

Data Mining

Dr. Raed Ibraheem Hamed

**University of Human Development,
College of Science and Technology
Department of Computer Science**

2016 – 2017



Road map

- The Apriori algorithm
 - Step 1: Mining all frequent itemsets
 - Definition of Apriori Algorithm
 - Definition (contd.)
 - Steps to Perform Apriori Algorithm
 - The Apriori Algorithm — Example-1
 - The Apriori Algorithm — Example-2
 - **Step 1:** Generating 1-itemset Frequent Pattern
 - **Step 2:** Generating 2-itemset Frequent Pattern
 - **Step 3:** Generating 3-itemset Frequent Pattern
 - **Step 4:** Generating 4-itemset Frequent Pattern
 - **Step 5:** Generating Association Rules from Frequent Itemsets
-

The Apriori algorithm Key Concepts :

1. **Frequent Itemsets:** The sets of item which has minimum support (denoted by **L_i** for i th-Itemset).
2. **Apriori Property:** Any subset of frequent itemset must be frequent.
3. **Join Operation:** To find **L_k** , a set of candidate k -itemsets is generated by joining L_{k-1} with itself.



Definition (contd.)

- Apriori uses a "bottom up" approach, where frequent subsets are extended one item at a time (a step known as **candidate generation**, and groups of candidates are tested against the data.
- The algorithm terminates when no further successful extensions are found.

Steps to Perform Apriori Algorithm

Apriori Algorithm

Step1

Scan the transaction database to get the support S of each 1-itemset, compare S with min_sup , and get a set of frequent 1-itemsets, L_1

Step2

Use L_{k-1} join L_{k-1} to generate a set of candidate k -itemsets. And use Apriori property to prune the unfrequented k -itemsets from this set

Step3

Scan the transaction database to get the support S of each candidate k -itemset in the final set, compare S with min_sup , and get a set of frequent k -itemsets, L_k

Step4:

The candidate set = Null

YES

NO

Step6

For every nonempty subset s of l , output the rule " $s \Rightarrow (l-s)$ " if confidence C of the rule " $s \Rightarrow (l-s)$ " ($= \text{support } S \text{ of } l / \text{support } S \text{ of } s$) $\geq \text{min_conf}$

Step5

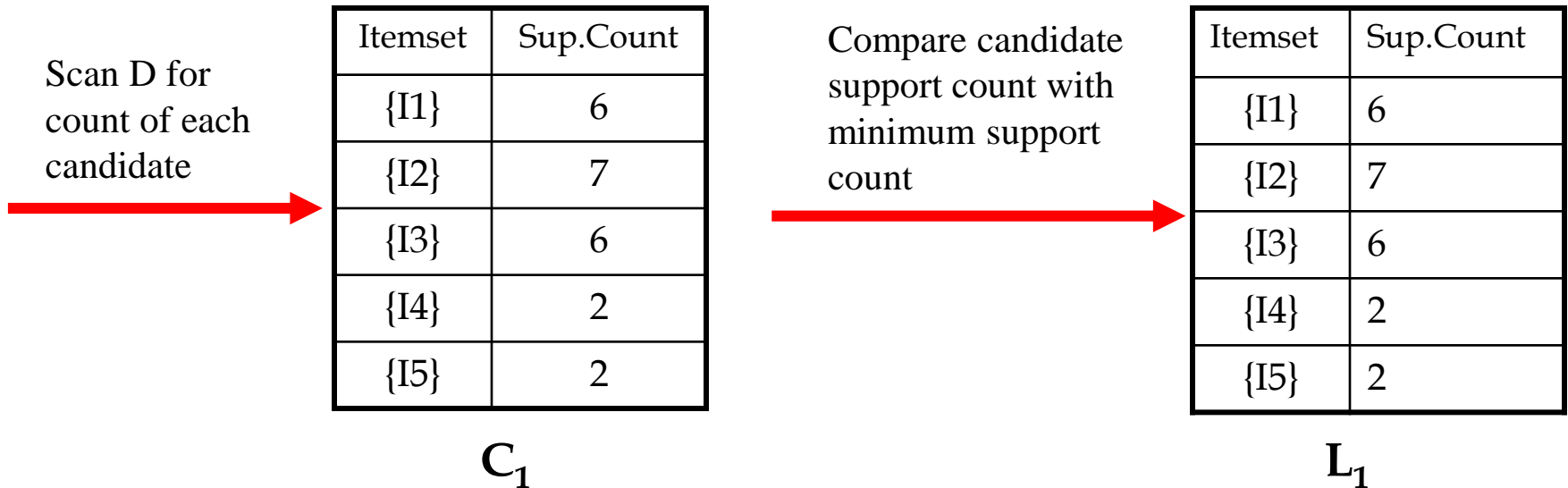
For each frequent itemset l , generate all nonempty subsets of l

The Apriori Algorithm: Example

| TID | List of Items |
|------|----------------|
| T100 | I1, I2, I5 |
| T200 | I2, I4 |
| T300 | I2, I3 |
| T400 | I1, I2, I4 |
| T500 | I1, I3 |
| T600 | I2, I3 |
| T700 | I1, I3 |
| T800 | I1, I2, I3, I5 |
| T900 | I1, I2, I3 |

- Consider a database, D , consisting of 9 transactions.
- Suppose min.support count required is 2 (i.e. $\text{min_sup} = 2/9 = 22\%$)
- Let **minimum confidence required is 70%**.
- We have to first find out the frequent itemset using Apriori algorithm.
- Then, Association rules will be generated using min. support & min. confidence.

Step 1: Generating 1-itemset Frequent Pattern



- In the first iteration of the algorithm, each item is a member of the set of candidate.
- The set of frequent 1-itemsets, L_1 , consists of the candidate 1-itemsets satisfying minimum support.

Step 2: Generating 2-itemset Frequent Pattern

Generate C_2 candidates from L_1

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

C_2

Scan D for count of each candidate

| Itemset | Sup. 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 |

C_2

Compare candidate support count with minimum support count

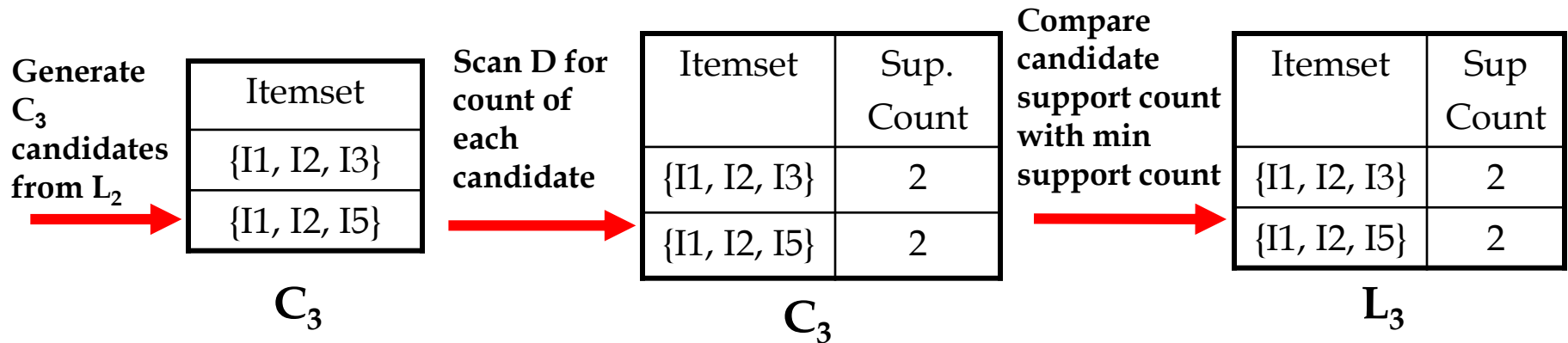
| Itemset | Sup Count |
|----------|-----------|
| {I1, I2} | 4 |
| {I1, I3} | 4 |
| {I1, I5} | 2 |
| {I2, I3} | 4 |
| {I2, I4} | 2 |
| {I2, I5} | 2 |

L_2

Step 2: Generating 2-itemset Frequent Pattern [Cont.]

- To discover the set of frequent 2-itemsets, L_2 , the algorithm uses $L_1 \text{ Join } L_1$ to generate a candidate set of 2-itemsets, C_2 .
- Next, the transactions in D are scanned and the support count for each candidate itemset in C_2 is accumulated (as shown in the middle table).
- The set of frequent 2-itemsets, L_2 , is then determined, consisting of those candidate 2-itemsets in C_2 having minimum support.
- **Note:** We haven't used Apriori Property yet.

Step 3: Generating 3-itemset Frequent Pattern



- The generation of the set of candidate 3-itemsets, C_3 , involves use of the Apriori Property.
- In order to find C_3 , we compute $L_2 \text{ Join } L_2$.
- $C_3 = L_2 \text{ Join } L_2 = \{\{I1, I2, I3\}, \{I1, I2, I5\}, \{I1, I3, I5\}, \{I2, I3, I4\}, \{I2, I3, I5\}, \{I2, I4, I5\}\}$.
- Now, Join step is complete and Prune step will be used to reduce the size of C_3 . Prune step helps to avoid heavy computation due to large C_k .

Step 3: Generating 3-itemset Frequent Pattern [Cont.]

- Based on the **Apriori property** that all subsets of a frequent itemset must also be frequent, we can determine that **four candidates cannot possibly be frequent**. How ?
- For example , lets take **{I1, I2, I3}**. The 2-item subsets of it are {I1, I2}, {I1, I3} & {I2, I3}. Since all 2-item subsets of {I1, I2, I3} are members of L_2 , We will keep {I1, I2, I3} in C_3 .
- Lets take another example of **{I2, I3, I5}** which shows how the pruning is performed. The 2-item subsets are {I2, I3}, {I2, I5} & {I3,I5}.
- BUT, {I3, I5} is not a member of L_2 and hence it is not frequent **violating Apriori Property**. Thus We will have to remove {I2, I3, I5} from C_3 .
- Therefore, $C_3 = \{\{I1, I2, I3\}, \{I1, I2, I5\}\}$ after checking for all members of **result of Join operation for Pruning**.
- Now, the transactions in D are scanned in order to determine L_3 , **consisting of those candidates 3-itemsets in C_3 having minimum support**.

Step 4: Generating 4-itemset Frequent Pattern

- The algorithm uses $L_3 \text{ Join } L_3$ to generate a candidate set of 4-itemsets, C_4 . Although the join results in $\{\{I1, I2, I3, I5\}\}$, this itemset is pruned since its subset $\{\{I2, I3, I5\}\}$ is not frequent.
- Thus, $C_4 = \varnothing$, and algorithm terminates, having found all of the frequent items. This completes our Apriori Algorithm.
- What's Next ?
These frequent itemsets will be used to generate strong association rules (where strong association rules satisfy both minimum support & minimum confidence).

The Apriori Algorithm — Example

Min support = 2

Database D

| TID | Items |
|-----|---------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |

Scan D

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {4} | 1 |
| {5} | 3 |

L_1

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {5} | 3 |

L_2

| itemset | sup |
|---------|-----|
| {1 3} | 2 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

C_2

| itemset | sup |
|---------|-----|
| {1 2} | 1 |
| {1 3} | 2 |
| {1 5} | 1 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

C_2

| itemset |
|---------|
| {1 2} |
| {1 3} |
| {1 5} |
| {2 3} |
| {2 5} |
| {3 5} |

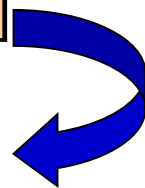
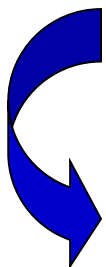
C_3

| itemset |
|---------|
| {2 3 5} |

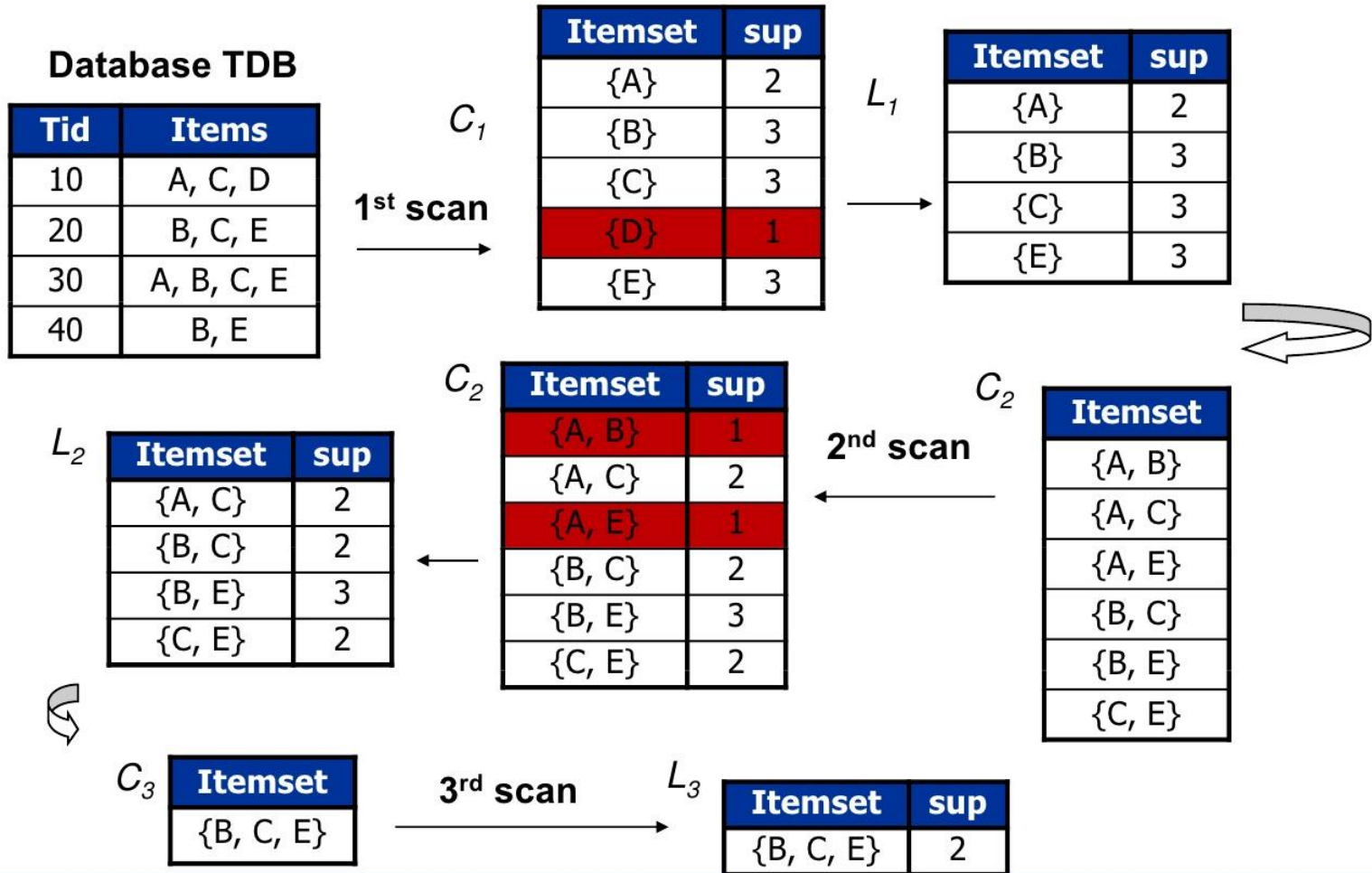
Scan D

| itemset | sup |
|---------|-----|
| {2 3 5} | 2 |

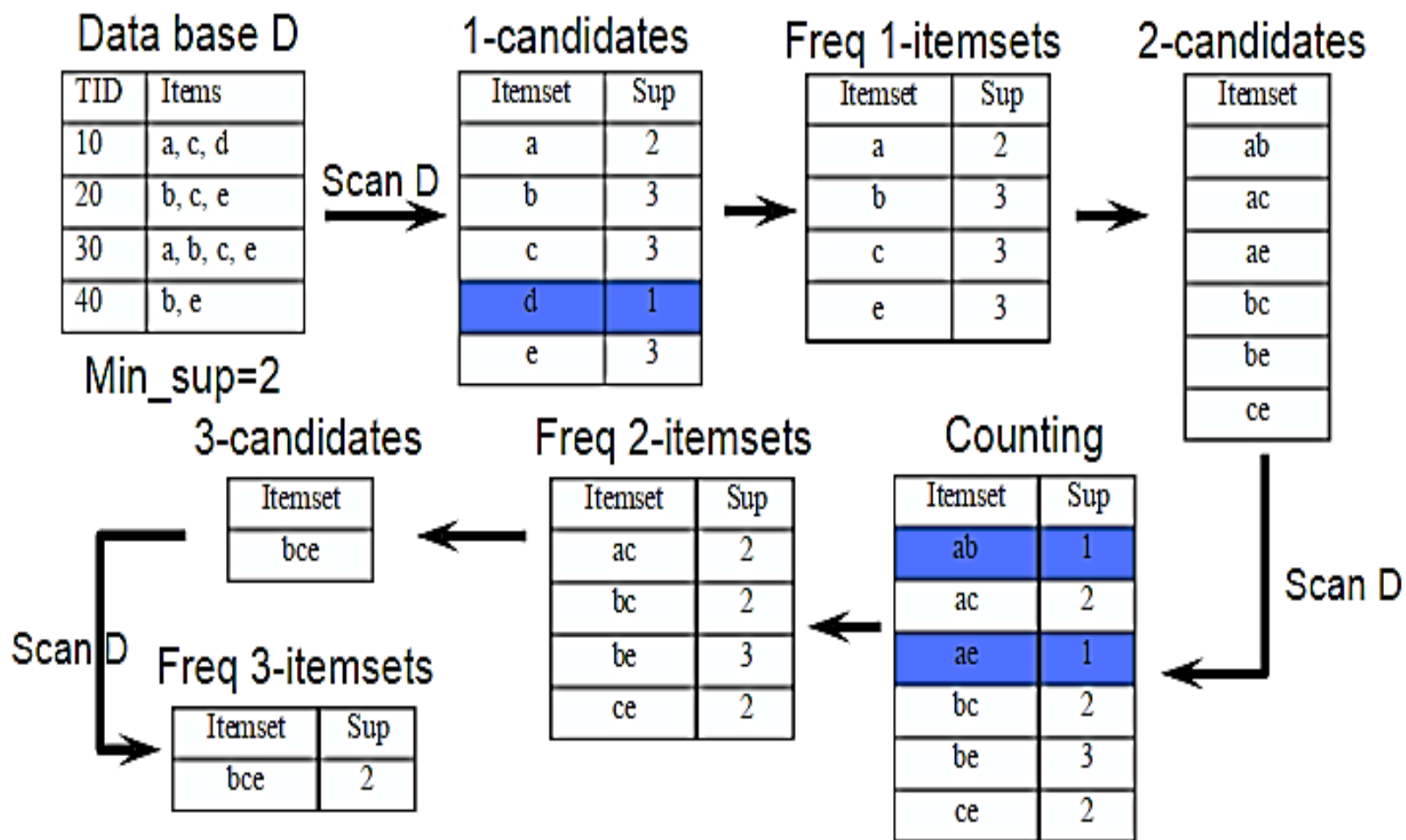
Note: {1,2,3} {1,2,5} and {1,3,5} not in C_3



Example of Apriori Run



Apriori algorithm example



Thank
you

