Apriori or FP Growth?	37
7.3	PrefixSpan:	38
7.3.1	How does it work?	38

7.3.2	Applying it via code:	38
7.3.3	Result and Analysis:	38
7.4	GSP: Generalized Sequential Patterns:	39
7.4.1	How does it work?	39
7.4.2	Applying it via code:	39
7.4.3	Result and Analysis:	39
7.5	Eclat Algorithm (Equivalence Class Transformation):	40
7.5.1	How does it work?	40
7.5.2	Applying it via code:	40
7.5.3	Result and Analysis:	40
7.6	Summarization of Key Differences	41
<b>8</b>	<b>Clustering</b>	<b>42</b>
8.1	Importance of feature selection	42
8.2	Selection of features and reasoning	42
8.3	Why a standardized dataset is required:	42
8.4	Applying Clustering	43
8.4.1	Kmeans	43
8.4.2	K-Median	43
8.4.3	ARI and NMI	44
8.4.4	DBSCAN:	44
<b>9</b>	<b>Classification</b>	<b>46</b>
9.1	Choosing the Target Variable:	46
9.2	Binning:	47
9.2.1	The concept of binning:	47
9.2.2	Choosing an ideal number of bins:	47
9.2.2.1	Square Root Rule:	47
9.2.2.2	Sturges' Formula:	48
9.2.2.3	Scott's Rule:	48
9.2.2.4	Freedman-Diaconis Rule:	48
9.2.2.5	Doane's Formula:	49
9.2.2.6	Results	49
9.2.3	Deciding on the bin edges:	50
9.2.4	Applying binning:	51
9.2.5	Analyzing 'Total' after binning:	52
9.3	Splitting data	53
9.4	Model Training and Testing:	53
<b>10</b>	<b>Regression</b>	<b>55</b>
10.1	Evaluation Metrics:	55
10.1.1	Why do we use different metrics for regression?	56
10.2	Applying Regression:	56
10.3	Analysis	56
10.4	Conclusion	57
<b>11</b>	<b>Outlier Detection and Statistical Validation</b>	<b>58</b>
11.1	Z-Score:	59

---

11.1.1	How does Z-Score work?	59
11.1.2	Applying Z-Score in code:	60
11.1.3	Result and Analysis:	60
11.2	Interquartile Range:	61
11.3	Isolation Forest:	62
11.4	Statistical Validation:	63
11.4.1	ANOVA:	63
11.4.2	Hypothesis Testing:	63
11.4.2.1	Two-Sample t-Test:	63
 <b>Bibliography</b>		 <b>65</b>
<b>References</b>		<b>65</b>

## Chapter 1

# Exploratory Data Analysis (EDA)

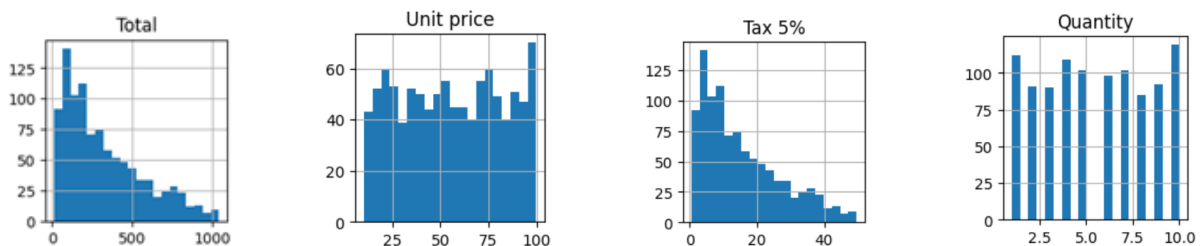
EDA: Making use of various Pandas functions, I conducted some initial analysis of our dataset

Invoice ID	750-67-8428
Branch	A
City	Yangon
Customer type	Member
Gender	Female
Product line	Health and beauty
Unit price	74.69
Quantity	7
Tax 5%	26.1415
Total	548.9715
Payment	Ewallet
Date	1/5/2019
Time	13:08
cogs	522.83
gross income	26.1415
Rating	9.1

Using Pandas' describe function, I was able to gain insight into the numerical features

Statistics	Unit price	Quantity	Tax 5%	Total	cogs	gross margin	gross income	Rating
count	1000.0000	1000.0000	1000.0000	1000.0000	1000.0000	1000.0000	1000.0000	1000.0000
mean	55.672130	5.510000	15.379369	322.966749	307.58738	4.761905	15.379369	6.97270
std	26.494628	2.923431	11.708825	245.885335	234.17651	0.000000	11.708825	1.71858
min	10.080000	1.000000	0.508500	10.678500	10.17000	4.761905	0.508500	4.00000
25%	32.875000	3.000000	5.924875	124.422375	118.49750	4.761905	5.924875	5.50000
50%	55.230000	5.000000	12.088000	253.848000	241.76000	4.761905	12.088000	7.00000
75%	77.935000	8.000000	22.445250	471.350250	448.90500	4.761905	22.445250	8.50000
max	99.960000	10.000000	49.650000	1042.650000	993.00000	4.761905	49.650000	10.00000

On average, customers spend approximately \$322.97 per transaction, with a standard deviation of \$245.89, indicating significant variability. The constant gross margin percentage of 4.76% suggests a consistent profit margin, while customer ratings, averaging 6.97 with a standard deviation of 1.72, reflect overall satisfaction. Distribution of numerical columns The visualizations below are to help us understand the distribution of some numerical columns. The x-axis of each histogram represents the range of values, and the y-axis represents the frequency (count) of data points within each range. The number of bins (20 in this case) determines the granularity of the distribution. The feature 'Unit Price' seems to be quite heavily and evenly distributed. In contrast, 'Quantity' sparsely populates its chart. The 'Total' and 'Tax' columns both show a central tendency towards higher values.

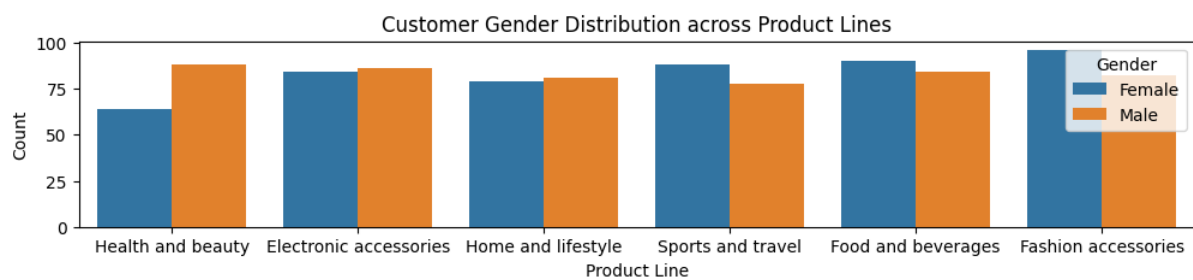




Initial correlation matrix for the dataset:

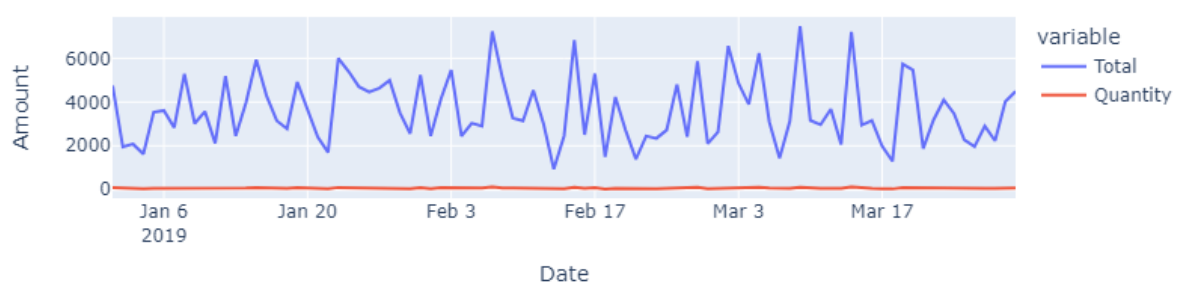
	Unit price	Quantity	Tax 5%	Total	cogs	gross margin percentage	gross income	Rating
Unit price	1.000000	0.010778	0.633962	0.633962	0.633962	NaN	0.633962	-0.008778
Quantity	0.010778	1.000000	0.705510	0.705510	0.705510	NaN	0.705510	-0.015815
Tax 5%	0.633962	0.705510	1.000000	1.000000	1.000000	NaN	1.000000	-0.036442
Total	0.633962	0.705510	1.000000	1.000000	1.000000	NaN	1.000000	-0.036442
cogs	0.633962	0.705510	1.000000	1.000000	1.000000	NaN	1.000000	-0.036442
gross margin percentage	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
gross income	0.633962	0.705510	1.000000	1.000000	1.000000	NaN	1.000000	-0.036442
Rating	-0.008778	-0.015815	-0.036442	-0.036442	-0.036442	NaN	-0.036442	1.000000

We can observe a strong positive correlation between Unit Price, Quantity, Tax 5%, Total, COGS, and Gross Income, suggesting that increases in one tend to correspond with increases in the others. However, the Gross Margin Percentage lacks variability in the given dataset, resulting in undefined correlations with other variables. Rating exhibits a very weak negative correlation with Unit Price, Quantity, Tax 5%, Total, COGS, and Gross Income, implying that changes in Rating are not strongly associated with changes in these variables.

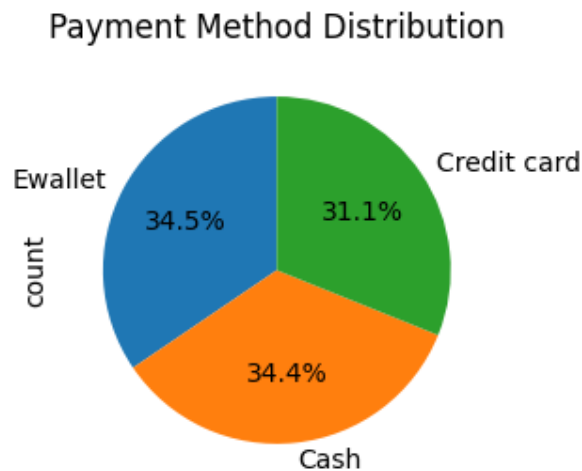


This helps understand how the purchase behavior varies between male and female customers for each product line. For example, it can reveal whether certain product lines are more popular among one gender, providing insights into potential marketing strategies or product positioning. If there's a significant gender bias in product preferences, the business can tailor its offerings accordingly.

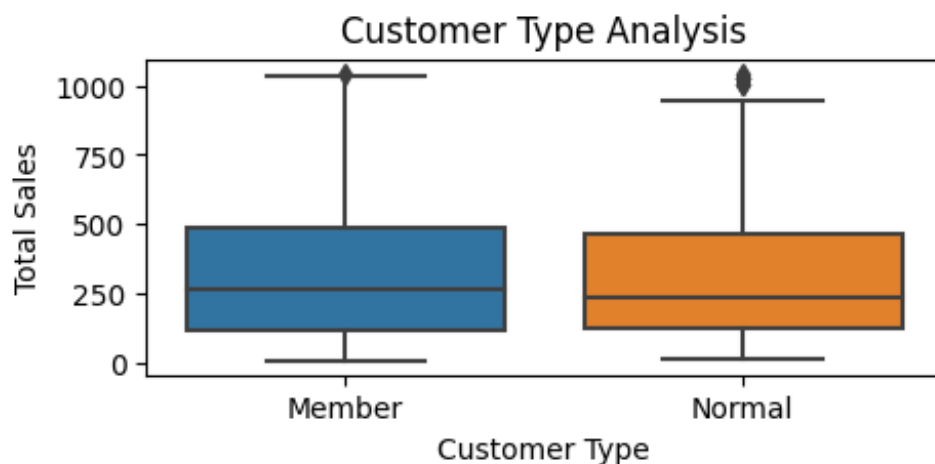
### Time-based Analysis



The line plot shows the overall trend of total sales and quantity sold over the given time. The sharp dips and spikes indicate the variance of the 'Total' feature. And the almost-straight, red line indicates, with its stagnancy, a low variation in the feature 'Quantity' throughout the months January to March.

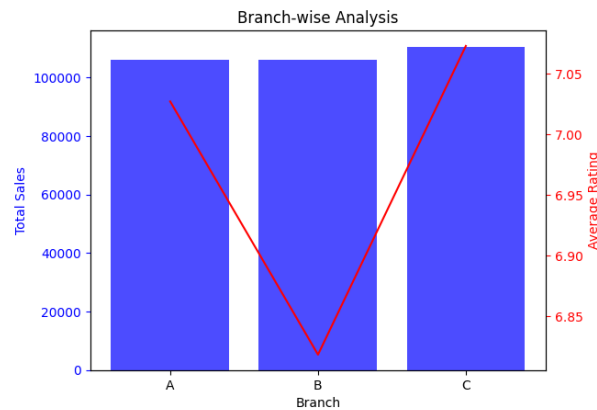


This Pie Chart visualization provides insights into the preferred payment methods of customers. Understanding the dominant payment method is crucial for optimizing payment processing systems and ensuring that the most convenient options are available. It also indicates trends in customer preferences and potentially influences decisions related to partnerships with specific payment providers. As can be seen by this chart, all 3 payment modes are near-to-identical and hence preferred almost equally by customers.



This provides a visual summary of the distribution of total sales for each customer type. There seem to be few outliers beyond the whiskers, indicating high sales for a few examples. Outliers for non-members could indicate specific cases where non-members made exceptionally high purchases.

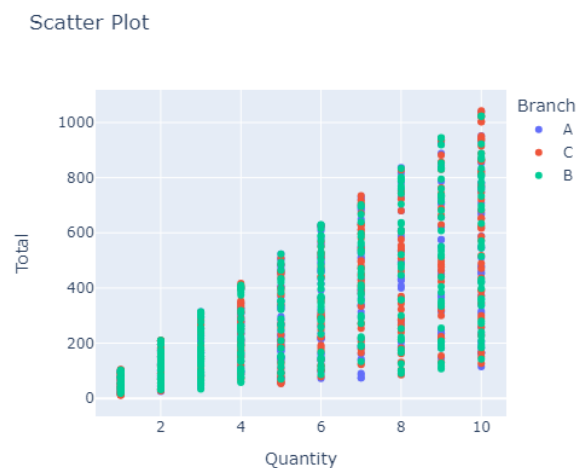
Exploring the performance of each branch based on total sales and average ratings:



The blue bars represent the total sales for each branch. The branches here seem to be almost identical in terms of Sales, with Branch C leading the other by just a bit. The red line represents the average rating for each branch. This helps evaluate the average customer satisfaction (rating) for each branch. An interesting thing to note is the steep descent of Average Rating for Branch B. This indicates that Branch A has the highest average rating, followed by Branch C. Relatively, Branch B has quite a low average rating. Grouping by Branch:

```
branch_analysis = df.groupby('Branch').agg({'Total': 'sum', 'Rating': 'mean'}).reset_index()
```

The code groups the dataset by the 'Branch' column, calculating the total sales and average rating for each branch.



This scatter plot represents the relationship between the 'Quantity' and 'Total' columns, with each point corresponding to a specific invoice. The color of the points is differentiated based on the 'Branch' column, and additional information about the 'Product line' is available when hovering over each point in the ipynb file. This interactive plot allows us to explore the relationship between the quantity and total sales.

## Chapter 2

# Data Preprocessing

### 2.1 Data Cleaning

I made use of relevant function in Pandas to check for missing values, duplicates, and outliers. However, the data was quite clean already, and barely required any initial cleaning. The table below shows no null values in any of the columns:

Column Name	Total Null Values
Invoice ID	0
Branch	0
City	0
Customer type	0
Gender	0
Product line	0
Unit price	0
Quantity	0
Tax 5%	0
Total	0
Date	0
Time	0
Payment	0
COGS	0
Gross margin percentage	0
Gross Income	0
Rating	0

TABLE 2.1: There were no duplicates or outliers

## 2.2 DateTime

The main task here was to make use of the existing Date and Time features and extract utilizable features. Other than this, I dropped a few irrelevant columns in the preprocessing phase. After doing so, I was able to get the following new features:

Date	Time
1/5/2019	13:08
3/8/2019	10:29
3/3/2019	13:23
1/27/2019	20:33
2/8/2019	10:37

TABLE 2.2: Original Features

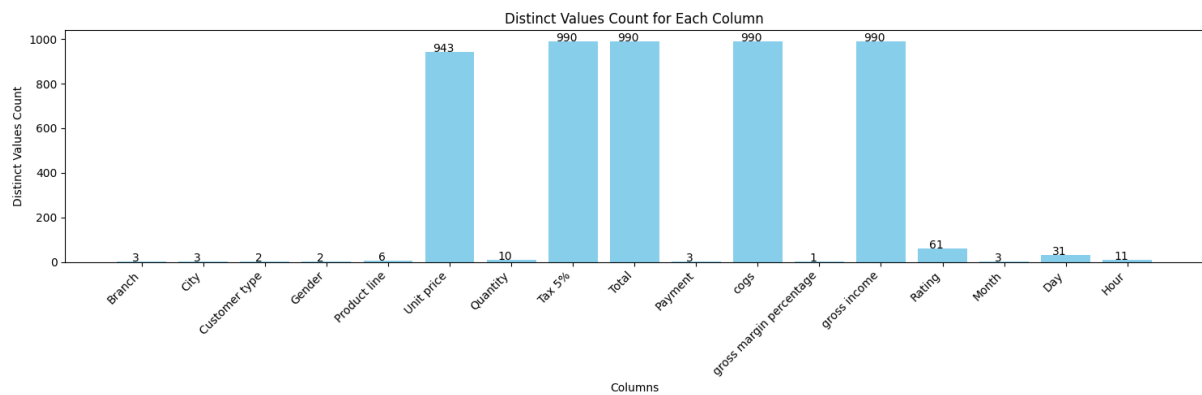
Date	Time	Duration
1	3	13
3	8	10
3	3	13
1	27	20
2	8	10

TABLE 2.3: New Features

## 2.3 Encoding

The main task here was to make use of the existing Date and Time features and extract utilizable features. Other than this, I dropped a few irrelevant columns in the preprocessing phase. After doing so, I was able to get the following new features:

Column Name	Total Null Values
Branch	3
City	3
Customer type	2
Gender	2
Product line	6
Unit price	943
Quantity	10
Tax 5%	990
Total	990
Payment	3
COGS	990
Gross margin percentage	1
Gross Income	990
Rating	61



From this we can identify the feature 'gross margin percentage' is not of any use, since it has only 1 distinct value. Which means it cannot have any correlation with other features, hence, it will not hold any importance in the variation of any target variable that we may choose when training a model.

	Day	Hour	Month	Quantity	Rating	Tax 5%	Total	Unit price	COGS	Gross income
Total	0.002515	0.00277	0.022301	0.70551	0.036442	1.0	1.0	0.633962	1.0	1.0

### 2.3.1 Label Encoding

suitable for ordinal variables, where there is an inherent order or ranking among the categories. For example, in our dataset, the days are ordinal because they have a clear and meaningful sequence. Each day follows the previous one in a systematic and sequential manner. Label encoding assigns integers based on this order, with the assumption that there is a linear relationship between the categories. This simplifies the representation of this variable while preserving the ordinal information. Hence, I'm using Label Encoding for the columns 'Day', 'Month' and 'Hour'.

### 2.3.2 One Hot Encoding

suitable for nominal variables, where there is no inherent order among the categories. Nominal variables represent categories with no implied order or ranking. In our dataset, columns like 'Branch', 'City', 'Gender', 'Product line', 'Customer type', and 'Payment' fall into this category.

The decision to use one-hot encoding is influenced by the understanding that these variables represent categories without any natural order. One-hot encoding creates binary columns for each category, providing a binary indicator (0 or 1) for the presence of each category. This is crucial to avoid introducing false ordinal relationships in the data.

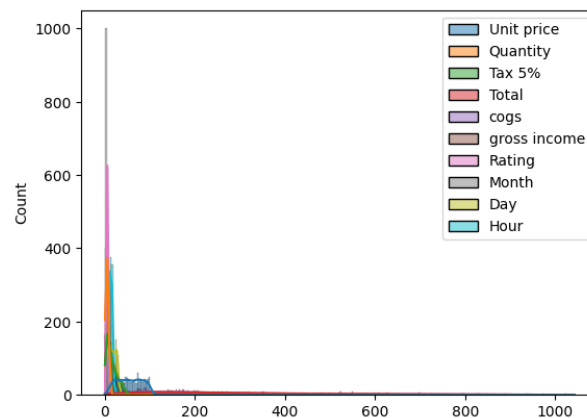
## 2.4 Standardization

### 2.4.1 Why was Standardization required here?

It is necessary to ensure that all features contribute equally to the analysis or model training, as many machine learning algorithms are sensitive to the scale of the input features. Standardization transforms the data so that it has a mean of 0 and a standard deviation of 1, making it easier for algorithms to converge and improving the interpretability of coefficients.

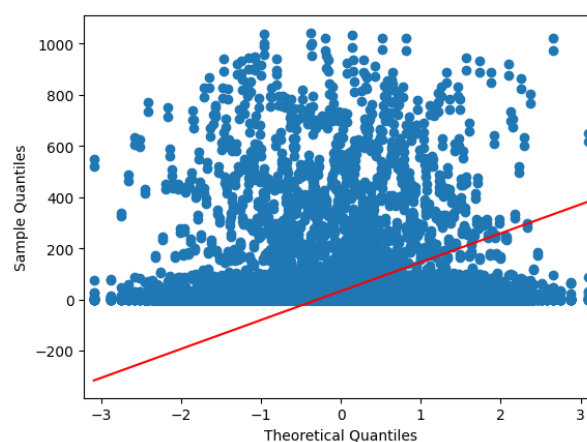
#### 2.4.1.1 Visual Insepection

##### a) Histogram:



Clearly, the histogram is not totally symmetric. Histograms depicting normal distribution show bell-shaped curves, and this is not bell-shaped. Hence, the data is not normally distributed.

##### b) Quantile-Quantile (Q-Q) Plot



In a Q-Q plot, if the points fall approximately along a straight line, the data is likely normally distributed. As can be seen clearly here, this is not the case. Deviations from the straight red line suggest departures from normality.

Hence, by both the visual inspections we have inferred that the dataset is not normally distributed. Another way of checking for data normality is by conducting statistical tests.

#### 2.4.1.2 Statistical Test (Shapiro-Wilk Test)

Upon conducting this test on our dataset, I got the following output:

Shapiro-Wilk Test: Statistic = 0.31571829319000244, p-value = 0.0

Since the p-value is lower than the chosen significance level (0.05), we can reject the null hypothesis; the data is not normally distributed.

Hence, both the visual and statistical tests have led us to believe that the data is not normally distributed and hence required normalization or standardization. But which one is suitable for our specific dataset and why?

#### 2.4.2 Which Features require Standardization and why?

Standardization becomes essential when examining features with numeric values, especially those displaying significant scale variations. The features that benefited from standardization include:

- Unit price: The prices of different products can vary widely, and standardizing this feature ensures that the magnitude of the price doesn't disproportionately influence other features during model training.
- Quantity: The number of items purchased can also exhibit considerable variability, and standardization aids in rendering this feature comparable with others.
- Tax 5%, Total, cogs, gross income: These features represent monetary values and may have different scales. Standardization is instrumental in ensuring that these monetary metrics are on a consistent scale, facilitating fair comparison.
- Rating: While this feature might not necessarily demand standardization if it already falls within a specific range, if the rating scale significantly differs from other features, standardization becomes beneficial in maintaining uniformity across all features.



### 2.4.3 Choosing between Normalization and Standardization

- **Magnitude Differences:** Standardization is robust to features with different scales, making it suitable for datasets where features like "Unit price," "Quantity," etc., may have different magnitudes.
- **Outliers:** Standardization is less sensitive to outliers compared to normalization, which can be beneficial in real-world datasets that might have some outliers.
- **Algorithms Sensitivity:** Many machine learning algorithms, including those used in data mining, can benefit from standardized features, especially when they involve calculations of distances or gradients.

### 2.4.4 Using StandardScaler

Before applying standardization, it is important to exclude the target column. Since we want to predict target values as they are in the original dataset, we do not apply any preprocessing like encoding or standardization to the target variable. The code involved in the actual standardization is basically 2 lines:

```
scaler = StandardScaler()
df_standardized[features] = scaler.fit_transform(df_standardized[features])
```

This code creates an instance of the StandardScaler class from scikit-learn. The StandardScaler is a transformer that standardizes features by removing the mean and scaling to unit variance. The `fit_transform` method calculates the mean and standard deviation of each selected column and then applies the standardization formula:

$$\text{StandardizedValue} = \text{OriginalValue} - \text{Mean} / \text{StandardDeviation} \quad (2.1)$$

	Customer type	Unit price	Quantity	Tax 5%	Total	cogs	gross income	Rating	Month	Day
0	-0.998002	0.711694	0.509930	0.951427	548	0.918142	0.951427	1.463146	-1.189455	-1.180314
1	1.002002	-1.517427	-0.174540	-1.013096	80	-0.987284	-1.013096	1.463146	1.206225	-0.835058
2	1.002002	-0.346194	0.509930	0.097287	340	0.072235	0.097287	0.276972	1.206225	-1.410484

FIGURE 2.1: Sample of a few features from the standardized dataset

### 2.4.5 How did this step contribute to effective data mining?

- **Improved Model Convergence:** Standardizing features helped machine learning models converge faster during training, especially for algorithms that relied on gradient descent optimization.

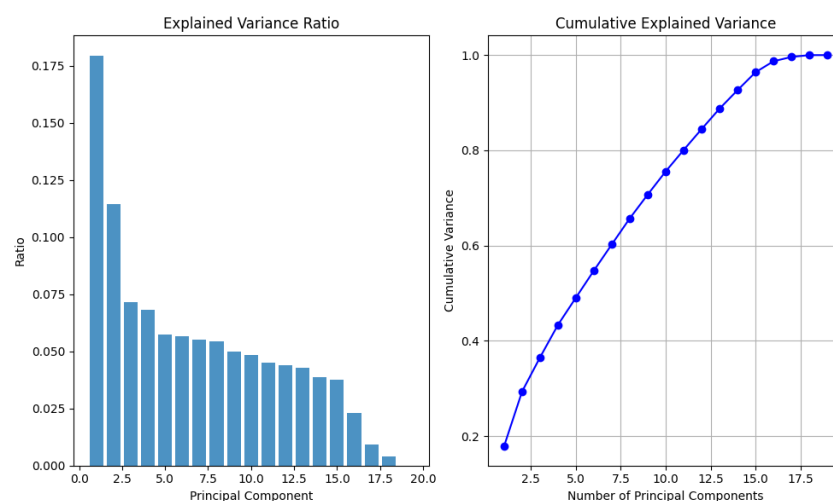
- **Enhanced Model Performance:** Standardization ensured that each feature contributed equally to the model, preventing features with larger scales from dominating the learning process. This led to better model performance and generalization.
- **Facilitated Interpretation:** Standardization made it easier to interpret the coefficients of the features in linear models. This was important for understanding the impact of each feature on the target variable.
- **Increased Robustness:** Scaling using StandardScaler made models more robust to variations in the data and helped avoid numerical instability issues during training data and helped avoid numerical instability issues during training.

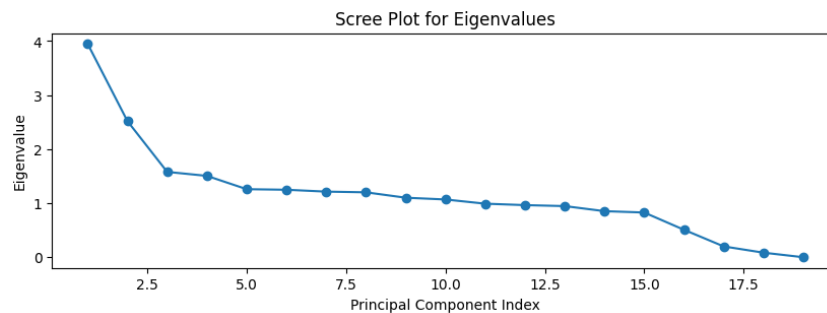
## 2.5 Feature Selection

I employed PCA to conduct feature selection. Technically, PCA makes use of both feature selection and creation. It basically ‘merges’ the top ‘k’ features together, where ‘k’ is user-defined. The motivation for using PCA specifically was simply because we had been taught that in depth, and I wanted to try it out. Choosing ‘k’ was the main factor.

### 2.5.1 Choosing K

As I mentioned above, ‘k’ is the number of principal components that we are choosing to work with. Choosing this value can be a challenging task, and I utilized various means to try and understand what the ideal value could be.





After analyzing these charts, I decided to initially start by experimenting with ‘k’ at 15. The most efficient way of finding the ideal value for ‘k’ is by trial and error. But using these means can help figure out a good starting point.

## 2.6 Splitting Data

### 2.6.1 Importance of splitting

- **Avoids Overfitting:** Imagine training your model on the same data you use to evaluate it. It would be like memorizing the answers in a test instead of truly understanding the concepts. The model might perform well on that specific data (high training accuracy), but it wouldn’t generalize well to new, unseen data (low testing accuracy). Splitting the data prevents this ”memorization” by reserving the testing set for unbiased evaluation.
- **Provides Unbiased Evaluation:** Training the model only on the training set eliminates any influence from the testing data. This ensures the evaluation metrics (e.g., R-squared, mean squared error) reflect the model’s true performance on unseen data, making it a reliable measure of its generalizability.
- **Helps Tune Hyperparameters:** Hyperparameters are knobs that influence the model’s learning process. Splitting the data allows you to experiment with different hyperparameter values on the training set, evaluating their impact on performance through the validation set, and finally selecting the best performing configuration for the final model.
- **Identifies Data Leakage:** Data leakage occurs when information from the testing set somehow influences the training process. For example, including future values in a time series prediction might give the model an unfair advantage in training. Splitting the data minimizes this risk and leads to a more accurate assessment of the model’s true capabilities.
- **Facilitates Cross-Validation:** Splitting the data can be extended to cross-validation techniques, where the dataset is further divided into multiple folds. Each fold is used for

training and testing in turn, providing a more robust and statistically sound evaluation of the model's performance across different data subsets.

### 2.6.2 Applying Train-Test Split

The `train_test_split` function from the `scikit-learn` library is used to split the data into training and testing sets. The parameters are as follows:

- `X`: The feature matrix (The preprocessed DataFrame)
- `y`: The target variable ('Total')
- `test_size`: The proportion of the dataset to include in the test split. In this case, it's set to 0.2, meaning 20% of the data will be used for testing.
- `random_state`: This is used to ensure reproducibility. Setting a seed (here, 42) means that if you run the code multiple times, you'll get the same train-test split each time.

The resulting dimensions of both sets: `X_train shape: (800, 22)`, `X_test shape: (200, 22)`, `y_train shape: (800,)`, `y_test shape: (200,)`

## Chapter 3

# Visualization

### 3.1 Variable Distributions and Patterns

#### 3.1.1 Numerical Variable Distributions:

For numerical variables like 'Unit price,' 'Quantity,' 'Tax 5%,' 'Total,' etc., their distributions can be visualized using histograms. Histograms show how the values are spread across different ranges or bins. Important aspects to consider include: Central Tendency: Is the distribution centered around a particular value (mean or median)? Spread: How spread out are the values? Are there outliers? Shape: Does the distribution have a particular shape, such as normal, skewed, or bimodal?

#### 3.1.2 Categorical Variable Distributions

For categorical variables like 'Branch,' 'City,' 'Customer type,' 'Payment,' etc., their distributions can be visualized using count plots or bar charts. These plots show the frequency of each category. Important considerations include:

- Frequency: Which categories are most common or rare?
- Imbalances: Are there significant imbalances in the distribution of categories?

Frequency: Which categories are most common or rare?

Imbalances: Are there significant imbalances in the distribution of categories?

### 3.1.3 Bivariate Relationships

Exploring relationships between variables, especially between numerical and categorical variables, helps understand how variables interact. Box plots, pair plots, or other bivariate visualizations provide insights into how the distribution of a numerical variable varies across different categories of a categorical variable.

### 3.1.4 Time Trends

For temporal variables like 'Date' in a time series dataset, understanding the distribution involves looking at how values change over time. Time series plots can reveal trends, seasonality, or patterns in the data.

## 3.2 Why is Visualization required here?

Visualizations are a cornerstone of data exploration and communication. They play a crucial role in this project by doing the following.

### 3.2.1 Unveiling Hidden Patterns and Relationships

- Revealing trends, correlations, and anomalies that might be obscured in raw data tables.
- Exposing patterns in customer behavior, product preferences, sales patterns, and seasonal influences.
- Highlighting relationships between variables, such as the impact of promotions on sales or the influence of customer demographics on purchasing habits.

### 3.2.2 Enhancing Understanding and Communication

- Transforming complex data into intuitive and accessible visual representations.
- Making insights easier to comprehend, interpret, and share among stakeholders with varying levels of technical expertise.
- Fostering collaboration and discussion by providing a visual platform for exploration and debate.

### 3.2.3 Guiding Further Analysis and Decision-Making:

- Informing subsequent steps in the data mining process, such as feature selection, model building, and evaluation.
- Supporting data-driven decision-making by providing evidence-based insights for marketing strategies, inventory management, customer segmentation, and other business initiatives.

## 3.3 Visualizations



FIGURE 3.1: Correlation Matrix Heatmap

Features like Tax, Total, Cost of Goods Sold, and Gross Income, can be seen having maximum positive correlation- which means they are highly directly correlated. Most other features are also highly correlated in a positive manner. Rating is slightly negatively correlated with all other features.

The distribution of the 'Total' column by Branch can be observed here. All Branches seem to have a higher distribution of values for 'Total' from the 50th to 75th percentile (middle to Q3). This means that all branches have higher totals per average transaction. How high though? As can be seen from the boxplots, the trend of higher distribution of values starts from around the \$200 mark.

The plot reveals the overall trend of total sales across all three branches over time. It allows us to identify seasonal patterns in sales, such as fluctuations due to holidays or periods of high

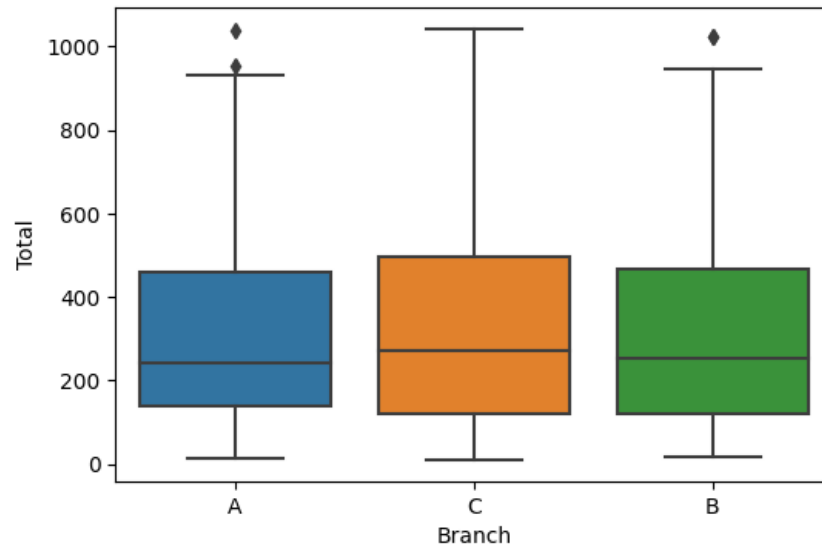


FIGURE 3.2: Box plot by branch



FIGURE 3.3: Time series graph

demand. It reveals higher trends in sales patterns, such as in mid-February and mid-March, that could indicate significant events or changes in store operations.

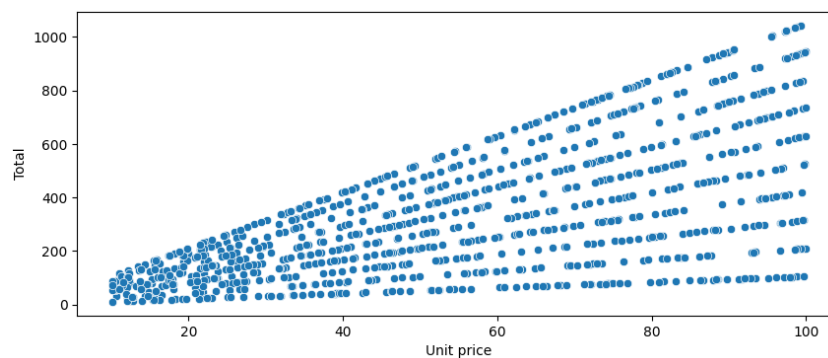


FIGURE 3.4: Scatter plot

Each point represents a single transaction, positioned according to its unit price and total value. Here we can observe a clear positive correlation, as values for 'Total' tend to increase with 'Unit



Price'. Not only this, but the density of points (transactions) also sees an upward trend as both features increase in their respective magnitudes.

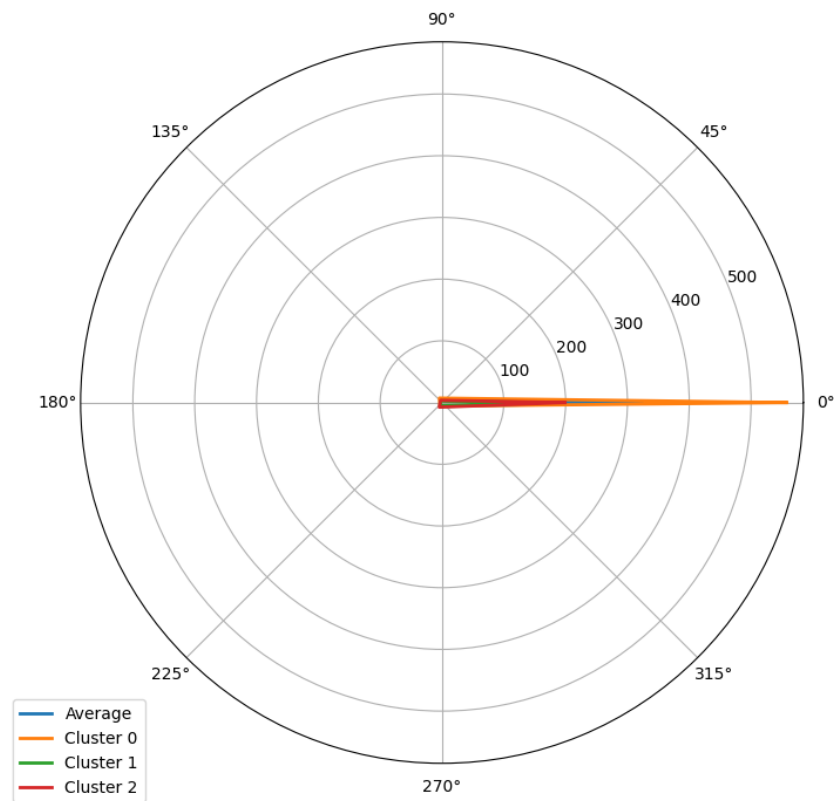
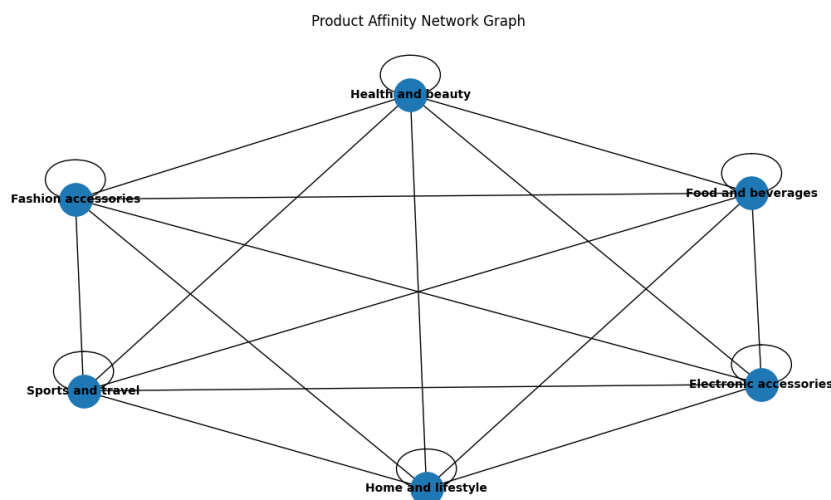


FIGURE 3.5: Radar graph

To make this, I selected specific features ('Total', 'Quantity', 'Rating') for segmentation. After this I standardized the features and then using K-Means, performed clustering (3 clusters). Each cluster represents a group of transactions with similar characteristics in terms of these variables. This helps identify distinct segments or groups within the data, allowing for a more nuanced understanding of how certain features contribute to the segmentation. Cluster 0 (represented by the outermost line in orange) covers most of the radar chart, occupying about 2/5ths of the outermost line. This suggests that a significant portion of the 3 features fall into Cluster 0. Cluster 1 is visible as a green segment for about 1/5th of the outer line. This indicates that there is a distinct group of data points (transactions) in Cluster 1 with specific characteristics that differentiate them from other clusters. Similarly, Cluster 2 is represented by a red segment for about 1/5th of the outer line. This suggests another distinctive group of data points with different characteristics compared to the rest of the dataset. The dominance of Cluster 0 suggests that it has a larger influence on the overall dataset. The cyclical appearance of different colored segments may indicate temporal patterns in customer behavior.



An undirected graph  $G$  is created using `nx.Graph()`. This graph represents relationships between products. The pairs of products represent the sequence of products in the 'Product line' column. Edges are added to the graph ( $G$ ) using `G.add_edges_from(pairs)`. Each pair of consecutive products in the 'Product line' forms an edge in the graph. The graph helps identify which products are often purchased together or follow each other in sales transactions. Each node here is connected to every other node which suggests a high level of interconnectedness. In the context of product affinity, it implies a strong association between all pairs of product lines. Hence, there are no specific clusters or groups of product lines; rather, all product lines are somewhat related to each other.

## Chapter 4

# Association Rule Mining

Apriori cannot be applied directly to the dataset. There are certain processes a dataset must go through before we can apply the algorithm. I chose to manually prepare the data to better understand and document the entire process.

### 4.1 Data preparation

#### 4.1.1 Creating a list of lists

First, I selected the ‘Product line’ feature from the dataset, as this was going to be the focus when implementing Apriori on the dataset. I split each element in the ‘Product line’ column into a list of substrings using the string ‘ and ’ as the separator. After this, I used the Pandas function ‘.explode’ to transform each element of a list-like to a row, replicating the index values. Finally, I converted the Pandas Series to a list. All these steps were performed in this single line of code:

```
exploded_list = apriori_df['Product line'].str.split(' and ').explode().tolist()
```

#### 4.1.2 Grouping products by count

I used the list of lists to create a Pandas DataFrame, and then added a new column ‘Count’ to this DataFrame. Using ‘groupby’ and ‘transform’, I was able to count the frequency of each product. Finally, I removed any duplicates to retain only the unique products and their corresponding counts.

### 4.1.3 Encoding the DataFrame using Transaction Encoder

The transaction encoder is also part of the `mlxtend` library and can be imported and directly used in Python to encode features based on their frequencies. Using the ‘fit’ method in ‘`tx.fit`’ along with ‘transform’, I was able to transform the data into a binary encoded array. 2 main changes made after this are to convert the array into a DataFrame, and to convert the binary values to an integer form (0s and 1s). Now, the dataset is finally ready as input to the Apriori algorithm.

## 4.2 Frequent Itemsets

### 4.2.1 Calculation

‘Apriori’ is a function from the `mlxtend.frequent_patterns` module, specifically designed for implementing the Apriori algorithm. I made use of it at this stage to apply the apriori algorithm to our dataset, with the objective of forming frequent itemsets. What exactly are frequent itemsets? Put simple, they are sets of items that appear together frequently in a dataset. Another term I’ll use here is minimum support, which measures the proportion of transactions in which an itemset occurs in a dataset. I set the ‘`min_support`’ parameter to 0.01 (1%), which specifies the minimum support threshold. Here, the ‘`min_support`’ parameter filters out infrequent itemsets, considering only those that occur with a frequency greater than or equal to 0.01.

### 4.2.2 Choosing Minimum Support

To answer this question, we must consider the ‘products by their counts’ DataFrame that I formed while preparing the data:

Product	Count
Health	152
Beauty	152
Electronic accessories	170
Home	160
Lifestyle	160
Sports	166
Travel	166
Food	174
Beverages	174
Fashion Accessories	178

The counts for each product line indicate that the dataset covers a variety of product lines, and there’s a reasonable distribution of counts across different categories. Setting a minsup of 0.01

means that we are considering itemsets that appear in at least 1% of the transactions. Since the entire dataset consists of a total of 1000 rows, this translates to a minimum count of 10 transactions. Considering that the lowest count is 152, a minsup of 0.01 seems reasonable as it is significantly above the minimum count of any individual product line, implying that we are capturing product lines which are relatively common in the dataset.

The choice of 0.01 strikes a balance between generality and specificity. It's low enough to capture patterns that occur frequently but high enough to avoid capturing extremely specific patterns. A lower minsup also allows for the discovery of a larger number of frequent itemsets, providing a more general view of patterns in the data.

### 4.2.3 The frequent itemset produced

The resultant DataFrame (`frequent_itemsets`) contains frequent itemsets along with their support values.

Support	Itemsets
0.166667	(Electronic accessories)
0.166667	(Fashion accessories)
0.166667	(Food)
0.166667	(Health)
0.166667	(Home)
0.166667	(Sports)
0.166667	(Beauty)
0.166667	(Beverages)
0.166667	(Lifestyle)
0.166667	(Beverages, Food)
0.166667	(Health, Beauty)
0.166667	(Lifestyle, Home)
0.166667	(Travel, Sports)

The itemsets of individual products (Electronic accessories), (Fashion accessories), (Food), (Health), (Home), (Sports), (beauty), (beverages), (lifestyle), and (travel) all have a support of 0.166667 (16.67%). This means that each of these product lines individually appears in approximately 16.67% of the transactions.

As for 2-itemsets: (beverages, Food), (Health, beauty), (lifestyle, Home), and (travel, Sports); each have a support of 0.166667 (16.67%), indicating that these pairs of product lines occur together in approximately 16.67% of the transactions.

## 4.3 Conducting Association Rule Mining

### 4.3.1 Applying association rule mining:

The ‘`association_rules`’ function is from the ‘`mlxtend.frequent_patterns`’ module, and it is used to generate association rules from frequent itemsets.

- `metric="confidence"`: Specifies the metric to be used for evaluating the generated rules. In this case, I chose confidence.
- `min_threshold=0.01`: Sets a minimum threshold for the chosen metric. Rules with a confidence value above this threshold will be considered.

### 4.3.2 Setting the confidence

#### 4.3.2.1 What is confidence?

Confidence is a measure used in association rule mining to quantify the likelihood that the occurrence of the antecedent (the left-hand side of the rule) implies the occurrence of the consequent (the right-hand side). It is defined as the ratio of the support of the combined antecedent and consequent to the support of the antecedent alone. Mathematically, it can be expressed as:

$$\text{Confidence}(A \rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)} \quad (4.1)$$

where A is the antecedent and B is the consequent.

#### 4.3.2.2 Choosing Minimum Confidence

Setting a low `min_confidence` threshold, such as 0.01, allows for the discovery of a broader range of association rules. It strikes a balance between generality and specificity. A lower `min_confidence` value may lead to the generation of more rules, providing a comprehensive view of potential associations. This can be valuable for exploratory data analysis and gaining insights into customer behavior. In association rule mining, there is often a trade-off between support and confidence. A lower `min_confidence` allows for the identification of rules with lower confidence but higher support, potentially revealing more common but weaker associations. Given the moderate size of our dataset (1000 transactions), a lower `min_confidence` allows for a more exploratory analysis without excluding potentially interesting rules.

### 4.3.3 The obtained association rules

The result is a DataFrame (rules) containing association rules with information such as antecedent, consequent, support, confidence, and lift. We have already talked about support, confidence, antecedents, and consequents. But what is lift?

#### 4.3.3.1 Lift

Lift is a measure of the strength of association between the antecedent and consequent in a rule. It quantifies how much more likely the occurrence of the antecedent and consequent together is compared to their independent occurrences. A lift value greater than 1 indicates a positive association, while a value less than 1 indicates a negative association.

$$\text{Lift}(A \rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A) \times \text{Support}(B)} \quad (4.2)$$

#### 4.3.3.2 Leverage

Leverage measures the deviation between the observed support of the combined antecedent and consequent and the expected support under independence. It assesses whether the co-occurrence of the antecedent and consequent is frequent than expected if they were independent.

$$\text{Leverage}(A \rightarrow B) = \text{Support}(A \cup B) - \text{Support}(A) \times \text{Support}(B) \quad (4.3)$$

#### 4.3.3.3 The obtained association rules:

Antecedents	Consequents	Antecedent Support	Consequents Support	Support	Confidence	Lift	Leverage
(beverages)	(Food)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(Food)	(beverages)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(Health)	(beauty)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(beauty)	(Health)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(Home)	(lifestyle)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(lifestyle)	(Home)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(travel)	(Sports)	0.166667	0.166667	0.166667	1.0	6.0	0.138889
(Sports)	(travel)	0.166667	0.166667	0.166667	1.0	6.0	0.138889

The provided association rules demonstrate compelling patterns within the dataset, showcasing a perfect confidence of 1.0 across various antecedent and consequent pairs, such as (beverages, Food), (Health, beauty), (Home, lifestyle), and (travel, Sports). The high lift values of 6.0 indicate positive associations, signifying that the occurrence of the antecedent significantly increases the likelihood of the consequent compared to their independent occurrences. The consistent leverage values of 0.138889 confirm a positive difference between observed and expected co-occurrences. These findings reveal strong relationships between specific product lines.



## Chapter 5

# Comparative Analysis with FP Growth Algorithm

### 5.1 Applying FP Growth Algorithm:

Since the objective of this step is comparison, not application of FP Growth, I will not go into detail about the application.

I simply used the same encoded `DataFrame` from the last step (Apriori algorithm application) to apply FP Growth. FP Growth is also available in the `mlxtend` Python library. The objective of FP Growth is the same as Apriori, hence it is a viable alternative. However, differences lie in the functionality, time complexity, memory usage and efficiency of both algorithms.

## 5.2 Comparison of Apriori and FP Growth Algorithms:

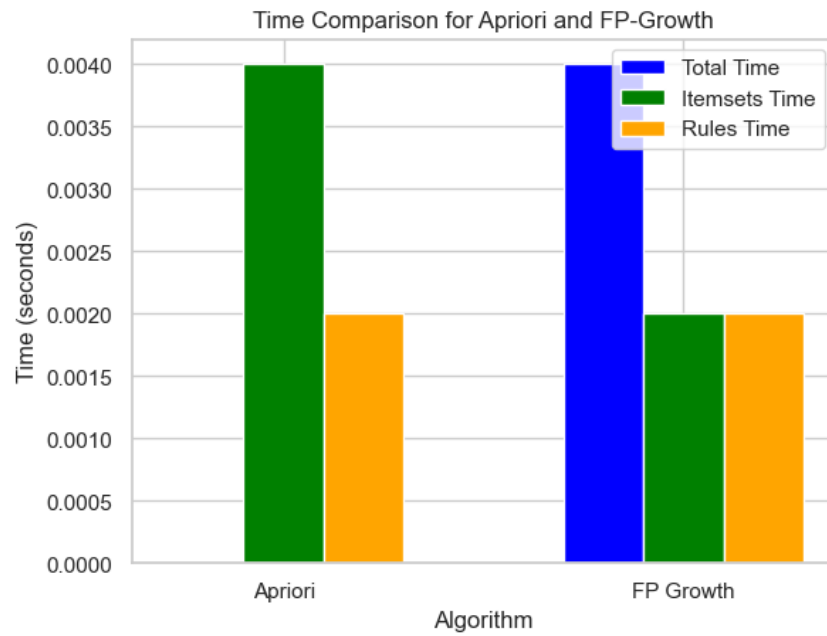
### 5.2.1 Functionality:

Feature/Aspect	Apriori	FP-Growth
Candidate Generation	Uses a candidate generation and test approach.	Uses a compact FP-tree structure, avoids explicit candidate generation.
Efficiency	Can be computationally expensive for large datasets.	Generally, more efficient, especially for large datasets.
Support Counting	Requires multiple passes through the dataset.	Performs support counting in a single pass with FP-tree.
Memory Usage	Can be memory-intensive due to storing candidates.	Generally, uses less memory with the compact FP-tree.
Handling Sparse Data	Less efficient with sparse datasets.	Well-suited for sparse datasets due to FP-tree structure.
Parallelization	Challenging to parallelize due to iterative nature.	Challenging to parallelize due to iterative nature. More amenable to parallelization, suitable for distributed computing.
Algorithmic Complexity	Exponential worst-case complexity.	Better average-case performance.

TABLE 5.1: There were no duplicates or outliers

### 5.2.2 Time Complexity:

Algorithm	Total Time	Frequent Itemsets	Association Rule
Apriori	0.044761	0.043750	0.008000
FP Growth	0.011839	0.008003	0.008027



This graph shows how the total time for both processes (generating frequent itemsets and rule mining), just the time for frequent itemset generation, and just the time for association rule mining, differ from one another.

### 5.2.3 Memory Usage:

Line#	Memory Usage	Increment	Occurances
41	142.1 MiB	142.1 MiB	1
42			
44	142.2 MiB	0.1 MiB	2
46	142.3 MiB	0.0 MiB	2

TABLE 5.2: Apriori

Line#	Memory Usage	Increment	Occurances
53	142.3 MiB	142.3 MiB	1
54			
56	142.3 MiB	0.0 MiB	2
58	142.3 MiB	0.0 MiB	2

TABLE 5.3: Apriori

Both algorithms show minimal memory usage variations during their execution. FP-Growth maintains a constant memory usage throughout its execution, indicating its ability to efficiently handle memory resources, especially for large datasets. Apriori exhibits a slightly higher memory increment during the discovery of frequent itemsets, but the additional memory is modest and may not be a significant concern.

## 5.2.4 Summarized Differences:

Aspect	Apriori	FP-Growth	Empirical Evidence
<b>Algorithmic Complexity</b>	Generally, has higher algorithmic complexity, especially in the frequent itemset generation phase.	Lower algorithmic complexity due to a single-pass tree construction.	Apriori's complexity can lead to longer execution times, as seen in your data.
<b>Scalability</b>	Tends to have scalability issues with large datasets and numerous itemsets.	More scalable, particularly for large datasets, as it requires fewer passes over the data.	FP-Growth may perform better as the dataset size increases, as suggested by your execution time results.
<b>Memory Requirements</b>	Typically, higher memory requirements, especially when dealing with many transactions and items.	Lower memory requirements, especially for large datasets, as it constructs a compact FP-tree.	FP-Growth is expected to be more memory-efficient based on the provided data.
<b>Empirical Evidence</b>	Apriori showed longer execution times, especially in the itemset generation phase.	FP-Growth demonstrated shorter execution times, indicating better performance.	The empirical evidence supports the idea that FP-Growth outperforms Apriori in terms of execution time.

TABLE 5.4: Table with 4 columns and 4 rows

## Chapter 6

# Tracking Patterns and Customer Behavior Analysis

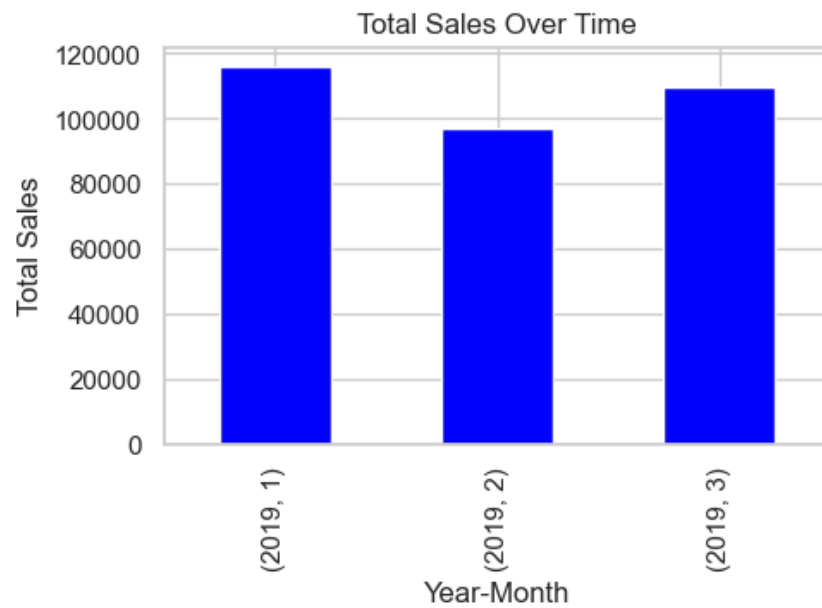
### 6.1 Creating Time-Series Data:

Time-series data is integral for tracking patterns and analyzing customer behavior evolution due to its ability to capture temporal nuances. By observing changes over time, businesses can uncover seasonal variations, identify trends, and forecast future behavior. This data-driven approach facilitates the detection of anomalies, adaptation of marketing strategies, and segmentation of customers based on temporal patterns. Utilizing time-series analysis improves decision-making, allowing businesses to make informed choices, allocate resources efficiently, and implement targeted campaigns, ultimately optimizing their strategies for dynamic customer preferences and market fluctuations.

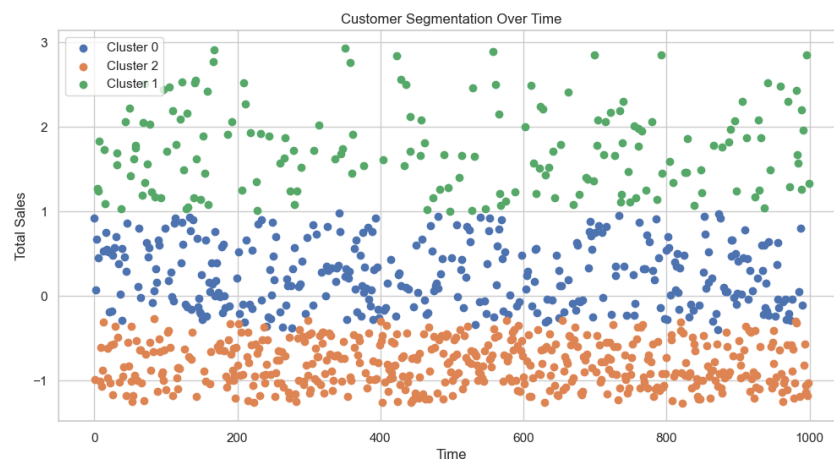
Hence, I utilized out existing 'Date' and 'Time' columns to create a 'Timestamp' feature for our dataset, which is evidently quite beneficial when conducting time-series analysis:

Timestamp Column:	
0	2019-01-05 13:08:00
1	2019-03-08 10:29:00
2	2019-03-03 13:23:00
3	2019-01-27 20:33:00
4	2019-02-08 10:37:00

## 6.2 Visualizations:

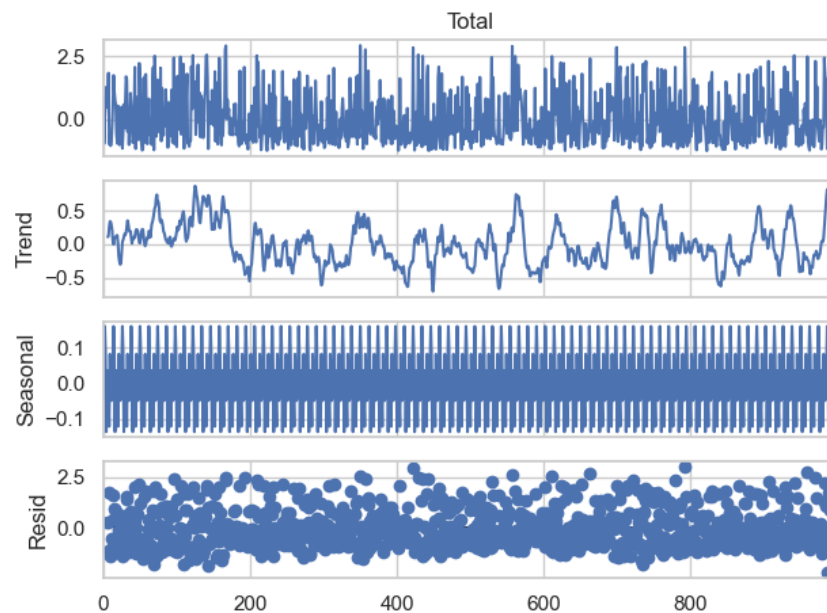


This chart depicts the total sales by year-month. As can be seen, February 2019 saw a fall in sales between January and March.



Using clustering, I was able to visualize customer segmentation over time.

### 6.2.1 Seasonal Decompose Graphs:



Using the `seasonal_decompose` function from `statsmodels.tsa`, I decomposed a time series into three components: trend, seasonality, and residual. The model is set to 'additive', which means the components are added together. The period is set to 12, indicating a seasonal pattern with a period of 12-time units.

- Observed Component (First Graph):

The top scatter plot represents the original time series, which is the observed component. As can be seen in the plot, the dense values with high variance suggest that the original time series has fluctuations, noise, or irregular patterns.

- Trend Component (Second Graph):

Sets a minimum threshold for the chosen metric. Rules with a confidence value above this threshold will be considered.

- Seasonal Component (Third Graph):

This graph captures the periodic patterns or seasonality in the data, showing how the data varies with a fixed interval (e.g., season).

- Residual Component (Fourth Graph):

This represents the remaining variability after removing the trend and seasonal components. It's essentially the difference between the observed data and the predicted values from the trend and seasonal components. The dense scatter plot suggests that the residuals still contain some variability or noise that is not explained by the trend and seasonal components.

### 6.2.2 Predictive Modeling using ARIMA:

ARIMA stands for AutoRegressive Integrated Moving Average. It is a popular time series forecasting model that combines three components: AutoRegressive (AR), Integrated (I), and Moving Average (MA).

#### 6.2.2.1 How is it beneficial here?

ARIMA helps identify and model patterns in the time series data, including seasonality, trends, and cyclical behaviors. It captures historical patterns to make predictions about future values.

Forecasting Customer Behavior: By applying ARIMA to historical customer behavior data (e.g., sales over time), the model can provide forecasts for future behavior. This can be valuable for planning marketing strategies, inventory management, and overall business decision-making.

Understanding Trends: The trend component of the ARIMA model helps in understanding the overall direction of the time series. This can be useful in identifying long-term trends in customer behavior.

ARIMA models, however, assume that the underlying patterns in the data are linear and stationary. For more complex and dynamic patterns, other time series models or machine learning approaches may be considered.

#### 6.2.2.2 Applying the model:

`ARIMA(df5['Total'], order=(1, 1, 1))`: This initializes an ARIMA model using the ARIMA class. `order=(1, 1, 1)` specifies the order of the ARIMA model. Here, (1, 1, 1) means:

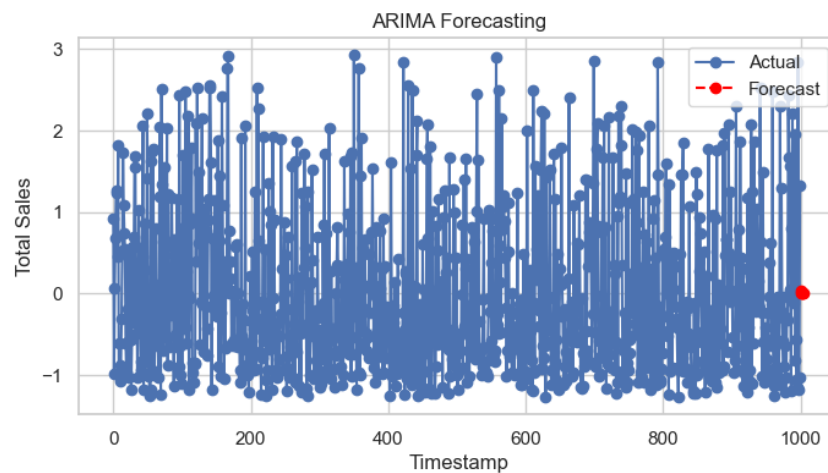
- `p=1`: One autoregressive term. It models the relationship between the current value and its immediate previous value (lag 1).
- `d=1`: One differencing. It indicates that the time series is differenced once to achieve stationarity. Differencing involves subtracting the previous value from the current value.
- `q=1`: One moving average term. It models the dependency between the current value and the residual error from the lagged moving average terms.

`model.fit()`: This fits the ARIMA model to the provided time series data (`df5['Total']`). `result.get_forecast(steps=forecast_steps)` generates forecasts for the next `forecast_steps` time points using the fitted ARIMA model.



**predicted\_mean:** This extracts the predicted mean values from the forecast result. The mean values represent the point estimates of the forecast.

### 6.2.2.3 Prediction and analysis

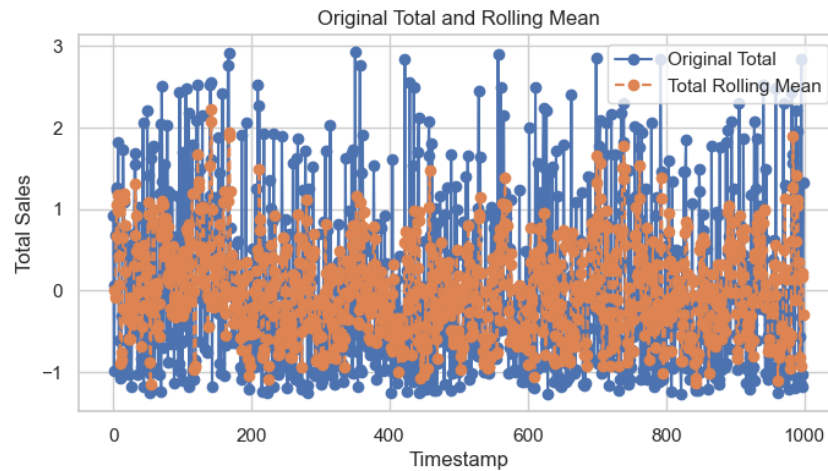


The dense distribution of blue points represents the actual "Total" sales values over time. The constant variation suggests that the sales data exhibit fluctuations or variability.

The sparse presence of red points (forecast) at the end of the x-axis indicates the predicted "Total" sales values for the next 5 time points. The small size of the red dots suggests that the model is forecasting relatively low variations or changes in sales over this short-term forecast horizon. Hence, there seems to be a relatively stable or slowly changing trend in future sales. The sparse and concentrated nature of the red points might also suggest a high level of confidence in the forecast, as the model is projecting a narrow range of potential outcomes. From all this we may conclude that there won't be much change in customer behavior soon.

### 6.2.3 Feature Engineering to calculate and visualize Rolling Mean:

### 6.2.3.1 The visualization:



### 6.2.3.2 How the graph works:

The purpose of calculating the rolling mean is to smooth out short-term fluctuations or noise in your data and highlight longer-term trends. In this case, we are using it to provide a clearer picture of how the "Total" sales values are evolving over time. We're using a rolling window of size 3. A rolling window refers to a fixed-size subset of consecutive data points that "rolls" or moves through the dataset as you analyze it. This helps in revealing trends and patterns.

### 6.2.3.3 Analyzing the visualization

The dense blue (Original Total) that can be seen throughout the graph, indicates frequent fluctuations or variations in the sales values. This could be due to regular transactions or short-term patterns in customer behavior.

In contrast, the orange points (Total Rolling Mean) are denser at the bottom half of the chart. Thus, the rolling mean is consistently lower than the actual sales value. This represents a downward trend or a smoothing effect that levels out short-term spikes or peaks in the original data. Hence, there are periods of lower sales that are being averaged out.

Based on all this, we can conclude that there are frequent fluctuations in total sales, which are potentially being influenced by daily transactions or short-term trends.

## Chapter 7

# Sequential Pattern Mining

### 7.1 Relevance to this project

Sequential Pattern Mining is a data mining technique that aims to discover patterns in sequential data, such as sequences of events or itemsets in a transactional dataset. In our project, sequential pattern mining can reveal interesting associations and orderings of itemsets across transactions and time periods. It may help identify frequent sequences of itemsets in the dataset. Patterns can reveal customer buying behavior, such as common sequences of products purchased together over time.

The goal of sequential pattern mining is to discover frequent sequences of items within a dataset. In the context of retail or supermarket sales, this could reveal common patterns of items that customers tend to purchase together, providing insights into customer behavior over time.

I decided to use a support value of 0.1 for all algorithms that I used.

### 7.2 How is this different from using Apriori or FP Growth?

Feature	Sequential Pattern Mining	Apriori Algorithm
Focus	Temporal order of events or itemsets	Frequent itemsets without order concern
Goal	Discover frequent sequences of items	Discover frequent sets of items
Order Sensitivity	Emphasizes order sensitivity	Not inherently order-sensitive
Data Structure	Utilizes sequence databases or similar	Utilizes itemsets and transaction data
Use Cases	Analyzing time-ordered data (e.g., transactions over time)	Finding associations in static datasets
Example Use Case	Identifying common sequences of products purchased over time	Identifying frequent sets of products irrespective of order

## 7.3 PrefixSpan:

### 7.3.1 How does it work?

PrefixSpan employs a depth-first search strategy to explore the space of possible sequential patterns. It uses a divide-and-conquer approach to efficiently find frequent sequences without generating all possible combinations.

The algorithm starts by considering each unique item in the dataset as a prefix and explores its extensions. For each prefix, it identifies the occurrences of that prefix in the dataset and extends it to find longer sequential patterns.

The recursion is a key feature of PrefixSpan. For each identified prefix, the algorithm recursively projects the remaining sequences and searches for frequent patterns in the projected sequences. This recursive approach allows PrefixSpan to efficiently discover sequential patterns of varying lengths.

### 7.3.2 Applying it via code:

I made use of the list of lists I had already computed in a previous task, as the input to this algorithm. Following is the main line in which I applied PrefixSpan:

```
patterns = PrefixSpan(transactions_list).frequent(0.1, closed=True)
```

The frequent method is applied to the PrefixSpan instance. It discovers frequent sequential patterns in the dataset based on a minimum support threshold of 0.1 (10%). The closed=True parameter indicates that closed sequential patterns should be included in the results.

The discovered patterns are then iterated over. This helps inspect and analyze the frequent sequential patterns in your dataset.

### 7.3.3 Result and Analysis:

**Output** (1, ['Health', 'beauty'])

(1, ['Home', 'lifestyle'])

(1, ['Sports', 'travel'])

(1, ['Electronic accessories'])

(1, ['Food', 'beverages'])

(1, ['Fashion accessories'])

These patterns signify sequences of itemsets that frequently occur together in transactions. The support value indicates how often each sequence appears in the dataset. Since all support values are 1 here, it indicates that each sequence occurred in exactly one transaction, meaning they are not very common or frequent.

## 7.4 GSP: Generalized Sequential Patterns:

### 7.4.1 How does it work?

GSP employs a sliding window approach to traverse the sequence of events systematically. This involves moving a fixed-size window or 'frame' across a sequence of data points, analyzing the data within the window at each position. The window slides through the sequence, typically one step at a time, capturing a subset of the data at each step.

GSP employs a breadth-first search strategy to explore the space of potential sequential patterns. GSP initially identifies frequent individual items and progressively extends them to discover longer patterns. The algorithm uses a candidate generation and pruning mechanism to efficiently identify patterns that meet the minimum support threshold. GSP iterates through the sequence, adjusting the size of the sliding window dynamically to capture patterns of varying lengths.

### 7.4.2 Applying it via code:

In the implementation of GSP, I utilized the seqmining library to discover frequent sequential patterns. The `freq_seq_nummethod` is applied to find sequences with a minimum support threshold of 0.1 (10%). The result is as follows:

### 7.4.3 Result and Analysis:

#### Output:

```
(( 'Food', ), 1)
(( 'Health', ), 1)
(( 'travel', ), 1)
(( 'Fashion accessories', ), 1)
(( 'Health', 'beauty', ), 1)
(( 'Food', 'beverages', ), 1)
(( 'Home', ), 1)
(( 'lifestyle', ), 1)
(( 'Electronic accessories', ), 1)
```

```
(('beverages',), 1)
(('Sports',), 1)
(('Sports', 'travel'), 1)
(('beauty',), 1)
(('Home', 'lifestyle'), 1)
```

The support count of 1 for each pattern indicates that the identified sequential patterns occur infrequently in the dataset. This might suggest that customers exhibit varied and unique purchasing patterns, and there isn't a dominant sequential pattern that applies to a large portion of the dataset.

## 7.5 Eclat Algorithm (Equivalence Class Transformation):

### 7.5.1 How does it work?

Eclat utilizes a vertical data representation, where each column corresponds to a unique item, and each row represents a transaction. It focuses on finding frequent itemsets by constructing equivalence classes for each item. Equivalence refers to a state of being equal or having the same value. These equivalence classes consist of sets of transactions containing the respective item. Eclat then intersects these equivalence classes to discover frequent itemsets of larger sizes. Its recursive approach ensures the discovery of progressively longer patterns. The vertical data representation and intersection-based technique contribute to Eclat's efficiency in mining frequent itemsets in datasets.

### 7.5.2 Applying it via code:

For Eclat, I employed the pyECLAT library, creating a DataFrame from the list of transactions. The ECLAT class is then used to find frequent itemsets with a minimum support of 0.1. The discovered itemsets are iterated over, facilitating an examination of the frequent itemsets in the dataset.

### 7.5.3 Result and Analysis:

**Output:** 'beverages': [4], 'Food': [4], 'Sports': [2], 'lifestyle': [1], 'Home': [1], 'Electronic accessories': [3], 'beauty': [0], 'Fashion accessories': [5], 'travel': [2], 'Health': [0], 'beverages Food': [4], 'Sports travel': [2], 'lifestyle Home': [1], 'beauty Health': [0] 'beverages': 0.16666666666666666, 'Food': 0.16666666666666666, 'Sports': 0.16666666666666666, 'lifestyle': 0.16666666666666666,

'Home': 0.16666666666666666, 'Electronic accessories': 0.16666666666666666, 'beauty': 0.16666666666666666,  
 'Fashion accessories': 0.16666666666666666, 'travel': 0.16666666666666666, 'Health': 0.16666666666666666,  
 'beverages Food': 0.16666666666666666, 'Sports travel': 0.16666666666666666, 'lifestyle Home':  
 0.16666666666666666, 'beauty Health': 0.16666666666666666

The diversity in item combinations and the uniform distribution of relative support values across different items imply that customers exhibit varied purchasing behaviors. There isn't a dominant sequential pattern or a specific set of items that consistently appear together in the dataset.

## 7.6 Summarization of Key Differences

Algorithm	Output (Sample)	Key Differences and Inferences
PrefixSpan	('Health', 'beauty')('Food', 'beverages')	Occurrence is infrequent, suggesting unique purchasing patterns; no dominant sequential pattern.
GSP	('Food')('Sports', 'travel')	Patterns occur infrequently; no clear dominant sequential pattern in the dataset.
Eclat	'beverages': [4], 'Food': [4], 'Sports': [2] 'beverages': 0.17, 'Food': 0.17, 'Sports': 0.17	Diversity in item combinations, uniform support distribution; no dominant sequential pattern identified.

## Chapter 8

# Clustering

### 8.1 Importance of feature selection

Feature selection in the context of clustering refers to the process of choosing a subset of relevant features from the original set of features in a dataset before applying a clustering algorithm. This selection is performed to enhance the quality of clustering results and improve the efficiency of the clustering process.

### 8.2 Selection of features and reasoning

I selected: 'Unit price', 'Quantity', 'Tax 5%', 'Total', 'cogs', 'gross income' and 'Rating'. Features like 'Unit price', 'Quantity', and 'Total', are likely to capture essential aspects of customer transactions. 'Rating' could provide insights into customer satisfaction. In addition to this, all the features I chose are commonly associated with sales transactions and customer behavior in retail. An alternative could've been employing some feature selection technique, like PCA. However, for this specific task, I chose to deal manually with feature selection. For all the plots, I've used the features 'Total' and 'Rating' for the axes.

### 8.3 Why a standardized dataset is required:

Standardization deals with issues related to scale disparities among features, which can significantly impact the performance of distance-based algorithms like KMeans. By transforming data to a common scale, standardization prevents biases toward features with larger magnitudes, enhances convergence, and facilitates meaningful comparisons. It also promotes model interpretability by ensuring that cluster centroids represent average values across features consistently.



## 8.4 Applying Clustering

### 8.4.1 Kmeans

K-Means is a popular partitioning algorithm that divides a dataset into K clusters. It aims to minimize the sum of squared distances between data points and the centroid of their assigned cluster. It assigns each data point to the cluster with the nearest centroid.

Metric Name	Score
Inertia (K-Means)	2705.273033683108
Silhouette Score	0.3037813725498845
Davies-Bouldin Index	1.1812055382786706
Calinski-Harabasz Index	791.3882872643744

The numbers indicate a moderate level of compactness (Inertia: 2705.27) and well-defined clusters, as indicated by the Silhouette Score (0.30), Davies-Bouldin Index (1.18), and Calinski-Harabasz Index (791.39). These metrics collectively suggest that the chosen number of clusters (3) is appropriate for grouping customers based on their purchase history. The clusters exhibit reasonable separation and distinctiveness, providing valuable insights for customer segmentation.

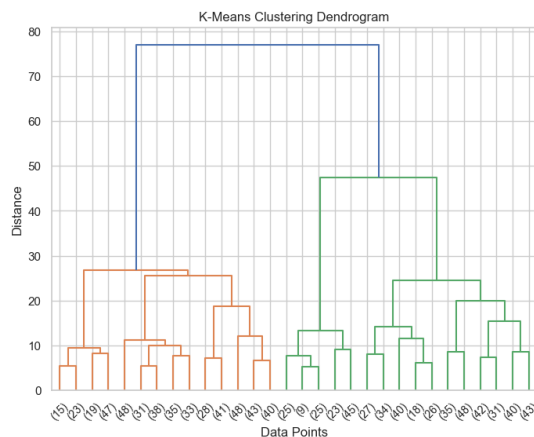


FIGURE 8.1: Dendrogram

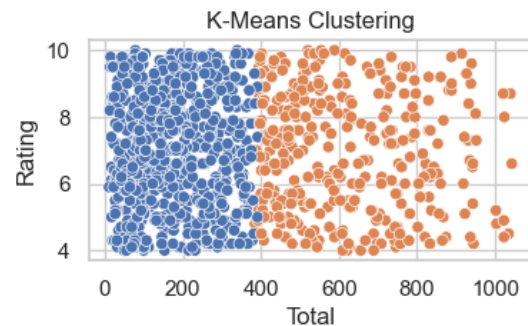


FIGURE 8.2: Scatterplot

### 8.4.2 K-Median

Like K-Means, K-Median is a partitioning algorithm that seeks to minimize the sum of absolute differences (medians) between data points and the centroid of their assigned cluster. K-Median is more robust to outliers compared to K-Means.

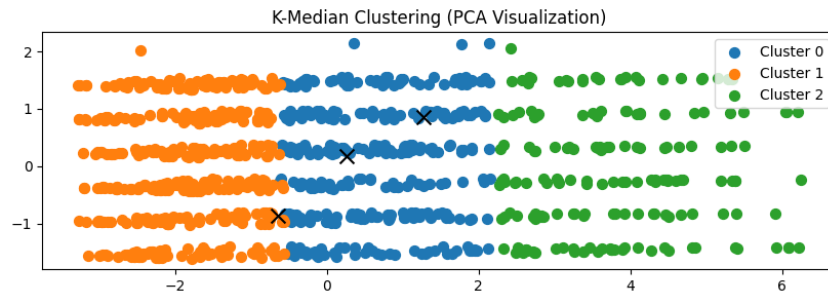


FIGURE 8.3: Scatter plot

Metric	Score
Adjusted Rand Index (ARI)	0.5451
Normalized Mutual Information (NMI)	0.5859

### 8.4.3 ARI and NMI

The Adjusted Rand Index (ARI) measures the similarity between clusters obtained by different algorithms. A score of 0.5451 indicates a moderate agreement between K-Means and Agglomerative Clustering.

Normalized Mutual Information (NMI) measures the mutual dependence between clusters. The score of 0.5859 suggests a reasonable degree of similarity between the two clustering methods.

The similarity in results between K-Means and Agglomerative Clustering suggests that both algorithms capture similar patterns in the data.

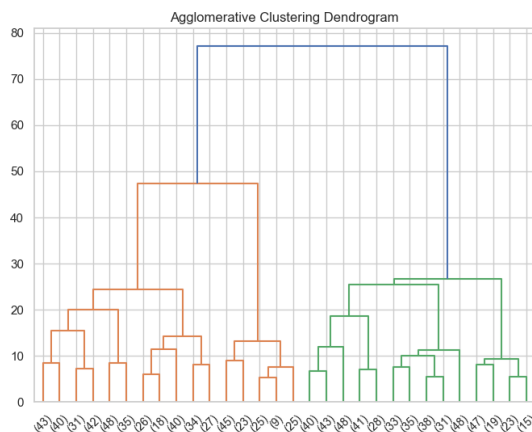


FIGURE 8.4: Dendrogram



FIGURE 8.5: Scatterplot

### 8.4.4 DBSCAN:

In applying DBSCAN to cluster customers based on purchase history, the algorithm identified 23 clusters. However, the negative Silhouette Score of approximately -0.125 indicates potential

Metric	Score
Silhouette Score (DBSCAN)	0.12528450383737091
Davies-Bouldin Index (DBSCAN)	1.5543217530912334
Number of Clusters (DBSCAN)	23

overlap between clusters, while the Davies-Bouldin Index of around 1.55 suggests moderate separation. These results may be indicative of the algorithm's sensitivity to data density and its tendency to partition the dataset into numerous, potentially small, dense regions.

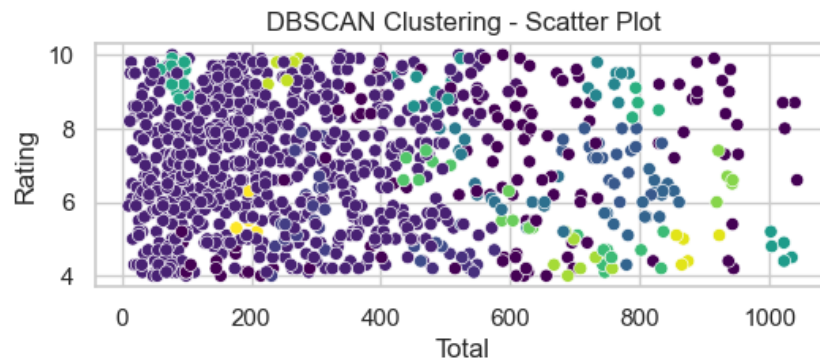


FIGURE 8.6: scatter plot

## Chapter 9

# Classification

Since most preprocessing required for classification here had already been performed in the first task of the project, I did not feel the need to conduct it again and simply re-used the dataset at the standardization stage. So, this is a sample of the DataFrame I started working with:

	Customer type	Unit price	Quantity	Tax 5%	Total	cogs	gross income	Rating	Month	Day	...
0	-0.998002	0.711694	0.509930	0.951427	548	0.918142	0.951427	1.463146	-1.189455	-1.180314	...
1	1.002002	-1.517427	-0.174540	-1.013096	80	-0.987284	-1.013096	1.463146	1.206225	-0.835058	...
2	1.002002	-0.346194	0.509930	0.097287	340	0.072235	0.097287	0.276972	1.206225	-1.410484	...
3	-0.998002	0.107187	0.852165	0.695185	489	0.674624	0.695185	0.870059	-1.189455	1.351560	...
4	1.002002	1.165074	0.509930	1.293083	634	1.268467	1.293083	-0.909202	0.008385	-0.835058	...

### 9.1 Choosing the Target Variable:

In the context of our problem (Customer Segmentation), the target variable was a choice between 2 features. Although other features could also be utilized, these 2 features seemed the best for performing customer segmentation. The 2 features are 'Total' and 'Customer type'.

I chose 'Total' due to its direct relevance to overall customer spending. It encapsulates the comprehensive financial transaction for each customer, making it a meaningful and practical variable for segmentation. The choice is justified by the potential predictive power of variations in spending patterns, the business relevance of understanding different customer segments, and the ability to employ classification metrics for model evaluation. Furthermore, the interpretability of 'Total' spending segments enhances the practical impact on business decision-making, aligning with the overarching goal of gaining actionable insights into customer behavior.

'Customer type', on the other hand, enables the model to differentiate between distinct customer types, such as "Member" and "Normal," offering valuable information for targeted marketing

strategies, loyalty programs, and operational decision-making. Predicting customer types can aid in customer retention efforts and enhance the overall customer experience by tailoring services to the specific needs and expectations of different segments.

Hence, from a business as well as a technical perspective (applying models), both features seem to make a lot of sense as the target variable. As I mentioned before, there are other features that could be used for this task as well. These include Product line, Payment Method, and Quantity. For customer segmentation specifically however, I felt 'Total' and 'Customer type' made more sense.

I chose 'Total' from these two because I wanted to explore the concept and application of binning. As 'Total' is a continuous feature and 'Customer type' is categorical.

By this point, I had cleaned, encoded, and standardized the data. However, binning was required before splitting the data, since my chosen target variable was 'Total', and this is a continuous feature in the dataset.

## 9.2 Binning:

### 9.2.1 The concept of binning:

Binning, also known as discretization, is a data preprocessing technique commonly used in data analysis and machine learning to handle numerical data. The general logic behind binning is to simplify the data and make it more manageable, especially when dealing with continuous variables. This process involves grouping a range of continuous or numerical data points into a smaller number of discrete "bins" or intervals.

### 9.2.2 Choosing an ideal number of bins:

For this step, I decided to use various rules which help indicate the ideal value for bins to use.

#### 9.2.2.1 Square Root Rule:

This is a simple method that calculates the number of bins by taking the square root of the total number of data points in the dataset. The formula is usually given as  $\sqrt{n}$ , where 'n' is the number of data points. It is easy to apply and works well for datasets of various sizes. It helps in preventing the creation of too many or too few bins.

**9.2.2.2 Sturges' Formula:**

$$k = 1 + \log^2(N) \quad (9.1)$$

where:

- $k$  is the number of bins.
- $N$  is the number of data points.

This formula assumes that the data follows a normal distribution. It is widely used and is suitable for datasets with a normal distribution. However, it may not perform well for datasets with outliers or non-normal distributions.

**9.2.2.3 Scott's Rule:**

$$3.5 * (S.D.) * (n^{1/3} - 1/3) \quad (9.2)$$

where:

- $h$  is the bin width.
- $N$  is the number of data points.

This rule considers the variability of the data. It is robust and performs well for datasets with varying degrees of variability. It adapts to the distribution of the data and is less sensitive to outliers.

**9.2.2.4 Freedman-Diaconis Rule:**

$$2 * IQR * (n^{1/3} - 1/3) \quad (9.3)$$

where:

- $h$  is the bin width.
- IQR is the interquartile range.
- $N$  is the number of data points.

This rule is particularly useful when dealing with datasets with outliers. It adjusts the bin width based on the spread of the data, making it less sensitive to extreme values.

### 9.2.2.5 Doane's Formula:

$$1 + \log_2(n) + \log_2(1 + |g_1|/\sigma_{g_1}) \quad (9.4)$$

where:

- k is the number of bins.
- N is the number of data points.
- $g_1$  is the skewness of the data.
- $\sigma(g_1)$  is the standard error of the skewness.

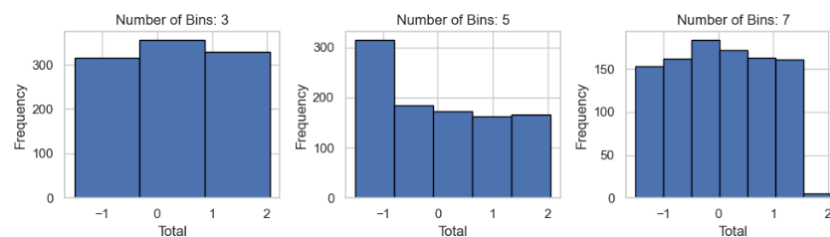
This Formula is an extension of Sturges' Formula that also considers the skewness of the data. It is suitable for datasets with non-normal distributions and can provide a more accurate estimate of the number of bins when skewness is present.

### 9.2.2.6 Results

Here are all the results summarized:

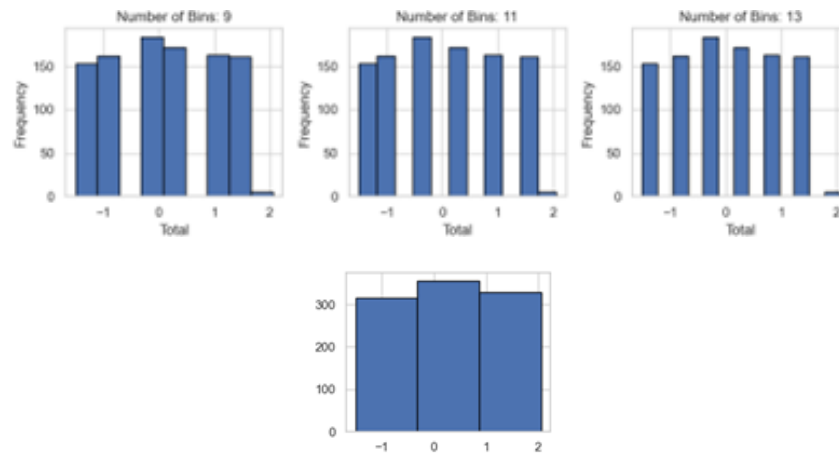
Method	Number of Bins
Square Root Rule	31
Sturges' Formula	11
Scott's Rule	11
Freedman-Diaconis Rule	14
Doane's Formula	13

Thus, from these results, the bin values were tending towards a value around 12. However, upon using such several bins, I faced issues. Hence, I went back to conduct further analysis:



As can be seen from these graphs, a higher value for the bins seems to form a sparser distribution. Thus, we might want to use a lower number of bins.

Judging from this, I decided to choose the number of bins to be 3, which seemed to be the ideal value as can be seen by this graph. All 3 bars, each representing a discrete 'Total' value, have an almost equal number of frequencies in our dataset.



### 9.2.3 Deciding on the bin edges:

Once I decided how many bins to use, the next step was figuring out what values were to be assigned to these bins. Let me explain. I chose the number of bins for 'Total' to be 3. Which means that the 'Total' column will be discretized into 3 categories. But on what basis? How will we decide what value the original 'Total' goes into which bin? For this, I explored our data some more:

As can be seen below, the count of the discrete values present in 'Total' is 546:

Total

125 6

175 6

..

721 1

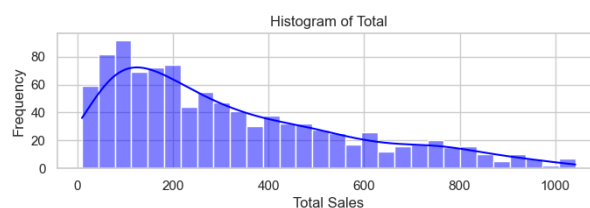
149 1

649 1

Name: count, Length: 546, dtype: int64

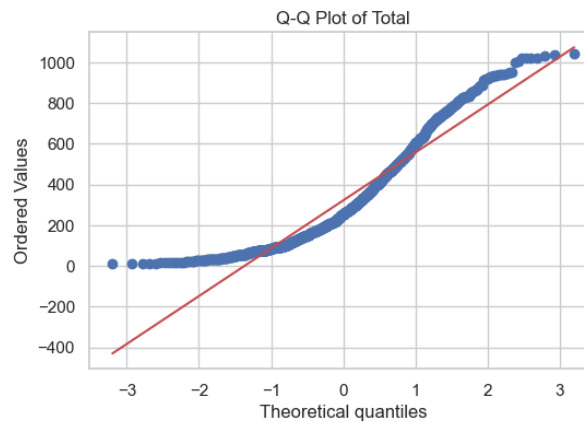
Also note, the Variance of the column 'Total' is 60457.8.

These values have given us an idea about the spread of the column 'Total'. To understand this spread more, its evenness and unevenness, I used the following visualizations: As can



be seen here, the blue line (actual values) does not follow a straight line, meaning the 'Total' column is not normally distributed. The points deviate from the line, especially towards the





tails, which indicates a departure from normality. The points curving upwards and downwards suggest skewness or light/heavy-tailed distribution.

I made use of the ‘Quantile’ method to calculate the bin edges. This method allows us to identify the values corresponding to specific percentiles within the distribution of the ‘Total’ column. I computed the 33rd and 66th percentiles, and the resulting values are utilized as reference points for defining the bin edges. Following is the code for this process:

```
percentile\_33 = df\_classif['Target'].quantile(0.33)
percentile\_66 = df\_classif['Target'].quantile(0.66)
bin\_edges = [df\_classif['Target'].min(),
, percentile\_33, percentile\_66, df\_classif['Target'].max()]
Bin Edges:
[10, 163.0, 381.680000000000006, 1042]
```

#### 9.2.4 Applying binning:

The last step before applying binning was to set the bin labels. I decided to set these to Low, Medium, and High. Since we’re applying binning on our ‘Total’ feature, which represents the total bill for each transaction, these labels seem to be good choices to categorize customer behavior.

After this, I applied binning using our decided labels, edges, and number of bins:

```
df\_classif['Total_bins'] = pd.cut(df\_classif[column_name], bins=bin_edges, labels=bin_labels, include_l
```

I made use of the powerful Pandas function ‘pd.cut’ to discretize the continuous feature ‘Total’. The first 3 parameters of the function are self-explanatory. The fourth parameter (include\_lowest) determines whether the leftmost edge of the bins should be considered as a valid bin edge or not. I’ve set this to true.

Here is the result following binning:

0 High

1 Low

2 Medium

3 High

4 High

...

995 Low

996 High

997 Low

998 Low

999 High

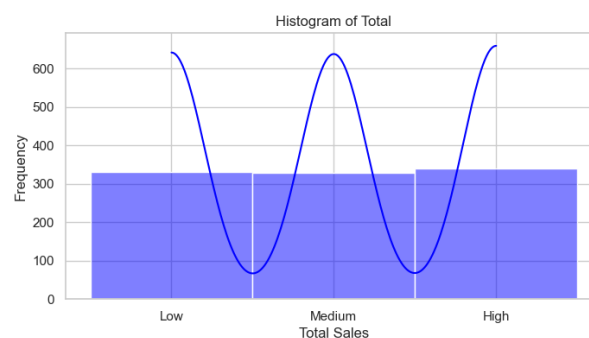
Name: Total, Length: 1000, dtype: category

Categories (3, object): ['Low' ; 'Medium' ; 'High']

### 9.2.5 Analyzing ‘Total’ after binning:

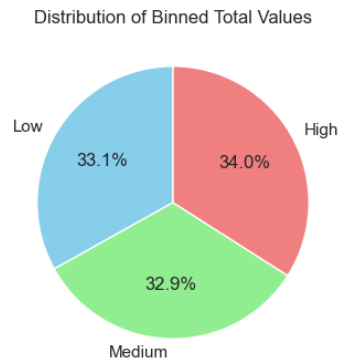
Category	Frequency
Low	331
Medium	329
High	340

In the above table, we can observe the frequencies of the new categories. Let’s visualize this using a histogram: All 3 categories reach approximately equal highs and lows. This indicates



that we have performed binning efficiently and can move onto the next step of the classification process.

This pie chart shows the spread of the categories High, Medium and Low as percentages of the ‘Total’ feature. It indicates that sampling is not required as these categories are roughly equal.



### 9.3 Splitting data

Using Train-Test split, I split the dataset. I already covered this step, in detail, during the preprocessing stage, so I won't go into detail here.

Training set shape: (800, 22) (800,)

Testing set shape: (200, 22) (200,)

### 9.4 Model Training and Testing:

Since the code after this is short, and applying the actual model is a small step, I decided to apply multiple models. Here are details about their obtained accuracy metrics:

Model	Accuracy	Precision	Recall	F1-Score
Random Forest	0.940	0.942260	0.940	0.940410
Gradient Boosting	0.925	0.927232	0.925	0.925458
Logistic Regression	0.940	0.940435	0.940	0.939632
Naive Bayes	0.920	0.921063	0.920	0.920297
K-Nearest Neighbors	0.810	0.816504	0.810	0.812382
Support Vector Machine	0.915	0.914824	0.915	0.914877

#### Key Findings:

- Random Forest and Logistic Regression outperform other models with the highest accuracy of 94.0
- Precision is consistently high across all models, indicating a low false positive rate in classifying customers.
- Recall values are generally balanced, suggesting good performance in capturing true positive instances.

- F1-Score combines precision and recall, showcasing the overall effectiveness of the model. Random Forest and Logistic Regression excel in this aspect.

**Conclusion** Considering the high accuracy, precision, recall, and F1-Score, Random Forest or Logistic Regression is recommended for predicting customer segments based on purchasing behavior.

## Chapter 10

# Regression

All previous preprocessing steps remain the same, the only difference is the exclusion of the process of binning since we need our target variable to be continuous for regression. I used the same target variable, 'Total'. As I did with classification, I applied multiple models here too, since the process after preprocessing is quite short and straightforward. Applying multiple models helps gain useful insight into the problem, and it becomes easier to choose a good model for predictions.

### 10.1 Evaluation Metrics:

Unlike classification problems, where metrics like accuracy, precision, and recall are common, regression problems involve predicting continuous numerical values. Here, we justify the selection of Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), R-squared (R<sup>2</sup>), Adjusted R-squared (Adjusted R<sup>2</sup>), and Explained Variance as evaluation metrics.

- Mean Squared Error (MSE):  
MSE measures the average of the squared differences between predicted and actual values. It amplifies the impact of larger errors, providing a comprehensive view of overall model accuracy.
- Root Mean Squared Error (RMSE):  
RMSE is the square root of MSE, providing an interpretable scale similar to the target variable. It helps in understanding the typical magnitude of errors in the predicted values.
- Mean Absolute Error (MAE):  
MAE calculates the average absolute differences between predicted and actual values. It offers a more straightforward interpretation, treating all errors equally.

- **R-squared (R2):**  
R2 quantifies the proportion of variance in the target variable explained by the model. A higher R2 indicates a better fit, with values closer to 1 indicating a stronger predictive capability.
- **Adjusted R-squared (Adjusted R2):**  
Adjusted R2 adjusts for the number of predictors in the model, providing a more accurate reflection of model fit. It penalizes the addition of irrelevant predictors, addressing potential overfitting.
- **Explained Variance:** Explained Variance measures the proportion of variance in the target variable that the model accounts for. It complements R2, offering another perspective on the model's explanatory power.

### 10.1.1 Why do we use different metrics for regression?

- **Continuous Prediction:**  
In regression, the goal is to predict numerical outcomes, making it inappropriate to use classification metrics designed for discrete class assignments.
- **Quantitative Error Assessment:**  
Regression metrics focus on quantifying the extent of errors between predicted and actual values, providing a nuanced understanding of model performance.

## 10.2 Applying Regression:

I made use of 5 different models, and 6 evaluation metrics, to help understand and evaluate the problem effectively. Here are the results:

Model	MSE	RMSE	MAE	R2	Adjusted R2	Explained Variance
Linear Regression	405.263702	20.131162	14.938315	0.993771	0.993263	0.993788
Ridge	405.359727	20.133547	14.947561	0.993769	0.993262	0.993787
ElasticNet	1366.993951	36.972881	27.411452	0.978989	0.977276	0.979146
Random Forest Regressor	302.036030	17.379184	13.506950	0.995358	0.994979	0.995358
Decision Tree Regressor	637.420000	25.247178	19.150000	0.990203	0.989404	0.990238

## 10.3 Analysis

- **Linear Regression and Ridge Regression:**  
Both linear models (Linear Regression and Ridge Regression) exhibit excellent performance with low MSE, RMSE, and MAE, indicating a strong fit to the data. The R<sup>2</sup> and Adjusted

$R^2$  values are close to 1, suggesting that these models explain a high proportion of the variance in the target variable.

- ElasticNet:

The ElasticNet model, while still performing reasonably well, shows higher errors compared to linear models. This could be due to the regularization applied, impacting the flexibility of the model.

- Random Forest Regression:

Random Forest Regression stands out as the top-performing model, with the lowest MSE, RMSE, and MAE. It demonstrates superior predictive accuracy and robustness, capturing complex relationships in the data.

- Decision Tree Regression:

Decision Tree Regression performs well, but not as accurately as Random Forest. It exhibits slightly higher errors and lower  $R^2$  values, indicating a lesser ability to explain the variance in the target variable.

## 10.4 Conclusion

Random Forest Regressor appears to be the most suitable for predicting sales and customer spending, providing a good balance between accuracy and generalization. The linear models (Linear Regression and Ridge Regression) also perform admirably.

## Chapter 11

# Outlier Detection and Statistical Validation

For this step, I'm going to be using encoded data, so that I'm able to apply the methods to all columns (since they are now numeric). Here is the DataFrame:

	Customer type	Unit price	Quantity	Tax 5%	Total	cogs	gross income	Rating	Month	Day	...
0	0	74	7	26	548	522	26	9	0	4	...
1	1	15	5	3	80	76	3	9	2	7	...
2	1	46	7	16	340	324	16	7	2	2	...
3	0	58	8	23	489	465	23	8	0	26	...
4	1	86	7	30	634	604	30	5	1	7	...

Let's consider the distinct value counts of all features, to get an idea of what we're working with:

All the encoded features have a distinct count of 2 since we used One-Hot encoding. Now, let's move onto applying different methods for outlier detection.



Customer type	Distinct Counts
Unit price	90
Quantity	10
Tax 5%	49
Total	546
COGS	527
Gross income	49
Rating	7
Month	3
Day	31
Hour	11
Branch_C	2
City_Naypyitaw	2
City_Yangon	2
Gender_Male	2
Product line_Fashion accesories	2
Product line_Food and beverages	2
Product line_Health and beauty	2
Product line_Home and lifestyle	2
Product line_Sports and travel	2
Payment_Credit card	2
Payment_Ewallet	2

## 11.1 Z-Score:

### 11.1.1 How does Z-Score work?

$$Z = (X - \mu) / \text{stddeviation} \quad (11.1)$$

The numerator  $X - \mu$  calculates the difference between the individual data point ( $X$ ) and the mean ( $\mu$ ). This part of the formula measures how far the data point is from the center of the distribution. The result of the numerator is then divided by the standard deviation ( $\sigma$ ). This division scales the difference in terms of standard deviations, providing a standardized measure of how far the data point is from the mean.

A Z-score of 0 indicates that the data point's value is exactly at the mean. A positive Z-score indicates that the data point is above the mean, and the magnitude of the Z-score tells us how many standard deviations above the mean the data point is. A negative Z-score indicates that the data point is below the mean, and again, the magnitude represents how many standard deviations below the mean.

### 11.1.2 Applying Z-Score in code:

In our code, we calculate the Z-score for each data point in selected columns of interest (e.g., Total, Unit price, Quantity) and identify outliers based on a user-defined threshold.

---

```
threshold_values = np.arange(1, 4, 0.5)
```

---

We must define a threshold value which will indicate the range for which we are considering the Z-score. The 'threshold\_values' is defined as an array ranging from 1 to 3. This array represents a range of threshold values that will be used in the outlier detection process.

As well as this, we make use of 'steps' to identify outliers. The concept of using steps refers to the systematic exploration of different threshold values in the outlier detection process. Here, I made use of the Numpy function 'arrange' to define the number of steps, which I set at 0.5. Basically, the code will consider the dataset in increments of 0.5 (Threshold Value).

The choice of starting from 1 and incrementing in steps of 0.5 allows for a nuanced exploration of outlier detection sensitivity. Lower thresholds may capture more points as outliers, while higher thresholds may only capture extreme anomalies.

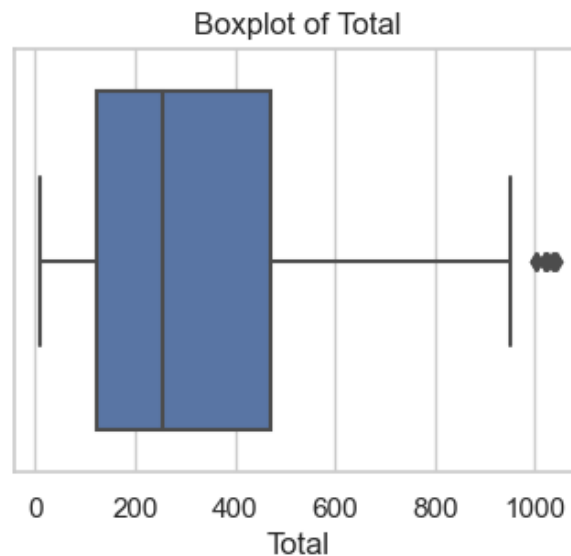
### 11.1.3 Result and Analysis:

Threshold	Total	Customer type	Unit price	Quantity	Tax 5%	Total	COGS	Gross Income
1.0	308	499	426	414	331	308	311	331
1.5	112	0	139	231	112	112	112	112
2.0	54	0	0	0	52	54	53	52
2.5	17	0	0	0	11	17	17	11
3.0	0	0	0	0	0	0	0	0
3.5	0	0	0	0	0	0	0	0

As expected, increasing the threshold leads to a decrease in the number of identified outliers. At the highest threshold (3.5), no outliers are detected in any column, indicating a more conservative approach to outlier identification.

- Total: The 'Total' column consistently has a high number of outliers across different thresholds, suggesting potential anomalies in the overall sales amount.
- Customer Type: The count of outliers for 'Customer type' is notably high at lower thresholds, indicating potential unusual patterns related to customer types.
- Unit Price, Quantity, Tax 5%: These columns show varying counts of outliers, suggesting potential discrepancies in pricing, purchase quantities, and tax amounts.
- Rating: Outliers in the 'Rating' column decrease with higher thresholds, indicating that extreme ratings become less significant as the threshold increases.

At certain thresholds, some columns have zero outliers. This is particularly evident for 'Hour' and 'Day,' which implies that these features might exhibit consistent patterns without significant outliers.



As can be observed, a certain number of outliers do seem to lie in the 'Total' column. I displayed all instances containing outliers in the column 'Total' and found these to be 17 rows in total.

## 11.2 Interquartile Range:

The Interquartile Range (IQR) is calculated as the difference between the third quartile (Q3) and the first quartile (Q1). This statistical measure is robust against extreme values and forms the basis for identifying outliers. The lower and upper bounds for potential outliers are established by extending 1.5 times the IQR below Q1 and above Q3, respectively.

Upon applying this method, I obtained 9 outlier rows based on the 'Total' feature.

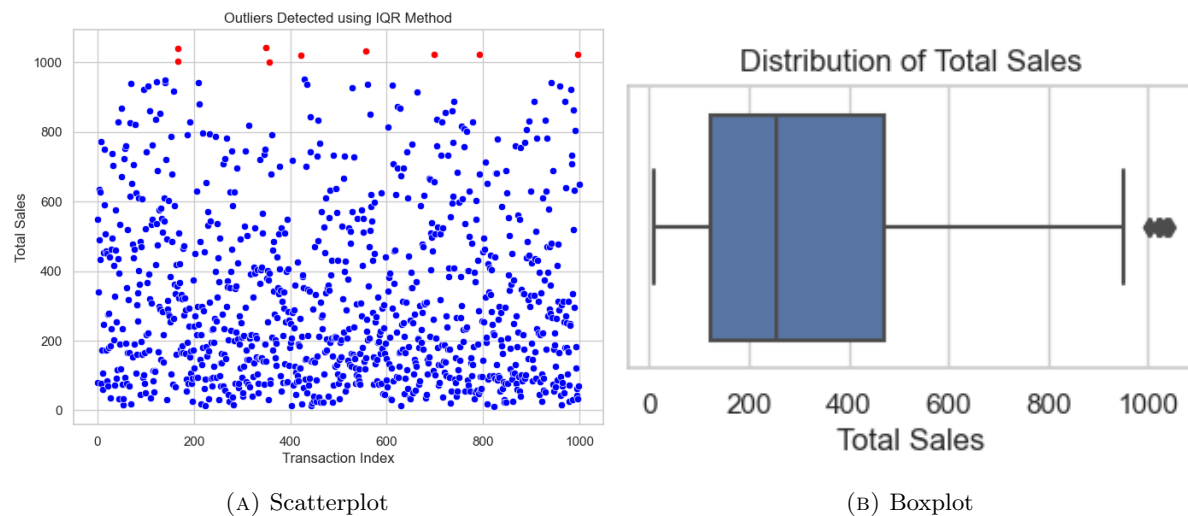


FIGURE 11.1: Two plots side by side

### 11.3 Isolation Forest:

The Isolation Forest algorithm is a robust technique for identifying outliers or anomalies within a dataset, particularly adept at handling high-dimensional data efficiently. Its core concept revolves around the construction of isolation trees, binary trees that isolate data points by recursively partitioning the feature space. Anomalies, being rare instances, are expected to have shorter path lengths in these trees. The algorithm builds an ensemble of such trees and derives anomaly scores based on the average path length. Lower scores signify points that are easier to isolate and are thus more likely to be outliers.

The focus is on key dimensions such as 'Total', 'Unit price', 'Quantity', and 'cogs' columns, which are deemed crucial for understanding customer.

```
model = IsolationForest(contamination=0.01)
```

In the code above, a contamination parameter of 0.01 is set, indicating the proportion of outliers expected in the dataset. The Isolation Forest algorithm assigns an anomaly score to each data point, and those with a score indicative of an outlier (equal to -1) are isolated and extracted into a separate dataframe.

A total of 10 rows are identified as outliers. Note that these are the top 1% outliers in our dataset based on the following features: 'Total', 'Unit price', 'Quantity' and 'cogs'.

## 11.4 Statistical Validation:

Statistical validation in the context of outlier detection involves applying formal statistical methods to assess and confirm the significance of identified outliers. It helps ensure that the patterns or anomalies observed in the data are not due to random chance and are, instead, statistically meaningful. The goal is to provide a robust foundation for drawing accurate conclusions from the detected outliers.

### 11.4.1 ANOVA:

Analysis of Variance is a statistical method used to compare means among multiple groups to determine if there are any statistically significant differences. It allows for the assessment of whether the variability within groups is comparable to the variability between groups. It is used to test the null hypothesis that there is no significant difference among the means of three or more groups. It also provides insights into whether any observed differences in means are likely due to genuine differences between groups or simply due to random variation.

After making use of ANOVA in my code, I was unable to detect any significant patterns. Upon testing for both ‘Total’ and ‘Customer type’ features, I found no significant difference amongst the different categories of both features. This was indicated to me by the values of T-statistic and P-value, which ANOVA calculates.

### 11.4.2 Hypothesis Testing:

In statistical validation, we often formulate hypotheses regarding the presence of outliers or patterns in the data.

- Null Hypothesis ( $H_0$ ): Assumes that there are no significant outliers or patterns.
- Alternative Hypothesis ( $H_1$ ): Assumes the presence of significant outliers or patterns.

Hence in our case, the Null Hypothesis would be that there is no significant difference in total sales between normal transactions and outliers. And the Alternative Hypothesis would be that there is a significant difference in total sales between normal transactions and outliers.

#### 11.4.2.1 Two-Sample t-Test:

- T-Statistic: This measures the difference between the means of the two groups relative to the variability within each group. A higher absolute t-statistic suggests a greater difference between the groups.

- **P-Value:** The p-value represents the probability of obtaining the observed results (or more extreme) under the assumption that the null hypothesis is true. A lower p-value indicates stronger evidence against the null hypothesis.

So, in our case, the computed t-statistic is a quantitative measure of the difference in total sales between normal transactions and outliers. And our P-Value will indicate whether there is a significant difference in total sales. If there is, we will reject the null hypothesis.

The significance level is denoted by alpha ( $\alpha$ ). This chosen level represents the threshold below which we would reject the null hypothesis. A common practice is to set alpha at 0.05, indicating a 5% probability of rejecting the null hypothesis when it is true. Hence, I also set it at 0.05.

<b>T-Statistic</b>	<b>8.916895886005445</b>
<b>P-Value</b>	<b>2.2421883966112505e-18</b>

With a p-value much smaller than the significance level, we reject the null hypothesis. Therefore, we conclude that there is a significant difference in total sales between normal transactions and outliers in our dataset.

The negative t-statistic indicates that outliers tend to have lower mean total sales compared to normal transactions.

# References

Here are some references to various documents and articles:

- Pandas Documentation [18]
- Scikit-Learn Documentation [22]
- NumPy Documentation [17]
- Matplotlib Documentation [16]
- Seaborn Documentation [24]
- OpenAI ChatGPT [26]
- Google Bard [2]
- Outliers Detection with IQR [12]
- Scipy Hierarchical Clustering [23]
- Understanding Boxplots [15]
- Pymining Documentation [21]
- Generalized Sequential Pattern (GSP) Algorithm [7]
- Mining Sequential Patterns [1]
- Sequential Pattern Mining [14]
- Agglomerative hierarchical clustering [3]
- K-means clustering [9]
- K-medians clustering [10]
- Interquartile Range (IQR) Method [8]
- Isolation Forest [25]

- DBSCAN [\[5\]](#)
- Mean shift [\[11\]](#)
- PrefixSpan Documentation [\[19\]](#)
- Principal Component Analysis (PCA) [\[13\]](#)
- Principal Component Analysis (PCA) - sklearn Documentation [\[20\]](#)
- Apriori Algorithm [\[4\]](#)
- FP-Growth Algorithm [\[6\]](#)



# Bibliography

- [1] Rakesh Agrawal and Ramakrishnan Srikant. “Mining Sequential Patterns”. In: (1995). URL: <https://dl.acm.org/doi/10.1145/223784.223791>.
- [2] Google AI. *Google Bard, Large language model*. not available. URL: <https://bard.google.com/>.
- [3] not available. *Article - "Agglomerative hierarchical clustering"*. not available. URL: [https://en.wikipedia.org/wiki/Hierarchical\\_clustering#Agglomerative\\_hierarchical\\_clustering](https://en.wikipedia.org/wiki/Hierarchical_clustering#Agglomerative_hierarchical_clustering).
- [4] not available. *Article - "Apriori Algorithm"*. not available. URL: [https://en.wikipedia.org/wiki/Apriori\\_algorithm](https://en.wikipedia.org/wiki/Apriori_algorithm).
- [5] not available. *Article - "DBSCAN"*. not available. URL: <https://en.wikipedia.org/wiki/DBSCAN>.
- [6] not available. *Article - "FP-Growth Algorithm"*. not available. URL: <https://en.wikipedia.org/wiki/FP-growth>.
- [7] not available. *Article - "Generalized Sequential Pattern (GSP) Algorithm"*. not available. URL: [https://en.wikipedia.org/wiki/Sequential\\_Pattern\\_Mining](https://en.wikipedia.org/wiki/Sequential_Pattern_Mining).
- [8] not available. *Article - "Interquartile Range (IQR) Method"*. not available. URL: [https://en.wikipedia.org/wiki/Interquartile\\_range](https://en.wikipedia.org/wiki/Interquartile_range).
- [9] not available. *Article - "K-means clustering"*. not available. URL: [https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering).
- [10] not available. *Article - "K-medians clustering"*. not available. URL: [https://en.wikipedia.org/wiki/K-medians\\_clustering](https://en.wikipedia.org/wiki/K-medians_clustering).
- [11] not available. *Article - "Mean shift"*. not available. URL: [https://en.wikipedia.org/wiki/Mean\\_shift](https://en.wikipedia.org/wiki/Mean_shift).
- [12] not available. *Article - "Outliers Detection with IQR"*. not available. URL: <https://www.geeksforgeeks.org/interquartile-range-to-detect-outliers-in-data/>.
- [13] not available. *Article - "Principal Component Analysis (PCA)"*. not available. URL: [https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis).

- [14] not available. *Article - "Sequential Pattern Mining"*. not available. URL: <https://www.geeksforgeeks.org/sequential-pattern-mining/>.
- [15] not available. *Article - "Understanding Boxplots"*. not available. URL: <https://towardsdatascience.com/understanding-boxplots-5e2df7bcbd51>.
- [16] not available. *Documentation - "Matplotlib Documentation"*. not available. URL: <https://matplotlib.org/stable/contents.html>.
- [17] not available. *Documentation - "NumPy Documentation"*. not available. URL: <https://numpy.org/doc/stable/>.
- [18] not available. *Documentation - "Pandas Documentation"*. not available. URL: <https://pandas.pydata.org/pandas-docs/stable/>.
- [19] not available. *Documentation - "PrefixSpan Documentation"*. not available. URL: <https://www.borgelt.net/prefixspan/prefixspan.pdf>.
- [20] not available. *Documentation - "Principal Component Analysis (PCA) - sklearn Documentation"*. not available. URL: <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>.
- [21] not available. *Documentation - "Pymining Documentation"*. not available. URL: <https://pymining.readthedocs.io/en/latest/index.html>.
- [22] not available. *Documentation - "Scikit-Learn Documentation"*. not available. URL: <https://scikit-learn.org/stable/>.
- [23] not available. *Documentation - "Scipy Hierarchical Clustering"*. not available. URL: <https://docs.scipy.org/doc/scipy/reference/generated/scipy.cluster.hierarchy.linkage.html>.
- [24] not available. *Documentation - "Seaborn Documentation"*. not available. URL: <https://seaborn.pydata.org/>.
- [25] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. *Isolation Forest*. 2008. URL: <https://cs.nju.edu.cn/zhoush/zhoush.files/publication/icdm08b.pdf>.
- [26] OpenAI. *ChatGPT, Version 3.5*. not available. URL: <https://openai.com/chatgpt>.