# Apriori Algorithm

## Techniques for DM Tasks

**Prediction Methods**

**Description Methods**

**Classification**

**Statistical**

**Clustering**

**Association Rules**

- ✓ Regression
- ✓ Naïve Bayes

- ✓ K-Mean

- ✓ Apriori Algorithm
- ✓ FP-Tree

- ✓ Decision Tree (ID3)
- ✓ Neural Net (Back propagation)
- ✓ Rule Base (ILA)
- ✓ Distance (K-nearset neighbor (K-NN))
- ✓ Support Vector Machines (SVM)

# Reference

- **Also called:**
  - ✓ **Finding Frequent Item-set Using Candidate Generation**
  - ✓ **Basket Problem (Data)**

Which items are frequently purchased together by my customers?

Shopping Baskets

Customer 1: milk, bread, cereal

Customer 2: milk, sugar, bread, eggs

Customer 3: milk, bread, butter

Customer n: sugar, eggs

Market Analyst

5

$I = \{i_1, i_2, \ldots, i_m\}$     *set of Items.*

$D = \{T_1, T_2, \ldots, T_i\}$    *set of Transactions.*

$T = \{i_k, i_j, \ldots\}$  , *where* $T \subseteq I$.

**TID** : is a unique identifier associated with each **T**

> **|D|** is # of transactions in D

- **Association Rule** of the form:    $X \Longrightarrow Y$

                *, where*

$$X \subset I$$
$$Y \subset I$$
$$X \cap Y = \varnothing$$

- **Given the Association Rule AR:  X $\Longrightarrow$ Y**

  ✓ **Support (S) of the AR :**

  $$s = \frac{Support\_count(X \cup Y)}{|D|}$$

  ✓ **Confidence (C)  of the AR :**

  $$C = \frac{Support\_count(X \cup Y)}{Support\_count(X)}$$

  - **min-conf  &  min-sup :** are user thresholds (i.e. KB)

- The **AR**:     $X \implies Y$  is generated if its:

  ✓    Support ($S$) ≥ **min_sup** threshold

  ✓    Confidence ($C$) ≥ **min_conf** threshold

- Apriori uses a "**bottom-up**" approach, where frequent subsets are extended one item at a time (a step known as **candidate generation**).

- The algorithm **terminates** when no further successful extensions are found.

- Apriori uses **breadth-first search** and a **hash tree structure** to **count candidate item sets efficiently**.

- The Apriori algorithm is based on the following two main steps:

  1. **Join** ⋈
  2. **Prune step**

# 1. The Join (⋈) step

- **Given two lists, the operation $L1 ⋈ L2$ will be executed if the following two conditions are satisfied :**

  ➤ **All the elements of the two lists are equal except the last elements**

  ➤ **(Last element of $L1$) < (Last element of $L2$)**

✓ **For example :** $L1 = \{1, 2, 3\}$

$L2 = \{1, 2, 4\}$

So, $L1 ⋈ L2 = \{1, 2, 3, 4\}$

# Example 1 of Join (⋈) step

- **Suppose we have the following set of item Lists**

**L**

| |
|---|
| { I1 , I2 } |
| { I1, I3 } |
| { I2, I3 } |

**L ⋈ L** →

| | Result |
|---|---|
| {I1 , I2} ⋈ {I1, I3} | = {I1 , I2, I3} |
| {I1 , I2} ⋈ {I2, I3} | no join |
| {I1 , I3} ⋈ {I2, I3} | no join |

# Example 2 of Join (⋈) step

- **Suppose** **L1** = **{I1}** and
  **L2** = **{I2}**

- ➢ **In this case we will consider**
  **L1** = **{∅ , I1}** and
  **L2** = **{∅ , I2}**

- ➢ So, **L1⋈L2** = **{I1, I2}**

- Let $L_k$ is a list of items

- $L_k \bowtie L_k = C_{k+1}$

- If $\{I_i, I_j, \ldots, I_r\}$ of $L_k \bowtie L_k \notin L_k$ Then this set of elements will be removed

✓ **Example** :

$L_2$

| |
|---|
| { I1 , I2 } |
| { I1 , I3 } |

$L_2 \bowtie L_2$

$C_3$

$= \{I1, I2, I3\}$

$\{I1, I2\} \in L_2$

$\{I1, I3\} \in L_2$

$\{I2, I3\} \notin L_2$

# Notation

$L_{k+1}$ : is generated from $C_{k+1}$ after pruning and all itemes are >= min-sup



$L_1$

| |
|---|
| {I1} |
| {I2} |
| {I3} |

$L_1 \bowtie L_1$

$C_2$

| |
|---|
| {I1, I2} |
| {I1, I3} |
| {I2, I3} |

$L_2$

| |
|---|
| {I1, I2} |
| {I1, I3} |

Suppose this itemset has sup_count=1 and the min_sup=2, so it will be removed

# of elements in $C_2 = L_1 \bowtie L_1 = \binom{|L_1|}{2}$

where,

$$\binom{n}{k} = \frac{n!}{(n-k)! \; k!}$$

- **Example:**

$L_1$

$C_2$

| $L_1$ |
|---|
| {I1} |
| {I2} |
| {I3} |
| {I4} |

$L_1 \bowtie L_1$

| $C_2$ |
|---|
| {I1, I2} |
| {I1, I3} |
| {I1, I4} |
| {I2, I3} |
| {I2, I4} |
| {I3, I4} |

- **# of elements in $C_2$**

$$= \binom{4}{2} = \frac{4!}{(2)! \; 2!} = 6$$

➢ Consider the following Transaction dataset

- Let the **min-conf.** = **70%**

- Let the **min-sup.** = **2**

| TID | Items |
|------|-------|
| T001 | { I1 , I2, I5 } |
| T002 | { I2, I4 } |
| T003 | { I2, I3 } |
| T004 | { I1, I2, I4 } |
| T005 | { I1, I3 } |
| T006 | { I2, I3 } |
| T007 | { I1, I3 } |
| T008 | { I1, I2, I3, I5 } |
| T009 | { I1, I2, I3 } |

17

| TID | Items |
|------|-------------------|
| T001 | { I1 , I2, I5 } |
| T002 | { I2, I4 } |
| T003 | { I2, I3 } |
| T004 | { I1, I2, I4 } |
| T005 | { I1, I3 } |
| T006 | { I2, I3 } |
| T007 | { I1, I3 } |
| T008 | { I1, I2, I3, I5 } |
| T009 | { I1, I2, I3 } |

**Scan D for count**

| ItemSet | Sup-Count |
|---------|-----------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

**Compare sup-count with min-sup**

| ItemSet | Sup-Count |
|---------|-----------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

18

## $L_1$

| ItemSet | Sup-Count |
|---------|-----------|
| {I1}    | 6         |
| {I2}    | 7         |
| {I3}    | 6         |
| {I4}    | 2         |
| {I5}    | 2         |

$L_1 \bowtie L_1$

## $C_2$

| Items      | count |
|------------|-------|
| {I1 , I2}  | 4     |
| {I1 , I3}  | 4     |
| {I1 , I4}  | 1     |
| {I1 , I5}  | 2     |
| {I2 , I3}  | 4     |
| {I2 , I4}  | 2     |
| {I2 , I5}  | 2     |
| {I3 , I4}  | 0     |
| {I3 , I5}  | 1     |
| {I4 , I5}  | 0     |

Compare sup-count with min-sup

## $L_2$

| Items      | count |
|------------|-------|
| {I1 , I2}  | 4     |
| {I1 , I3}  | 4     |
| {I1 , I5}  | 2     |
| {I2 , I3}  | 4     |
| {I2 , I4}  | 2     |
| {I2 , I5}  | 2     |

$L_2$

| Items | count |
|---|---|
| {I1 , I2} | 4 |
| {I1 , I3} | 4 |
| {I1 , I5} | 2 |
| {I2 , I3} | 4 |
| {I2 , I4} | 2 |
| {I2 , I5} | 2 |

$L_2 \bowtie L_2$

$C_3$

| Items |
|---|
| {I1, I2, I3} |
| {I1, I2, I5} |
| {I1, I3, I5} |
| {I2, I3, I4} |
| {I2, I3, I5} |
| {I2, I4, I5} |

**Pruning step**

Removed bcz {I3, I5} $\notin L_2$

Removed bcz {I3, I4} $\notin L_2$

Removed bcz {I3, I5} $\notin L_2$

Removed bcz {I4, I5} $\notin L_2$

$C_3$

| Items | Sup.count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

**Compare sup-count with min-sup**

$L_3$

| Items | Sup.count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

$L_3$

| Items | Sup.count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

$L_3 \bowtie L_3$

$C_4$

| Items |
|---|
| {I1, I2, I3, I5} |

**Pruning step**

$\{I2, I3, I5\} \notin L_3$

So, $C_4 = \emptyset$

**So, the Frequent item sets are in $L_3$**

| Items | Sup.count |
|---|---|
| {I1, I2, I3} | 2 |
| {I1, I2, I5} | 2 |

- **Discovering ARs From** $L_3$   <span style="background:yellow">(√) means **Interesting**, (X) means **Rejected**</span>

| Items |
|:---:|
| **{I1, I2, I3}** |
| **{I1, I2, I5}** |

**Rule1**: I1 ^ I2 $\Longrightarrow$ I3  confidence = 2/4 = 50% (X)

**Rule2**: I1 ^ I3 $\Longrightarrow$ I2  confidence = 2/4 = 50% (X)

**Rule3**: I2 ^ I3 $\Longrightarrow$ I1  confidence = 2/4 = 50% (X)

**Rule4**: I1 $\Longrightarrow$ I2 ^ I3  confidence = 2/6 = 33% (X)

**Rule5**: I2 $\Longrightarrow$ I1 ^ I3  confidence = 2/7 = 29% (X)

**Rule6**: I3 $\Longrightarrow$ I1 ^ I2  confidence = 2/6 = 33% (X)

**Rule7**: I1 ^ I2 $\Longrightarrow$ I5  confid. = 2/4 = 50% (X)

**Rule8**: I1 ^ I5 $\Longrightarrow$ I2  confid. = 2/2 = 100% (√)

**Rule9**: I2 ^ I5 $\Longrightarrow$ I1  confid. = 2/2 = 100% (√)

**Rule10**: I1 $\Longrightarrow$ I2 ^ I5  confid. = 2/6 = 33% (X)

**Rule11**: I2 $\Longrightarrow$ I1 ^ I5  confid. = 2/7 = 29% (X)

**Rule12**: I5 $\Longrightarrow$ I1 ^ I2  confid. = 2/2 = 100% (√)

Prof.Fadl Ba-Alwi

# ➢ Apriori Algorithm

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$, a database of transactions;

- $min\_sup$, the minimum support count threshold.

**Output:** $L$, frequent itemsets in $D$.

**Method:**

```
(1)      L₁ = find_frequent_1-itemsets(D);
(2)      for (k = 2; Lₖ₋₁ ≠ φ; k++) {
(3)          Cₖ = apriori_gen(Lₖ₋₁);
(4)          for each transaction t ∈ D { // scan D for counts
(5)              Cₜ = subset(Cₖ, t); // get the subsets of t that are candidates
(6)              for each candidate c ∈ Cₜ
(7)                  c.count++;
(8)          }
(9)          Lₖ = {c ∈ Cₖ | c.count ≥ min_sup}
(10)     }
(11)     return L = ∪ₖLₖ;
```

# ➢ Apriori Algorithm

procedure apriori_gen($L_{k-1}$:frequent $(k-1)$-itemsets)
(1)      for each itemset $l_1 \in L_{k-1}$
(2)          for each itemset $l_2 \in L_{k-1}$
(3)              if $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge ... \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$ then {
(4)                  $c = l_1 \bowtie l_2$; // join step: generate candidates
(5)                  if has_infrequent_subset$(c, L_{k-1})$ then
(6)                      delete $c$; // prune step: remove unfruitful candidate
(7)                  else add $c$ to $C_k$;
(8)              }
(9)      return $C_k$;

procedure has_infrequent_subset($c$: candidate $k$-itemset;
        $L_{k-1}$: frequent $(k-1)$-itemsets); // use prior knowledge
(1)      for each $(k-1)$-subset $s$ of $c$
(2)          if $s \notin L_{k-1}$ then
(3)              return TRUE;
(4)      return FALSE;

**DM**

Prof. Fadl Ba-Alwi

$$\sum^{Prof.} \sqrt{\dfrac{Fadl}{Ba - Alwi}}$$

thank you