

1/2017

Efficient algorithm: Given  $T$ , generate all possible frequent itemsets.

Apriori - too many repetitions - bottom up approach.  
fp growth - 2 scans.

Last step of Apriori is  $C_{k+1} = \emptyset$ .

- Upward closure - if an itemset is infrequent then its immediate supersets are infrequent.
- Downward closure - exact opp of upward closure.

No of scans in Apriori =  $|C_1| + |C_2| + \dots + |C_k|$   
Time complexity of apriori.

- How is fp-growth better?  
2nd scan is to generate fp-tree.
- Sequence pattern mining in Apriori -> Generalised Sequence pattern mine.

Apriori is adjustable in SPM domain.  
fp growth is not suitable for GSP.

fp growth was developed by Han + Team (2010).  
fp growth is not for sequence pattern mining.

- fp growth for sequence pattern -> PrefixSpan.  
SPM used in web patterns. Generic to specific websites.
- fp growth, does it require only 2 overall scans?  
Prere is referred as the projection of the database.  
Projection is to support efficient counting.

Apriori given by Rakesh Agrawal.

DIC was about parallel counting.

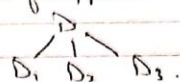
Where do more counting occur in fp tree?  
 $D_1, D_2$  might occur in many paths. So the subset merging requires smaller scans (not the entire DB). So it suffers from repeated scans.  
Subset merging: looking at sequences where  $D_1$  is present.

Both Apriori & fp-growth are time complex.  
DIC is space complex.  
fp growth is a good balance b/w Apriori and DIC.

Prere Search (Maximal) ; A-close (Closed FS).

- Apriori is good for sparse dataset. That means when there are many infrequent itemsets.
- Can we try improving the efficiency of Apriori?  
Apriori is a balance b/w space & time complexity.

- Partitioning approach of Apriori.



if  $10F = 15$  (minsup = 3)  
 $|D_1| = 5 = |D_2| = |D_3|$ .

Local mining approach is done (Local Apriori).  
To arrive at the global count we need merging of  $D_1, D_2$  &  $D_3$  (Union all three).

7/1/2017

In partitioning, something that is locally infrequent might be globally frequent.  
Do union.

Hash pruning approach. → in early levels for counting. Let in the  $C_2$  level.

→ Why not in higher levels? Occurrence of the itemset with higher length will come down.

1, 2, 5	
2, 4	
2, 3	Min sup = 2
1, 2, 4	
1, 3	
2, 3	
1, 3	
1, 2, 3, 5	
1, 2, 3	

If a b c d are there, then  
 a → 1  
 b → 2  
 c → 3  
 d → 4

$L_1 = \{1, 2, 3, 4, 5\}$

$C_2 = L_1 \times L_1$

$C_2$   
 $L_2$   $C_2$   
 freq infreq

Instead of going thru the whole dataset, we look at the hash table. In this way Apriori with hashing is better.

$$H(x, y) = (\text{order of } x \times 10^3 + \text{order of } y) \% 7$$

Buckets →	0	1	2	3	4	5	6	for (1, 4)
		(1, 4) (1, 5) (2, 3)	(2, 4) (2, 5)	(1, 2) (1, 3)				
Contents →	(3, 5) (1, 5) (2, 3)	(2, 4)	(1, 2) (1, 3)	(1, 2) (1, 3)	(1, 2) (1, 3)			
	(2, 3)		(1, 2) (1, 3)	(1, 2) (1, 3)				
		2	2	4	2	2	4	4

Repeat for all. Collision happens only for infrequent.  $C_2$  is better.

In case of collision, revise the hash function

### Transaction Reduction:-

If  $B(1, 2)$  is infrequent, then remove the transaction from the Database.

1, 3, 7	MS=3	1 (5)
2, 3, 7		2 (7)
1, 2, 3		3 (7)
2, 3		4 (3)
2, 3, 4, 5		
2, 3		$L_1 \neq C_1$
1, 2, 3, 4, 6		
2, 3, 4, 6		
1, 3		
1		

If an itemset is infrequent in the earlier levels remove it from the later levels.

5, 7, 6 are infrequent. Remove them from the Dataset.

So new DB is,

1, 3
2, 3
1, 2, 3
2, 3
2, 3, 4
2, 3
1, 2, 3, 4
2, 3, 4
1, 3
1

Atleast one character is being thrown out.

(Time complexity)

The order might remain same.

Apriori's complexity is exponential.



Transaction reduction is also level based.  
Transaction reduction is not efficient as an algo, but implementation wise it is efficient.

ECLAT → IFS  
Apriori → CF  
Dives → MF

\* ECLAT - Equivalence Class Lattice Traversal

Apriori is RF traversal

ECLAT → vertical representation of DB → focus on item  
Horizontal focuses on Transactions

$I_1 \rightarrow \{T_1, T_4, T_5, T_8, T_9\}$   
 $I_2 \rightarrow \{T_1, T_2, T_3, T_4, T_6, T_8, T_9\}$   
 $I_3 \rightarrow \{T_2, T_6, T_7, T_8, T_9\}$   
 $I_4 \rightarrow \{T_2, T_6\}$   
 $I_5 \rightarrow \{T_1, T_9\}$

MS = 2

'n' in order (Intersection in order)

✓  $I_1 I_2 = \{T_1, T_4, T_8, T_9\}$   
 ✓  $I_1 I_3 = \{T_1, T_4, T_8, T_9\}$   
 ✓  $I_1 I_4 = \{T_1, T_4, T_8, T_9\}$   
 ✓  $I_1 I_5 = \{T_1, T_4, T_8, T_9\}$   
 ✓  $I_2 I_3 = \{T_2, T_6, T_8, T_9\}$   
 ✓  $I_2 I_4 = \{T_2, T_6\}$   
 ✓  $I_2 I_5 = \{T_1, T_9\}$   
 ✓  $I_3 I_4 = \{T_2, T_6\}$   
 ✓  $I_3 I_5 = \{T_1, T_9\}$   
 ✓  $I_4 I_5 = \{T_1, T_9\}$   
 ✓  $I_1 I_6 = \emptyset$   
 ✓  $I_2 I_6 = \{T_2, T_6\}$   
 ✓  $I_3 I_6 = \emptyset$

→  $C_2$  in terms of Apriori

ECLAT is better than Apriori as it ignores on counting.  
No need to go to the dataset at all.

→  $T_1, T_2, T_3$

$I_1, I_2, I_3 = \{T_1, T_2, T_3\}$   
 $I_4, I_5, I_6 = \{T_1, T_2, T_3\}$

Cost function following

CHART / Read about matrix

What was the need for CF, MF and LF?  
Rare is not frequent entirely. Sensitive & high utility items are other types of items.  
Finding sensitive items → Privacy preservation data mining.

11/1/2019

A-close → Apriori-close

Why do we need CF, MF & LF?

The application is not happy with so many FI's. So think it down.



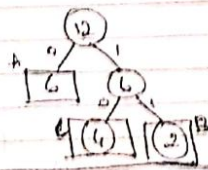
Closed is from closure. The result should belong to the universe of Discourse. Here  $I_2$  is the universe of Discourse. CF is lowest, MF is low. Lowest → from CF I can use reconstruct  $I_2$ ? Yes.  $q_2, p_2, r_2$ , etc. (low length encoding). Huffman tree was taught.

AAAAA B A/BAAAAA CCCCC

A → 0 B → 1 C → 1

descending order of frequency

$\frac{6}{12} \frac{1}{6} \frac{1}{6}$



D → rest

for A → 0 ; B → 1 ; C → 10  
The item with more frequency gets the more or lesser code length.

Biology + Data Mining = Bioinformatics.

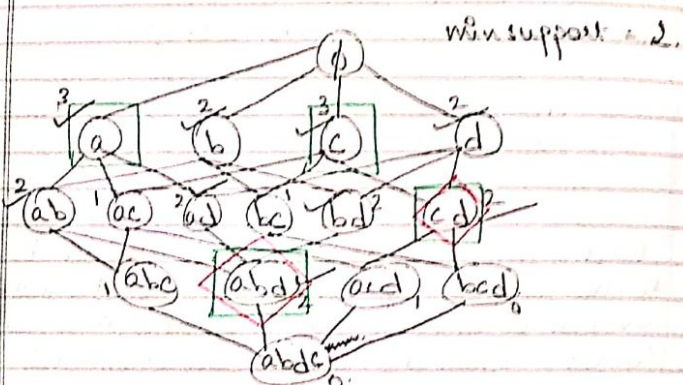
Decoding should be unambiguous. (Prefix free encoding)

CFS, MFS find applications in compression.

Vide and image compression can be long.

CFS =  $A(s) \cdot AC(s)$ .

1/2019



- ✓ → f.i. (1)
- → cfi's (4)
- ◇ → mfi (2)

Cf is lossless, mfi is lossy.

are we getting the entire fi from cfi's.

MFS is more about IFI (cardinality)

CFS is a tighter constraint on MFS.

MFS also comments on the subset.

We will not be able to say the subset count from MFS. Confidence calculation won't be possible.

With CFS we can tell the support count of subsets.

$x \rightarrow$  if there is no superset of  $x$  with the same support count as  $x$  (or less support)

CFS

$\text{support}(x) \leq \text{sc}(x)$

$x \rightarrow$  if any superset of  $x$  which is frequent, MFS.

Difference between closed itemset and closed frequent itemset.

\* A "close" Algorithm:-

S items (A, C, T, D, W).

all the transactions that contain

$i(x)$  itemset  $x$ .  $\# \{T \in I : x \subseteq T\}$

$i(y) : y \subseteq T$

1 A C T W

2 C D W

3 A C T W

4 A C D W

5 A C D T W

6 C D T

$i(C) = \{1, 2, 3, 4, 5, 6\}$

$i(C) = \text{undefined}$

$i(T_1, T_2) = \text{Intersection of } 1 \text{ \& } 2 = C, W$

$i(C) = C$

$i(C, D) = \{2, 4, 5, 6\}$

$i(C, D) =$

$i(A) = \{1, 3, 4, 5\}$  }  $i(A) \neq A$  unlike the examples before this.

$i(A) = \{A, C, W\}$  } it is equal to a greater set

Do,  $i(A, C, W) = \{A, C, W\}$ .

Do  $i(x)$  then  $i(i(x))$ . If  $i(x) \neq x$  but  $i$  equal to  $y$ , Do  $i(i(y))$ . Here it has to stop.

1 level we have elements gaining.



- (i)  $i(t(x)) = x$ .  
 (ii)  $i(t(x)) = y$  ;  $x \subset y$  → elements gained  
 $i(t(y)) = y$ .

### Generators :-

Identify those elements that act as generators.

$x$  is gen(y) if  $x^+ = y$ .  
 generator of y      closure of x yields y.

In the example above, A is a generator of ACW.

In CFP we are interested about minimal generators.  
 There is no proper subset of x generates y.

What all generates ACW.  
 ACW and A.

We cannot have more than 1 minimal generator.

1/20/19

### A - Close Trace.

1	ACD		A
2	BEF	ms=3	A B
3	ABCE		B C
4	BE		C E
5	ABCE		E

$$C_1 = \{A, B, C, D, E\}$$

$$L_1 = \{A, B, C, E\}$$

$$C_2 = \{AB, AC, AE, BC, BE, CE\}$$

With CF1 we reconstruct the FI's but with Closed Itemset we can traceback to the transactions.

$$L_2 = \{AC, BC, BE, CE\}$$

Is any item having the same support as it's subset. Here AC  $\neq$  A and BE  $\neq$  B. Throw them away.

$$L_2 = \{BC, CE\}$$

$$C_3 = \emptyset$$

Since we left out AC & BE we must prove that AC can be reached from some other minimal set.

### (x) Generators

	$t(x)$	$i(t(x))$
A	1, 3, 5	AC
B	2, 3, 4, 5	BE
C	1, 2, 3, 5	C
E	2, 3, 4, 5	x BE
BC	2, 3, 5	BCE
CE	2, 3, 5	x BCE

Whatever is left out is able to be reached from  $i(t(x))$ . Reaching out in mathematical terms is closure ( $x^+$ ).

Remove the duplicates in order.

$$CFI's = \{AC, C, BCE, BE\}$$

$$AC_3 \quad BCE_3$$

$$C_4 \quad BE_4$$

with this as a start pt we should be able to trace back all the FI's.

By not retaining the higher cardinality sets, we are saving space. The cardinality of the Generators should be minimum.  
 One CFI alone cannot give all the FI's. Go for the CFI with the larger count.

bottom up.

\* Primer Search.  $\uparrow \downarrow$  "Dataset is given"

support = 20%, relative support = 15.  
 $\therefore$  min sup = 3

$\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$  9 symbols.

Move in both bottom up and topdown.

Maximal frequent candidate set.  
 $MFCs = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$

$S_i \rightarrow$  infrequent at  $i^{th}$  level.

Let  $L_1 = \{2, 3, 4, 5, 6, 7, 8\}$ .

$S_1 = \{1, 9\}$

We keep record of infrequent to see that they are not part of MFC.

Keep resizing MFCs.

Trimming down MFCs is topdown approach.

MFCs and  $S_i$  are helping to bring in the top down flavour to the trace.

$S_1 = \{1, 9\}$ .  $MFCs = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$   
 for  $\{1\}$  in  $S_1$  :  $\{2, 3, 4, 5, 6, 7, 8, 9\} = MFCs$   
 $MFCs \setminus S_1$  minus in set theory.  
 for  $\{9\}$  in  $S_1$  :  $\{2, 3, 4, 5, 6, 7, 8\} = MFCs$ .  
 $MFCs \setminus S_1$ .

$C_2 = L_1 \times L_1$

$L_2 = \{23, 24, 35, 37, 56, 57, 67\}$

11/2019  
 $S_2 = \{25, 26, 27, 28, 34, 36, 38, 45, 46, 47, 48, 58, 68, 78\}$

Exhaust  $S_i$

smf18@iitd.ac.in.

Communique

2 & 5 are infrequent together

'26' is occurring together only here.

(25)  $S_2 \rightarrow MFCs = \{345678, 234678\}$

(26)  $S_2 \rightarrow MFCs = \{345678, 34678, 23478\}$ . trim down only if they are occurring together

it's already in

So do not consider '34678'

updated MFCs =  $\{23478, 345678\}$

(29)  $S_2 \rightarrow MFCs = \{3478, 2348, 345678\}$

this is a subset of

$MFCs = \{2348, 345678\}$

(28)  $S_2 \rightarrow MFCs = \{348, 234, 345678\}$

subset of

$MFCs = \{234, 345678\}$

(30)  $S_2 \rightarrow MFCs = \{23, 24, 45678, 35678\}$

continue for the other  $S_i$ 's

(31)  $S_2 \rightarrow MFCs = \{23, 24, 357, 5678, 8\}$

$C_3 = \{234, 3578, 567\}$

Topdown will stop if  $S_i$  is null set. There is no more trimming down of MFCs

Here  $S_3$  is a null set.

MFC is a long reconstruction of FF. We do not have the sup count. So, we cannot have the confidence.



## LAB-2.

Implement A close and Pincer Search.  
Two move algs. MAFIA & CHARM.

2019

- Difference b/w BigData and Data Mining.  
How well may how well the storage systems can handle the distributed data. Incremental learning (with added data, the system must be intelligent enough to relearn the patterns). Storage issues according to scaling up of data must be addressed. Scaling up is done across networks.

Online DB  $\equiv$  Incremental learning.

Bucket Brigade  $\rightarrow$  Genetic Algorithms

↓  
evolutionary computation.  
↓  
evolve for the good.

GA is one evolutionary computation strategy.  
Tabu search, PSO, Ant colony are others.

↓  
particle  
swarm  
optimization.

Reinforcement learning thrives on penalty & awards.  
Most evolutionary computation algs work on Initial population.

All the evolutionary computation methods are basically for optimization (on the computers front).

Bother for optimum after seeing if the solution is feasible or not.

David Goldberg  $\rightarrow$  GA book (1989)  
James Holland  $\rightarrow$  inventor of GA.

Communiqué

- Enumerative search  $\rightarrow$  point by point search.  
Also called try and error. (but it is random)
- Operations Research  $\rightarrow$  mathematical optimization methods.
- Derivative based  $\rightarrow$  Gradient ascent or descent.  
Hill climbing is also another word. A slope should always exist. But real life data is not always continuous. Also, the data can be multimodal, the hill climbing gets stuck in local optimum.

In the search for optimal value, the path towards optima is also important.

Where will GA fit in?  $\rightarrow$  Classification.

If we are not able to give unique class labels in Dtree, it is non feasible.

Convergence  $\rightarrow$  how quickly will it reach the best solution?

Efficiency vs efficacy.

28/1/2019

- Why EVOLUTIONARY?

Can my overall objective keep bettering. GA, Tabu Search etc are search algs that search for optimal.

Fitness of a Solution  $\rightarrow$  how good is a solution.

iteration  $\equiv$  Generation.

As we do more generations, the fitness should improve.