# Introduction to Data Mining
## Pang-Ning Tan, Michael Steinbach, Vipin Kumar

## HW 1

# Chapter 6.10 Exercises

1. For each of the following questions, provide an example of an association rule from the market basket domain that satisfies the following conditions. Also, describe whether such rules are subjectively interesting.

   (a) A rule that has high support and high confidence.

   (b) A rule that has reasonably high support but low confidence.

   (c) A rule that has low support and low confidence.

   (d) A rule that has low support and high confidence.

**Table 6.1.** Example of market basket transactions.

| Customer ID | Transaction ID | Items Bought |
|---|---|---|
| 1 | 0001 | $\{a, d, e\}$ |
| 1 | 0024 | $\{a, b, c, e\}$ |
| 2 | 0012 | $\{a, b, d, e\}$ |
| 2 | 0031 | $\{a, c, d, e\}$ |
| 3 | 0015 | $\{b, c, e\}$ |
| 3 | 0022 | $\{b, d, e\}$ |
| 4 | 0029 | $\{c, d\}$ |
| 4 | 0040 | $\{a, b, c\}$ |
| 5 | 0033 | $\{a, d, e\}$ |
| 5 | 0038 | $\{a, b, e\}$ |

2. Consider the data set shown in Table 6.1.

   (a) Compute the support for itemsets $\{e\}$, $\{b, d\}$, and $\{b, d, e\}$ by treating each transaction ID as a market basket.

   (b) Use the results in part (a) to compute the confidence for the association rules $\{b, d\} \longrightarrow \{e\}$ and $\{e\} \longrightarrow \{b, d\}$. Is confidence a symmetric measure?

   (c) Repeat part (a) by treating each customer ID as a market basket. Each item should be treated as a binary variable (1 if an item appears in at least one transaction bought by the customer, and 0 otherwise.)

   (d) Use the results in part (c) to compute the confidence for the association rules $\{b, d\} \longrightarrow \{e\}$ and $\{e\} \longrightarrow \{b, d\}$.

   (e) Suppose $s_1$ and $c_1$ are the support and confidence values of an association rule $r$ when treating each transaction ID as a market basket. Also, let $s_2$ and $c_2$ be the support and confidence values of $r$ when treating each customer ID as a market basket. Discuss whether there are any relationships between $s_1$ and $s_2$ or $c_1$ and $c_2$.

8. The *Apriori* algorithm uses a generate-and-count strategy for deriving frequent itemsets. Candidate itemsets of size $k+1$ are created by joining a pair of frequent itemsets of size $k$ (this is known as the candidate generation step). A candidate is discarded if any one of its subsets is found to be infrequent during the candidate pruning step. Suppose the *Apriori* algorithm is applied to the data set shown in Table 6.3 with *minsup* $= 30\%$, i.e., any itemset occurring in less than 3 transactions is considered to be infrequent.

Table 6.3. Example of market basket transactions.

| Transaction ID | Items Bought |
|----------------|--------------|
| 1 | $\{a, b, d, e\}$ |
| 2 | $\{b, c, d\}$ |
| 3 | $\{a, b, d, e\}$ |
| 4 | $\{a, c, d, e\}$ |
| 5 | $\{b, c, d, e\}$ |
| 6 | $\{b, d, e\}$ |
| 7 | $\{c, d\}$ |
| 8 | $\{a, b, c\}$ |
| 9 | $\{a, d, e\}$ |
| 10 | $\{b, d\}$ |

(a) Draw an itemset lattice representing the data set given in Table 6.3. Label each node in the lattice with the following letter(s):

- **N:** If the itemset is not considered to be a candidate itemset by the *Apriori* algorithm. There are two reasons for an itemset not to be considered as a candidate itemset: (1) it is not generated at all during the candidate generation step, or (2) it is generated during the candidate generation step but is subsequently removed during the candidate pruning step because one of its subsets is found to be infrequent.
- **F:** If the candidate itemset is found to be frequent by the *Apriori* algorithm.
- **I:** If the candidate itemset is found to be infrequent after support counting.

(b) What is the percentage of frequent itemsets (with respect to all itemsets in the lattice)?

(c) What is the pruning ratio of the *Apriori* algorithm on this data set? (Pruning ratio is defined as the percentage of itemsets not considered to be a candidate because (1) they are not generated during candidate generation or (2) they are pruned during the candidate pruning step.)

(d) What is the false alarm rate (i.e, percentage of candidate itemsets that are found to be infrequent after performing support counting)?

9. The *Apriori* algorithm uses a hash tree data structure to efficiently count the support of candidate itemsets. Consider the hash tree for candidate 3-itemsets shown in Figure 6.2.

(a) Given a transaction that contains items $\{1, 3, 4, 5, 8\}$, which of the hash tree leaf nodes will be visited when finding the candidates of the trans-

(b) Use the visited leaf nodes in part (b) to determine the candidate item-
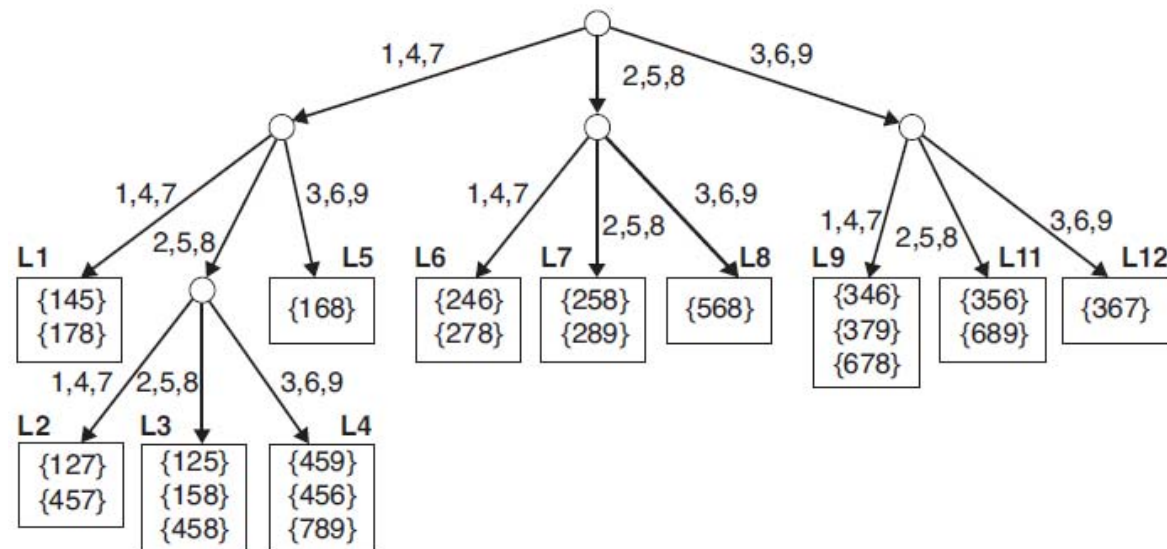


Figure 6.2. An example of a hash tree structure.

11. Given the lattice structure shown in Figure 6.4 and the transactions given in Table 6.3, label each node with the following letter(s):

- $M$ if the node is a maximal frequent itemset,
- $C$ if it is a closed frequent itemset,
- $N$ if it is frequent but neither maximal nor closed, and
- $I$ if it is infrequent.

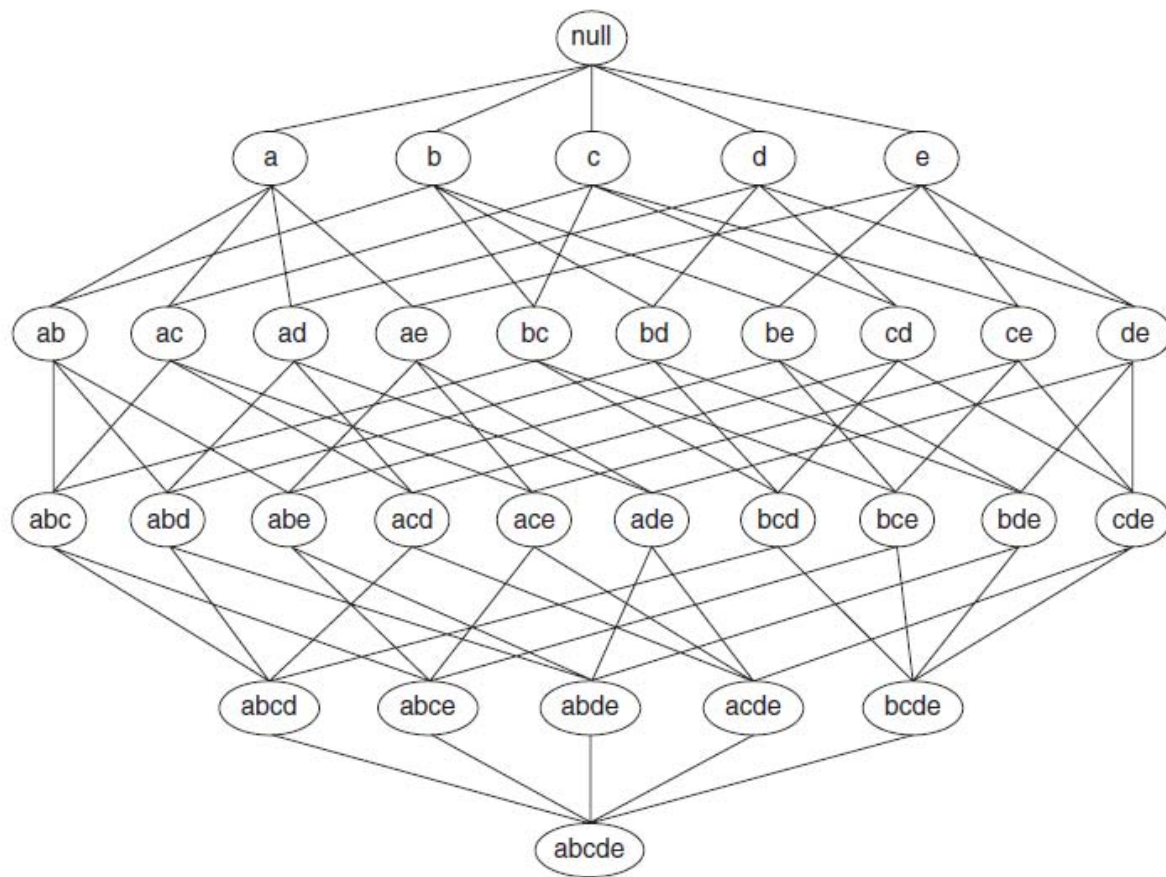Assume that the support threshold is equal to 30%.

**Figure 6.4.** An itemset lattice