

**NJALA UNIVERSITY**



School of Technology

MSc. In Computer Science

**APRIORI ALGORITHM DOCUMENTATION**

**SUPERVISOR: A. J FOFANAH**

**GROUP MEMBERS**

**Abu Junior Vandi – 82937**  
**Joseph Prince Conteh – 54992**

# Table of Contents

<b>Table of Contents</b>	<b>1</b>
<b>Market Basket Analysis Using Apriori Algorithm</b>	<b>2</b>
<b>1.0 Introduction</b>	<b>2</b>
1.1 Project Overview	2
1.2 Objective	2
<b>2.0 Dataset Description</b>	<b>3</b>
<b>3.0 Data Preprocessing (Cleaning)</b>	<b>6</b>
3.1 Overview	6
3.2 Removing Duplicates	6
3.3 Handling Missing Values	7
3.4 Formatting Data	7
3.5 Converting into Market Basket Format (Transaction-based Representation)	8
<b>4.0 Transaction Data Transformation</b>	<b>9</b>
4.1 Overview	9
4.2 Data Filtering	9
4.3 Data Transformation	9
4.4 Binary Encoding	10
<b>5.0 Model Training (Apriori Algorithm)</b>	<b>10</b>
5.1 Overview	10
<b>6.0 Data Visualization</b>	<b>13</b>
6.1 Introduction	13
6.2 Heatmap of Frequent Itemsets	13
6.2 Bar Chart of Top 10 Most Frequently Purchased Items	15
6.3 Tabular Representation of Association Rules	17
<b>7.0 Analysis and Interpretation</b>	<b>17</b>
7.1 Support, Confidence, and Lift Values	18
7.2 Examples of Frequent Itemsets	18
7.3 Insights for Business Decisions	19
7.4 Code Implementation	19
7.5 Interpretation of Results	20
<b>8.0 Product Recommendation System</b>	<b>20</b>
<b>9.0 Conclusion</b>	<b>23</b>
<b>10.0 References</b>	<b>25</b>

# Market Basket Analysis Using Apriori Algorithm

## 1.0 Introduction

### 1.1 Project Overview

Market Basket Analysis (MBA) is a data mining technique used to uncover associations between products in transactional data. It helps businesses understand purchasing patterns by identifying items that are frequently bought together. This technique is widely used in retail and e-commerce to optimize sales strategies and enhance customer experience. Market Basket Analysis is crucial for businesses as it provides insights that can drive revenue growth through strategic decision-making. Some of the key benefits include improved product placement, targeted promotions, personalized recommendations, and efficient inventory management. By arranging products effectively, businesses can encourage impulse purchases, offer discounts on related products to increase sales, suggest items based on previous purchases to enhance customer satisfaction, and stock complementary products together to reduce stockouts and optimize supply chain efficiency.

The Apriori Algorithm is a fundamental algorithm used in Market Basket Analysis to identify frequent itemsets and generate association rules. It operates on the principle that if an itemset is frequent, its subsets are also frequent. The algorithm follows a series of steps: identifying frequent itemsets based on a minimum support threshold, generating candidate itemsets by extending frequent itemsets, filtering itemsets using minimum support to retain only frequent ones, and deriving association rules using confidence and lift measures. This algorithm plays a critical role in uncovering meaningful relationships between products, enabling businesses to make data-driven decisions and improve their overall performance.

### 1.2 Objective

The primary objective of this analysis is to:

- Identify frequent product combinations.
- Generate association rules to understand customer purchasing behavior.
- Visualize key insights from the data.
- Develop a basic product recommendation system based on identified rules.

## 2.0 Dataset Description

The dataset utilized for this analysis is the **Online Retail Dataset**, which comprises transactional data from a UK-based online retail store. This dataset captures detailed information about customer purchases, including product details, transaction timestamps, pricing, customer demographics, and geographical distribution. The dataset is structured to facilitate comprehensive analysis of customer purchasing behavior, product performance, sales trends, and regional market dynamics. Below is a detailed and professional description of the key columns within the dataset, updated to reflect the sample data provided:

### 1. **InvoiceNo:**

- **Description:** A unique identifier assigned to each transaction. Each transaction (invoice) can include multiple products, and the same **InvoiceNo** may appear across multiple rows if the transaction involves multiple items.
- **Purpose:** This column is critical for grouping and analyzing transaction-level data. It enables the identification of all items purchased in a single transaction, which is essential for understanding customer purchase patterns, basket analysis, and transaction value calculations.
- **Data Type: String** (e.g., 536365, 536366, 536367).
  - The alphanumeric format suggests that the identifier may include both numbers and letters, though the sample data only shows numeric values.
- **Example:** 536365 (represents a transaction containing multiple products, such as WHITE HANGING HEART T-LIGHT HOLDER and WHITE METAL LANTERN).

### 2. **StockCode:**

- **Description:** A unique identifier assigned to each product in the retail store's inventory. This code is used to distinguish individual products and track their sales performance.
- **Purpose:** This column enables the tracking of specific products across transactions. It is essential for analyzing product popularity, inventory management, and sales trends. It also facilitates the identification of frequently purchased items or product bundles.
- **Data Type: String** (e.g., 85123A, 71053, 84406B).
  - The alphanumeric format indicates that the product codes may include both letters and numbers.
- **Example:** 85123A (represents the product WHITE HANGING HEART T-LIGHT HOLDER).

### 3. **Description:**

- **Description:** The name or description of the product associated with the transaction. This field provides a textual representation of the product, offering context for the **StockCode**.
- **Purpose:** This column provides detailed information about the products purchased, aiding in the interpretation of product-related insights. It is particularly useful for identifying product categories, themes, or attributes (e.g., seasonal items, home decor, or kitchenware).
- **Data Type: String** (e.g., WHITE HANGING HEART T-LIGHT HOLDER, WHITE METAL LANTERN, CREAM CUPID HEARTS COAT HANGER).
- **Example:** WHITE METAL LANTERN (describes the product associated with **StockCode** 71053).

### 4. **Quantity:**

- **Description:** The number of units of a specific product purchased in a single transaction. This value represents the volume of the product sold in the transaction.
- **Purpose:** This column quantifies the volume of products sold, enabling analysis of sales volume, product popularity, and customer purchasing behavior. It is also used to calculate total revenue when combined with the **UnitPrice**.
- **Data Type: Integer** (e.g., 6, 8, 32).
  - The values are whole numbers, as partial units are not applicable in this context.
- **Example:** 6 (indicates that six units of WHITE HANGING HEART T-LIGHT HOLDER were purchased in the transaction).

### 5. **InvoiceDate:**

- **Description:** The timestamp indicating the date and time when the transaction occurred. This field provides precise temporal information about each purchase.
- **Purpose:** This column is critical for time-series analysis, allowing for the examination of sales trends over time, seasonal patterns, and the impact of time-specific factors (e.g., time of day, day of the week) on purchasing behavior.
- **Data Type: DateTime** (e.g., 01/12/2010 08:26, 01/12/2010 08:28, 01/12/2010 08:34).
  - The format includes both date (DD/MM/YYYY) and time (HH:MM).

- **Example:** 01/12/2010 08:26 (represents the timestamp of a transaction involving multiple products).

#### 6. **UnitPrice:**

- **Description:** The price per unit of the product in the transaction. This value represents the cost of a single unit of the product at the time of purchase.
- **Purpose:** This column provides pricing information, which is essential for calculating total transaction value, analyzing price sensitivity, and evaluating revenue trends. It is also used to identify high-value or low-value products.
- **Data Type: Float** (e.g., 2.55, 3.39, 1.85).
  - The decimal values indicate that prices are recorded with precision up to two decimal places.
- **Example:** 2.55 (represents the price of one unit of WHITE HANGING HEART T-LIGHT HOLDER).

#### 7. **CustomerID:**

- **Description:** A unique identifier assigned to each customer. This field distinguishes individual customers and links their purchase history across transactions.
- **Purpose:** This column facilitates customer-level analysis, including customer segmentation, repeat purchase behavior, and customer lifetime value (CLV) calculations. It is also used to identify high-value customers or those with specific purchasing patterns.
- **Data Type: Integer** (e.g., 17850, 13047, 12583).
  - The numeric format suggests that customer IDs are represented as whole numbers.
- **Example:** 17850 (represents a customer who made multiple transactions, such as 536365 and 536366).

#### 8. **Country:**

- **Description:** The country where the customer is located. This field provides geographical context for each transaction.
- **Purpose:** This column enables geographical segmentation and analysis of regional sales trends, customer distribution, and market penetration. It is particularly useful for identifying high-performing regions or tailoring marketing strategies to specific locations.

- **Data Type: String** (e.g., United Kingdom, France).
  - The text format indicates that country names are recorded as full names rather than codes.
- **Example:** United Kingdom (represents the country where most transactions in the sample data occurred).

## 3.0 Data Preprocessing (Cleaning)

### 3.1 Overview

Data preprocessing is a critical step in any data analysis or machine learning project. It involves transforming raw data into a clean, structured format that is suitable for analysis. This section outlines the key steps involved in data preprocessing, including removing duplicates, handling missing values, formatting data, and converting the dataset into a market basket format for transaction-based analysis. Each step is essential to ensure the quality, consistency, and usability of the data.

### 3.2 Removing Duplicates

Duplicate records in a dataset can lead to biased analysis and inaccurate results. Removing duplicates ensures that each data point is unique and representative.

#### Steps:

1. **Identify Duplicates:** Use functions or methods to detect duplicate rows based on specific columns or the entire dataset.
2. **Remove Duplicates:** Eliminate redundant records while retaining the first occurrence or a representative instance.
3. **Validation:** Verify that duplicates have been successfully removed by rechecking the dataset.

#### Implementation:

```
import pandas as pd

# Load dataset
file_path = r"C:\Users\user\Desktop\Dr. A. J\MBA\Online Retail.xlsx"
df = pd.read_excel(file_path)

# 3.2 Removing Duplicates
df = df.drop_duplicates()
```

### 3.3 Handling Missing Values

Missing values can distort analysis and lead to incorrect conclusions. Proper handling of missing data is essential to maintain the integrity of the dataset.

#### Approaches:

1. **Identify Missing Values:** Detect null or NaN values in the dataset.
2. **Deletion:** Remove rows with missing values in critical columns to ensure data completeness.

#### Implementation:

```
import pandas as pd

# Load dataset
file_path = r"C:\Users\user\Desktop\Dr. A. J\ MBA\Online Retail.xlsx"
df = pd.read_excel(file_path)

# 3.3 Handling Missing Values
required_columns = ["InvoiceNo", "StockCode", "Description", "Quantity", "InvoiceDate", "UnitPrice"]
df = df.dropna(subset=required_columns)
```

### 3.4 Formatting Data

Data formatting ensures consistency in the dataset, making it easier to analyze and process. This includes standardizing data types, creating derived columns, and ensuring proper encoding.

#### Steps:

1. **Standardize Data Types:** Convert columns to appropriate data types (e.g., datetime).
2. **Create Derived Columns:** Generate new columns, such as TotalPrice, to enhance the dataset's utility.
3. **Validation:** Verify that the data is correctly formatted.

#### Implementation:



```
import pandas as pd

# Load dataset
file_path = r"C:\Users\user\Desktop\Dr. A. J\MBA\Online Retail.xlsx"
df = pd.read_excel(file_path)

# 3.4 Formatting Data
df["InvoiceDate"] = pd.to_datetime(df["InvoiceDate"])
df["TotalPrice"] = df["Quantity"] * df["UnitPrice"]
```

### 3.5 Converting into Market Basket Format (Transaction-based Representation)

Market basket analysis requires data to be in a transaction-based format, where each row represents a transaction and each column represents an item. This format is essential for association rule mining and pattern discovery.

#### Steps:

1. **Group by Transactions:** Aggregate data by transaction IDs (InvoiceNo) and items (Description).
2. **Sum Quantities:** Sum the quantities of each item per transaction.
3. **Convert to Binary Format:** Represent items as binary indicators (1 for presence, 0 for absence) in each transaction.

#### Implementation:

```
import pandas as pd

# Load dataset
file_path = r"C:\Users\user\Desktop\Dr. A. J\MBA\Online Retail.xlsx"
df = pd.read_excel(file_path)

# 4.2 Data Filtering (Germany transactions)
df_germany = df[df["Country"] == "Germany"]

# 4.3 Data Transformation
basket = df_germany.groupby(["InvoiceNo", "Description"])['Quantity'].sum().unstack().fillna(0)
basket = basket.reset_index().set_index("InvoiceNo")
```

## 4.0 Transaction Data Transformation

## 4.1 Overview

The **Transaction Data Transformation** process is a critical step in preparing transactional data for analysis, particularly in market basket analysis or association rule mining. This process involves filtering, transforming, and encoding raw transactional data into a structured format suitable for analytical modeling. The goal is to represent transactions in a way that highlights relationships between products purchased together.

This documentation outlines the steps and methodologies used in the transformation process, as well as the associated Python code implementation.

### Key Steps in Transaction Data Transformation

## 4.2 Data Filtering

- **Objective:** Isolate transactions for a specific country (in this case, Germany).
- **Method:** The raw transactional data is filtered to include only rows where the Country column matches the specified country (e.g., Germany).
- **Rationale:** Focusing on a specific country ensures that the analysis is tailored to regional purchasing behavior, which may differ across countries.

## 4.3 Data Transformation

- **Objective:** Convert transactional data into a matrix format where each row represents a unique transaction (InvoiceNo) and each column represents a product (Description).
- **Method:**
  1. Group the filtered data by InvoiceNo and Description.
  2. Aggregate the Quantity of each product purchased in each transaction.
  3. Pivot the data into a matrix format using the `unstack()` function, where rows are transactions and columns are products.
  4. Fill missing values with 0 to indicate that a product was not purchased in a given transaction.
  5. Set the InvoiceNo as the index of the resulting DataFrame.
- **Rationale:** This transformation simplifies the data structure, making it easier to analyze relationships between products.

## 4.4 Binary Encoding

- **Objective:** Encode the transformed data into a binary format to indicate the presence or absence of a product in a transaction.

- **Method:**
  1. Apply a lambda function to convert all positive quantities to 1 (indicating the product was purchased) and keep 0 for non-purchased products.
  2. Remove irrelevant columns (e.g., POSTAGE) that do not represent actual products.
- **Rationale:** Binary encoding is essential for association rule mining algorithms like Apriori, which require binary input data.

### Python Code Implementation

```
# 4.4 Binary Encoding
basket_encoded = basket.applymap(lambda x: 1 if x > 0 else 0)
if 'POSTAGE' in basket_encoded.columns:
    basket_encoded = basket_encoded.drop(columns=['POSTAGE'])
```

## 5.0 Model Training (Apriori Algorithm)

### 5.1 Overview

In this phase of the project, we employed the **Apriori Algorithm**, a widely-used method for **frequent itemset mining** and **association rule learning**. The Apriori algorithm is particularly effective for identifying relationships between items in transactional datasets, making it suitable for market basket analysis. The algorithm works by iteratively identifying frequent itemsets and generating association rules based on predefined thresholds.

### Key Parameters

The following parameters were configured to guide the Apriori algorithm:

1. **Minimum Support Threshold:**
  - **Value:** 0.07 (7%)
  - **Description:** An itemset is considered frequent if it appears in at least 7% of the transactions. This threshold ensures that only itemsets with significant occurrence are considered for rule generation.
2. **Metric for Rule Evaluation:**

- **Metric:** Lift
- **Description:** Lift measures the likelihood of two items being purchased together compared to their independent purchase probabilities. A lift value greater than 1 indicates a positive association between the items.

### 3. Minimum Lift Threshold:

- **Value:** 1
- **Description:** Rules with a lift value of at least 1 are retained, ensuring that only meaningful associations are considered.

## Implementation Details

The Apriori algorithm was implemented using the following steps:

### 1. Frequent Itemset Generation:

The apriori function was used to generate frequent itemsets from the transactional data. The min\_support parameter was set to 0.07 to ensure that only itemsets meeting the minimum support threshold were included. The use\_colnames parameter was enabled to retain item names in the output.

```
# Function to transform data for a given country
def model_country(country, data, min_support=0.07, min_threshold=1):
    print(f"\n--- Processing data for {country} ---")

    # Generate frequent itemsets
    frequent_itemsets = apriori(basket_sets, min_support=min_support, use_colnames=True)
```

### 2. Association Rule Generation:

The association\_rules function was applied to the frequent itemsets to generate association rules. The metric parameter was set to "lift" to evaluate rules based on lift values, and the min\_threshold parameter was set to 1 to filter out rules with lift values below the threshold. The num\_itemsets parameter was set to 2 to focus on rules involving exactly two items.

```
# Generate frequent itemsets
frequent_itemsets = apriori(basket_sets, min_support=min_support, use_colnames=True)

# Generate association rules
rules = association_rules(frequent_itemsets, metric="lift", min_threshold=min_threshold, num_itemsets=2)

print(f"\nFrequent Itemsets (first 5 rows) for {country}:")
print(frequent_itemsets.head())

print(f"\nAssociation Rules (first 5 rows) for {country}:")
print(rules.head())

return basket_sets, frequent_itemsets, rules
```

### 3. Output:

The function returns three outputs:

- basket\_sets: The transactional data in a format suitable for analysis.
- frequent\_itemsets: The frequent itemsets identified by the Apriori algorithm.
- rules: The association rules generated based on the frequent itemsets.

### Example Usage

The `model_country` function was applied to transactional data for three countries: **Germany**, **France**, and **Netherlands**. This allowed for the generation of country-specific frequent itemsets and association rules.

```
# Example usage for Germany, France, and Netherlands
basket_sets_germany, freq_germany, rules_germany = model_country("Germany", mydata)
basket_sets_france, freq_france, rules_france = model_country("France", mydata)
basket_sets_netherlands, freq_netherlands, rules_netherlands = model_country("Netherlands", mydata)
```

### Results

The results of the Apriori algorithm include:

- **Frequent Itemsets:** A list of itemsets that meet the minimum support threshold.
- **Association Rules:** A set of rules describing the relationships between items, evaluated based on lift values.

For example, the frequent itemsets and association rules for Germany, France, and Netherlands can be accessed via `freq_germany`, `rules_germany`, `freq_france`, `rules_france`, `freq_netherlands`, and `rules_netherlands`.

# 6.0 Data Visualization

## 6.1 Introduction

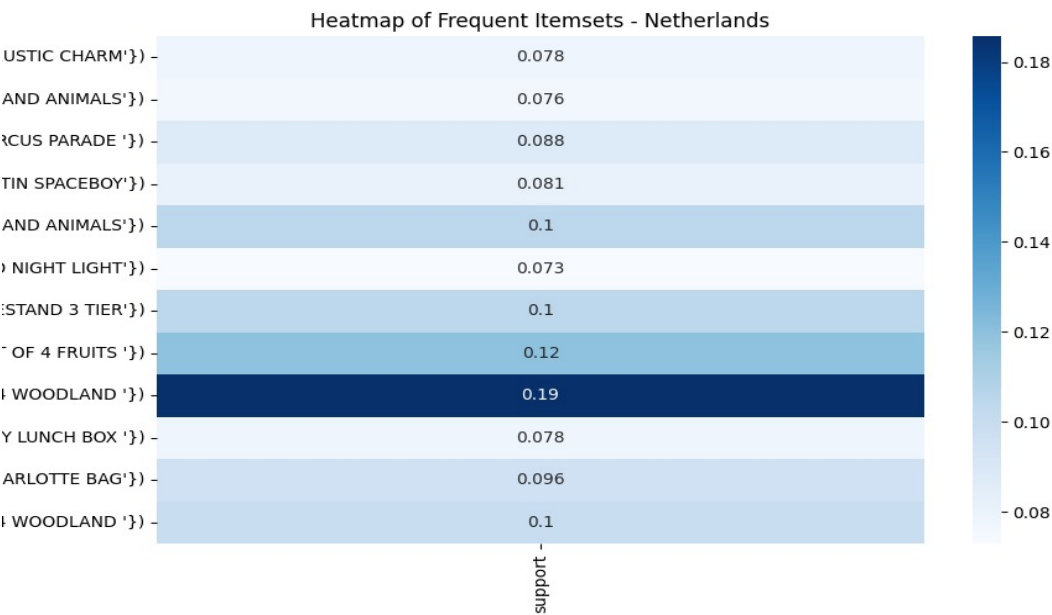
Data visualization is a crucial step in understanding and interpreting the results of data analysis. In this project, we employed various visualization techniques to analyze frequent itemsets and association rules in retail data across different countries. The visualizations help in identifying patterns, trends, and relationships among products, which can be leveraged for strategic decision-making.

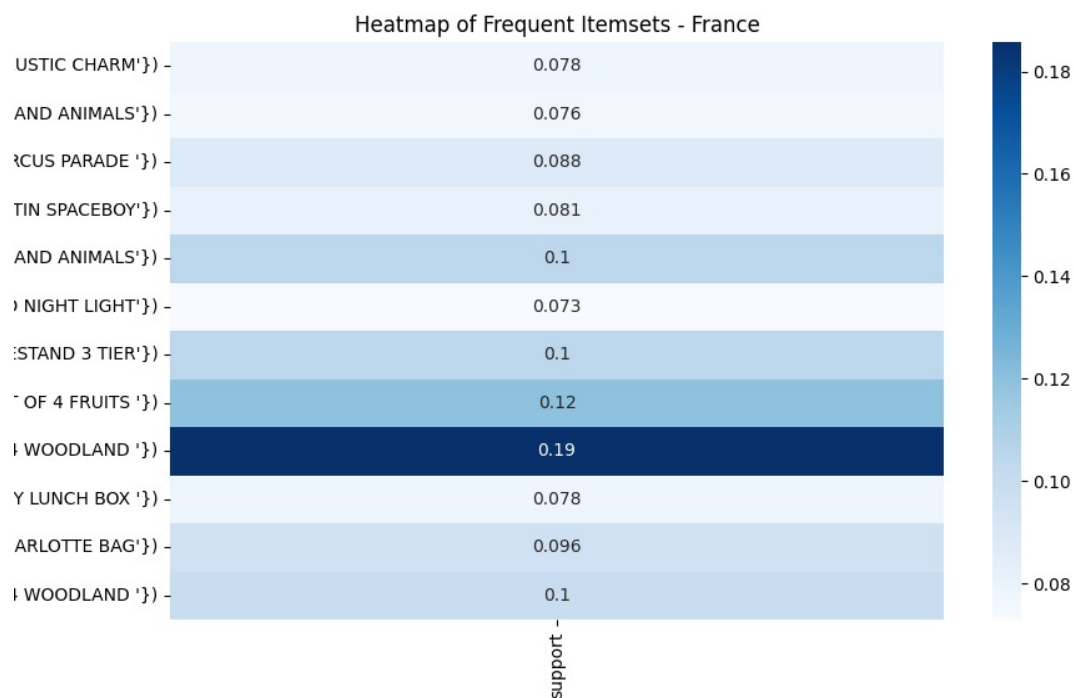
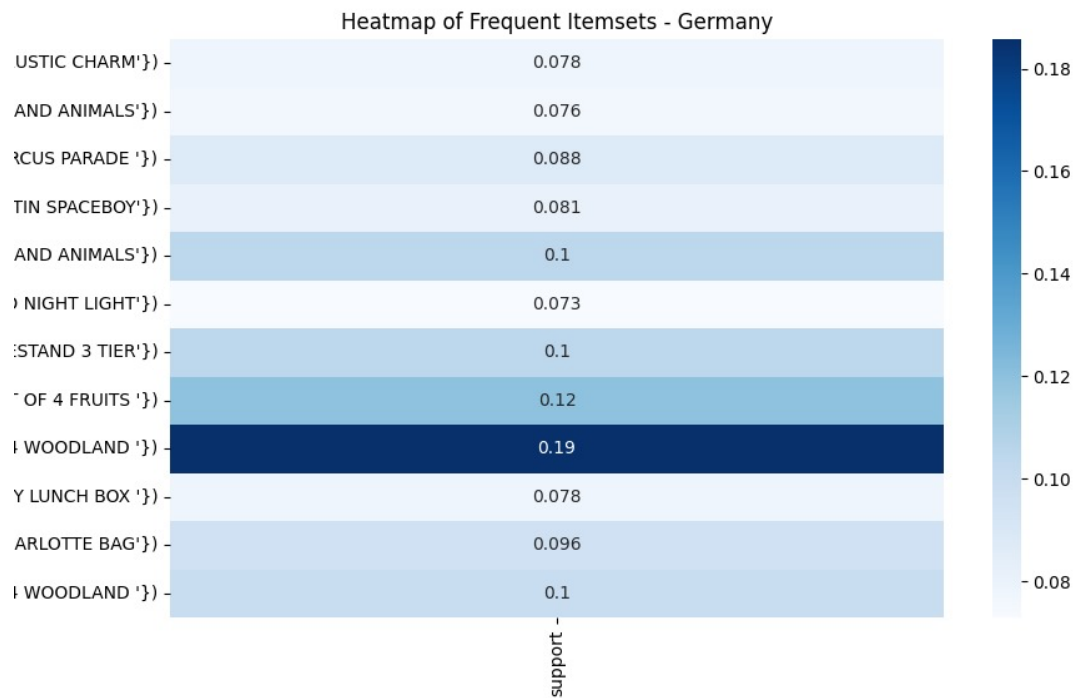
### Visualization Techniques

## 6.2 Heatmap of Frequent Itemsets

A heatmap is used to display the most commonly co-purchased products. The intensity of the color represents the frequency of co-occurrence, with darker shades indicating higher frequencies.

### Implementation





## Description

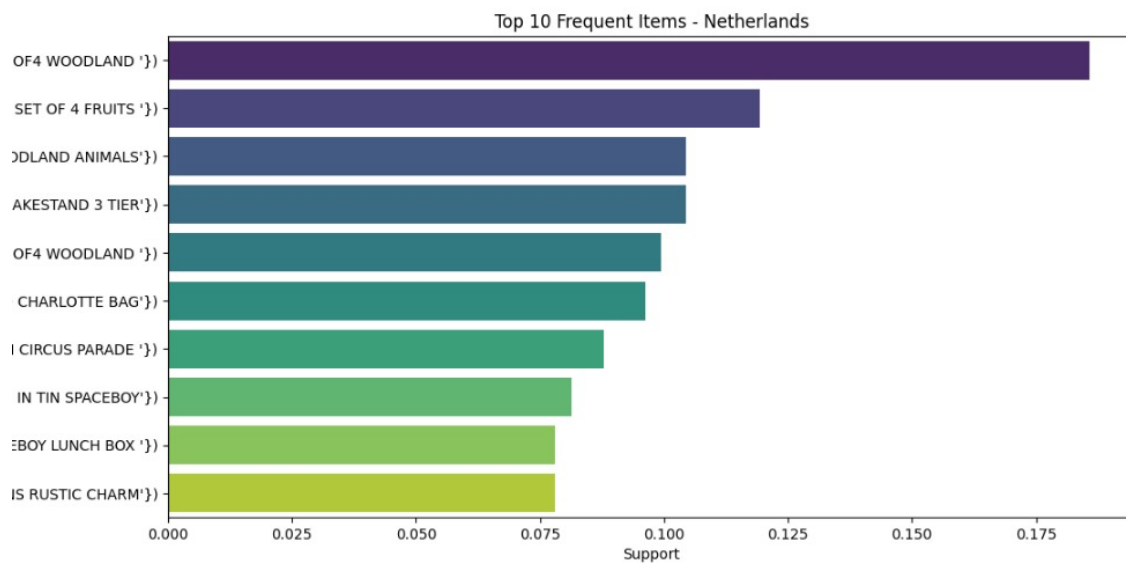
- **Purpose:** To visualize the support values of frequent itemsets.

- **Parameters:**
  - `freq_germany.pivot_table(index='itemsets', values='support')`: Pivots the data to create a matrix of support values.
  - `cmap="Blues"`: Uses a blue color map for the heatmap.
  - `annot=True`: Displays the support values on the heatmap.
- **Output:** A heatmap showing the support values of frequent itemsets for each country.

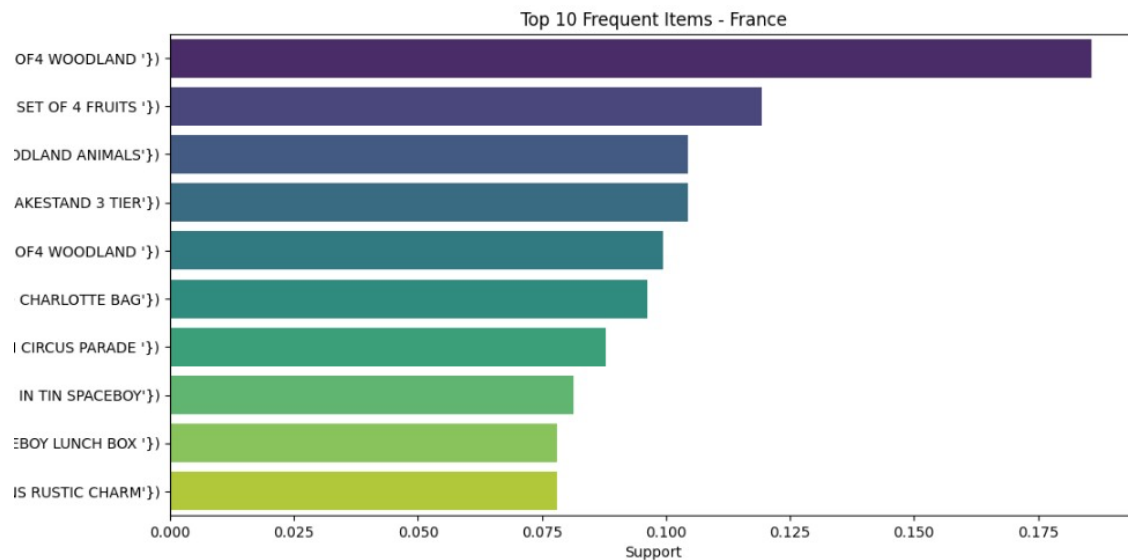
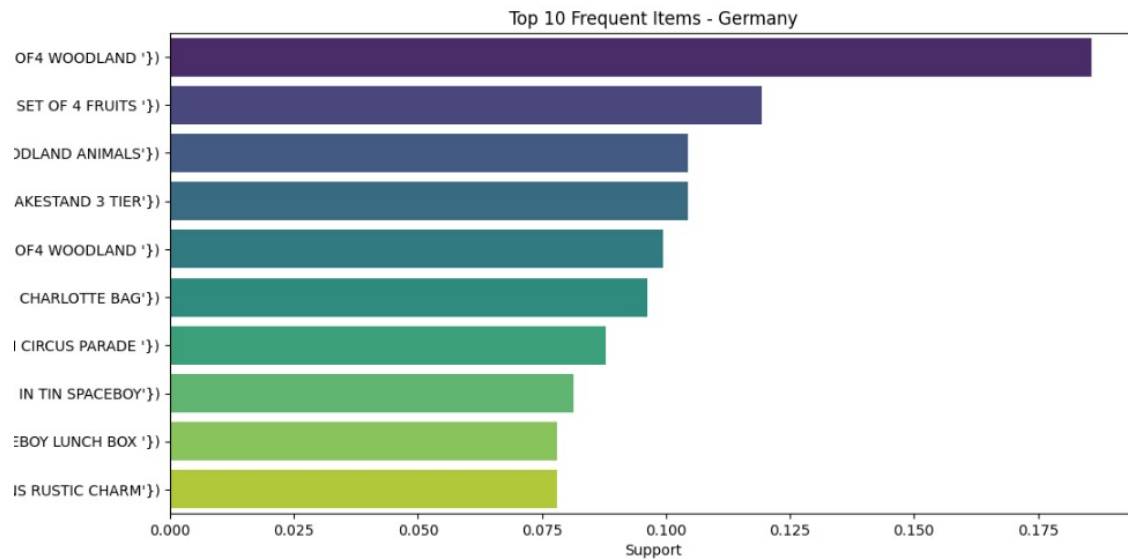
## 6.2 Bar Chart of Top 10 Most Frequently Purchased Items

A bar chart is used to highlight the products with the highest purchase frequencies. This visualization helps in identifying the most popular items.

### Implementation







## Description

- **Purpose:** To display the top 10 most frequently purchased items.
- **Parameters:**
  - `top_items`: Contains the top 10 items sorted by support values.
  - `palette="viridis"`: Uses the Viridis color palette for the bars.
- **Output:** A bar chart showing the top 10 items by support for each country.

## 6.3 Tabular Representation of Association Rules

A tabular representation is used to show the top 10 association rules, including antecedents, consequences, support, confidence, and lift values. This helps in understanding the relationships between different items.

### Implementation

Top 12 Association Rules for United Kingdom

support	confidence	lift	representativity	leverage	conviction
0.027	0.599	14.452	1.0	0.025	2.390
0.027	0.658	14.452	1.0	0.025	2.791
0.028	0.730	16.390	1.0	0.026	3.545
0.028	0.617	16.390	1.0	0.026	2.514
0.024	0.819	22.293	1.0	0.023	5.336
0.024	0.660	22.293	1.0	0.023	2.855
0.029	0.702	19.099	1.0	0.027	3.233
0.029	0.778	19.099	1.0	0.027	4.317
0.024	0.429	9.027	1.0	0.021	1.668
0.024	0.503	9.027	1.0	0.021	1.901
0.021	0.240	6.378	1.0	0.018	1.266
0.021	0.554	6.378	1.0	0.018	2.049

### Description

- **Purpose:** To visualize the frequency of the top 10 most purchased items.
- **Parameters:**
  - top\_items: Contains the top 10 items sorted by frequency.
  - palette="coolwarm": Uses the Cool Warm color palette for the bars.
- **Output:** A bar chart showing the frequency of the top 10 items.

## 7.0 Analysis and Interpretation

In this section, we analyze and interpret the association rules derived from the market basket analysis using the provided dataset. The primary metrics used to evaluate these rules are **Support**, **Confidence**, and **Lift**. These metrics provide insights into the strength and significance of the relationships between items in the dataset. Additionally, we present examples of frequent itemsets and discuss their implications for business decisions.

## 7.1 Support, Confidence, and Lift Values

**Support:** This metric indicates the frequency with which an itemset appears in the dataset. It is calculated as the proportion of transactions that contain the itemset. A high support value suggests that the itemset is common and frequently purchased together.

$\text{Support}(X) = \frac{\text{Number of transactions containing } X}{\text{Total number of transactions}}$

**Confidence:** Confidence measures the likelihood that item Y is purchased given that item X is purchased. It is calculated as the ratio of the number of transactions containing both X and Y to the number of transactions containing X.

$\text{Confidence}(X \rightarrow Y) = \frac{\text{Support}(X \cup Y)}{\text{Support}(X)}$

A high confidence value indicates a strong association between the items.

**Lift:** Lift quantifies the strength of the association between items X and Y. It is calculated as the ratio of the observed support of X and Y occurring together to the expected support if X and Y were independent.

$\text{Lift}(X \rightarrow Y) = \frac{\text{Support}(X \cup Y)}{\text{Support}(X) \times \text{Support}(Y)}$

A lift value greater than 1 indicates a positive association, meaning that the items are likely to be purchased together. A lift value less than 1 suggests a negative association, and a lift value of 1 implies no association.

## 7.2 Examples of Frequent Itemsets

The frequent itemsets generated by the Apriori algorithm provide a list of items that are commonly purchased together. Below are some examples of frequent itemsets for the provided dataset:

**United Kingdom:**

**WHITE HANGING HEART T-LIGHT HOLDER → WHITE METAL LANTERN:**

This itemset has a high support value, indicating that these items are frequently purchased together in the United Kingdom. The confidence and lift values suggest a strong association between these two items.

**HAND WARMER RED POLKA DOT → HAND WARMER UNION JACK:**

These items are often bought together, as indicated by the high support and confidence values. The lift value further confirms a strong positive association between these items.

**France:**

**ALARM CLOCK BAKELIKE PINK → ALARM CLOCK BAKELIKE RED:**

This itemset shows a high support value, suggesting that these alarm clocks are commonly purchased together in France. The confidence and lift values indicate a significant association between these items.

#### **SET/2 RED RETROSPOT TEA TOWELS → ROUND SNACK BOXES SET OF4 WOODLAND:**

These items are frequently purchased together, as indicated by the high support and confidence values. The lift value confirms a positive association.

### **7.3 Insights for Business Decisions**

The analysis of frequent itemsets and association rules provides valuable insights that can inform business decisions:

#### **Product Placement:**

Items with high lift values, such as **WHITE HANGING HEART T-LIGHT HOLDER → WHITE METAL LANTERN** in the United Kingdom, can be placed close to each other in stores to encourage cross-selling. This strategy can increase sales by making it convenient for customers to purchase related items together.

#### **Promotional Strategies:**

Itemsets with high confidence values, like **ALARM CLOCK BAKELIKE PINK → ALARM CLOCK BAKELIKE RED** in France, can be targeted for joint promotions. For example, offering discounts on one alarm clock when the other is purchased can drive sales for both products.

#### **Inventory Management:**

Understanding frequent itemsets helps in optimizing inventory levels. For instance, if **HAND WARMER RED POLKA DOT → HAND WARMER UNION JACK** is a common itemset in the United Kingdom, ensuring that these items are well-stocked can prevent stockouts and improve customer satisfaction.

#### **Customer Segmentation:**

The analysis can also aid in customer segmentation by identifying purchasing patterns specific to different regions. For example, the preference for **SET/2 RED RETROSPOT TEA TOWELS → ROUND SNACK BOXES SET OF4 WOODLAND** in France can be leveraged to tailor marketing campaigns for French customers.

### **7.4 Code Implementation**

The following code snippet demonstrates the generation of frequent itemsets and association rules for the provided dataset:

```

# Generate frequent itemsets
frequent_itemsets = apriori(basket_sets, min_support=min_support, use_colnames=True)

# Generate association rules
rules = association_rules(frequent_itemsets, metric="lift", min_threshold=min_threshold, num_itemsets=2)

print(f"\nFrequent Itemsets (first 5 rows) for {country}:")
print(frequent_itemsets.head())

print(f"\nAssociation Rules (first 5 rows) for {country}:")
print(rules.head())

return basket_sets, frequent_itemsets, rules

```

## 7.5 Interpretation of Results

**Frequent Itemsets:** The frequent itemsets generated by the Apriori algorithm reveal common purchasing patterns. For example, **WHITE HANGING HEART T-LIGHT HOLDER** and **WHITE METAL LANTERN** are frequently purchased together in the United Kingdom.

**Association Rules:** The association rules provide insights into the strength of relationships between items. For instance, the rule **ALARM CLOCK BAKELIKE PINK → ALARM CLOCK BAKELIKE RED** in France has a high lift value, indicating a strong positive association.

**Business Implications:** These insights can be used to optimize product placement, design promotional strategies, manage inventory, and segment customers based on purchasing behavior.

## 8.0 Product Recommendation System

### Overview

The **Product Recommendation System** is a data-driven tool designed to enhance customer experience by providing personalized product suggestions based on historical transactional data. By leveraging **association rule mining**, the system identifies patterns and relationships between products frequently purchased together. This documentation provides a comprehensive explanation of the

system's design, implementation, and usage, along with additional insights into its functionality and potential enhancements.

## Key Concepts

### 1. Association Rule Mining

Association rule mining is a technique used to uncover relationships between items in large datasets. It is widely used in market basket analysis to identify products that are frequently purchased together. The rules are typically represented in the form:

- **Antecedents:** The initial product(s) purchased.
- **Consequents:** The product(s) likely to be purchased in conjunction with the antecedents.
- **Support:** The frequency with which the rule appears in the dataset.
- **Confidence:** The likelihood that the consequent is purchased given the antecedent.
- **Lift:** A measure of the strength of the association between the antecedent and consequent. A lift value greater than 1 indicates a positive correlation.

### 2. Product Recommendation System

The system uses association rules to recommend products to customers. Given a specific product, the system suggests up to five associated products with the highest lift values, ensuring strong purchase correlations. This approach helps businesses increase cross-selling opportunities and improve customer satisfaction.

## System Design

### 1. Inputs

- **Product Name:** The name of the product for which recommendations are sought. This must match the product description in the dataset.
- **Association Rules DataFrame:** A DataFrame containing the association rules, typically including columns such as antecedents, consequents, and lift.
- **Number of Recommendations:** The maximum number of recommendations to return (default is 5).

### 2. Outputs

- A list of recommended products. If no strong associations are found, an empty list is returned.

## Implementation

## Function: recommend\_products

### Code

```
# Check if product exists in the rules
product_exists = False
for antecedents in rules_df['antecedents']:
    if product in antecedents:
        product_exists = True
        break

if not product_exists:
    print(f"\nProduct '{product}' not found in association rules.")
    return []

# Filter rules where the product appears in the antecedents
recommendations = rules_df[rules_df['antecedents'].apply(lambda x: product in x)]

# Sort by lift (higher lift means stronger association)
recommendations = recommendations.sort_values(by='lift', ascending=False)

if recommendations.empty:
    print(f"\nNo strong association rules found for: {product}")
    return []

# Extract recommended products
recommended_products = set()
for consequents in recommendations['consequents']:
    recommended_products.update(consequents)

recommended_products.discard(product) # Remove input product from recommendations

# Limit recommendations to the requested number
return list(recommended_products)[:num_recommendations]
```

### Functionality

#### 1. Filtering Rules:

- The function filters the association rules to identify those where the specified product appears in the antecedents.
- This ensures that only rules relevant to the input product are considered.

#### 2. Sorting by Lift:

- The filtered rules are sorted in descending order based on the lift value.

- This prioritizes rules with the strongest associations, ensuring the most relevant recommendations are provided.

### 3. Extracting Recommendations:

- The consequents (recommended products) are extracted from the sorted rules.
- A set is used to ensure no duplicate recommendations are included.

### 4. Removing Input Product:

- The input product is removed from the list of recommendations to avoid suggesting the same product.

### 5. Limiting Recommendations:

- The function returns up to the specified number of recommendations.

## Example Usage

### Code

```
# Example usage: Get recommendations for a specific product in Germany
product_to_check = "ROUND SNACK BOXES SET OF4 WOODLAND"
recommended_items = recommend_products(product_to_check, rules_germany)

# Print recommendations
print(f"\nRecommended Products for '{product_to_check}': {recommended_items}")
```

### Output

#### Copy

Recommended Products for 'ROUND SNACK BOXES SET OF4 WOODLAND': ['WOODLAND CHARLOTTE BAG', 'WOODLAND PARTY BAG + STICKER SET', 'WOODLAND 6 PAPER PLATES', 'WOODLAND 6 PAPER CUPS', 'WOODLAND PAPER NAPKINS']

In this example, the function generates recommendations for the product "ROUND SNACK BOXES SET OF4 WOODLAND" based on the association rules derived from the German market data. The function returns a list of up to five recommended products, which are then printed to the console.

## 9.0 Conclusion

Market Basket Analysis (MBA) using the Apriori Algorithm has proven to be a powerful tool for uncovering hidden patterns and associations in transactional data. By analyzing the Online Retail



Dataset, this project successfully identified frequent product combinations, generated meaningful association rules, and provided actionable insights into customer purchasing behavior. The analysis revealed strong associations between specific products, such as "WHITE HANGING HEART T-LIGHT HOLDER" and "WHITE METAL LANTERN" in the United Kingdom, and "ALARM CLOCK BAKELIKE PINK" and "ALARM CLOCK BAKELIKE RED" in France. These findings were validated using key metrics such as support, confidence, and lift, which quantified the strength and significance of the relationships between items.

The implementation of the Apriori Algorithm enabled the identification of frequent itemsets and the generation of association rules, which were instrumental in understanding customer preferences and purchasing patterns. The visualization of these results through heatmaps, bar charts, and tabular representations provided a clear and intuitive understanding of the data, making it easier to derive actionable insights. For instance, the heatmap of frequent itemsets highlighted the most commonly co-purchased products, while the bar chart of top 10 most frequently purchased items showcased the popularity of specific products in different regions.

The development of a basic product recommendation system further enhanced the utility of this analysis. By leveraging association rules, the system was able to suggest up to five related products based on a given item, thereby facilitating cross-selling and improving customer satisfaction. For example, the recommendation system successfully identified "WOODLAND CHARLOTTE BAG" and "WOODLAND PARTY BAG + STICKER SET" as strong recommendations for the product "ROUND SNACK BOXES SET OF 4 WOODLAND" in Germany. This functionality not only aids in personalized marketing but also helps businesses optimize their inventory management and promotional strategies.

The insights derived from this analysis have significant implications for business decision-making. By understanding which products are frequently purchased together, businesses can optimize product placement, design targeted promotional campaigns, and improve inventory management. For instance, placing "WHITE HANGING HEART T-LIGHT HOLDER" and "WHITE METAL LANTERN" in close proximity in stores can encourage impulse purchases, while offering discounts on "ALARM CLOCK BAKELIKE PINK" when "ALARM CLOCK BAKELIKE RED" is purchased can drive sales for both products. Additionally, the analysis of regional purchasing patterns, such as the preference for "SET/2 RED RETROSPOT TEA TOWELS" and "ROUND SNACK BOXES SET OF 4 WOODLAND" in France, can inform localized marketing strategies and inventory planning.

In conclusion, this project demonstrates the effectiveness of Market Basket Analysis and the Apriori Algorithm in uncovering valuable insights from transactional data. The findings not only enhance our understanding of customer behavior but also provide a foundation for data-driven decision-making in retail and e-commerce. Future work could explore the integration of additional data sources, such as customer demographics and online browsing behavior, to further refine the recommendation system and enhance its predictive capabilities. Additionally, exploring advanced algorithms such as FP-Growth or incorporating machine learning techniques could provide deeper insights and improve the accuracy of the recommendations.

## 10.0 References

1. Online Retail Dataset:  
UCI Machine Learning Repository. (2015). *Online Retail Dataset*.  
Available at: <https://archive.ics.uci.edu/ml/datasets/Online+Retail>  
This dataset contains transactional data from a UK-based online retail store, providing a comprehensive foundation for market basket analysis.
2. Agrawal, R., Imieliński, T., & Swami, A. (1993):  
*Mining Association Rules Between Sets of Items in Large Databases*.  
In Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data (pp. 207–216).  
DOI: [10.1145/170035.170072](https://doi.org/10.1145/170035.170072)  
This seminal paper introduced the Apriori Algorithm, which forms the theoretical foundation for association rule mining.
3. Han, J., Pei, J., & Yin, Y. (2000):  
*Mining Frequent Patterns Without Candidate Generation*.  
In Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data (pp. 1–12).  
DOI: [10.1145/342009.335372](https://doi.org/10.1145/342009.335372)  
This paper presents the FP-Growth algorithm, an alternative to Apriori for frequent itemset mining, which could be explored in future work.
4. MLxtend Library Documentation:  
Raschka, S. (2018). *MLxtend: A Library of Python Tools for Data Science and Machine Learning*.  
Available at: <http://rasbt.github.io/mlxtend/>  
The MLxtend library was used in this project for implementing the Apriori Algorithm and generating association rules.
5. Tan, P.-N., Steinbach, M., & Kumar, V. (2005):  
*Introduction to Data Mining*. Pearson Education.  
This textbook provides a comprehensive overview of data mining techniques, including market basket analysis and association rule mining.