

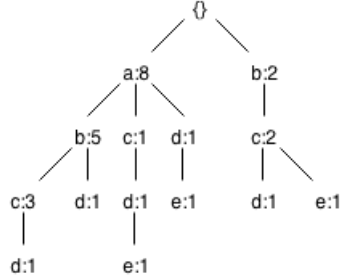
CS6220 Data Mining Fall 2014 Homework 3, Wei Luo

1. Frequent Pattern Mining for Set Data

(a) Scan the Database once, we get:

a:8 b:7 c:6 d:5 e:3

Sort them and build the FP-tree:

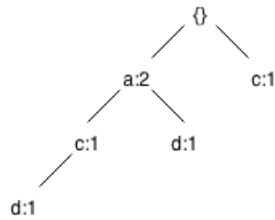


(b) e's conditional pattern base is:

acd:1, ad:1, bc:1

So for e's conditional FP-tree, we have:

a:2 b:1 c:2 d:2 , remove b:1 since it doesn't reach the min_support=2, e's conditional FP-tree is:



Frequent patterns based on e's conditional FP-tree: e, ae, ce, de, ade.

2. Correction Analysis

(a) Based on the observed values given, we can calculate the probabilities and get the observed values table:

	Beer	No Beer	Total
Nuts	17	833	850
No Nuts	183	8967	9150
Total	200	9800	10000

$$\text{confidence}(Beer \Rightarrow Nuts) = 50/200 = 0.25$$

$$\text{lift}(Beer, Nuts) = \frac{P(Beer \cup Nuts)}{P(Beer)P(Nuts)} = \frac{50/10000}{200/10000 \cdot 850/10000} = 2.9412$$

$$\chi^2 = \frac{(50-17)^2}{17} + \frac{(800-833)^2}{833} + \frac{(150-183)^2}{183} + \frac{(9000-8967)^2}{8967} = 71.4384$$

$$\text{all_confidence} = \min(P(Beer|Nuts), P(Nuts|Beer)) = \min(50/850, 50/200) = 0.0588$$

(b) Since $\text{lift}(Beer, Nuts) = 2.9412 > 1$. Buying beer and buying nuts are positively correlated.

3. Sequential Pattern Mining (GSP Algorithm)

(a) For a sequence $s = \langle (ab)(cd)ef \rangle$. s contains 4 elements. The length of s is 6. It contains 63 non-empty subsequences.

(b) From L_3 drop the first and last element of each sequence, we can get:

ID	s	drop first	drop last
1.	$\langle (ab)c \rangle$	$\langle (b)c \rangle$	$\langle (ab) \rangle$
2.	$\langle (ab)d \rangle$	$\langle (b)d \rangle$	$\langle (ab) \rangle$
3.	$\langle a(cd) \rangle$	$\langle (cd) \rangle$	$\langle a(c) \rangle$
4.	$\langle (ac)e \rangle$	$\langle (c)e \rangle$	$\langle (ac) \rangle$
5.	$\langle b(cd) \rangle$	$\langle (cd) \rangle$	$\langle b(c) \rangle$
6.	$\langle bce \rangle$	$\langle ce \rangle$	$\langle bc \rangle$

We can see that 1 and 5, 1 and 6 can be joined. Join them we get:

$\langle (ab)(cd) \rangle$, $\langle (ab)ce \rangle$

For $\langle (ab)(cd) \rangle$, all its length 3 subsequences are in L_3 , keep it.

For $\langle (ab)ce \rangle$, its subsequence $\langle (ab)e \rangle$ is not in L_3 , so we prune it.

So C_4 is $\langle (ab)(cd) \rangle$.

4. Application

First, get the titles of papers in the 20 conferences in time periods 2001-2005, 2008-2012. Then, use tools from nltk to get tokens from each title. Ignore stop words, punctuations and numbers. After that, use PrefixSpan algorithm to find the frequent sequence patterns.

The results are:

For year 2001-2005:

Top 20 most frequent patterns with length 1

$\langle \text{data} \rangle :1231$	$\langle \text{information} \rangle :418$	$\langle \text{efficient} \rangle :289$	$\langle \text{queries} \rangle :226$
$\langle \text{based} \rangle :764$	$\langle \text{xml} \rangle :383$	$\langle \text{approach} \rangle :276$	$\langle \text{systems} \rangle :225$
$\langle \text{mining} \rangle :665$	$\langle \text{retrieval} \rangle :316$	$\langle \text{classification} \rangle :264$	
$\langle \text{using} \rangle :607$	$\langle \text{search} \rangle :314$	$\langle \text{system} \rangle :243$	
$\langle \text{web} \rangle :591$	$\langle \text{query} \rangle :313$	$\langle \text{database} \rangle :238$	
$\langle \text{learning} \rangle :521$	$\langle \text{clustering} \rangle :307$	$\langle \text{model} \rangle :231$	

Top 20 most frequent patterns with length 2

$\langle \text{data mining} \rangle :297$	$\langle \text{xml data} \rangle :82$	$\langle \text{search web} \rangle :72$
$\langle \text{information retrieval} \rangle :151$	$\langle \text{based clustering} \rangle :75$	$\langle \text{web based} \rangle :70$
$\langle \text{based data} \rangle :105$	$\langle \text{proceedings conference} \rangle :74$	$\langle \text{dimensional high} \rangle :69$
$\langle \text{data streams} \rangle :95$	$\langle \text{web data} \rangle :74$	$\langle \text{frequent mining} \rangle :67$
$\langle \text{data using} \rangle :88$	$\langle \text{web mining} \rangle :74$	$\langle \text{report workshop} \rangle :66$
$\langle \text{time series} \rangle :83$	$\langle \text{association rules} \rangle :72$	$\langle \text{processing query} \rangle :64$
$\langle \text{data clustering} \rangle :82$	$\langle \text{management data} \rangle :72$	

Top 20 most frequent patterns with length 3

$\langle \text{proceedings international conference} \rangle :47$	$\langle \text{association rules mining} \rangle :25$
$\langle \text{high dimensional data} \rangle :39$	$\langle \text{data mining workshop} \rangle :25$
$\langle \text{support vector machines} \rangle :39$	$\langle \text{discovery knowledge data} \rangle :24$
$\langle \text{proceedings conference data} \rangle :32$	$\langle \text{proceedings international data} \rangle :24$
$\langle \text{language information retrieval} \rangle :28$	$\langle \text{proceedings data mining} \rangle :24$
$\langle \text{data mining knowledge} \rangle :26$	$\langle \text{international conference data} \rangle :23$

⟨proceedings conference mining⟩ :23	⟨preserving privacy data⟩ :22
⟨report workshop data⟩ :23	⟨workshop retrieval information⟩ :22
⟨rule association mining⟩ :23	⟨series time data⟩ :21
⟨conference data mining⟩ :22	
⟨discovery knowledge mining⟩ :22	

Top 20 most frequent patterns with length 4

⟨proceedings international conference data⟩ :22	⟨proceedings discovery data mining⟩ :16
⟨proceedings conference data mining⟩ :21	⟨conference discovery data knowledge⟩ :15
⟨discovery knowledge data mining⟩ :19	⟨conference discovery data mining⟩ :15
⟨proceedings conference discovery knowledge⟩ :18	⟨proceedings conference knowledge data⟩ :15
⟨proceedings discovery knowledge mining⟩ :18	⟨acm proceedings international conference⟩ :14
⟨conference discovery mining knowledge⟩ :17	⟨cross language retrieval information⟩ :14
⟨proceedings conference knowledge mining⟩ :17	⟨dimensional high clustering data⟩ :14
⟨proceedings knowledge data mining⟩ :17	⟨proceedings international conference mining⟩ :14
⟨conference knowledge data mining⟩ :16	⟨proceedings conference discovery data⟩ :14
⟨proceedings conference discovery mining⟩ :16	
⟨proceedings discovery knowledge data⟩ :16	

For year 2008-2012:

Top 20 most frequent patterns with length 1

⟨data⟩ :1856	⟨information⟩ :637	⟨efficient⟩ :511	⟨multi⟩ :474
⟨based⟩ :1783	⟨analysis⟩ :617	⟨clustering⟩ :505	⟨time⟩ :386
⟨using⟩ :1129	⟨web⟩ :596	⟨retrieval⟩ :493	
⟨learning⟩ :1099	⟨system⟩ :569	⟨model⟩ :484	
⟨mining⟩ :1004	⟨classification⟩ :546	⟨networks⟩ :480	
⟨search⟩ :738	⟨query⟩ :542	⟨approach⟩ :474	

Top 20 most frequent patterns with length 2

⟨data mining⟩ :421	⟨feature selection⟩ :131	⟨semi supervised⟩ :122
⟨information retrieval⟩ :233	⟨web search⟩ :128	⟨clustering based⟩ :119
⟨based data⟩ :179	⟨data streams⟩ :127	⟨using based⟩ :117
⟨social networks⟩ :155	⟨using data⟩ :127	⟨learning using⟩ :115
⟨system based⟩ :152	⟨approach based⟩ :123	⟨algorithm based⟩ :113
⟨scale large⟩ :139	⟨machine learning⟩ :123	⟨mining based⟩ :111
⟨model based⟩ :136	⟨time series⟩ :123	

Top 20 most frequent patterns with length 3

⟨proceedings international conference⟩ :63	⟨high dimensional data⟩ :43
⟨proceedings data mining⟩ :61	⟨discovery knowledge data⟩ :43
⟨international data mining⟩ :53	⟨proceedings conference mining⟩ :42
⟨proceedings international data⟩ :53	⟨proceedings international mining⟩ :42
⟨proceedings conference data⟩ :52	⟨international workshop data⟩ :42
⟨semi supervised learning⟩ :51	⟨data mining based⟩ :42
⟨data mining workshop⟩ :51	⟨proceedings discovery knowledge⟩ :41
⟨conference data mining⟩ :45	⟨machines vector support⟩ :38

⟨conference discovery knowledge⟩ :37	⟨series time data⟩ :34
⟨international mining workshop⟩ :35	
⟨conference international data⟩ :34	

Top 20 most frequent patterns with length 4

⟨proceedings conference data mining⟩ :41	⟨part conference discovery knowledge⟩ :22
⟨proceedings international data mining⟩ :38	⟨preface workshop data mining⟩ :22
⟨proceedings conference international data⟩ :32	⟨conference discovery data knowledge⟩ :21
⟨proceedings conference discovery knowledge⟩ :32	⟨proceedings international workshop data⟩ :21
⟨discovery knowledge data mining⟩ :27	⟨proceedings discovery knowledge mining⟩ :21
⟨international workshop data mining⟩ :27	⟨proceedings discovery data mining⟩ :21
⟨proceedings discovery knowledge data⟩ :25	⟨conference discovery data mining⟩ :20
⟨conference international data mining⟩ :24	⟨conference discovery mining knowledge⟩ :20
⟨proceedings conference international mining⟩ :23	⟨conference knowledge data mining⟩ :20
⟨proceedings conference machine learning⟩ :22	⟨part proceedings conference discovery⟩ :20