

COMP 345: Data Mining

Mining Frequent Patterns, Associations and Correlations

Slides Adapted From : Jiawei Han, Micheline Kamber & Jian Pei
Data Mining: Concepts and Techniques, 3rd ed.



1

Reminders

- Assignment 3 – due Wed. Sept. 19th /Thurs. Sept. 20th
- Group Project Proposals – due Monday, Oct. 1st /Tues. Oct. 2nd

2

What Is Frequent Pattern Analysis?

- **Frequent pattern**: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set
- First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of **frequent itemsets** and **association rule mining**
- Motivation: Finding inherent regularities in data
 - What products were often purchased together?— Beer and diapers?!
 - What are the subsequent purchases after buying a PC?
 - What kinds of DNA are sensitive to this new drug?
 - Can we automatically classify web documents?
- Applications
 - Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis.

3

Apriori Example

There's a haunted corn maze every Fall that sells lots of fun things. Items that can be purchased are:

1. Hot Cider, 2. Pumpkin, 3. Gourd, 4. Hayride, 5. Maze Tour

You are given the transaction data for a morning of sales. (Items referred to by #).

Sales ID	List of Item IDs	Sales ID	List of item IDs
Order 1	1, 2, 5	Order 6	2, 3
Order 2	2, 4	Order 7	1, 3
Order 3	2, 3	Order 8	1, 2, 3, 5
Order 4	1, 2, 4	Order 9	1, 2, 3
Order 5	1, 3		

Assuming that minimum support = $2/9$ (.222) and minimum confidence is $7/9$ (.777)

1. Apply the Apriori algorithm to the dataset and identify **all** frequent k-itemsets
2. Find all **strong** association rules of the form $X \wedge Y \rightarrow Z$

4

Pattern-Growth Approach: Mining Frequent Patterns Without Candidate Generation

- Bottlenecks of the Apriori approach
 - Breadth-first (i.e., level-wise) search
 - Candidate generation and test
 - Often generates a huge number of candidates
- The FPGrowth Approach (J. Han, J. Pei, and Y. Yin, SIGMOD' 00)
 - Depth-first search
 - Avoid explicit candidate generation
- Major philosophy: Grow long patterns from short ones using local frequent items only
 - “abc” is a frequent pattern
 - Get all transactions having “abc”, i.e., project DB on abc: DB|abc
 - “d” is a local frequent item in DB|abc → abcd is a frequent pattern

5

Construct FP-tree from a Transaction Database

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o, w}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

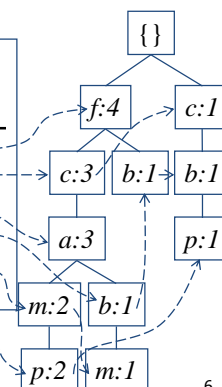
min_support = 3

1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Sort frequent items in frequency descending order, f-list
3. Scan DB again, construct FP-tree

Header Table

<i>Item</i>	<i>frequency</i>	<i>head</i>
f	4	
c	4	
a	3	
b	3	
m	3	
p	3	

F-list = f-c-a-b-m-p



6

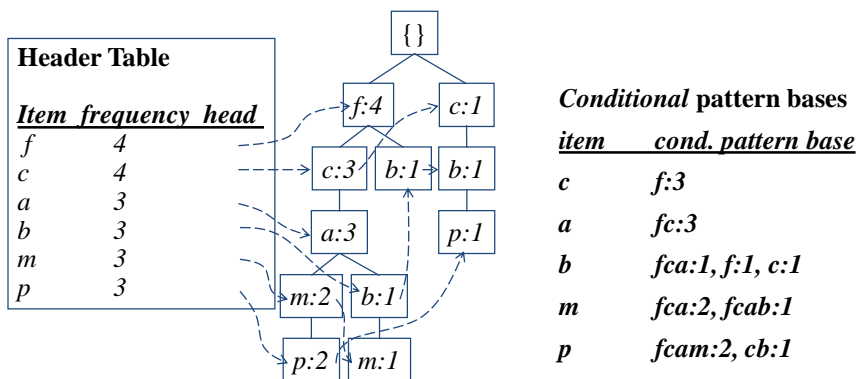
Partition Patterns and Databases

- Frequent patterns can be partitioned into subsets according to f-list
 - F-list = f-c-a-b-m-p
 - Patterns containing p
 - Patterns having m but no p
 - ...
 - Patterns having c but no a nor b, m, p
 - Pattern f
- Completeness and non-redundancy

7

Find Patterns Having P From P-conditional Database

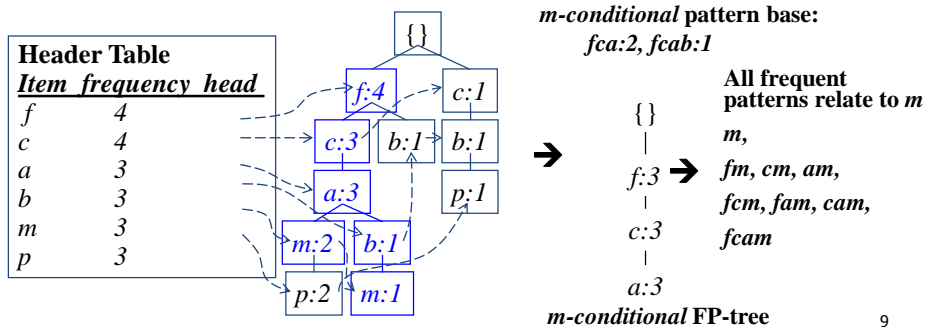
- Starting at the frequent item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item p
- Accumulate all of *transformed prefix paths* of item p to form p 's conditional pattern base



8

From Conditional Pattern-bases to Conditional FP-trees

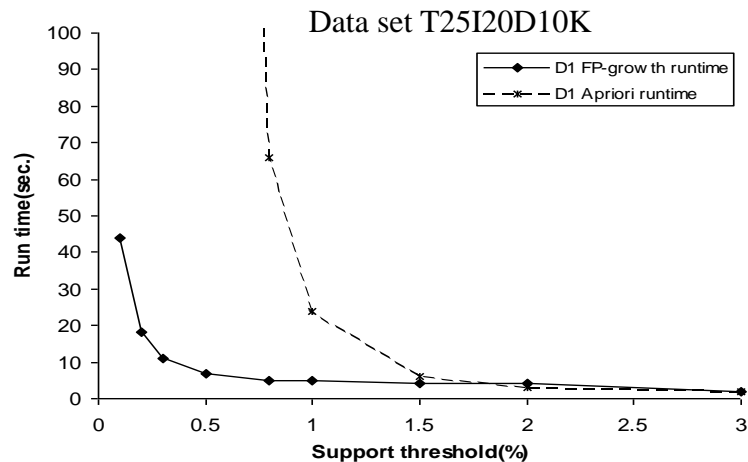
- For each pattern-base
 - Accumulate the count for each item in the base
 - Construct the FP-tree for the frequent items of the pattern base



Benefits of the FP-tree Structure

- Completeness
 - Preserve complete information for frequent pattern mining
 - Never break a long pattern of any transaction
- Compactness
 - Reduce irrelevant info—infrequent items are gone
 - Items in frequency descending order: the more frequently occurring, the more likely to be shared
 - Never be larger than the original database (not counting node-links and the *count* field)

FP-Growth vs. Apriori: Scalability With the Support Threshold



11

FP-Growth Example

We're having a picnic and you're asked to pick up a couple of items to bring along. Items are Hot Dogs, Buns, Ketchup, Chips, and Coke. Here are the transactions from each purchase

Transaction ID	List of Item IDs
T1	HotDogs, Buns, Ketchup
T2	HotDogs, Buns
T3	HotDogs, Coke, Chips
T4	Chips, Coke
T5	Chips, Ketchup
T6	HotDogs, Coke, Chips

Assuming that minimum support = 33.3% (2) and minimum confidence is 60%

1. Apply the FP-growth algorithm to the dataset and identify **all** frequent k-itemsets
2. Find all **strong** association rules.

12

Interestingness Measure: Correlations (Lift)

- *play basketball* \Rightarrow *eat cereal* [40%, 66.7%] is misleading
 - The overall % of students eating cereal is 75% > 66.7%.
- *play basketball* \Rightarrow *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence
- Measure of dependent/correlated events: **lift**

$$lift = \frac{P(A \cup B)}{P(A)P(B)}$$

$$lift(B, C) = \frac{2000 / 5000}{3000 / 5000 * 3750 / 5000} = 0.89$$

$$lift(B, \neg C) = \frac{1000 / 5000}{3000 / 5000 * 1250 / 5000} = 1.33$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

13

Are lift and χ^2 Good Measures of Correlation?

- “Buy walnuts \Rightarrow buy milk [1%, 80%]” is misleading if 85% of customers buy milk
- Support and confidence are not good to indicate correlations
- Over 20 interestingness measures have been proposed (see Tan, Kumar, Sritastava @KDD'02)
- Which are good ones?

symbol	measure	range	formula
ϕ	ϕ -coefficient	-1 ... 1	$\frac{P(A, B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
Q	Yule's Q	-1 ... 1	$\frac{P(A, B)P(\bar{A}, \bar{B}) - P(A, \bar{B})P(\bar{A}, B)}{P(A, B)P(\bar{A}, \bar{B}) + P(A, \bar{B})P(\bar{A}, B)}$
Y	Yule's Y	-1 ... 1	$\frac{\sqrt{P(A, B)P(\bar{A}, \bar{B})} - \sqrt{P(A, \bar{B})P(\bar{A}, B)}}{\sqrt{P(A, B)P(\bar{A}, \bar{B})} + \sqrt{P(A, \bar{B})P(\bar{A}, B)}}$
k	Cohen's	-1 ... 1	$\frac{P(A, B) + P(\bar{A}, \bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{P(A, B) + P(\bar{A}, \bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}$
PS	Piatetsky-Shapiro's	-0.25 ... 0.25	$P(A, B) - P(A)P(B)$
F	Certainty factor	-1 ... 1	$\max(\frac{P(A, B) - P(A)P(B)}{1 - P(B)}, \frac{P(A, B) - P(A)}{1 - P(A)})$
AV	added value	-0.5 ... 1	$\max(P(B A) - P(B), P(A B) - P(A))$
K	Klogsen's Q	-0.33 ... 0.38	$\frac{\sqrt{P(A, B) \max(P(B A) - P(B), P(A B) - P(A))}}{\sum_j \max_k P(A_j, B_k) + \sum_j \max_k P(A_j, \bar{B}_k) - \max_k P(A_j) - \max_k P(\bar{B}_k)}$
g	Goodman-kruskal's	0 ... 1	$\frac{\sum_j \max_k P(A_j, B_k) - \max_k P(B_k)}{2 - \max_k P(A_j) - \max_k P(\bar{B}_k)}$
M	Mutual Information	0 ... 1	$\frac{\sum_j \sum_k P(A_j, B_k) \log \frac{P(A_j, B_k)}{P(A_j)P(B_k)}}{\sum_j \sum_k P(A_j, B_k) \log \frac{P(A_j, B_k)}{P(A_j)P(B_k)} + \sum_j \sum_k P(A_j, \bar{B}_k) \log \frac{P(A_j, \bar{B}_k)}{P(A_j)P(\bar{B}_k)}}$
J	J-Measure	0 ... 1	$\max(P(A, B) \log(\frac{P(A, B)}{P(A)P(B)}) + P(\bar{A}, \bar{B}) \log(\frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})}))$
G	Gini index	0 ... 1	$\frac{P(A, B) \log(\frac{P(A, B)}{P(A)P(B)}) + P(\bar{A}, \bar{B}) \log(\frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})})}{\max(P(A)P(B)(A)^2 + P(\bar{B})(B)^2 + P(\bar{A})P(B)(\bar{A})^2 + P(B)(\bar{B})^2 - P(B)^2 - P(\bar{B})^2, P(B)(P(A)^2 + P(\bar{A})B)^2 + P(\bar{B})(P(A)\bar{B})^2 + P(\bar{A})P(\bar{B})^2 - P(A)^2 - P(\bar{A})^2)}$
s	support	0 ... 1	$P(A, B)$
c	confidence	0 ... 1	$\max(P(B A), P(A B))$
L	Laplace	0 ... 1	$\max(\frac{N P(A, B) + 1}{N P(A) + 2}, \frac{N P(A, B) + 1}{N P(B) + 2})$
IS	Cosine	0 ... 1	$\frac{P(A, B)}{\sqrt{P(A)P(B)}}$
γ	coherence(Jaccard)	0 ... 1	$\frac{P(A, B)}{P(A) + P(B) - P(A, B)}$
α	all confidence	0 ... 1	$\frac{\max(P(A), P(B))}{P(A, B)}$
o	odds ratio	0 ... ∞	$\frac{P(A, B)P(\bar{A}, \bar{B})}{P(A, \bar{B})P(\bar{A}, B)}$
V	Conviction	0.5 ... ∞	$\max(\frac{P(A)P(\bar{B})}{P(A, B)}, \frac{P(B)P(\bar{A})}{P(A, B)})$
λ	lift	0 ... ∞	$\frac{P(A, B)}{P(A)P(B)}$
S	Collective strength	0 ... ∞	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A, B) - P(\bar{A}, \bar{B})}$
χ^2	χ^2	0 ... ∞	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{\sum_j (\frac{P(A_j, B_k)^2}{P(A_j)P(B_k)})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A, B) - P(\bar{A}, \bar{B})}$

14

Summary

- Basic concepts: association rules, support-confident framework, closed and max-patterns
- Scalable frequent pattern mining methods
 - Apriori (Candidate generation & test)
 - Projection-based (FP-growth)
- Which patterns are interesting?
 - Pattern evaluation methods