# Recitation Problems

**Q2.** Consider the data set shown in Table 6.22.

**Table 6.22.** 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 | $\{a, b, e\}$ |

(a) Compute the support for itemsets {e}, {b,d}, and {b,d,e} by treating each transaction ID as a market basket.
**Answer:**
$S(\{e\}) = \sigma(e)/N = 8/10 = 0.8$
$S(\{b,d\}) = \sigma(bUd)/N = 2/10 = 0.2$
$S(\{b,d,e\}) = \sigma(bUdUe)/N = 2/10 = 0.2$

(b) Use the results in part (a) to compute the confidence for the association rules {b,d} → {e} and {e} →{b ,d}. Is confidence a symmetric measure?
**Answer:**
$C(b,d \to e) = \sigma(b,d \, U \, e)/\sigma(b,d) = 0.2/0.2 = 100\%$
$C(e \to b,d) = \sigma(b,d \, U \, e)/\sigma(e) = 0.2/0.8 = 25\%$
So, by above example we know that confidence is not symmetric.

(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 atleast one transaction bought by the customer, and 0 otherwise.)
**Answer:**
$S(\{e\}) = \sigma(e)/N = 4/5 = 0.8$
$S(\{b,d\}) = \sigma(bUd)/N = 5/5 = 1$
$S(\{b,d,e\}) = \sigma(bUdUe)/N = 4/5 = 0.8$

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

**Answer:**

$C(b,d \rightarrow e) = \sigma(b,d \cup e)/\sigma(b,d) = 0.8/1 = 80\%$

$C(e \rightarrow b,d) = \sigma(b,d \cup e)/\sigma(e) = 0.8/0.8 = 100\%$

(e) Suppose s1 and c1 are the support and confidence values of an association rule r when treating each transaction ID as a market basket. Also, let s2 and c2 be the support and confidence values of r when treating each customer ID as a market basket. Discuss whether there are any relationships between s1 and s2 or c1 and c2.

**Answer:**

No, there is no relationship between s1,s2,c1 and c2.

**Q6.** Consider the market basket transactions shown in Table 6.23.

Table 6.23. Market basket transactions.

| Transaction ID | Items Bought |
|---|---|
| 1 | {Milk, Beer, Diapers} |
| 2 | {Bread, Butter, Milk} |
| 3 | {Milk, Diapers, Cookies} |
| 4 | {Bread, Butter, Cookies} |
| 5 | {Beer, Cookies, Diapers} |
| 6 | {Milk, Diapers, Bread, Butter} |
| 7 | {Bread, Butter, Diapers} |
| 8 | {Beer, Diapers} |
| 9 | {Milk, Diapers, Bread, Butter} |
| 10 | {Beer, Cookies} |

(a) What is the maximum number of association rules that can be extracted from this data (including rules that have zero support)?

**Answer:** The maximum number of association rules that can be extracted is calculated by given formula:

$$R = 3^n - 2^{n+1} + 1$$
$$= 3^6 - 2^{6+1} + 1$$
$$= 602$$

(b) What is the maximum size of frequent itemsets that can be extracted (assuming minsup > 0)?
**Answer:** The maximum size of frequent itemsets that can be extracted is 4 because the longest transaction is of size 4.

(c) Write an expression for the maximum number of size-3 itemsets that can be derived from this data set.
**Answer:** $\binom{6}{3} = 20$.

(d) Find an itemset (of size 2 or larger) that has the largest support.
**Answer:** The itemset which has maximum support is {Bread, Butter} with support =5.

(e) Find a pair of items ,a and b, such that the rules {a} → {b} and {b} →{a} have the same confidence.
**Answer:** {Beer, Cookies}=2/4, {Butter, Cookies}=1/4, {Milk, Diaper}=4/5,{Milk, Bread}=3/5,{Diaper, Bread}=3/5.

**Q7.** Consider the following set of frequent 3-itemsets:
{1, 2, 3}, {1, 2, 4), {r, 2, 5}, {r, 3, 4},{ 1, 3, 5}, {2, 3, 4}, {2, 3, 5}, {3, 4, 5 } .
Assume that there are only five items in the data set.

(a) List all candidate 4-itemsets obtained by a candidate generation procedure using the $F_{k-1}$ x $F_k$ merging strategy.
**Answer:** {1,2,3,4},{1,2,3,5},{1,2,4,5},{1,3,4,5},{2,3,4,5}.

(b) List all candidate 4-itemsets obtained by the candidate generation procedure in Apriori.
**Answer:** {1,2,3,5},{1,2,3,4},{1,2,4,5},{1,3,4,5},{2,3,4,5}.

(c) List all candidate 4-itemsets that survive the candidate pruning step of the Apriori, algorithm.
**Answer:** {1,2,3,4}

**Q9.** 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 6.32.
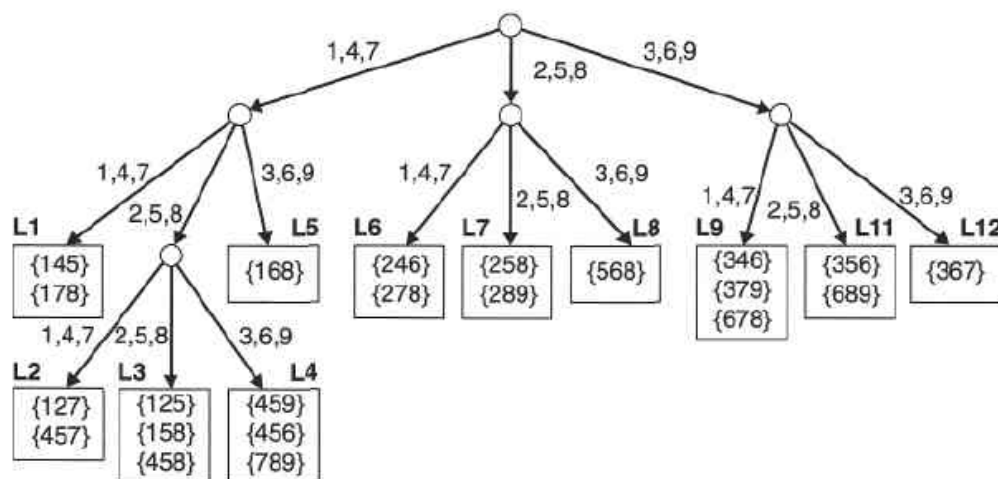


**Figure 6.32.** An example of a hash tree structure.

(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 transaction?
**Answer:** The leaf nodes visited that will be visited are L1, L3, L5, L9, and L11.

(b) Use the visited leaf nodes in part (b) to determine the candidate itemsets that are contained in the transaction {1,3,4,5,8}.
**Answer:** The candidates contained in the transaction are {1, 4, 5}, {1, 5, 8}, and {4, 5, 8}.

**Q11.** Given the lattice structure shown in Figure 6.33 and the transactions given in Table 6.24, label each node with the following letter(s):

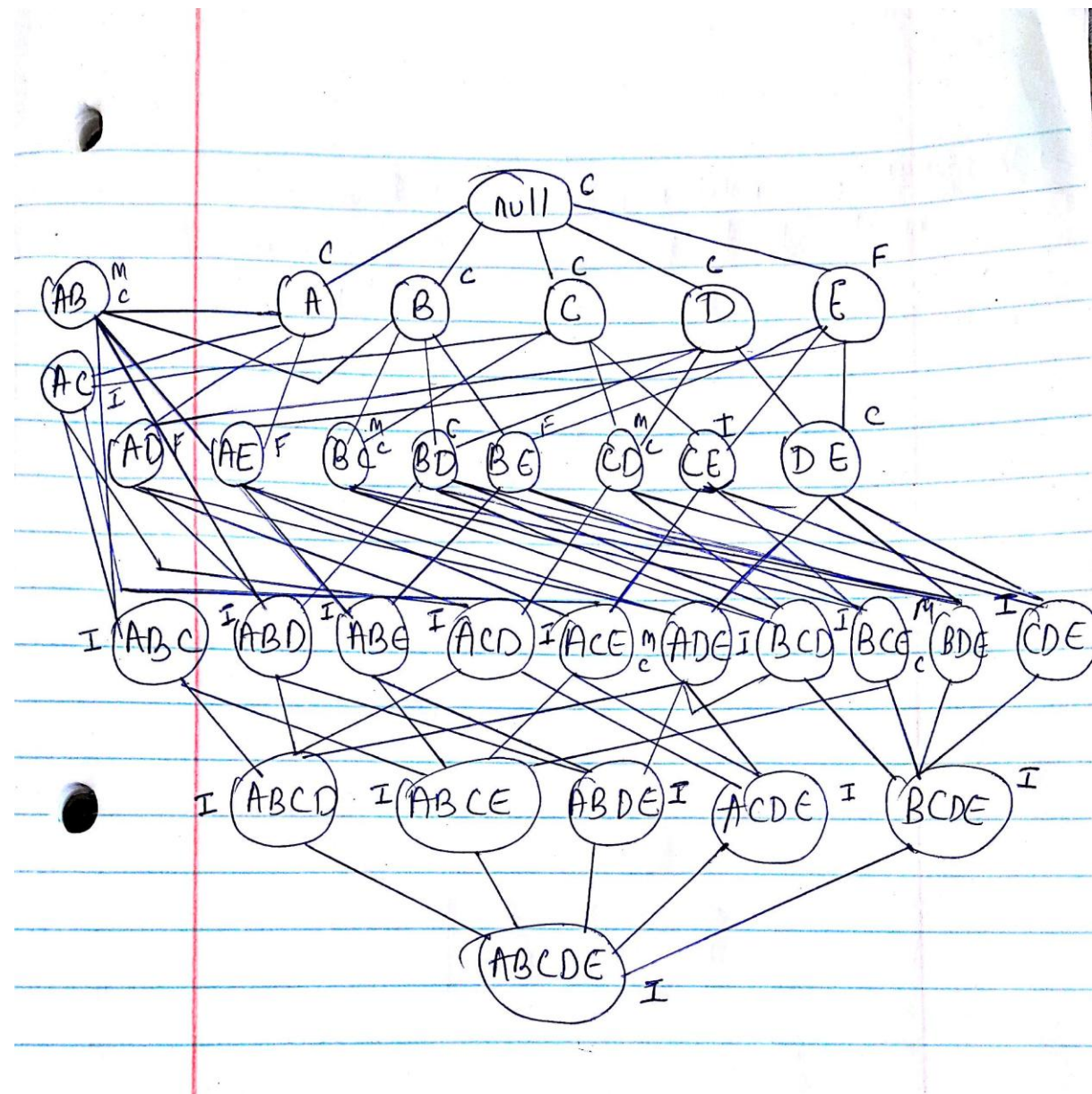M if the node is a maximal frequent itemset,
C if it is a closed frequent itemset,
A if it is frequent but neither maximal nor closed,
I if it is infrequent.

Assume that the support threshold is equal to 30

**Answer:**

**Q12.** The original association rule mining formulation uses the support and confidence measures to prune uninteresting rules.

**Table 6.25.** 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 a contingency table for each of the following rules using the transactions shown in Table 6.25.

Rules: $\{b\} \rightarrow \{c\}$, $\{a\} \rightarrow \{d\}$, $\{b\} \rightarrow \{d\}$, $\{e\} \rightarrow \{c\}$, $\{c\} \rightarrow \{a\}$.

**Answer:**

|  | b | $c^-$ |
|:---:|:---:|:---:|
| b | 3 | 4 |
| $c^-$ | 2 | 1 |

|  | a | $d^-$ |
|:---:|:---:|:---:|
| a | 4 | 1 |
| $d^-$ | 5 | 0 |

|  | b | $d^-$ |
|:---:|:---:|:---:|
| b | 6 | 1 |
| $d^-$ | 3 | 0 |

|  | e | $c^-$ |
|---|---|---|
| e | 2 | 4 |
| $c^-$ | 3 | 1 |

|  | c | $a^-$ |
|---|---|---|
| c | 2 | 3 |
| $a^-$ | 3 | 2 |

(b) Use the contingency tables in part (a) to compute and rank the rules in decreasing order according to the following measures.

i. Support.

**Answer:**

| Rules | Support | Rank |
|---|---|---|
| {b}→{c} | 3/10 | 3 |
| {a} → {d} | 4/10 | 2 |
| {b} → {d} | 6/10 | 1 |
| {e} → {c} | 2/10 | 4 |
| {c}→{a} | 2/10 | 4 |

ii. Confidence.

**Answer:**

| Rules | Confidence | Rank |
|---|---|---|
| {b}→{c} | 3/7 | 3 |
| {a} → {d} | 4/5 | 2 |
| {b} → {d} | 6/7 | 1 |
| {e} → {c} | 2/6 | 5 |
| {c}→{a} | 2/5 | 4 |

iii. Interest. P(X,Y)/P(X) * P(Y)

**Answer:**

| Rules | Interest | Rank |
|---|---|---|
| {b}→{c} | (0.3/0.7)*0.5=0.214 | 3 |
| {a} → {d} | (0.4/0.5)*0.9=0.72 | 2 |
| {b} → {d} | (0.6/0.7)*0.9=0.771 | 1 |
| {e} → {c} | (0.2/0.6)*0.5=0.167 | 5 |
| {c}→{a} | (0.2/0.5)*0.5=0.2 | 4 |

iv. IS $P(X,Y)/\sqrt{(P(X)*P(Y))}$

**Answer:**

| Rules | IS | Rank |
|---|---|---|
| {b}→{c} | 0.507 | 3 |
| {a} → {d} | 0.596 | 2 |
| {b} → {d} | 0.756 | 1 |
| {e} → {c} | 0.365 | 5 |
| {c}→{a} | 0.4 | 4 |

v. Klosgen. $\sqrt{(P(X,Y))}*(P(Y|X)-P(Y))$

**Answer:**

| Rules | Klosgen | Rank |
|---|---|---|
| {b}→{c} | -0.039 | 2 |
| {a} → {d} | -0.063 | 4 |
| {b} → {d} | -0.33 | 1 |
| {e} → {c} | -0.075 | 5 |
| {c}→{a} | -0.045 | 3 |

vi. Odds ratio. $P(X,Y)\ P(\overline{X},\overline{Y})/P(\overline{X},Y)P(X,\overline{Y})$

**Answer:**

| Rules | Odds ratio | Rank |
|---|---|---|
| {b}→{c} | 0.375 | 2 |
| {a} → {d} | 0 | 4 |
| {b} → {d} | 0 | 4 |
| {e} → {c} | 0.167 | 3 |
| {c}→{a} | 0.44 | 1 |

**Q18.** Table 6.26 shows a 2 x 2 x 2 contingency table for the binary variables A and B at different values of the control variable C.

**Table 6.26.** A Contingency Table.

| | | | A | |
|---|---|---|---|---|
| | | | 1 | 0 |
| C = 0 | B | 1 | 0 | 15 |
| | | 0 | 15 | 30 |
| C = 1 | B | 1 | 5 | 0 |
| | | 0 | 0 | 15 |

(a) Compute the phi coefficient for A and B when C=0, C =1, and C= 0 or 1.
**Answer:**

When C=0, we have

Phi =  P(A,B)-P(A)P(B)/√(P(A)P(B)(1-P(A)(1-P(B))

= 0-0.25*0.25/√ (0.25*0.25*0.75*0.75)

= -0.33

When C=1, we have

Phi =  P(A,B)-P(A)P(B)/√(P(A)P(B)(1-P(A)(1-P(B))

= 0.25-0.0625/√ (0.0625*0.5625)

= 1

When C=0 or 1

Phi = 0

(b) What conclusions can you draw from the above result?

**Answer:** If confounding factors are not considered then the result may show discrepancies.

**Q19.** Consider the contingency tables shown in Table 6.27.

**Table 6.27.** Contingency tables for Exercise 19.

|   | $B$ | $\overline{B}$ |
|---|-----|-----|
| $A$ | 9 | 1 |
| $\overline{A}$ | 1 | 89 |

(a) Table I.

|   | $B$ | $\overline{B}$ |
|---|-----|-----|
| $A$ | 89 | 1 |
| $\overline{A}$ | 1 | 9 |

(b) Table II.

(a) For table I, compute support, the interest measure, and the correlation coefficient for the association pattern {A, B}. Also, compute the confidence of rules A → B and B → A.
**Answer:**

s(A) = 0.1, s(B) = 0.9, s(A,B) = 0.09.
I(A,B) = 9, φ(A,B) = 0.89.
c(A → B) = 0.9, c(B → A) = 0.9.

(b) For table II, compute support, the interest measure, and the correlation coefficient for the association pattern {A, B}. Also, compute the confidence of rules A → B and B → A.
**Answer:**
s(A) = 0.9, s(B) = 0.9, s(A,B) = 0.89.
I(A,B) = 1.09, φ(A,B) = 0.89.
c(A → B) = 0.98, c(B → A) = 0.98.

(c) What conclusions can you draw from the results of (a) and (b)?

**Answer:** We can conclude that phi co-efficient is invariant as it takes into consideration both absence and presence of items.

**Q20.** Consider the relationship between customers who buy high-definition televisions and exercise machines as shown in Tables 6.19 and 6.20.
(a) Compute the odds ratios for both tables.
**Answer:**
For Table 6.19, odds ratio = 1.4938.
For Table 6.20, the odds ratios are 0.8333 and 0.98.

(b) Compute the coefficient for both tables.
**Answer:**
For table 6.19, $\varphi = 0.098$.
For Table 6.20, the $\varphi$-coefficients are -0.0233 and -0.0047.

(c) Compute the interest factor for both tables.
**Answer:**
For table 6.19, $\varphi = 1.0784$.
For Table 6.20, the $\varphi$-coefficients are 0.88 and 0.9971.

For each of the measures given above, describe how the direction of association changes when data is pooled together instead of being stratified.
**Answer:**
When data is pooled together, association direction changes from negative to positive.