

## **DATA WARE HOUSE AND DATA MINING**

### **1. What can be extracted using market basket analysis? Who uses information extracted from market basket analysis and what for they use.**

Market basket analysis aims to identify associations or relationships between items that are frequently purchased together. By analyzing customer transaction data, market basket analysis can extract the following information:

1. Association Rules: Market basket analysis can uncover association rules that indicate the likelihood of certain items being purchased together. For example, it can identify that customers who buy bread are also likely to buy butter.
2. Support: Support measures the frequency of a particular item or itemset occurring together in transactions. It helps determine the popularity or occurrence of associations.
3. Confidence: Confidence measures the reliability or strength of the association rule. It represents the conditional probability that an item will be purchased given that another item has already been purchased. High-confidence rules indicate strong associations.
4. Lift: Lift measures the strength of the association between items and represents how much more likely an item is to be purchased when another item is already in the basket. Lift values greater than 1 indicate a positive association, while values less than 1 indicate a negative association.

The information extracted from market basket analysis is utilized by various stakeholders for different purposes:

1. Retailers: Retailers use market basket analysis to improve their marketing strategies, optimize product placement, and enhance cross-selling opportunities. By understanding which items are frequently purchased together, retailers can optimize store layouts, design effective promotional campaigns, and develop personalized recommendations for customers.
2. Manufacturers and Suppliers: Manufacturers and suppliers can utilize market basket analysis to identify complementary products or optimize their product offerings. This information helps them make informed decisions regarding product bundling, inventory management, and supply chain optimization.
3. E-commerce Companies: Online retailers leverage market basket analysis to enhance their recommendation systems and provide personalized product suggestions to customers. By

understanding customer purchase patterns, e-commerce companies can improve customer engagement and increase sales.

4. Marketers: Market basket analysis provides valuable insights for marketers to design targeted cross-selling and upselling campaigns. By identifying associations between products, marketers can develop customized promotions and offers based on customers' purchasing behavior.

5. Business Analysts: Business analysts use market basket analysis to gain insights into customer behavior, identify trends, and make data-driven business decisions. It helps them understand customer preferences, optimize pricing strategies, and develop effective sales and marketing plans.

Overall, market basket analysis enables businesses to enhance customer satisfaction, increase sales, and improve operational efficiency by understanding and leveraging associations between products or services.

## **2. What are two measures of rule interestingness and explain the role of these measures in mining association rules?**

Two commonly used measures of rule interestingness in association rule mining are support and confidence. These measures play a crucial role in determining the significance and usefulness of association rules:

**Support:** Support is a measure of the frequency or occurrence of an itemset in a dataset. It quantifies how often a particular itemset appears in the transactions being analyzed. Support is calculated by dividing the number of transactions containing the itemset by the total number of transactions in the dataset.

The role of support in mining association rules is to identify itemsets that occur frequently enough to be considered significant. Association rules with high support indicate that the rule's antecedent and consequent items co-occur together frequently. Higher support values indicate stronger associations between items.

**Confidence:** Confidence measures the reliability or strength of an association rule. It quantifies the conditional probability that an itemset or a set of items will be present in a transaction given the presence of another itemset or set of items. Confidence is calculated by dividing the number of transactions containing both the antecedent and consequent by the number of transactions containing only the antecedent.

The role of confidence in mining association rules is to determine the strength of the relationship between items. Rules with high confidence indicate a strong association between the antecedent and consequent. A high-confidence rule suggests that when the antecedent is present in a transaction, there is a high probability that the consequent will also be present.

I apologize for any confusion caused by the numbering issue. Please let me know if you have any further questions!

**3. Explain how closed frequent item set and maximal frequent items set concepts are used, with example.**

In frequent itemset mining, two important concepts are closed frequent itemsets and maximal frequent itemsets. These concepts help in summarizing and reducing the number of discovered frequent itemsets while preserving their interestingness. Here's an explanation of each concept with an example:

**1. Closed Frequent Itemsets:**

A closed frequent itemset is an itemset that does not have any supersets with the same support. In other words, if an itemset is closed, it means that no other itemset contains the same items with the same or higher support. Closed frequent itemsets provide a concise representation of frequent itemsets by eliminating redundant itemsets that provide no additional information.

For example, let's consider a transaction dataset with the following transactions:

Transaction 1: {A, B, C}

Transaction 2: {A, B, D}

Transaction 3: {B, C}

Transaction 4: {A, B, C, D}

Let's assume the minimum support threshold is 2 (i.e., an itemset should appear in at least 2 transactions to be considered frequent). Based on this threshold, we can identify the following frequent itemsets:

Frequent itemsets:

{A}: 3

{B}: 4

{C}: 2

{D}: 2

{A, B}: 3

{A, C}: 2

{B, C}: 2

{A, B, C}: 2

{A, D}: 2

{B, D}: 2

{A, B, D}: 2

From the above list, we can observe that {A, C} and {B, C} have the same support as {A, B, C}. However, {A, B, C} is considered a closed frequent itemset because there is no other superset with the same support ({A, B, C, D} has a higher support). Thus, {A, C} and {B, C} are redundant and can be eliminated.

## 2. Maximal Frequent Itemsets:

A maximal frequent itemset is an itemset that is frequent but does not have any proper supersets that are also frequent. In other words, a maximal frequent itemset cannot be extended further while maintaining the required minimum support. Maximal frequent itemsets provide a higher level of abstraction by eliminating subsets of frequent itemsets that do not add any additional insights.

Continuing with the previous example, we can identify the maximal frequent itemsets as follows:

Maximal frequent itemsets:

{A, B}: 3

{A, B, C}: 2

{A, B, D}: 2

In this case, {A, B} is a maximal frequent itemset because it does not have any proper superset that is also frequent. {A, B, C} and {A, B, D} are not maximal frequent itemsets because they have proper supersets ({A, B, C, D}) that are also frequent.

By considering closed frequent itemsets and maximal frequent itemsets, we can summarize the frequent itemsets and reduce the number of itemsets for further analysis. This helps in simplifying the interpretation and understanding of the discovered associations in large datasets.

## **Q.no 4::**

Apologies for the confusion in my previous response. While Apriori-based methods and pattern-growth methods are two commonly used approaches for frequent pattern mining, there are other methods as well. Here are a few additional methods used in frequent pattern mining:

3. Closed Pattern Mining: This method focuses on discovering closed frequent patterns, which are patterns that have no super-pattern with the same frequency. It aims to reduce redundancy in the discovered patterns.

4. Maximal Pattern Mining: Similar to closed pattern mining, this method aims to find maximal frequent patterns, which are patterns that cannot be extended to form larger frequent patterns. It helps in reducing the number of patterns generated and focuses on the most significant ones.

5. Constraint-based Pattern Mining: This method incorporates constraints or user-defined measures to guide the mining process. It allows users to specify interestingness measures or constraints that the patterns should satisfy, enabling more customized pattern discovery.

6. Pattern Mining with Constraints: This method involves incorporating constraints on various aspects, such as item constraints, sequence constraints, or structural constraints, to guide the pattern mining process. It helps in discovering patterns that meet specific requirements or have certain characteristics.

These are just a few examples of additional methods used in frequent pattern mining. Different methods may be suitable depending on the nature of the data, the specific goals of the mining task, and the desired patterns to be discovered.

## Q.no 5:

Certainly! Here's a simplified version of the Apriori algorithm for discovering frequent itemsets for mining boolean association rules:

1. Initialize a list of candidate itemsets with single items from the dataset.
2. Scan the dataset to count the occurrences of each candidate itemset.
3. Prune the candidate itemsets that do not meet the minimum support threshold.
4. Generate new candidate itemsets by combining the frequent (k-1)-itemsets.
  - Join Step: Take the frequent (k-1)-itemsets and generate new candidate k-itemsets by joining them.
  - Prune Step: Eliminate candidate k-itemsets that contain subsets not present in the frequent (k-1)-itemsets.
5. Repeat steps 2 to 4 until no new frequent itemsets can be generated.
6. Extract the frequent itemsets that satisfy the minimum support threshold.

Here's a Python-like pseudocode representation of the algorithm:

```
...
def apriori(dataset, min_support):
    # Initialize frequent itemsets
    frequent_itemsets = []

    # Generate frequent 1-itemsets
    frequent_1_itemsets = find_frequent_1_itemsets(dataset, min_support)
    frequent_itemsets.append(frequent_1_itemsets)

    k = 2
    while frequent_itemsets[k-2]:
        # Generate candidate itemsets
        candidate_itemsets = generate_candidate_itemsets(frequent_itemsets[k-2])

        # Count the occurrences of candidate itemsets
        frequent_k_itemsets = find_frequent_k_itemsets(candidate_itemsets, dataset, min_support)

        # Append frequent k-itemsets to the list
        frequent_itemsets.append(frequent_k_itemsets)

        k += 1

    return frequent_itemsets

def find_frequent_1_itemsets(dataset, min_support):
    frequent_1_itemsets = []
    item_counts = {}

    # Count the occurrences of each item
    for transaction in dataset:
        for item in transaction:
            item_counts[item] = item_counts.get(item, 0) + 1

    # Prune infrequent items
    for item, count in item_counts.items():
        if count >= min_support:
            frequent_1_itemsets.append({item})

    return frequent_1_itemsets

def generate_candidate_itemsets(frequent_itemsets):
    candidate_itemsets = []

    for i in range(len(frequent_itemsets)):
        for j in range(i+1, len(frequent_itemsets)):
            itemset1 = frequent_itemsets[i]
            itemset2 = frequent_itemsets[j]
```

```

    # Join Step
    new_itemset = itemset1.union(itemset2)

    # Prune Step
    if len(new_itemset) == len(itemset1) + 1:
        candidate_itemsets.append(new_itemset)

return candidate_itemsets

def find_frequent_k_itemsets(candidate_itemsets, dataset, min_support):
    frequent_k_itemsets = []
    item_counts = {}

    # Count the occurrences of candidate itemsets in the dataset
    for transaction in dataset:
        for itemset in candidate_itemsets:
            if itemset.issubset(transaction):
                item_counts[itemset] = item_counts.get(itemset, 0) + 1

    # Prune infrequent itemsets
    for itemset, count in item_counts.items():
        if count < min_support:
            frequent_k_itemsets.append(itemset)

    return frequent_k_itemsets

# Usage example
dataset = [
    {'A', 'B', 'C'},
    {'A', 'B', 'D'},
    {'B', 'C', 'D'},
    {'A', 'C', 'D'},
    {'A', 'B', 'C', 'D'},
]

min_support = 3

frequent_itemsets = apriori(dataset, min_support)

# Print the frequent itemsets
for k, itemsets in enumerate(frequent_itemsets):
    print(f"Frequent {k+1}-itemsets:")
    for itemset in itemsets:
        print(itemset)
...

```

This implementation assumes that the dataset is represented as a list of sets, where each set represents a transaction and contains the items present in that transaction. The `min\_support` parameter specifies the minimum support threshold for an itemset to be considered frequent.

Please note that this is a simplified version of the Apriori algorithm, and there are more efficient optimizations and variations available.

## **Q.no 6**

Association rules are generated from frequent itemsets to discover interesting relationships or patterns in the data. The rules express associations between items or itemsets in a transactional dataset. They are typically represented in the form of "if-then" statements, where the "if" part represents the antecedent (premise) and the "then" part represents the consequent (conclusion).

To generate association rules from frequent itemsets, the following steps are typically followed:

1. Generate Frequent Itemsets: Use a frequent itemset mining algorithm (such as the Apriori algorithm) to discover all frequent itemsets in the dataset based on a specified minimum support threshold.

2. Generate Candidate Rules: For each frequent itemset, generate candidate association rules by considering different combinations of items or itemsets within the frequent itemset. The candidate rules will have one itemset as the antecedent and the remaining items as the consequent.

3. Calculate Rule Measures: For each candidate rule, calculate various measures that assess the interestingness or significance of the rule. Common measures include support, confidence, lift, and conviction. These measures quantify the statistical relationships between items and help evaluate the strength of the association rule.

- Support: The proportion of transactions in the dataset that contain both the antecedent and the consequent of the rule.

- Confidence: The conditional probability that the consequent occurs in a transaction given that the antecedent is present.

- Lift: The ratio of the observed support for the rule to the expected support if the antecedent and the consequent were independent.

- Conviction: Measures the degree of implication of the consequent by the antecedent, considering the independence between them.

4. Prune Rules: Apply a minimum support threshold, minimum confidence threshold, or other interestingness measures to prune or filter out less significant rules. This step helps focus on the most relevant and interesting association rules.

5. Select High-Quality Rules: Select the association rules that meet the desired threshold or criteria for support, confidence, lift, or other interestingness measures. These rules represent the discovered associations between items or itemsets in the dataset.

Here's an example to illustrate the process:

Suppose we have a transactional dataset containing information about customer purchases at a grocery store. Let's assume we have already discovered the frequent itemsets using a frequent itemset mining algorithm:

Frequent Itemsets:

{Milk}, {Bread}, {Butter}, {Milk, Bread}, {Milk, Butter}

Now, let's generate association rules:

1. Generate Candidate Rules:

- {Milk} => {Bread}
- {Milk} => {Butter}
- {Bread} => {Milk}
- {Butter} => {Milk}

2. Calculate Rule Measures:

For each candidate rule, calculate measures such as support, confidence, lift, and conviction. Let's say we calculate the following measures:

- {Milk} => {Bread}: support=0.3, confidence=0.6, lift=1.5, conviction=1.25
- {Milk} => {Butter}: support=0.4, confidence=0.8, lift=1.33, conviction=1.67
- {Bread} => {Milk}: support=0.3, confidence=0.75, lift=1.5, conviction=1.67
- {Butter} => {Milk}: support=0.4, confidence=0.67, lift=1.33, conviction=1.25

3. Prune Rules:

Apply a minimum support threshold and minimum confidence threshold to filter out less significant rules. Let's assume we set a minimum support of 0.3 and a minimum confidence of 0.6. Based on these thresholds, we keep the following rules:

- {Milk} => {Bread}: support=0.3, confidence=0.6, lift=1.5, conviction=1.25
- {Milk} => {Butter}: support=0.4, confidence=0.8, lift=1.33, conviction=1.67

4. Select High-Quality Rules:

The remaining rules meet the specified thresholds and are considered high-quality association rules. These rules suggest that customers who buy milk are likely to buy bread and butter as well, based on the support, confidence, lift, and conviction values.

Association Rules:

- {Milk} => {Bread}: support=0.3, confidence=0.6, lift=1.5, conviction=1.25
- {Milk} => {Butter}: support=0.4, confidence=0.8, lift=1.33, conviction=1.67



These association rules provide insights into the purchasing behavior of customers at the grocery store, helping businesses understand and potentially optimize product placement, promotions, or recommendations.

Note that in practice, there can be additional steps or variations in the process of generating association rules, such as considering higher-order itemsets, applying multiple thresholds, or using different interestingness measures depending on the specific requirements and objectives of the analysis.

### **Q.no 7**

Mining frequent itemsets without candidate generation is achieved through an algorithm called FP-Growth (Frequent Pattern Growth). The FP-Growth algorithm avoids the costly process of generating and testing candidate itemsets by utilizing a tree-based structure called the FP-Tree. This approach allows for efficient and scalable mining of frequent itemsets. Here's an explanation of how FP-Growth works, along with an example:

#### **1. Construct the FP-Tree:**

- Scan the dataset to collect the frequency of each item and sort them in descending order.
- Build the FP-Tree structure based on the sorted items. Each item represents a node in the tree, and the frequency of the item determines the count associated with the node.
- For each transaction, insert the items into the FP-Tree, ensuring that the tree branches represent the order of items in the transaction.

#### **2. Mine Frequent Itemsets from the FP-Tree:**

- Starting with the infrequent items, traverse the FP-Tree in a bottom-up manner.
- For each item, generate conditional pattern bases by tracing its paths in the tree. A conditional pattern base consists of all the paths leading to the item in question.
- From each conditional pattern base, extract the frequent itemsets by considering the combinations of the item and its ancestor nodes. These combinations form frequent itemsets.
- Recursively repeat the process on the conditional pattern bases to mine all frequent itemsets.

Let's consider an example with a transactional dataset:

Dataset:

Transaction 1: {A, B, D}

Transaction 2: {B, C, E}

Transaction 3: {A, B, C, E}

Transaction 4: {B, E}

1. Construct the FP-Tree:

- Scan the dataset and count the frequency of each item:

- A: 2

- B: 4

- C: 2

- D: 1

- E: 3

- Sort the items in descending order based on frequency:  $B > E > A > C > D$ .

- Build the FP-Tree based on the sorted items:

...

B(4)

/ \

E(3) A(2)

/ / \

A(1) C(2) D(1)

...

2. Mine Frequent Itemsets from the FP-Tree:

- Starting with the infrequent item D, trace its paths in the tree to obtain conditional pattern bases. In this case, there is only one path:  $B \rightarrow A \rightarrow D$ .

- From the conditional pattern base, generate the frequent itemset: {B, A, D}.

- Recursively repeat the process on the conditional pattern base of B:

- Generate the frequent itemset: {B, A}.

- For the conditional pattern base of A, generate the frequent itemset: {A}.

- Move to the next infrequent item C:

- Trace its paths in the tree to obtain the conditional pattern base:  $B \rightarrow A \rightarrow C$ ,  $B \rightarrow C$ .

- Generate the frequent itemsets: {B, A, C} and {B, C}.

- Finally, move to the infrequent item E:

- Trace its paths in the tree to obtain the conditional pattern bases:  $B$ ,  $B \rightarrow C$ .

- Generate the frequent itemsets: {B, E}, {B, C, E}, and {E}.

The frequent itemsets mined from the FP-Tree are:

- {B, A, D}
- {B, A}
- {A}
- {B, A, C}
- {B, C}
- {B, E}
- {B, C, E}
- {E}

The FP-Growth algorithm effectively mines frequent itemsets without the need for candidate generation, making it more efficient for large datasets compared to traditional Apriori-based algorithms. It takes advantage of the compact representation of frequent itemsets in the FP-Tree structure, which allows for faster and more scalable mining.

#### **Q.no 8**

Constraint-based mining in data mining refers to incorporating user-defined constraints or interestingness measures during the mining process. These constraints help guide the mining algorithm to discover patterns that satisfy specific requirements or possess certain characteristics. Here are some commonly used constraints in constraint-based mining:

1. Minimum/Maximum Support Constraint: This constraint sets a threshold for the minimum or maximum support of itemsets or rules. It ensures that only itemsets or rules with sufficient or limited support are considered. It helps filter out patterns that are either too common or too rare.
2. Minimum/Maximum Confidence Constraint: This constraint specifies a threshold for the minimum or maximum confidence of association rules. It ensures that only rules with a certain level of confidence (probability) are considered. It helps filter out rules that are either too weak or too strong.
3. Minimum/Maximum Lift Constraint: This constraint sets a threshold for the minimum or maximum lift of association rules. Lift measures the ratio of observed support to expected support, indicating the strength of the association. Applying this constraint helps identify rules with a desired level of impact or dependency between items.

4. Minimum/Maximum Conviction Constraint: Conviction is a measure that quantifies the degree of implication of the consequent by the antecedent, considering the independence between them. This constraint sets a threshold for the minimum or maximum conviction of association rules. It helps filter out rules that have weak or strong implication.

5. Length Constraint: This constraint limits the length of the discovered patterns or rules. It specifies a minimum or maximum number of items or itemsets in the pattern or rule. This constraint helps focus on patterns of a specific length, allowing users to explore different levels of granularity in the discovered patterns.

6. Item Constraints: Item constraints restrict the presence or absence of specific items or itemsets in the discovered patterns or rules. These constraints enforce specific item relationships or exclusions. For example, specifying that item A must be present or item B must not be present in the pattern.

7. Sequence Constraints: Sequence constraints are specific to sequential pattern mining. They define constraints on the order of items in sequential patterns. For example, enforcing that item A must precede item B in the sequence.

8. Structural Constraints: Structural constraints apply to patterns or rules with a specific structural property. For example, constraints may be imposed to discover only closed patterns, maximal patterns, or patterns with a specific hierarchical or nested structure.

By incorporating these constraints, constraint-based mining allows users to customize the mining process and focus on patterns or rules that align with their specific interests, requirements, or domain knowledge. These constraints help refine the mining results and uncover patterns of higher relevance and significance.

### **Q.no 9**

The Apriori algorithm is a classic approach for mining frequent itemsets in a transactional dataset. It works by iteratively generating candidate itemsets and counting their support to identify the frequent itemsets. Here's a step-by-step explanation of how frequent itemsets are mined using the Apriori algorithm:

1. Initialization:

- Scan the dataset to collect the support count for each individual item.

- Set a minimum support threshold, which determines the minimum frequency required for an itemset to be considered frequent.

2. Generating Frequent 1-Itemsets:

- Identify the frequent 1-itemsets by selecting the items with support counts above the minimum support threshold.

3. Generating Candidate Itemsets:

- Generate candidate itemsets of size  $k$  ( $k > 1$ ) by joining the frequent  $(k-1)$ -itemsets.
- To avoid generating duplicate candidates, ensure that all  $(k-1)$ -item subsets of a candidate itemset are themselves frequent.

4. Pruning Infrequent Itemsets:

- Scan the dataset to count the support of each candidate itemset.
- Prune the candidate itemsets that have support below the minimum support threshold.
- The remaining candidate itemsets are considered frequent  $k$ -itemsets.

5. Iteratively Generating Frequent Itemsets:

- Repeat steps 3 and 4 until no more frequent itemsets can be generated.
- Increase the value of  $k$  with each iteration to generate larger itemsets.

6. Output:

- Collect all the frequent itemsets obtained from each iteration.
- The frequent itemsets represent the sets of items that satisfy the minimum support threshold.

The Apriori algorithm leverages the "apriori" property, which states that if an itemset is infrequent, all its supersets must also be infrequent. By exploiting this property, the algorithm reduces the search space by eliminating the need to count the support for all possible itemsets. Instead, it focuses on generating and counting only the candidate itemsets that have the potential to be frequent.

By following these steps, the Apriori algorithm efficiently discovers the frequent itemsets in a transactional dataset, which can then be further utilized for association rule mining or other data analysis tasks.

