# NAANMUTHALVAN

# ARTIFICIAL INTELLIGENCE

# PROJECT TITLE

# MARKET BASKET INSIGHTS

# PHASE 4

# Performing association analysis Generating insights

# INTRODUCTION

The purpose of this project is to conduct a comprehensive market basket analysis to gain insights into customer purchasing behavior. Market basket analysis is a data-driven approach that involves examining transaction data to identify patterns and associations between products that customers tend to buy together. By understanding customer purchasing behavior through this analysis, the project aims to achieve the following objectives:

❖ Product Recommendations: To provide personalized product recommendations to customers based on their past purchase history and preferences.

❖ Inventory Management: To optimize inventory management by identifying frequently co-purchased items and ensuring that these items are well-stocked together.

❖ Pricing Strategies: To develop effective pricing strategies, such as bundling or discount offers, that encourage customers to purchase complementary products.

❖ Customer Segmentation: To segment customers based on their purchasing behavior, enabling targeted marketing campaigns and better understanding of different customer segments.

❖ Cross-Selling and Upselling: To identify opportunities for cross-selling and upselling, encouraging customers to explore related or higher-value products.

❖ Enhanced Customer Experience: Personalized product recommendations can significantly improve the customer experience, making it more relevant and tailored to individual preferences.

❖ In summary, market basket analysis is a valuable tool that empowers businesses to better understand their customers, optimize their operations, and ultimately, drive growth and success in today's dynamic and customer-centric marketplace

# 1.DATA COLLECTION

Collecting transaction data is a crucial step in market basket analysis, as the quality and quantity of data directly impact the accuracy and usefulness of the analysis. Here's an overview of the process of collecting transaction data and the importance of quality data and data sources:

## PROCESS OF COLLECTING TRANSACTION DATA:

1. **Identify Data Sources:** The first step is to identify the sources of transaction data. These sources can include point-of-sale (POS) systems, e-commerce platforms, mobile apps, and other customer touchpoints where purchases are recorded.

2. **Data Capture:** Transaction data is captured in real-time or at regular intervals, depending on the nature of the business. This data typically includes information about the items purchased, the transaction time and date, customer details (if available), and transaction value.

3. **Data Storage:** The collected data is stored in a structured format, often in databases or data warehouses, to ensure it's easily accessible and can be analyzed efficiently.

4. **Data Cleaning and Preprocessing:** Data cleaning involves identifying and rectifying errors, inconsistencies, or missing values in the transaction data. This step is essential for maintaining data quality.

5. **Data Integration:** In cases where data is collected from multiple sources or locations, data integration may be necessary to consolidate all the transaction data into a single, comprehensive dataset.

6. **Pattern Recognition:** To identify meaningful patterns and associations in customer purchasing behavior, the data needs to be reliable. Poor-quality data can obscure important relationships between products.

7. **Customer Insights:** High-quality customer data, including demographic and behavioral information, can provide additional context for understanding purchasing behavior. This enriched data can lead to more targeted and effective marketing strategies.

Let's assume we have a fictitious e-commerce database with transaction data, and we want to collect data for analysis using Python:

```python
import pandas as pd

data_file = 'Market_basket_insights.csv'

def capture_transaction_data(file_path):
    try:
        data = pd.read_csv(file_path)
        return data
    except FileNotFoundError:
        print("The data file was not found.")
        return None
    except Exception as e:
        print("An error occurred while reading the data:", str(e))
        return None
```

# 2.DATA  PREPROCESSING

Cleaning and formatting data is a crucial step in the data preprocessing phase. It involves identifying and addressing issues in the dataset to ensure that the data is of high quality and ready for analysis. Here are the typical steps involved in cleaning and formatting data, along with some common challenges and issues encountered during this phase:

## Steps Involved in Cleaning and Formatting Data:

1. **Handling Missing Data:**
   - **Challenge:** Missing data can lead to inaccurate analysis. Dealing with missing values is a critical step. Common methods include imputation (filling in missing values with reasonable estimates) or removing rows or columns with too many missing values.

2. **Removing Duplicates:**

- **Challenge:** Duplicate records can skew analysis and lead to incorrect insights. Identifying and removing duplicates is important to ensure the data is not biased.

3. **Handling Outliers:**
   - **Challenge:** Outliers can significantly impact statistical analysis. Identifying and deciding how to handle outliers (e.g., removing, transforming, or treating them as special cases) is essential.

4. **Data Transformation:**
   - **Challenge:** Data may need to be transformed to meet the assumptions of the analysis methods. This can include standardization, normalization, or log transformation, depending on the data and analysis requirements.

5. **Dealing with Inconsistent Data:**
   - **Challenge:** Inconsistent data formats, naming conventions, or units of measurement can lead to errors in analysis. Standardizing data is necessary for accurate results.

6. **Addressing Data Quality Issues:**
   - **Challenge:** Data quality issues can include incorrect entries, typos, or data entry errors. These need to be corrected or validated.

7. **Handling Categorical Data:**
   - **Challenge:** Categorical data (e.g., product categories, customer segments) may need to be encoded or one-hot encoded to be used in machine learning models.

8. **Data Scaling and Normalization:**
   - **Challenge:** Data with varying scales can affect the performance of some machine learning algorithms. Scaling and normalization may be necessary to ensure that all features have equal weight in the analysis.

9. **Dealing with Data Imbalance:**
   - **Challenge:** In classification tasks, data imbalance can occur when one class significantly outnumbers the others. Techniques like oversampling, undersampling, or using different evaluation metrics need to be employed to address this issue.

Let's assume we have a fictitious e-commerce database with transaction data and data preprocessing :

```
import pandas as pd
data_file = 'Market_based_insights.csv'
DataFrame
def capture_transaction_data(file_path):
    try:
        data = pd.read_csv(file_path)
        return data
    except FileNotFoundError:
        print("The data file was not found.")
        return None
    except Exception as e:
        print("An error occurred while reading the data:", str(e))
        return None

# Data cleaning and preprocessing (simulated)
def clean_and_preprocess(data):
    # Perform data cleaning and preprocessing steps here
    # For example, removing duplicates, handling missing values, etc.
    # You can customize this part based on your data and requirements.
    cleaned_data = data.drop_duplicates()  # Remove duplicates as an example
    return cleaned_data

# Data storage: Save the cleaned data to a new CSV file
def store_data(data, output_file):
    data.to_csv(output_file, index=False)
    print("Cleaned data saved to", output_file)
```

# 3.ASSOCIATION ANALYSIS

## Introduction to Association Analysis:

Association analysis is a data mining technique used to discover relationships, patterns, and associations within large datasets. It is primarily applied in market basket analysis, where the goal is to identify associations among items frequently purchased together. The insights gained from association analysis are valuable for making data-driven decisions in various domains, such as retail, e-commerce, and recommendation systems.

## Apriori Algorithm:

The Apriori algorithm is one of the fundamental and widely used techniques in association analysis. It's designed to find frequent itemsets and generate association rules. Here's an explanation of how the Apriori algorithm works:

1. **Frequent Itemset:** An itemset is considered frequent if it appears in the dataset with a frequency greater than or equal to a predefined minimum support threshold.

2. **Candidate Generation:** The Apriori algorithm employs an iterative approach. It begins with 1-item itemsets, scans the dataset to find frequent itemsets, and then generates candidate itemsets of larger sizes based on the frequent itemsets discovered in the previous step.

3. **Support Calculation:** For each candidate itemset, the algorithm scans the dataset again to calculate the support, which is the proportion of transactions containing the candidate itemset.

4. **Pruning:** Candidate itemsets that do not meet the minimum support threshold are pruned, as they cannot be frequent themselves or contribute to the generation of larger frequent itemsets.

5. **Repeat:** Steps 2-4 are repeated iteratively until no new frequent itemsets can be found.

6. **Generate Association Rules:** Once all frequent itemsets are identified, the Apriori algorithm generates association rules based on these frequent itemsets. Association rules have two parts: an antecedent (the items on the left-hand side) and a consequent (the items on the right-hand side). The strength of the rule is determined by support and confidence.

Let's consider a simple example using grocery store transaction data:

Transaction 1: {bread, milk, eggs}
Transaction 2: {bread, milk}
Transaction 3: {milk, eggs}
Transaction 4: {bread, butter}
Transaction 5: {bread, milk, eggs, butter}

**Iteration 1:**

1. Find 1-item frequent itemsets: {bread}, {milk}, {eggs}, {butter}.

**Iteration 2:**

2. Generate 2-item candidate itemsets: {bread, milk}, {bread, eggs}, {bread, butter}, {milk, eggs}, {milk, butter}, {eggs, butter}.
3. Calculate support for the candidate itemsets:
   - {bread, milk}: 3 (Transaction 1, Transaction 2, Transaction 5)
   - {bread, eggs}: 2 (Transaction 1, Transaction 5)
   - {bread, butter}: 1 (Transaction 4)
   - {milk, eggs}: 2 (Transaction 1, Transaction 3)
   - {milk, butter}: 1 (Transaction 5)
   - {eggs, butter}: 1 (Transaction 5)
4. Prune infrequent itemsets:
   - {bread, butter}, {milk, butter}, and {eggs, butter} are pruned because their support is less than 2.

**Iteration 3:**

5. Generate 3-item candidate itemsets (e.g., {bread, milk, eggs}) based on the frequent 2-item itemsets found in the previous step.
6. Calculate support and prune as before.

# 4.FREQUENT ITEMSETS

The identification of frequent itemsets in a dataset is a critical step in association analysis, often performed using algorithms like Apriori or FP-growth. Frequent itemsets are sets of items that appear together in transactions with a frequency greater than or equal to a predefined minimum support threshold. These frequent itemsets provide valuable insights into patterns, associations, and co-occurrences of items in the dataset.

## Relevant Statistics and Findings:

1. **Frequent Itemset Statistics:**
   - For each frequent itemset, you can record statistics such as support (frequency of occurrence), itemset size (number of items in the set), and itemset itself (the items that constitute the frequent itemset). These statistics provide insights into which itemsets are popular among customers.

2. **Item Co-occurrence Patterns:**
   - Frequent itemsets reveal which items tend to be purchased together. For example, you might find that "bread" and "milk" are frequently bought together, which can be used for product placement or bundling strategies.

3. **Market Basket Insights:**
   - Understanding frequent itemsets enables businesses to make data-driven decisions for product recommendations, cross-selling, upselling, and targeted marketing campaigns. For instance, if "bread" and "butter" frequently appear together, you can create promotions around this pair.

4. **Segmentation Opportunities:**
   - Frequent itemsets can be used to segment customers based on their purchasing patterns. For instance, you may identify one group of customers who buy "bread," "milk," and "eggs" and another group that buys "bread" and "butter." This segmentation can guide personalized marketing strategies.

5. **Inventory Management:**
   - Retailers can optimize inventory by ensuring that frequently co-purchased items are stocked together. This can help reduce stockouts and improve the shopping experience.

6. **Association Rules:**
   - Frequent itemsets serve as the basis for generating association rules, which provide actionable insights. For example, an association rule might reveal that if a customer buys "bread" and "milk," there's a high probability they'll also buy "eggs."
7. **Threshold Sensitivity:**
   - Be aware of how changing the minimum support threshold affects the number and nature of frequent itemsets. A lower threshold yields more frequent itemsets, while a higher threshold yields fewer, more significant associations.

```
from mlxtend.frequent_patterns import apriori

from mlxtend.frequent_patterns import association_rules

import pandas as pd

# Sample transaction dataset

data = {'TransactionID': [1, 2, 3, 4, 5],

     'Items': [['bread', 'milk', 'eggs'],

            ['bread', 'milk'],

            ['milk', 'eggs'],

            ['bread', 'butter'],

            ['bread', 'milk', 'eggs', 'butter']]}

df = pd.DataFrame(data)

# Convert the transaction data into a one-hot encoded format

oht = pd.get_dummies(pd.DataFrame(df['Items'].tolist()), prefix='', prefix_sep='')

# Find frequent itemsets using the Apriori algorithm

min_support = 0.4  # Minimum support threshold

frequent_itemsets = apriori(oht, min_support=min_support, use_colnames=True)

# Print frequent itemsets

print("Frequent Itemsets:")
```

```
print(frequent_itemsets)

# Generate association rules

min_confidence = 0.7  # Minimum confidence threshold

rules = association_rules(frequent_itemsets, metric="confidence",
min_threshold=min_confidence)

# Print association rules

print("\nAssociation Rules:")

print(rules)
```

# 5.GENERATING INSIGHTS

Association rules, generated through techniques like Apriori or FP-growth in association analysis, provide valuable insights by revealing patterns and relationships among items in a dataset. These insights can be used for various purposes. Here's how association rules are used to generate insights:

1. **Product Recommendations:** Association rules can be used to recommend products to customers based on their current or past purchases. For example, if a customer adds "bread" and "milk" to their cart, the association rule might suggest adding "eggs" due to a high confidence that these items are frequently purchased together.

2. **Cross-Selling:** Retailers can use association rules to identify related or complementary products and cross-sell them to customers. If a customer buys a camera, a rule might suggest purchasing memory cards or camera bags.

3. **Upselling:** By identifying items often purchased together, businesses can upsell customers to higher-value or premium versions of products. For example, if a customer selects a basic smartphone, a rule might suggest considering a more advanced model.

4. **Inventory Management:** Retailers can optimize inventory by ensuring that items often purchased together are stocked together. This helps reduce stockouts, improve shelf organization, and enhance the shopping experience.
5. **Pricing Strategies:** Understanding which items are frequently bought together can inform pricing strategies. For instance, offering a discount when customers purchase a bundle of frequently associated products can increase sales.

6. **Customer Segmentation:** Association rules can help segment customers based on their purchasing behavior. For example, one group of customers might frequently buy "coffee" and "milk," while another group buys "tea" and "sugar." This segmentation is valuable for targeted marketing campaigns.

7. **Market Basket Analysis:** Association rules provide a comprehensive view of market basket analysis. Retailers can gain insights into which combinations of items are popular and adjust their strategies accordingly.


**Practical Applications of Association Rule Insights in Retail:**

In a retail setting, the insights generated from association rules have several practical applications:

1. **Personalized Recommendations:** Retailers can create personalized shopping experiences by recommending products based on the customer's current or historical cart contents. This increases the likelihood of additional purchases.
2. **Store Layout and Visual Merchandising:** By knowing which items are often bought together, retailers can optimize store layouts and product placements. Items that frequently co-occur can be placed near each other to encourage cross-sales.
3. **Dynamic Pricing:** Retailers can implement dynamic pricing strategies based on association rules. For example, if "laptop" and "laptop bag" are frequently purchased together, a discount can be offered when both items are added to the cart.
4. **Inventory Management and Stock Control:** Efficient inventory management based on association rules reduces the chances of understocking or overstocking. It ensures that products are available when customers want them.
5. **Targeted Marketing:** Segmenting customers based on their preferences allows for more targeted marketing efforts. Retailers can send tailored promotions, offers, and product recommendations.

## Practical Applications in a Retail Setting :

```python
from mlxtend.frequent_patterns import apriori
from mlxtend.frequent_patterns import association_rules

import pandas as pd

# Sample transaction dataset

data = {'TransactionID': [1, 2, 3, 4, 5],

    'Items': [['bread', 'milk', 'eggs'],

        ['bread', 'milk'],

        ['milk', 'eggs'],

        ['bread', 'butter'],

        ['bread', 'milk', 'eggs', 'butter']]}

df = pd.DataFrame(data)

# Convert the transaction data into a one-hot encoded format

oht = pd.get_dummies(pd.DataFrame(df['Items'].tolist()), prefix='', prefix_sep='')

# Find frequent itemsets using the Apriori algorithm

min_support = 0.4  # Minimum support threshold

frequent_itemsets = apriori(oht, min_support=min_support, use_colnames=True)

# Generate association rules

min_confidence = 0.7  # Minimum confidence threshold

rules = association_rules(frequent_itemsets, metric="confidence",
min_threshold=min_confidence)


# Print association rules

print("Association Rules:")

print(rules)
```

# 6.VISUALIZATION

Creating visual representations of insights is a powerful way to convey complex data patterns and findings in a more understandable and interpretable manner. In the context of association analysis, visualizations can help retail businesses and analysts gain a deeper understanding of the data and the relationships between items. Below, I'll explain how specific types of visualizations can assist in understanding association rule insights and provide examples:

## 1. Support vs. Itemsets (Bar Chart):

- A bar chart can show the support of different itemsets. Each bar represents an itemset, and its height corresponds to the support value.

- This visualization helps in identifying the most and least frequent itemsets.

```python
import matplotlib.pyplot as plt

# Assuming you have a DataFrame of frequent itemsets with 'itemset' and 'support' columns
frequent_itemsets.plot(kind='bar', x='itemset', y='support')
plt.xlabel('Frequent Itemsets')
plt.ylabel('Support')
plt.title('Support of Frequent Itemsets')
plt.show()
```

## 2. Confidence vs. Lift (Scatter Plot):

- A scatter plot can illustrate the relationship between confidence and lift for different association rules. Each point represents an association rule, with confidence on one axis and lift on the other.

- This visualization helps in identifying rules with high confidence and high lift, which are strong associations.

```python
import matplotlib.pyplot as plt
# Assuming you have a DataFrame of association rules with 'confidence' and 'lift' columns
plt.scatter(rules['confidence'], rules['lift'])
```

```
plt.xlabel('Confidence')
plt.ylabel('Lift')
plt.title('Confidence vs. Lift for Association Rules')
plt.grid()
plt.show()
```

## 3. Network Graph (Network Visualization):

- A network graph can represent itemsets and their associations as nodes and edges. Nodes represent items or itemsets, and edges represent the association between them.

- This visualization provides an intuitive view of how items are connected.

```
import networkx as nx
import matplotlib.pyplot as plt

G = nx.Graph()

# Assuming you have a list of itemsets and their associations
for itemset, association in itemset_associations:
    G.add_node(itemset)
    for assoc_itemset in association:
        G.add_edge(itemset, assoc_itemset)

pos = nx.spring_layout(G)
nx.draw(G, pos, with_labels=True, node_size=3000, node_color='lightblue')
plt.title('Itemset Associations Network')
plt.show()
```

# CONCLUSION

Market-based insights play a pivotal role in today's data-driven world, offering invaluable advantages to businesses across various industries. The ability to understand and harness the relationships and patterns hidden within customer purchasing behavior is a powerful tool for making informed decisions, optimizing operations, and enhancing customer experiences.

Through techniques like association analysis, including the Apriori and FP-growth algorithms, businesses can unlock a wealth of knowledge from transaction data. Key takeaways and achievements of market basket analysis projects include the discovery of item associations, the generation of association rules, segmentation opportunities, optimized inventory management, enhanced customer experiences, and informed decision-making. These insights have a direct and profound impact on various facets of a business, from boosting sales and profitability to nurturing customer loyalty and refining marketing strategies.

As technology continues to advance, and as data becomes increasingly accessible, the importance of market basket insights is more pronounced than ever. Retailers, e-commerce platforms, and businesses across multiple sectors rely on these insights to remain competitive, adapt to evolving customer preferences, and deliver products and services that resonate with their target audiences.