CS 521 Data Mining Techniques Instructor: Abdullah Mueen

LECTURE 4: FREQUENT PATTERN MINING



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.

Why Is Freq. Pattern Mining Important?

Freq. pattern: An intrinsic and important property of datasets

Foundation for many essential data mining tasks

- Association, correlation, and causality analysis
- Sequential, structural (e.g., sub-graph) patterns
- Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
- Classification: discriminative, frequent pattern analysis
- Cluster analysis: frequent pattern-based clustering
- Data warehousing: iceberg cube and cube-gradient
- Semantic data compression: fascicles
- Broad applications

Basic Concepts: Frequent Patterns

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



itemset: A set of one or more items

k-itemset
$$X = \{x_1, ..., x_k\}$$

(absolute) support, or, *support count* of X: Frequency or occurrence of an itemset X

(relative) support, s, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)

An itemset X is *frequent* if X's support is no less than a *minsup* threshold

Basic Concepts: Association Rules

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



Find all the rules $X \rightarrow Y$ with minimum support and confidence

- support, s, probability that a transaction contains $X \cup Y$
- confidence, c, conditional probability that a transaction having X also contains Y

Let minsup = 50%, minconf = 50%

Freq. Pat.: Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

- Association rules: (many more!)
 - Beer \rightarrow Diaper (60%, 100%)
 - Diaper \rightarrow Beer (60%, 75%)

Closed Patterns and Max-Patterns

A long pattern contains a combinatorial number of sub-patterns, e.g., $\{a_1, ..., a_{100}\}$ contains $\binom{100}{1} + \binom{100}{2} + ... + \binom{1}{1} \binom{0}{0} = 2^{100} - 1 = 1.27 \times 10^{30}$ sub-patterns!

Solution: *Mine closed patterns and max-patterns instead*

An itemset X is closed if X is *frequent* and there exists *no super-pattern* Y > X, *with the same support* as X (proposed by Pasquier, et al. @ ICDT'99)

An itemset X is a max-pattern if X is frequent and there exists no frequent super-pattern Y > X (proposed by Bayardo @ SIGMOD'98)

Closed pattern is a lossless compression of freq. patterns

• Reducing the # of patterns and rules

Closed Patterns and Max-Patterns

Exercise: Suppose a DB contains only two transactions

- <a₁, ..., a₁₀₀>, <a₁, ..., a₅₀>
- Let min_sup = 1

What is the set of closed itemset?

- {a₁, ..., a₁₀₀}: 1
- {a₁, ..., a₅₀}: 2

What is the set of max-pattern?

• {a₁, ..., a₁₀₀}: 1

What is the set of all patterns?

- {a₁}: 2, ..., {a₁, a₂}: 2, ..., {a₁, a₅₁}: 1, ..., {a₁, a₂, ..., a₁₀₀}: 1
- A big number: 2¹⁰⁰ 1? Why?

Chapter 5: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

Basic Concepts



Frequent Itemset Mining Methods

Which Patterns Are Interesting?—Pattern Evaluation Methods

Summary

Scalable Frequent Itemset Mining Methods

Apriori: A Candidate Generation-and-Test Approach

Improving the Efficiency of Apriori

FPGrowth: A Frequent Pattern-Growth Approach

ECLAT: Frequent Pattern Mining with Vertical Data Format

\wedge	
∇	

The Downward Closure Property and Scalable Mining Methods

The downward closure property of frequent patterns

- Any subset of a frequent itemset must be frequent
- If {beer, diaper, nuts} is frequent, so is {beer, diaper}
- i.e., every transaction having {beer, diaper, nuts} also contains {beer, diaper}

Scalable mining methods: Three major approaches

- Apriori (Agrawal & Srikant@VLDB'94)
- Freq. pattern growth (FPgrowth—Han, Pei & Yin @SIGMOD'00)
- Vertical data format approach (Charm—Zaki & Hsiao @SDM'02)

Apriori: A Candidate Generation & Test Approach

Apriori pruning principle: If there is any itemset which is infrequent, its superset should not be generated/tested! (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)

Method:

- Initially, scan DB once to get frequent 1-itemset
- Generate length (k+1) candidate itemsets from length k frequent itemsets
- Test the candidates against DB
- Terminate when no frequent or candidate set can be generated

The Apriori Algorithm—An Example



The Apriori Algorithm (Pseudo-Code)

 C_k : Candidate itemset of size k L_k : frequent itemset of size k

 $\begin{array}{l} \mathcal{L}_{1} = \{ \text{frequent items} \}; \\ \text{for } (k = 1; \mathcal{L}_{k} \mid = \varnothing; k + +) \text{ do begin} \\ \mathcal{C}_{k+1} = \text{candidates generated from } \mathcal{L}_{k}; \\ \text{for each transaction } t \text{ in database do} \\ \text{increment the count of all candidates in } \mathcal{C}_{k+1} \text{ that are contained in } t \\ \mathcal{L}_{k+1} = \text{candidates in } \mathcal{C}_{k+1} \text{ with min_support} \\ \text{end} \end{array}$

return $\cup_k L_k$;

Implementation of Apriori

How to generate candidates?

- Step 1: self-joining *L*_k
- Step 2: pruning

Example of Candidate-generation

- L₃={abc, abd, acd, ace, bcd}
- Self-joining: L₃*L₃
 - *abcd* from *abc* and *abd*
 - acde from acd and ace
- Pruning:
 - *acde* is removed because *ade* is not in L_3
- $C_4 = \{abcd\}$

Further Improvement of the Apriori Method

Major computational challenges

- Multiple scans of transaction database
- Huge number of candidates
- Tedious workload of support counting for candidates

Improving Apriori: general ideas

- Reduce passes of transaction database scans
- Shrink number of candidates
- Facilitate support counting of candidates

Partition: Scan Database Only Twice

Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB

- Scan 1: partition database and find local frequent patterns
- Scan 2: consolidate global frequent patterns

A. Savasere, E. Omiecinski and S. Navathe, VLDB'95



DHP: Reduce the Number of Candidates

A *k*-itemset whose corresponding hashing bucket count is below the threshold cannot be frequent

- Candidates: a, b, c, d, e
- Hash entries
 - {ab, ad, ae}
 - {bd, be, de}
 - ...
- Frequent 1-itemset: a, b, d, e
- ab is not a candidate 2-itemset if the sum of count of {ab, ad, ae} is below support threshold

J. Park, M. Chen, and P. Yu. An effective hash-based algorithm for mining association rules. *SIGMOD'95*

count	itemsets	
35	{ab, ad, ae}	
88	{bd, be, de}	
•		
102	{yz, qs, wt}	
Hach Tablo		

Sampling for Frequent Patterns

Select a sample of original database, mine frequent patterns within sample using Apriori

Scan database once to verify frequent itemsets found in sample, only *borders* of closure of frequent patterns are checked

• Example: check *abcd* instead of *ab, ac, ..., etc.*

Scan database again to find missed frequent patterns

H. Toivonen. Sampling large databases for association rules. In VLDB'96

DIC: Reduce Number of Scans



Dynamic itemset counting and implication rules for market basket data. In SIGMOD'97 Once both A and D are determined frequent, the counting of AD begins

Once all length-2 subsets of BCD are determined frequent, the counting of BCD begins



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 \rightarrow abcd is a frequent pattern

Construct FP-tree from a Transaction Database



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-redundency

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



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





Cond. pattern base of "cam": (f:3)
$$\begin{cases} \\ f:3 \\ f:3 \end{cases}$$

cam-conditional **FP-tree**

A Special Case: Single Prefix Path in FP-tree

Suppose a (conditional) FP-tree T has a shared single prefixpath P

Mining can be decomposed into two parts

Reduction of the single prefix path into one node

 $a_2:n_2$

{ }

 $a_1:n_1$

Concatenation of the mining results of the two parts



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 count node-links and the count field)

The Frequent Pattern Growth Mining Method

Idea: Frequent pattern growth

Recursively grow frequent patterns by pattern and database partition

Method

- For each frequent item, construct its conditional pattern-base, and then its conditional FP-tree
- Repeat the process on each newly created conditional FP-tree
- Until the resulting FP-tree is empty, or it contains only one path—single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

Advantages of the Pattern Growth Approach

Divide-and-conquer:

- Decompose both the mining task and DB according to the frequent patterns obtained so far
- Lead to focused search of smaller databases

Other factors

- No candidate generation, no candidate test
- Compressed database: FP-tree structure
- No repeated scan of entire database
- Basic ops: counting local freq items and building sub FP-tree, no pattern search and matching

A good open-source implementation and refinement of FPGrowth

• FPGrowth+ (Grahne and J. Zhu, FIMI'03)

Mining Frequent Closed Patterns: CLOSET

Flist: list of all frequent items in support ascending order

• Flist: d-a-f-e-c

Divide search space

- Patterns having d
- Patterns having d but no a, etc.

Find frequent closed pattern recursively

• Every transaction having d also has $cfa \rightarrow cfad$ is a frequent closed pattern

J. Pei, J. Han & R. Mao. "CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets", DMKD'00.

Min_sup=2			
TID	Items		
10	a, c, d, e, f		
20	a, b, e		
30	c, e, f		
40	a, c, d, f		
50	c, e, f		

CLOSET+: Mining Closed Itemsets by Pattern-Growth

Itemset merging: if Y appears in every occurrence of X, then Y is merged with X

Sub-itemset pruning: if $Y \supset X$, and sup(X) = sup(Y), X and all of X's descendants in the set enumeration tree can be pruned

Hybrid tree projection

- Bottom-up physical tree-projection
- Top-down pseudo tree-projection

Item skipping: if a local frequent item has the same support in several header tables at different levels, one can prune it from the header table at higher levels

Efficient subset checking

MaxMiner: Mining Max-Patterns



Since BCDE is a max-pattern, no need to check BCD, BDE, CDE in later scan

R. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98

Computational Complexity of Frequent Itemset Mining

How many itemsets are potentially to be generated in the worst case?

- The number of frequent itemsets to be generated is senstive to the minsup threshold
- When minsup is low, there exist potentially an exponential number of frequent itemsets
- The worst case: M^N where M: # distinct items, and N: max length of transactions

The worst case complexty vs. the expected probability

- Ex. Suppose Walmart has 10⁴ kinds of products
 - The chance to pick up one product 10⁻⁴
 - The chance to pick up a particular set of 10 products: ~10⁻⁴⁰
 - What is the chance this particular set of 10 products to be frequent 10³ times in 10⁹ transactions?

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

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

 $lift(B,\neg C) = \frac{1000/5000}{3000/5000*1250/5000} = 1.33$ $lift(B,C) = \frac{2000/5000}{3000/5000*3750/5000} = 0.89$

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

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	-11	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1 - P(A))(1 - P(B))}}$
Q	Yule's Q	-11	$\frac{\dot{P}(A,B)P(\overline{A},\overline{B}) - P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{A},\overline{B}) + P(A,\overline{B})P(\overline{A},B)}$
Y	Yule's Y	-11	$\frac{\sqrt{P(A,B)P(\overline{A},\overline{B})} - \sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{A},\overline{B})} + \sqrt{P(A,\overline{B})P(\overline{A},B)}}$
k	Cohen's	-11	$\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$
PS	Piatetsky-Shapiro's	-0.250.25	P(A, B) - P(A)P(B)
F	Certainty factor	-11	$\max(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)})$
AV	added value	$-0.5 \dots 1$	$\max(P(B A) - P(B), P(A B) - P(A))$
K	Klosgen's Q	-0.330.38	$\sqrt{P(A,B)}\max(P(B A) - P(B), P(A B) - P(A))$
g	Goodman-kruskal's	$0 \dots 1$	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
M	Mutual Information	$0 \dots 1$	$\frac{\sum_{i} \sum_{j} P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i) P(B_j)}}{\min(-\sum_{i} P(A_i) \log P(A_i) \log P(A_i) - \sum_{i} P(B_i) \log P(B_i) \log P(B_i))}$
J	J-Measure	$0 \dots 1$	$\max(P(A, B)\log(\frac{P(B A)}{P(B)}) + P(A\overline{B})\log(\frac{P(B A)}{P(B)}))$
			$P(A, B) \log(\frac{P(A B)}{P(A)}) + P(\overline{A}B) \log(\frac{P(\overline{A} B)}{P(\overline{A})})$
G	Gini index	$0 \dots 1$	$\max(P(A)[P(B A)^2 + P(\overline{B} A)^2] + P(\overline{A}[P(B \overline{A})^2 + P(\overline{B} \overline{A})^2] - P(B)^2 - P(\overline{B})^2,$
			$P(B)[P(A B)^{2} + P(\overline{A} B)^{2}] + P(\overline{B}[P(A \overline{B})^{2} + P(\overline{A} \overline{B})^{2}] - P(A)^{2} - P(\overline{A})^{2})$
s	support	01	P(A,B)
	confidence	01	$\max(P(B A), P(A B))$ $(NP(A,B)+1, NP(A,B)+1)$
	Laplace	01	$\max(\frac{-NP(A)+2}{NP(A)+2}, \frac{-NP(B)+2}{NP(B)+2})$
IS	Cosine	$0 \dots 1$	$\frac{\Gamma(A,B)}{\sqrt{P(A)P(B)}}$
γ	$\operatorname{coherence}(\operatorname{Jaccard})$	$0 \dots 1$	$\frac{P(A,B)}{P(A) \pm P(B) - P(A,B)}$
α	$all_confidence$	$0 \dots 1$	$\frac{P(A,B)}{\max(P(A),P(B))}$
0	odds ratio	$0 \dots \infty$	$\frac{P(A,B)P(\overline{A},\overline{B})}{P(\overline{A},B)P(A,\overline{B})}$
V	Conviction	$0.5 \ldots \infty$	$\max\left(\frac{P(A)P(\overline{B})}{P(A\overline{B})}, \frac{P(B)P(\overline{A})}{P(B\overline{A})}\right)$
λ	lift	$0 \dots \infty$	$\frac{P(A,B)}{P(A)P(B)}$
S	Collective strength	$0 \dots \infty$	$\frac{P(A,B)+P(AB)}{P(A)P(B)+P(\overline{A})P(\overline{B})} \times \frac{1-P(A)P(B)-P(A)P(B)}{1-P(A,B)-P(\overline{AB})}$
χ^2	χ^2	$0 \dots \infty$	$\sum_{i} \frac{(P(A_i) - E_i)^2}{E_i}$

Summary

Basic concepts: association rules, support-confident framework, closed and max-patterns

Scalable frequent pattern mining methods

- Apriori (Candidate generation & test)
- Projection-based (FPgrowth, CLOSET+, ...)
- Which patterns are interesting?
 - Pattern evaluation methods