

Data-240 sec 12 Data Mining / AnalysisHW-4Prayag Nikul Puroani (017416737)Question 1 (frequent item set & Associative rules)

T10 — 1, 3, 4

T20 — 2, 3, 5

T30 — 1, 2, 3, 5

T40 — 2, 5

min-sup = 40%

min-conf = 70%

a)

	items	sup	min-sup %
<u>Scan 1</u> \Rightarrow	1	2	50%
	2	3	75%
	3	3	75%
	4	1	25%
	5	3	75%

Pruning 1 \Rightarrow {1}, {2}, {3}, {5}frequent 1-item set (L_1) \Rightarrow sup > 40%
{1}, {2}, {3}, {5}

Scan 2 \Rightarrow

$\{1, 2\}$	1	25%
$\{1, 3\}$	2	50%
$\{1, 5\}$	1	25%
$\{2, 3\}$	2	50%
$\{2, 5\}$	3	75%
$\{3, 5\}$	2	50%

Pruning 2 $\Rightarrow \{1, 3\}, \{2, 3\}, \{2, 5\}, \{3, 5\}$

Frequent 2-itemset (L_2) $\text{sup} \geq 40\%$

Scan 3 \Rightarrow

$\{1, 2, 3\}$	1	25%
$\{1, 2, 5\}$	1	25%
$\{2, 3, 5\}$	2	50%

Frequent 3-itemset (L_3) $\text{sup} \geq 40\%$

$\{2, 3, 5\}$

b) Association rules

from $\{2, 3, 5\}$

we will generate all the combinations of association rule and then

will pruning them based on the $\text{min-conf} = 70\%$.

2

	$\sup_{\mathcal{L}}(X \cup Y)$	$\sup_{\mathcal{L}}(X)$	conf
$\{2, 3\} \rightarrow \{5\}$	2	2	100%
$\{2, 5\} \rightarrow \{3\}$	2	3	66.67%
$\{3, 5\} \rightarrow \{2\}$	2	2	100%
$\{2\} \rightarrow \{3\}$	2	3	66.67%
$\{3\} \rightarrow \{2\}$	2	3	66.67%
$\{2\} \rightarrow \{5\}$	3	3	100%
$\{5\} \rightarrow \{2\}$	3	3	100%
$\{3\} \rightarrow \{5\}$	2	3	66.67%
$\{5\} \rightarrow \{3\}$	2	3	66.67%

final strong association rules

$\{2, 3\} \rightarrow \{5\}$ conf = 100%

$\{3, 5\} \rightarrow \{2\}$ conf = 100%

$\{2\} \rightarrow \{5\}$ conf = 100%

$\{5\} \rightarrow \{2\}$ conf = 100%

Question 2 (Association rule mining)

a) $\{ \text{hot dogs} \} \rightarrow \{ \text{hamburgers} \}$

$$\text{sup}(\text{hot dog} \cup \text{hamburgers}) = \frac{2000}{5000} = 0.4$$

40%

$$\Rightarrow \text{since } \text{sup}(\text{hot dog} \cup \text{hamburger}) > 25\%$$

$$\text{conf}(\text{hot dogs} \rightarrow \text{hamburgers})$$

$$= \frac{\text{sup}(\text{hot dog} \cup \text{hamburgers})}{\text{sup}(\text{hot dog})}$$

$$= \frac{\frac{2000}{5000}}{\frac{3000}{5000}} = \frac{2000}{3000} = 66.67\%$$

$$\text{since confidence } (66.67\%) > 50\%$$

So, hot dogs \Rightarrow hamburger is a strong association rule as both

$$\text{sup}(40\%) > \text{min} - \text{sup}(25\%)$$

$$\text{conf}(66.67\%) > \text{min} - \text{conf}(50\%)$$

meets the threshold.

b)

To check for if the purchases of hot dogs & hamburgers are independent, we compare expected & observed value

$$\text{Expected count} = \frac{3000 \times 2500}{5000} = 1500$$

⇒ Since the observed value for both is 2000 is higher than expected value (1500) so purchases of hot dogs & hamburgers are independent.

Conclusion → $2000 > 1500$, there is a +ve correlation between hot dog and hamburgers.

c)

$$1) \text{Leff} = \frac{P(HD \cup H)}{P(HD) \times P(H)}$$

$$= \frac{0.4}{0.6 \times 0.5} = 1.33$$

$1.33 > 1$ so there is a +ve association.

2) Correlation

Observed joint frequency of HD & H is higher than expected

so +ve correlation b/w hobday & hamburgers

3) All confidence (min conf)

$$\begin{aligned}
 &= \min(P(\text{hobday}), P(\text{hamburger})) \\
 &= \min(0.6, 0.5) \\
 &= \underline{\underline{0.5}}
 \end{aligned}$$

as all confidence (0.5) is lower than

confidence (0.667) showing lower bound of support across both items.

4) Max confidence

$$\begin{aligned}
 &= \max(P(\text{hobday}), P(\text{hamburger})) \\
 &= \max(0.6, 0.5) \\
 &= 0.6
 \end{aligned}$$

Max conf higher than All conf but still lower than actual conf. This measure shows the potential max sup. for either of items

5) Gini measure

$$= \sqrt{P(A|B) \times P(B|A)}$$

$$= \sqrt{\frac{P(A \cup B)}{P(B)} \times \frac{P(B \cup A)}{P(A)}}$$

$$= \sqrt{\frac{P(H \cup H)}{P(H)} \times \frac{P(H \cup H)}{P(H)}}$$

$$= \sqrt{\frac{2000}{2500} \times \frac{2000}{3000}} = 0.73$$

0.73 indicates a strong relationship
b/t H0 & H by balancing both
conditional probabilities based towards
either item.

All conf \Rightarrow conservative lower bound
useful when stability across both
items is desired.

Max conf \Rightarrow highest likelihood \rightarrow overestimation

Gini measure \Rightarrow balances both item set \Rightarrow
no biases without favouring
either item.

Lift \Rightarrow a direct measure of strength, by comparing joint probability with independence.

correlation \Rightarrow Offers directional insight if relationship is +ve/-ve.

Question 3 Sequential pattern mining with prefix span

S_1	$\langle (ab)ca \rangle$	min_sup = 3
S'_1	$\langle (ab)bc \rangle$	
S_2	$\langle bcd \rangle$	
S_3	$\langle b(ab) \rangle$	
S_4		

a) Length 1 sequential patterns

$\langle a \rangle$ 3

$\langle b \rangle$ 4

$\langle c \rangle$ 4

~~$\langle d \rangle$ 1~~

so length 1 pattern
 $\Rightarrow \langle a \rangle, \langle b \rangle, \langle c \rangle$

b) Projection of database $\langle a \rangle$

S_1	$\langle (-b)ca \rangle$
S'_1	$\langle (-b)bc \rangle$
S_2	$\langle (-b) \rangle$
S_4	

Projection of database $\langle b \rangle$

S_1	$\langle ca \rangle$
S'_1	$\langle bc \rangle$
S_2	$\langle cd \rangle$
S_3	
S_4	$\langle (ab) \rangle$

Projection of $\langle c \rangle$

$S_1 \quad \langle a \rangle$
 $S_2 \quad \langle \rangle$
 $S_3^2 \quad \langle d \rangle$

c) projection of $\langle a \rangle$

$S_1 \quad \langle (-b)ca \rangle$

$S_2 \quad \langle (-b)bc \rangle$

$S_3^2 \quad \langle (-b) \rangle$

S_4

~~$\langle aa \rangle - 1$~~

~~$\langle ab \rangle - 1$~~

~~$\langle ac \rangle - 2$~~

~~$\langle ab \rangle - 3$~~

min - sup = 3

project of $\langle b \rangle$

$S_1 \quad \langle ca \rangle$

$S_2 \quad \langle bc \rangle$

$S_3^2 \quad \langle cd \rangle$

$S_4 \quad \langle (ab) \rangle$

~~$\langle ba \rangle - 2$~~

~~$\langle bb \rangle - 2$~~

~~$\langle bc \rangle - 2$~~

~~$\langle bd \rangle - 1$~~

projection of $\langle c \rangle$

$S_1 \quad \langle a \rangle$

$S_2 \quad \langle \rangle$

$S_3^2 \quad \langle d \rangle$

~~$\langle ca \rangle - 1$~~

~~$\langle cd \rangle - 1$~~

So the only frequent sequence of length-2 is $\langle (ab) \rangle$

d) Projection of $\langle ab \rangle$.

$S_1 \langle ca \rangle$

$S_1' \langle bc \rangle$

$S_2' \langle \rangle$

$\langle (ab) a \rangle - 1$

$\langle (ab) b \rangle - 1$

$\langle (ab) c \rangle - 2$

as $\text{min_sup} = 3$ and we will

not have any more frequent
to proceed with.

Depth-first search

