



# 4. Data Mining

- Exercise 1: Multi MinSup

- $M = \{3, 2, 5, 4, 6, 1\}$

- Read transactions:

Item	Count	SUP %	MIS %
1	8	80	70
2	2	20	17
3	2	20	15
4	2	20	30
5	3	30	30
6	6	60	35

- $L = \{3, 2, 5, 4, 6, 1\}$

- $F_1 = \{3, 2, 5, 6, 1\}$

Transactions	Item	MIS %
1, 4, 6	1	70
1	2	17
1, 5, 6	3	15
1, 6	4	30
4, 6	5	30
1, 2, 3, 5	6	35
1, 2, 3, 5		
6		
1		
1, 6		

F	Item	SUP %	MIS %
F1	3	20	15
	2	20	17
	5	30	30
	6	60	35
	1	80	70



# 4. Data Mining

–  $L = \{3, 2, 5, 4, 6, 1\}$

– Candidate gen.,  $K=2$

- $\{3, 2\}$  :  $\text{sup}(2) = 20\%$   
 $20\% > \text{MIS}(3) = 15$  and  
 $|\text{sup}(3) - \text{sup}(2)| = |20 - 20| = 0 < \varphi = 20\%$   
so  $\{3, 2\}$  is a good candidate
- $\{3, 5\}$ : is a good candidate
- $\{3, 4\}$ : is a good candidate
- $\{3, 6\}$ : is NOT a good candidate ( $> \varphi$ )
- $\{3, 1\}$ : is NOT a good candidate ( $> \varphi$ )

Item	Count	SUP %	MIS %	Transactions
1	8	80	70	1, 4, 6
2	2	20	17	1
3	2	20	15	1, 5, 6
4	2	20	30	1, 6
5	3	30	30	4, 6
6	6	60	35	1, 2, 3, 5
				1, 2, 3, 5
				6
				1
				1, 6

$\varphi = 20\%$



# 4. Data Mining

–  $L = \{3, 2, 5, 4, 6, 1\}$

- $\{2, 5\}$ : is a good candidate
- $\{2, 4\}$ : is a good candidate
- $\{2, 6\}$ : is NOT a good candidate ( $> \varphi$ )
- $\{2, 1\}$ : is NOT a good candidate ( $> \varphi$ )

$$\varphi = 20\%$$

Item	Count	SUP %	MIS %
1	8	80	70
2	2	20	17
3	2	20	15
4	2	20	30
5	3	30	30
6	6	60	35



# 4. Data Mining

$\varphi = 20\%$

Transactions
1, 4, 6
1
1, 5, 6
1, 6
4, 6
1, 2, 3, 5
1, 2, 3, 5
6
1
1, 6

- $L = \{3, 2, 5, 4, 6, 1\}$ 
  - $\{5, 4\}$ :  $\text{sup}(4) = 20\% < \text{MIS}(5) = 30\%$   
so  $\{5, 4\}$  is NOT a good candidate
  - $\{5, 6\}$ : is NOT a good candidate
  - $\{5, 1\}$ : is NOT a good candidate ( $> \varphi$ )
  - 4 can't be used as seed since  $\text{sup}(4) < \text{MIS}(4)$
  - $\{6, 1\}$ : is a good candidate
- $C2 = \{\{3, 2\}, \{3, 5\}, \{3, 4\}, \{2, 5\}, \{2, 4\}, \{6, 1\}\}$

Item	Count	SUP %	MIS %
1	8	80	70
2	2	20	17
3	2	20	15
4	2	20	30
5	3	30	30
6	6	60	35



# 4. Data Mining

- $C2 = \{\{3, 2\}, \{3, 5\}, \{3, 4\}, \{2, 5\}, \{2, 4\}, \{6, 1\}\}$
- Read Transactions to calculate F2
  - $F2 = \{\{3, 2\}, \{3, 5\}, \{2, 5\}, \{6, 1\}\}$

F	Item	SUP %	MIS %
F1	3	20	15
	2	20	17
	5	30	30
	6	60	35
	1	80	70
F2	{3, 2}	20	15
	{3, 5}	20	15
	{2, 5}	20	17
	{6, 1}	40	35

Transactions
1, 4, 6
1
1, 5, 6
1, 6
4, 6
1, 2, 3, 5
1, 2, 3, 5
6
1
1, 6



# 4. Data Mining

- $F2 = \{\{3, 2\}, \{3, 5\}, \{2, 5\}, \{6, 1\}\}; k = 3$
- Join:
  - $\{3, 2, 5\}$ :  $MIS(2) < MIS(5)$  and  $|\text{sup}(2) - \text{sup}(5)| = 10 < \varphi$ , so it can be joined
  - Nothing else can be joined
- Prune
  - $\{3, 2\}$  and  $\{3, 5\} \in F2$
  - Since  $\{2, 5\} \in F2$  the head problem is avoided otherwise we should have recorded also  $\text{sup}(\{2, 5\})$
- $C3 = \{3, 2, 5\}$

Transactions
1, 4, 6
1
1, 5, 6
1, 6
4, 6
1, 2, 3, 5
1, 2, 3, 5
6
1
1, 6

Item	Count	SUP %	MIS %
1	8	80	70
2	2	20	17
3	2	20	15
4	2	20	30
5	3	30	30
6	6	60	35



# 4. Data Mining

minconf = 60%

- Scan transactions,  $F3 = \{3, 2, 5\}$ 
  - $Sup(\{3, 2, 5\}) = 20\% > MIS(3) = 15$
- Step 2: rule generation from  $F3 = \{3, 2, 5\}$ 
  - Non-empty subsets:  $\{3, 2\}, \{3, 5\}, \{2, 5\}, \{3\}, \{2\}, \{5\}$
  - Possible rules derived from  $F_3$ :
    - $\{3, 2\} \rightarrow \{5\}, [sup = 20\%, conf = 100\%]$
    - $\{3, 5\} \rightarrow \{2\}, [sup = 20\%, conf = 100\%]$
    - $\{2, 5\} \rightarrow \{3\}, [sup = 20\%, conf = 100\%]$
    - $\{3\} \rightarrow \{2, 5\}, [sup = 20\%, conf = 100\%]$
    - $\{2\} \rightarrow \{3, 5\}, [sup = 20\%, conf = 100\%]$
    - $\{5\} \rightarrow \{3, 2\}, [sup = 20\%, conf = 67\%]$
  - All are valid since minconf = 60%

F	Item	SUP %	MIS %
F1	3	20	15
	2	20	17
	5	30	30
	6	60	35
	1	80	70
F2	{3, 2}	20	15
	{3, 5}	20	15
	{2, 5}	20	17
	{6, 1}	40	35
F3	{3, 2, 5}	20	15



# 4. Data Mining

- Possible rules derived from  $F_2$ :
  - $\{3\} \rightarrow \{2\}$ , [sup = 20%, conf = 100%]
  - $\{2\} \rightarrow \{3\}$ , [sup = 20%, conf = 100%]
  - $\{3\} \rightarrow \{5\}$ , [sup = 20%, conf = 100%]
  - $\{5\} \rightarrow \{3\}$ , [sup = 20%, conf = 67%]
  - $\{2\} \rightarrow \{5\}$ , [sup = 20%, conf = 100%]
  - $\{5\} \rightarrow \{2\}$ , [sup = 20%, conf = 67%]
  - $\{6\} \rightarrow \{1\}$ , [sup = 40%, conf = 67%]
  - $\{1\} \rightarrow \{6\}$ , [sup = 40%, conf = 50%]
- Except  $\{1\} \rightarrow \{6\}$ , all are valid

minconf = 60%

F	Item	SUP %	MIS %
F1	3	20	15
	2	20	17
	5	30	30
	6	60	35
	1	80	70
F2	{3, 2}	20	15
	{3, 5}	20	15
	{2, 5}	20	17
	{6, 1}	40	35
F3	{3, 2, 5}	20	15





# 4. Data Mining

- **Exercise 2: GSP**

- Initial step

- All singleton sequences are  $\langle a \rangle$ ,  $\langle b \rangle$ ,  $\langle c \rangle$ ,  $\langle d \rangle$

- General step,  $k = 1$

- $\langle d \rangle$  can't form patterns so it can be left out

SID	Sequence
1	$\langle (dc)b(ac) \rangle$
2	$\langle bc(bac) \rangle$
3	$\langle (ab)a \rangle$

Cand	Support
$\langle a \rangle$	3
$\langle b \rangle$	3
$\langle c \rangle$	2
$\langle d \rangle$	1



# 4. Data Mining

- General step,  $k = 1$ , generate length 2 candidates
  - First generate 2 event candidates

	<a>	<b>	<c>
<a>	<aa>	<ab>	<ac>
<b>	<ba>	<bb>	<bc>
<c>	<ca>	<cb>	<cc>

- Then generate 1 sequence candidates, each event with 2 items

	<a>	<b>	<c>
<a>		<(ab)>	<(ac)>
<b>			<(bc)>
<c>			



# 4. Data Mining

–  $k = 2$ , we have 12 2-length candidates

- After the second table scan we remain with 7 2-patterns:  
<ba>, <bc>, <ca>, <cb>, <cc>, <(ab)>, <(ac)>

SID	Sequence
1	<(dc)b(ac)>
2	<bc(bac)>
3	<(ab)a>

Candidate	Support	SIDs
<aa>	1	3
<ab>	0	-
<ac>	0	-
<ba>	3	1, 2, 3
<bb>	1	2
<bc>	2	1, 2
<ca>	2	1, 2
<cb>	2	1, 2
<cc>	2	1, 2
<(ab)>	2	2, 3
<(ac)>	2	1, 2
<(bc)>	1	2



## – Generalization:

- Join

- Joining  $k-1$  elements together to obtain  $k$ -length candidates
- Idea by join is that two sequences,  $s_1$  and  $s_2$  can be joined if after dropping the first item from  $s_1$  and the last item from  $s_2$ , we obtain the same sequence
- E.g.:
  - »  $\langle bc \rangle$  and  $\langle ca \rangle$  can be joined since by dropping  $b$  from  $\langle bc \rangle$  and  $a$  from  $\langle ca \rangle$  we obtain  $\langle c \rangle$ . The joined result is  $\langle bca \rangle$
  - »  $\langle ba \rangle$  and  $\langle (ab) \rangle$  can also be joined and we obtain  $\langle b(ab) \rangle$

- Prune

- Is similar to the apriori algorithm
- $\langle bca \rangle$  passes pruning only if  $\langle bc \rangle$ ,  $\langle ba \rangle$  and  $\langle ca \rangle \in F_2$
- $\langle b(ab) \rangle$  passes pruning only if  $\langle ba \rangle$ ,  $\langle bb \rangle$  and  $\langle (ab) \rangle \in F_2$



# 4. Data Mining

- $k = 2$ , generate length 3 candidates
  - $\langle ba \rangle, \langle bc \rangle, \langle ca \rangle, \langle cb \rangle, \langle cc \rangle, \langle (ab) \rangle, \langle (ac) \rangle$

	$\langle ba \rangle$	$\langle bc \rangle$	$\langle ca \rangle$	$\langle cb \rangle$	$\langle cc \rangle$	$\langle (ab) \rangle$	$\langle (ac) \rangle$
$\langle ba \rangle$	-	-	-	-	-	$\langle b(ab) \rangle$	$\langle b(ac) \rangle$
$\langle bc \rangle$	-	-	$\langle bca \rangle$	$\langle bcb \rangle$	$\langle bcc \rangle$		
$\langle ca \rangle$	-	-	-	-	-	$\langle c(ab) \rangle$	$\langle c(ac) \rangle$
$\langle cb \rangle$	$\langle cba \rangle$	$\langle cbc \rangle$	-	-	-	-	-
$\langle cc \rangle$	-	-	$\langle cca \rangle$	$\langle ccb \rangle$	-	-	-
$\langle (ab) \rangle$	$\langle (ab)a \rangle$	$\langle (ab)c \rangle$	-	-	-	-	-
$\langle (ac) \rangle$	-	-	$\langle (ac)a \rangle$	$\langle (ac)b \rangle$	$\langle (ac)c \rangle$	-	-

- Now perform pruning
  - $\langle bc \rangle, \langle ba \rangle$  and  $\langle ca \rangle \in F_2$  so  $\langle bca \rangle$  is a good candidate
  - $\langle bcb \rangle$  is not, because  $\langle bb \rangle \notin F_2$
  - ...
- After pruning
  - $C_3 = \langle b(ac) \rangle, \langle bca \rangle, \langle bcc \rangle, \langle c(ab) \rangle, \langle c(ac) \rangle, \langle cba \rangle, \langle cbc \rangle, \langle cca \rangle, \langle ccb \rangle$



# 4. Data Mining

–  $k = 3$ , we have 9 3-length candidates

- $C_3 = \langle b(ac) \rangle, \langle bca \rangle, \langle bcc \rangle, \langle c(ab) \rangle, \langle c(ac) \rangle, \langle cba \rangle, \langle cbc \rangle, \langle cca \rangle, \langle ccb \rangle$
- After table scan  
 $F_3 = \langle b(ac) \rangle, \langle c(ac) \rangle$

Candidate	Support	SIDs
$\langle b(ac) \rangle$	2	1, 2
$\langle bca \rangle$	1	2
$\langle bcc \rangle$	0	-
$\langle c(ab) \rangle$	1	2
$\langle c(ac) \rangle$	2	1, 2
$\langle cba \rangle$	1	1
$\langle cbc \rangle$	1	1
$\langle cca \rangle$	0	-
$\langle ccb \rangle$	0	-

SID	Sequence
1	$\langle (dc)b(ac) \rangle$
2	$\langle bc(bac) \rangle$
3	$\langle (ab)a \rangle$



# 4. Data Mining

– Build  $C_4$  from  $F_3 = \langle b(ac) \rangle, \langle c(ac) \rangle$

	$\langle b(ac) \rangle$	$\langle c(ac) \rangle$
$\langle b(ac) \rangle$	-	-
$\langle c(ac) \rangle$	-	-

- We can't build any 4 length candidate so we remain with  $\langle b(ac) \rangle, \langle c(ac) \rangle$  as 3-patterns



# 4. Data Mining

- Exercise 3: MA(4)

Data	MA(4)
4,38	
4,19	
4,65	
6,40	4,905
6,26	5,375
13,51	7,705
4,19	7,59
8,41	8,0925
6,50	8,1525
8,43	6,8825
9,87	8,3025
9,56	8,59
6,57	8,6075
9,03	8,7575
10,18	8,835

