



10 Data Mining

Solutions

- Exercise 1: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



10 Data Mining

Solutions

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

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

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

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



10 Data Mining

Solutions

- $k = 2$, 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 $\langle (ab) \rangle \in F_2$



10 Data Mining

Solutions

- $k = 2$, 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
 - $C_3 = \langle b(ac) \rangle, \langle bca \rangle, \langle bcc \rangle, \langle c(ab) \rangle, \langle c(ac) \rangle, \langle cba \rangle, \langle cbc \rangle, \langle cca \rangle, \langle ccb \rangle$



10 Data Mining

Solutions

– $k = 3$, 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	$\langle (ab)a \rangle$



10 Data Mining

Solutions

– 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



10 Data Mining

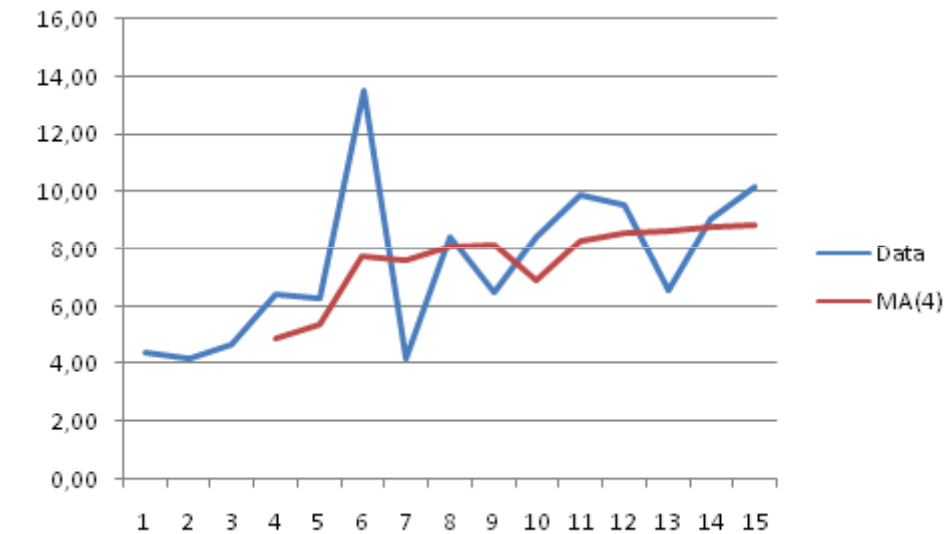
Solutions

- Exercise 2.1: time-series
 - A sequences of values or events **changing with time**
 - Data is recorded at **regular intervals**



- Exercise 2.2: 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





- Exercise 2.3: whole matching method
 - Index building
 - Obtain the DFT coefficients of each sequence in the database
 - Build a $2k$ -dimensional index using the first k Fourier coefficients ($2k$ -dimensions are needed because Fourier coefficients are complex numbers)
 - Query processing
 - Obtain the DFT coefficients of the query sequence
 - Use the $2k$ -dimensional index to filter out such sequences that are at most ε distance away from the query sequence
 - Discards false alarms by computing the actual distance between two sequences