

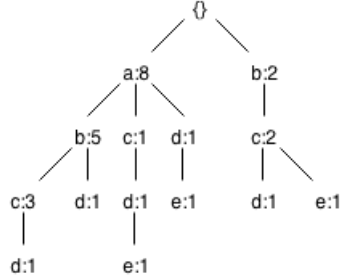
CS6220 Data Mining Fall 2014 Homework 3, Wei Luo

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:

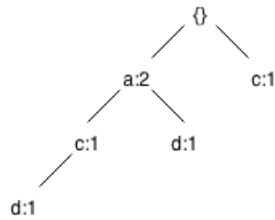


(b) e's conditional pattern base is:

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, cde, ace, acde.

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 8. 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 PyFIM with APRIORI 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

$\langle \text{mining data} \rangle :297$	$\langle \text{web data} \rangle :74$	$\langle \text{association mining} \rangle :55$
$\langle \text{based data} \rangle :105$	$\langle \text{web mining} \rangle :74$	$\langle \text{learning data} \rangle :54$
$\langle \text{streams data} \rangle :95$	$\langle \text{management data} \rangle :72$	$\langle \text{efficient data} \rangle :53$
$\langle \text{using data} \rangle :88$	$\langle \text{web based} \rangle :70$	$\langle \text{model based} \rangle :53$
$\langle \text{clustering data} \rangle :82$	$\langle \text{frequent mining} \rangle :67$	$\langle \text{dimensional data} \rangle :51$
$\langle \text{xml data} \rangle :82$	$\langle \text{patterns mining} \rangle :62$	$\langle \text{learning based} \rangle :50$
$\langle \text{clustering based} \rangle :75$	$\langle \text{approach based} \rangle :58$	

Top 20 most frequent patterns with length 3

$\langle \text{asia mining data} \rangle :6$	$\langle \text{cliques data mining} \rangle :2$	$\langle \text{ssp mining data} \rangle :2$
$\langle \text{mdm mining data} \rangle :6$	$\langle \text{drifting mining data} \rangle :2$	$\langle \text{ugly mining data} \rangle :2$
$\langle \text{pacific mining data} \rangle :6$	$\langle \text{fractals mining data} \rangle :2$	$\langle \text{warehouse based data} \rangle :2$
$\langle \text{dm mining data} \rangle :3$	$\langle \text{grids mining data} \rangle :2$	$\langle \text{7th data mining} \rangle :1$
$\langle \text{pakdd mining data} \rangle :3$	$\langle \text{medicine mining data} \rangle :2$	$\langle \text{aaand mining using} \rangle :1$
$\langle \text{academy mining data} \rangle :2$	$\langle \text{ole mining data} \rangle :2$	$\langle \text{aboutness using based} \rangle :1$
$\langle \text{bad mining data} \rangle :2$	$\langle \text{peculiarity data mining} \rangle :2$	

Top 20 most frequent patterns with length 4

$\langle \text{adherence using based data} \rangle :1$	$\langle \text{ids mining based data} \rangle :1$
$\langle \text{admit mining based data} \rangle :1$	$\langle \text{infer mining using data} \rangle :1$
$\langle \text{bibfinder mining using data} \rangle :1$	$\langle \text{lead mining using data} \rangle :1$
$\langle \text{columbia mining based data} \rangle :1$	$\langle \text{meningitis mining using data} \rangle :1$
$\langle \text{cooperatively mining using data} \rangle :1$	$\langle \text{rs mining using data} \rangle :1$
$\langle \text{deployment mining based data} \rangle :1$	$\langle \text{rsbr mining using data} \rangle :1$
$\langle \text{divisive mining using data} \rangle :1$	$\langle \text{simplicial mining using data} \rangle :1$
$\langle \text{effectively mining using data} \rangle :1$	$\langle \text{statminer mining using data} \rangle :1$
$\langle \text{gdt mining using data} \rangle :1$	$\langle \text{ubdm mining based data} \rangle :1$
$\langle \text{generalised mining using data} \rangle :1$	
$\langle \text{ibl mining using data} \rangle :1$	

For year 2008-2012:

Top 20 most frequent patterns with length 1

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

Top 20 most frequent patterns with length 2

$\langle \text{mining data} \rangle :421$	$\langle \text{machine learning} \rangle :123$	$\langle \text{learning based} \rangle :108$
$\langle \text{data based} \rangle :179$	$\langle \text{clustering based} \rangle :119$	$\langle \text{research based} \rangle :100$
$\langle \text{system based} \rangle :152$	$\langle \text{using based} \rangle :117$	$\langle \text{analysis based} \rangle :95$
$\langle \text{model based} \rangle :136$	$\langle \text{learning using} \rangle :115$	$\langle \text{learning data} \rangle :94$
$\langle \text{streams data} \rangle :127$	$\langle \text{algorithm based} \rangle :113$	$\langle \text{multi learning} \rangle :94$
$\langle \text{using data} \rangle :127$	$\langle \text{mining based} \rangle :111$	$\langle \text{management data} \rangle :93$
$\langle \text{approach based} \rangle :123$	$\langle \text{analysis data} \rangle :110$	

Top 20 most frequent patterns with length 3

$\langle \text{warehouses based data} \rangle :3$	$\langle \text{albatross using data} \rangle :1$
$\langle \text{modis using data} \rangle :2$	$\langle \text{ale using based} \rangle :1$
$\langle \text{pathway data based} \rangle :2$	$\langle \text{alias learning based} \rangle :1$
$\langle \text{abnormalities learning data} \rangle :1$	$\langle \text{allow based data} \rangle :1$
$\langle \text{abnormalities learning using} \rangle :1$	$\langle \text{alphabets based using} \rangle :1$
$\langle \text{abnormalities using data} \rangle :1$	$\langle \text{alpos learning data} \rangle :1$
$\langle \text{abundant learning based} \rangle :1$	$\langle \text{american learning using} \rangle :1$
$\langle \text{accents learning using} \rangle :1$	$\langle \text{analyst data based} \rangle :1$
$\langle \text{aco using based} \rangle :1$	$\langle \text{ancheng data based} \rangle :1$
$\langle \text{adverse using data} \rangle :1$	
$\langle \text{affected using data} \rangle :1$	

Top 20 most frequent patterns with length 4

⟨abnormalities learning using data⟩ :1	⟨remaining using data based⟩ :1
⟨bee using data based⟩ :1	⟨reverible learning using data⟩ :1
⟨ciphertext using data based⟩ :1	⟨rotation learning data based⟩ :1
⟨comet learning using data⟩ :1	⟨sigma using data based⟩ :1
⟨froc learning using data⟩ :1	⟨subdivision learning data based⟩ :1
⟨hazards using data based⟩ :1	⟨tailed using data based⟩ :1
⟨homogenous learning using data⟩ :1	⟨topographic learning using data⟩ :1
⟨pathway using data based⟩ :1	⟨vibratory using data based⟩ :1
⟨periodical learning data based⟩ :1	
⟨recipe learning using data⟩ :1	