

Data Mining Techniques: Frequent Patterns in Sets and Sequences

Mirek Riedewald

Some slides based on presentations by
Han/Kamber and Tan/Steinbach/Kumar

Frequent Pattern Mining Overview

- Basic Concepts and Challenges
- Efficient and Scalable Methods for Frequent Itemsets and Association Rules
- Pattern Interestingness Measures
- Sequence Mining

2

What Is Frequent Pattern Analysis?

- Find patterns (itemset, sequence, structure, etc.) that occur frequently in a data set
- First proposed for frequent itemsets and association rule mining
- Motivation: Find inherent regularities in data
 - What products were often purchased together?
 - What are the subsequent purchases after buying a PC?
 - What kinds of DNA are sensitive to a new drug?
- Applications
 - Market basket analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, DNA sequence analysis

3

Association Rule Mining

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Association Rules

(Diaper) → {Beer},
{Milk, Bread} → {Eggs,Coke},
{Beer, Bread} → {Milk},

Implication means co-occurrence,
not causality!

4

Definition: Frequent Itemset

- Itemset
 - A collection of one or more items
 - Example: {Milk, Bread, Diaper}
 - k-itemset: itemset that contains k items
- Support count (σ)
 - Frequency of occurrence of an itemset
 - E.g., $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- Support (s)
 - Fraction of transactions that contain an itemset
 - E.g., $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$
- Frequent Itemset
 - An itemset whose support is greater than or equal to a **minsup** threshold

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

5

Definition: Association Rule

- Association Rule = implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
 - Ex.: {Milk, Diaper} → {Beer}

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

- Rule Evaluation Metrics

- Support (s) = $P(X \cup Y)$
 - Estimated by fraction of transactions that contain both X and Y
- Confidence (c) = $P(Y | X)$
 - Estimated by fraction of transactions that contain X and Y among all transactions containing X

Example: {Milk, Diaper} → Beer

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|D|} = \frac{2}{5}$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3}$$

6

Association Rule Mining Task

- Given a transaction database DB, find all rules having support $\geq \text{minsup}$ and confidence $\geq \text{minconf}$
- Brute-force approach:
 - List all possible association rules
 - Compute support and confidence for each rule
 - Remove rules that fail the minsup or minconf thresholds
 - Computationally prohibitive!

7

Mining Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example rules:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$ ($s=0.4, c=0.67$)
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$ ($s=0.4, c=0.5$)
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$ ($s=0.4, c=0.5$)

Observations:

- All the above rules are binary partitions of the same itemset {Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

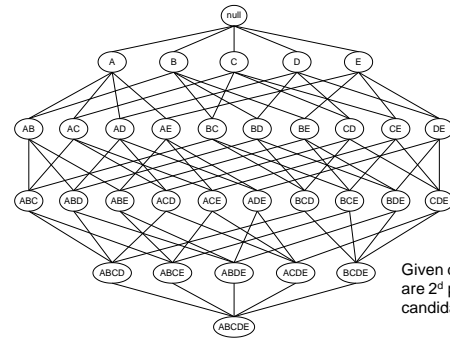
8

Mining Association Rules

- Two-step approach:
 - Frequent Itemset Generation
 - Generate all itemsets that have support $\geq \text{minsup}$
 - Rule Generation
 - Generate high-confidence rules from each frequent itemset, where each rule is a binary partitioning of the frequent itemset
- Frequent itemset generation is still computationally expensive

9

Frequent Itemset Generation

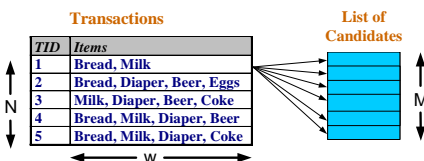


Given d items, there are 2^d possible candidate itemsets

10

Frequent Itemset Generation

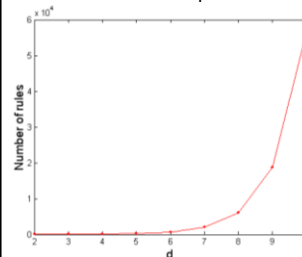
- Brute-force approach:
 - Each itemset in the lattice is a candidate frequent itemset
 - Count the support of each candidate by scanning the database
 - Match each transaction against every candidate
 - Complexity $\approx O(N * M * w) \Rightarrow$ expensive since $M=2^d$



11

Computational Complexity

- Given d unique items, total number of itemsets = 2^d
- Total number of possible association rules?



$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \cdot \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$

$$= 3^d - 2^{d+1} + 1$$

If $d=6$, $R = 602$ possible rules

12

Frequent Pattern Mining Overview

- Basic Concepts and Challenges
- Efficient and Scalable Methods for Frequent Itemsets and Association Rules
- Pattern Interestingness Measures
- Sequence Mining

13

Frequent Itemset Generation Strategies

- Reduce the **number of candidates** (M)
 - Complete search: $M=2^d$
 - Use pruning techniques to reduce M
- Reduce the **number of transactions** (N)
 - Skip short transactions as size of itemset increases
- Reduce the **number of comparisons** ($N*M$)
 - Use efficient data structures to store the candidates or transactions
 - No need to match every candidate against every transaction

14

Reducing Number of Candidates

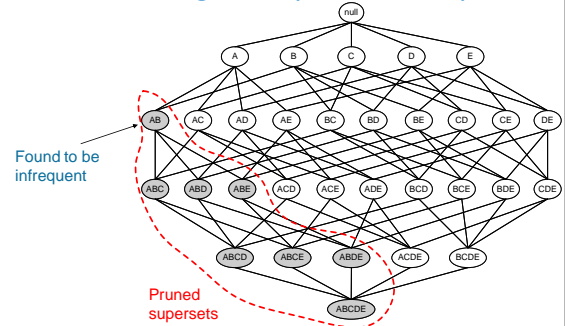
- **Apriori principle:**
 - If an itemset is frequent, then all of its subsets must also be frequent
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the **anti-monotone** property of support

15

Illustrating the Apriori Principle



16

Illustrating the Apriori Principle

Item	Count
Bread	4
Coke	4
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)

Itemset	Count
{Bread, Milk}	3
{Bread, Beer}	2
{Bread, Diaper}	3
{Milk, Beer}	2
{Milk, Diaper}	3
{Beer, Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$
 With support-based pruning,
 $6 + 6 + 1 = 13$

Itemset	Count
{Bread, Milk, Diaper}	3

Triples (3-itemsets)

17

Apriori Algorithm

- Generate L_1 = frequent itemsets of length $k=1$
- Repeat until no new frequent itemsets are found
 - Generate C_{k+1} , the length- $(k+1)$ candidate itemsets, from L_k
 - Prune candidate itemsets in C_{k+1} containing subsets of length k that are not in L_k (and hence infrequent)
 - Count support of each remaining candidate by scanning DB; eliminate infrequent ones from C_{k+1}
 - $L_{k+1} = C_{k+1}$; $k = k+1$

18

Important Details of Apriori

- How to generate candidates?
 - Step 1: self-joining L_k
 - Step 2: pruning
- How to count support of candidates?
- Example of Candidate-generation for $L_3 = \{ \{a,b,c\}, \{a,b,d\}, \{a,c,d\}, \{a,c,e\}, \{b,c,d\} \}$
 - Self-joining L_3
 - $\{a,b,c,d\}$ from $\{a,b,c\}$ and $\{a,b,d\}$
 - $\{a,c,d,e\}$ from $\{a,c,d\}$ and $\{a,c,e\}$
 - Pruning:
 - $\{a,c,d,e\}$ is removed because $\{a,d,e\}$ is not in L_3
 - $C_4 = \{ \{a,b,c,d\} \}$

19

How to Generate Candidates?

- Step 1: self-joining L_{k-1}

```

insert into Ck
select p.item1, p.item2, ..., p.itemk-1, q.itemk-1
from Lk-1 p, Lk-1 q
where p.item1=q.item1 AND ... AND p.itemk-2=q.itemk-2
AND p.itemk-1 < q.itemk-1
            
```
- Step 2: pruning
 - for all itemsets c in C_k do
 - for all $(k-1)$ -subsets s of c do
 - if s is not in L_{k-1} then delete c from C_k

20

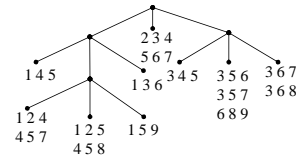
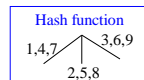
How to Count Supports of Candidates?

- Why is counting supports of candidates a problem?
 - Total number of candidates can be very large
 - One transaction may contain many candidates
- Method:
 - Candidate itemsets stored in a hash-tree
 - Leaf node contains list of itemsets
 - Interior node contains a hash table
 - Subset function finds all candidates contained in a transaction

21

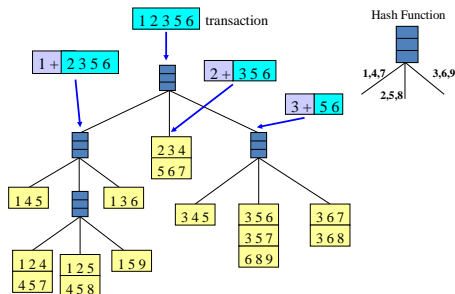
Generate Hash Tree

- Suppose we have 15 candidate itemsets of length 3:
 - $\{1,4,5\}, \{1,2,4\}, \{4,5,7\}, \{1,2,5\}, \{4,5,8\}, \{1,5,9\}, \{1,3,6\}, \{2,3,4\}, \{5,6,7\}, \{3,4,5\}, \{3,5,6\}, \{3,5,7\}, \{6,8,9\}, \{3,6,7\}, \{3,6,8\}$
- We need:
 - Hash function
 - Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)



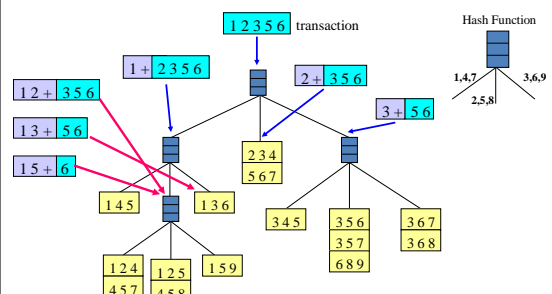
22

Subset Operation Using Hash Tree



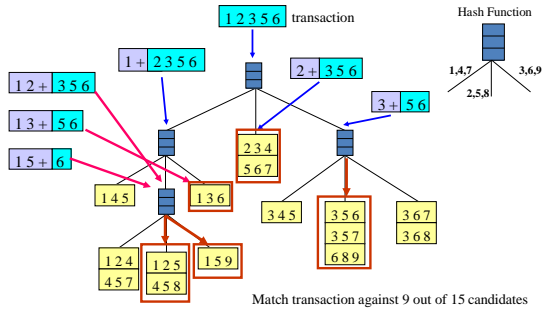
23

Subset Operation Using Hash Tree



24

Subset Operation Using Hash Tree



25

Association Rule Generation

- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the minimum confidence requirement
 - If $\{A, B, C, D\}$ is a frequent itemset, candidate rules are:
 - $ABC \rightarrow D$, $ABD \rightarrow C$, $ACD \rightarrow B$, $BCD \rightarrow A$,
 - $A \rightarrow BCD$, $B \rightarrow ACD$, $C \rightarrow ABD$, $D \rightarrow ABC$,
 - $AB \rightarrow CD$, $AC \rightarrow BD$, $AD \rightarrow BC$, $BC \rightarrow AD$,
 - $BD \rightarrow AC$, $CD \rightarrow AB$
- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

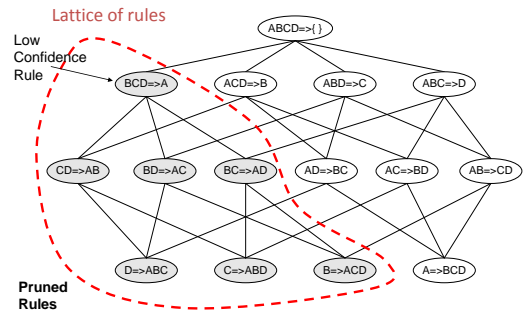
26

Rule Generation

- How do we efficiently generate association rules from frequent itemsets?
 - In general, confidence does not have an anti-monotone property
 - $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$
 - But confidence of rules **generated from the same itemset** has an anti-monotone property
 - For $\{A, B, C, D\}$, $c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$
 - Confidence is anti-monotone w.r.t. number of items on the right-hand side of the rule

27

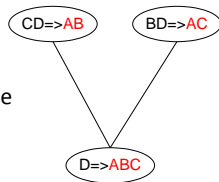
Rule Generation for Apriori Algorithm



28

Rule Generation for Apriori Algorithm

- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
 - Join($CD \rightarrow AB$, $BD \rightarrow AC$) would produce the candidate rule $D \rightarrow ABC$
- Prune rule $D \rightarrow ABC$ if its subset $AD \rightarrow BC$ does not have high confidence



29

Improving Apriori

- Challenges
 - Multiple scans of transaction database
 - Huge number of candidates
 - Tedious workload of support counting for candidates
- General ideas
 - Reduce passes of transaction database scans
 - Further shrink number of candidates
 - Facilitate support counting of candidates

30

Bottleneck of Frequent-Pattern Mining

- Apriori generates a very large number of candidates
 - 10^4 frequent 1-itemsets can result in more than 10^7 candidate 2-itemsets
 - Many candidates might have low support, or do not even exist in the database
- Apriori scans entire transaction database for every round of support counting
- Bottleneck:** candidate-generation-and-test
- Can we avoid candidate generation?

31

How to Avoid Candidate Generation

- Grow long patterns from short ones using **local frequent items**
 - Assume $\{a,b,c\}$ is a frequent pattern in transaction database DB
 - Get all transactions containing $\{a,b,c\}$
 - Notation: $DB|\{a,b,c\}$
 - $\{d\}$ is a local frequent item in $DB|\{a,b,c\}$, if and only if $\{a,b,c,d\}$ is a frequent pattern in DB

32

Construct FP-tree from a Transaction Database

TID	Items bought	(ordered) frequent items
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}

$min_support = 3$

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

Header Table	
Item	frequency head
f	4
c	4
a	3
b	3
m	3
p	3

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

33

Construct FP-tree from a Transaction Database

TID	Items bought	(ordered) frequent items
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}

$min_support = 3$

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

Header Table	
Item	frequency head
f	4
c	4
a	3
b	3
m	3
p	3

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

34

Construct FP-tree from a Transaction Database

TID	Items bought	(ordered) frequent items
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}

$min_support = 3$

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

Header Table	
Item	frequency head
f	4
c	4
a	3
b	3
m	3
p	3

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

35

Construct FP-tree from a Transaction Database

TID	Items bought	(ordered) frequent items
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}

$min_support = 3$

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

Header Table	
Item	frequency head
f	4
c	4
a	3
b	3
m	3
p	3

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

36

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 larger than the original database (if we do not count node-links and the count field)
 - For some example DBs, compression ratio over 100

37

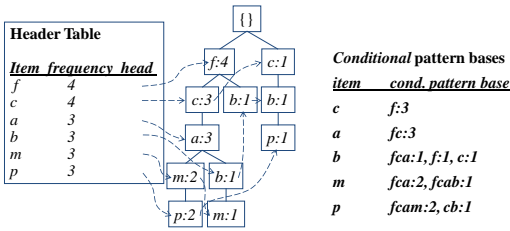
Partition Patterns and Databases

- 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 b, but neither m nor p
 - ...
 - Patterns having c, but neither a, b, m, nor p
 - Pattern f
- This partitioning is **complete** and **non-redundant**

38

Construct Conditional Pattern Base For Item X

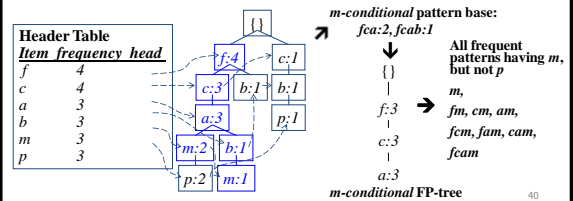
- Conditional pattern base = set of prefix paths in FP-tree that co-occur with x
- Traverse FP-tree by following link of frequent item x in header table
- Accumulate paths with their frequency counts



39

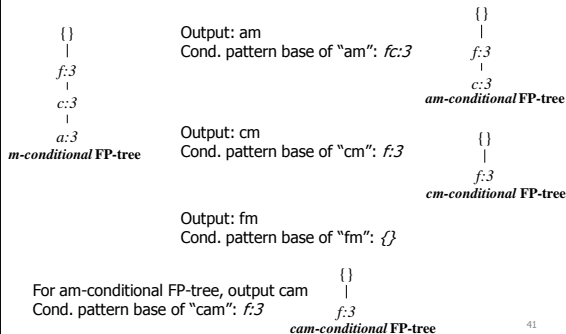
From Conditional Pattern Bases to Conditional FP-Trees

- 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



40

Recursion: Mining Conditional FP-Trees



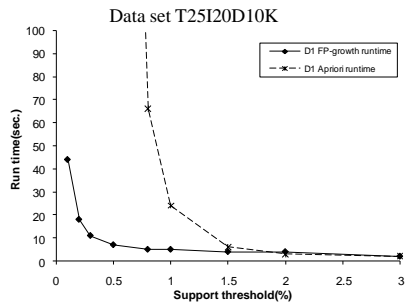
41

FP-Tree Algorithm Summary

- **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 recursively on each newly created conditional FP-tree
 - Stop recursion when resulting FP-tree is empty
 - Optimization if tree contains only one path: single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

42

FP-Growth vs. Apriori: Scalability With Support Threshold



43

Why Is FP-Growth the Winner?

- Divide-and-conquer
 - Decompose both the mining task and DB according to the frequent patterns obtained so far
 - Leads 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 operations: counting local frequent single items and building sub FP-tree
 - No pattern search and matching

44

Factors Affecting Mining Cost

- Choice of minimum support threshold
 - Lower support threshold => more frequent itemsets
 - More candidates, longer frequent itemsets
- Dimensionality (number of items) of the data set
 - More space needed to store support count of each item
 - If number of frequent items also increases, both computation and I/O costs may increase
- Size of database
 - Each pass over DB is more expensive
- Average transaction width
 - May increase max. length of frequent itemsets and traversals of hash tree (more subsets supported by transaction)
- How can we further reduce some of these costs?

45

Compact Representation of Frequent Itemsets

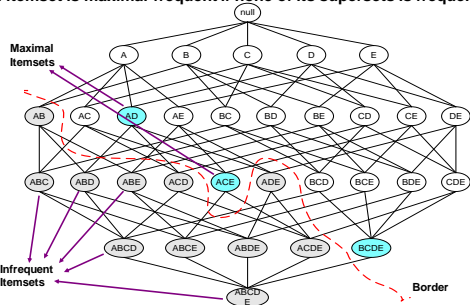
- Some itemsets are redundant because they have identical support as their supersets
- Number of frequent itemsets = $3 \times \sum_{k=1}^{10} \binom{10}{k}$
- Need a compact representation

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	

46

Maximal Frequent Itemset

An itemset is maximal-frequent if none of its supersets is frequent



47

Closed Itemset

- A frequent itemset is closed if none of its supersets has the **same support**
 - Lossless compression of the set of all frequent itemsets

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

min_sup = 2

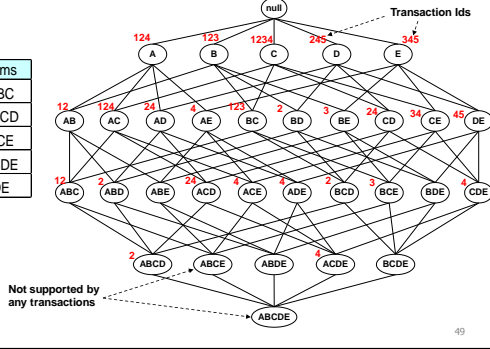
Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	3
{A,B,C,D}	2

48

Maximal vs Closed Frequent Itemsets

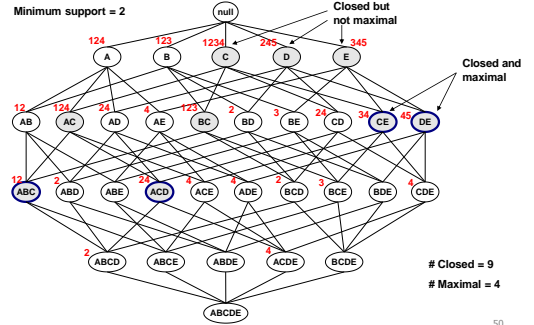
TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



49

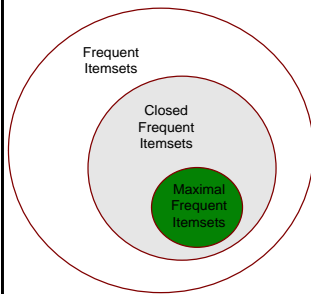
Maximal vs Closed Frequent Itemsets

Minimum support = 2



50

Maximal vs Closed Frequent Itemsets

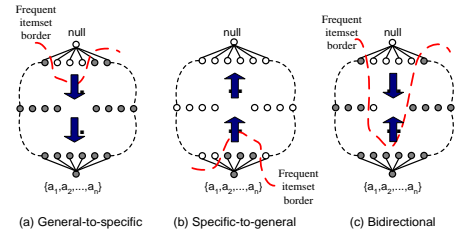


- How to efficiently find maximal frequent itemsets? (similar for closed ones)
 - Naive: first find all frequent itemsets, then remove non-maximal ones
 - Better: use maximality property for pruning
- Effectiveness depends on itemset generation strategy
- See book for details

51

Methods for Frequent Itemset Generation

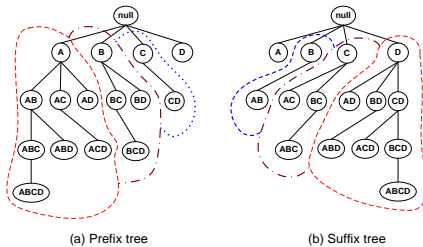
- Traversal of itemset lattice
 - General-to-specific: Apriori
 - Specific-to-general: good for pruning for maximal frequent itemsets



54

Alternative Methods for Frequent Itemset Generation

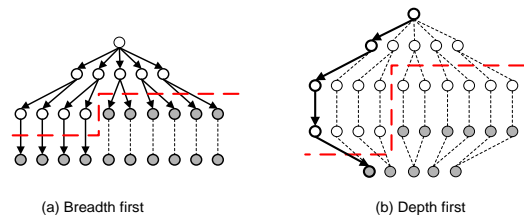
- Traversal of itemset lattice
 - Equivalence Classes: search one class first, before moving on to the next one



55

Alternative Methods for Frequent Itemset Generation

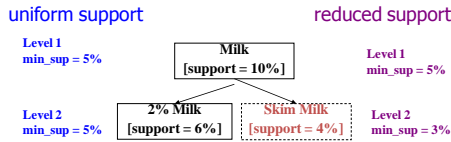
- Traversal of Itemset Lattice
 - Breadth-first vs Depth-first
 - Apriori is breadth-first (good for pruning)
 - Depth-first often good for maximal frequent itemsets: discover large frequent itemsets quickly, use for pruning



56

Extension: Mining Multiple-Level Association Rules

- Items often form hierarchies
 - Most relevant pattern might only show at the right granularity
- Flexible support settings
 - Items at the lower level are expected to have lower support



57

Extension: Mining Multi-Dimensional Associations

- Single-dimensional rules: one type of predicate
 - $\text{buys}(X, \text{"milk"}) \rightarrow \text{buys}(X, \text{"bread"})$
- Multi-dimensional rules: ≥ 2 types of predicates
 - Interdimensional association rules (no repeated predicates)
 - $\text{age}(X, \text{"19-25"}) \wedge \text{occupation}(X, \text{"student"}) \rightarrow \text{buys}(X, \text{"coke"})$
 - Hybrid-dimensional association rules (repeated predicates)
 - $\text{age}(X, \text{"19-25"}) \wedge \text{buys}(X, \text{"popcorn"}) \rightarrow \text{buys}(X, \text{"coke"})$
- See book for efficient mining algorithms

58

Frequent Pattern Mining Overview

- Basic Concepts and Challenges
- Efficient and Scalable Methods for Frequent Itemsets and Association Rules
- Pattern Interestingness Measures
- Sequence Mining

59

Lift

- Rule found: $\text{Basketball} \rightarrow \text{Cereal}$ [40%, 66.7%]
 - Misleading, because overall % of students eating cereal is 75% (which is $> 66.7\%$)
- $\text{Basketball} \rightarrow \text{Not_cereal}$ [20%, 33.3%] is more useful, although with lower support and confidence
- Measure of dependent/correlated events: **lift**

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

A, B are itemsets

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

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

60

Lift vs. Other Correlation Measures

- Intuition: Are milk and coffee usually bought together?
 - $(m, c) > (\sim m, c) + (m, \sim c)$
- m and c are...
 - bought together in A's
 - independent in B
 - not bought together in C's
- All measures good for B
- Lift, χ^2 bad for A's, C's
 - Reason: strongly affected by number of null-transactions (those without m, c)
- all_conf, cosine good for A's, C's
 - Not affected by number of null-transactions

		Milk	No Milk
Coffee	m, c	$\sim m, c$	
No Coffee	m, $\sim c$	$\sim m, \sim c$	

$$\text{all_conf}(A) = \frac{\text{sup}(A)}{\max_item_sup(A)}$$

$$\text{cosine}(A, B) = \frac{P(A \cup B)}{\sqrt{P(A)P(B)}} \quad \text{Lift vs. cosine: cosine does not depend on size of DB}$$

Data Set	mc	m \bar{c}	$\bar{m}c$	$\bar{m}\bar{c}$	all_conf	cosine	lift	χ^2
A ₁	1,000	100	100	100,000	0.91	0.91	83.64	83,452.6
A ₂	1,000	100	100	10,000	0.91	0.91	9.26	9,055.7
A ₃	1,000	100	100	1,000	0.91	0.91	1.82	1,472.7
A ₄	1,000	100	100	0	0.91	0.91	0.99	9.9
B ₁	1,000	1,000	1,000	1,000	0.50	0.50	1.00	0.0
C ₁	100	1,000	1,000	100,000	0.09	0.09	8.44	670.0
C ₂	1,000	100	10,000	100,000	0.09	0.29	9.18	8,172.8
C ₃	1	1	100	10,000	0.01	0.07	50.0	48.8

61

Which Measure Is Best?

- Does it identify the right patterns?
- Does it result in an efficient mining algorithm?

symbol	measure	range	formula
ϕ	co-efficient	-1 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B)}$
Q	Yule's Q	-1 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
Y	Yule's Y	-1 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
k	Chen's k	-1 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
PS	Platt-Roberts-Shapiro's F	-0.25 ... 0.25	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
F	Certainty factor	-1 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
AV	added value	-0.5 ... 0.5	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
K	Klingman's Q	-0.25 ... 0.25	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
g	Goodman-Kruskal's g	0 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
M	Mutual Information	0 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
J	J-Measure	0 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
G	Gini index	0 ... 1	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B) + P(B)P(A)}$
s	support	0 ... 1	$P(A \cup B)$
c	confidence	0 ... 1	$\frac{P(A \cup B)}{P(A)}$
L	Laplace	0 ... 1	$\frac{P(A \cup B)}{P(A) + P(B)}$
IS	Conid	0 ... 1	$\frac{P(A \cup B)}{P(A) + P(B)}$
r	odds ratios (hazard)	0 ... 1	$\frac{P(A \cup B)}{P(A) + P(B)}$
o	all_confidence	0 ... 1	$\frac{P(A \cup B)}{P(A) + P(B)}$
o	odds ratio	0 ... ∞	$\frac{P(A \cup B)}{P(A) + P(B)}$
V	Conviction	0.5 ... ∞	$\frac{P(A \cup B)}{P(A) + P(B)}$
X	lift	0 ... ∞	$\frac{P(A \cup B)}{P(A)P(B)}$
S	Collective strength	0 ... ∞	$\frac{P(A \cup B)}{P(A)P(B)}$
χ^2		0 ... ∞	$\frac{P(A \cup B) - P(A)P(B)}{P(A)P(B)}$

62

Symbol	Measure	Range	P1	P2	P3	O1	O2	O3	O3'	O4
Φ	Correlation	-1 ... 0 ... 1	Yes	Yes	Yes	Yes	No	Yes	Yes	No
λ	Lambda	0 ... 1	Yes	No	No	Yes	No	No*	Yes	No
α	Odds ratio	0 ... 1 ... ∞	Yes*	Yes	Yes	Yes	Yes	Yes*	Yes	No
Q	Yule's Q	-1 ... 0 ... 1	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
Y	Yule's Y	-1 ... 0 ... 1	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
k	Cohen's k	-1 ... 0 ... 1	Yes	Yes	Yes	Yes	No	No	Yes	No
M	Mutual Information	0 ... 1	Yes	Yes	Yes*	Yes	No	No*	Yes	No
J	J-Measure	0 ... 1	Yes	No	No	No	No	No	No	No
G	Gini Index	0 ... 1	Yes	No	No	No	No	No*	Yes	No
s	Support	0 ... 1	No	Yes	No	Yes	No	No	No	No
c	Confidence	0 ... 1	No	Yes	No	Yes	No	No	No	Yes
L	Laplace	0 ... 1	No	Yes	No	Yes	No	No	No	No
V	Conviction	0.5 ... 1 ... ∞	No	Yes	No	Yes**	No	No	Yes	No
I	Interest	0 ... 1 ... ∞	Yes*	Yes	Yes	Yes	No	No	No	No
IS	IS (cosine)	0 ... 1	No	Yes	Yes	Yes	No	No	No	Yes
PS	Ratelsky-Shapiro's	-0.25 ... 0 ... 0.25	Yes	Yes	Yes	Yes	No	Yes	Yes	No
F	Certainty Factor	-1 ... 0 ... 1	Yes	Yes	Yes	No	No	No	Yes	No
AV	Added value	0.5 ... 1 ... 1	Yes	Yes	Yes	No	No	No	No	No
S	Collective strength	0 ... 1 ... ∞	No	Yes	Yes	Yes	No	Yes*	Yes	No
ζ	Jaccard	0 ... 1	No	Yes	Yes	Yes	No	No	No	Yes
K	Klosgen's	$\left(\frac{2}{\sqrt{2P}-1}\right) \left(2 - \sqrt{2 - \frac{1}{2P}}\right) \cdot 0 \dots \frac{2}{\sqrt{2P}}$	Yes	Yes	Yes	No	No	No	No	No

The P's and O's are various desirable properties, e.g., symmetry under variable permutation (O1), which we do not cover in this class. Take-away message: **no interestingness measure has all the desirable properties.**

Frequent Pattern Mining Overview

- Basic Concepts and Challenges
- Efficient and Scalable Methods for Frequent Itemsets and Association Rules
- Pattern Interestingness Measures
- Sequence Mining

Introduction

- Sequence mining is relevant for transaction databases, time-series databases, and sequence databases
- Applications of sequential pattern mining
 - Customer shopping sequences:
 - First buy computer, then CD-ROM, and then digital camera, within 3 months
 - Medical treatments, natural disasters (e.g., earthquakes), science & engineering processes, stocks and markets
 - Telephone calling patterns, Weblog click streams
 - DNA sequences and gene structures

What Is Sequential Pattern Mining?

- Given a set of sequences, find all frequent subsequences

A sequence: $\langle (ef)(ab)(df)cb \rangle$

A sequence database

SID	sequence
10	$\langle a(abc)(ac)d(cf) \rangle$
20	$\langle (ad)c(bc)(ae) \rangle$
30	$\langle (ef)(ab)(df)cb \rangle$
40	$\langle eg(af)cbc \rangle$

An element may contain a set of items. Items within an element are unordered and we list them alphabetically

$\langle a(bc)dc \rangle$ is a **subsequence** of $\langle a(abc)(ac)d(cf) \rangle$

Given support threshold $min_sup = 2$, $\langle (ab)c \rangle$ is a **sequential pattern**

Challenges of Sequential Pattern Mining

- Huge number of possible patterns
- A mining algorithm should
 - find all patterns satisfying the minimum support threshold
 - be highly efficient and scalable
 - be able to incorporate user-specific constraints

Apriori Property of Sequential Patterns

- If a sequence S is not frequent, then none of the super-sequences of S is frequent
 - E.g, if $\langle hb \rangle$ is infrequent, then so are $\langle hab \rangle$ and $\langle (ah)b \rangle$

Seq. ID	Sequence
10	$\langle (bd)cb(ac) \rangle$
20	$\langle (bf)(ce)b(fg) \rangle$
30	$\langle (ah)(bf)abf \rangle$
40	$\langle (be)(ce)d \rangle$
50	$\langle a(bd)bcb(ade) \rangle$

Given support threshold $min_sup = 2$, find all frequent subsequences

GSP: Generalized Sequential Pattern Mining

- Initially, every item in DB is a candidate of length $k=1$
- For each level (i.e., sequences of length k) do
 - Scan database to collect support count for each candidate sequence
 - Generate candidate length- $(k+1)$ sequences from length- k frequent sequences
 - Join phase: sequences s_1 and s_2 join, if s_1 without its first item is identical to s_2 without its last item
 - Prune phase: delete candidates that contain a length- k subsequence that is not among the frequent ones
- Repeat until no frequent sequence or no candidate can be found
- Major strength: Candidate pruning by Apriori

80

Finding Length-1 Sequential Patterns

- Initial candidates: all singleton sequences
 - $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle, \langle g \rangle, \langle h \rangle$
- Scan database once, count support for candidates

Seq. ID	Sequence
10	$\langle (bd)cb(ac) \rangle$
20	$\langle (bf)(ce)b(fg) \rangle$
30	$\langle (ah)(bf)abf \rangle$
40	$\langle (be)(ce)d \rangle$
50	$\langle a(bd)bcb(ade) \rangle$

$min_sup = 2$

Cand	Sup
$\langle a \rangle$	3
$\langle b \rangle$	5
$\langle c \rangle$	4
$\langle d \rangle$	3
$\langle e \rangle$	3
$\langle f \rangle$	2
$\langle g \rangle$	1
$\langle h \rangle$	1

81

GSP: Generating Length-2 Candidates

51 length-2 Candidates

	$\langle a \rangle$	$\langle b \rangle$	$\langle c \rangle$	$\langle d \rangle$	$\langle e \rangle$	$\langle f \rangle$
$\langle a \rangle$	$\langle aa \rangle$	$\langle ab \rangle$	$\langle ac \rangle$	$\langle ad \rangle$	$\langle ae \rangle$	$\langle af \rangle$
$\langle b \rangle$	$\langle ba \rangle$	$\langle bb \rangle$	$\langle bc \rangle$	$\langle bd \rangle$	$\langle be \rangle$	$\langle bf \rangle$
$\langle c \rangle$	$\langle ca \rangle$	$\langle cb \rangle$	$\langle cc \rangle$	$\langle cd \rangle$	$\langle ce \rangle$	$\langle cf \rangle$
$\langle d \rangle$	$\langle da \rangle$	$\langle db \rangle$	$\langle dc \rangle$	$\langle dd \rangle$	$\langle de \rangle$	$\langle df \rangle$
$\langle e \rangle$	$\langle ea \rangle$	$\langle eb \rangle$	$\langle ec \rangle$	$\langle ed \rangle$	$\langle ee \rangle$	$\langle ef \rangle$
$\langle f \rangle$	$\langle fa \rangle$	$\langle fb \rangle$	$\langle fc \rangle$	$\langle fd \rangle$	$\langle fe \rangle$	$\langle ff \rangle$

	$\langle a \rangle$	$\langle b \rangle$	$\langle c \rangle$	$\langle d \rangle$	$\langle e \rangle$	$\langle f \rangle$
$\langle a \rangle$		$\langle (ab) \rangle$	$\langle (ac) \rangle$	$\langle (ad) \rangle$	$\langle (ae) \rangle$	$\langle (af) \rangle$
$\langle b \rangle$			$\langle (bc) \rangle$	$\langle (bd) \rangle$	$\langle (be) \rangle$	$\langle (bf) \rangle$
$\langle c \rangle$				$\langle (cd) \rangle$	$\langle (ce) \rangle$	$\langle (cf) \rangle$
$\langle d \rangle$					$\langle (de) \rangle$	$\langle (df) \rangle$
$\langle e \rangle$						$\langle (ef) \rangle$
$\langle f \rangle$						

Without Apriori property,
 $8 \times 8 + 8 \times 7 / 2 = 92$ candidates

Apriori prunes
 44.57% candidates

82

The GSP Mining Process

- Scan 5: 1 candidate, 1 length-5 seq. pattern
 $\langle (bd)cb(a) \rangle$ (Cand. does not pass support threshold)
- Scan 4: 8 candidates, 6 length-4 seq. patterns
 $\langle abba \rangle, \langle (bd)bc \rangle, \dots$ (Cand. not in DB at all)
- Scan 3: 47 candidates, 19 length-3 seq. patterns, 20 candidates not in DB at all
 $\langle abb \rangle, \langle aab \rangle, \langle aba \rangle, \langle baa \rangle, \langle bab \rangle, \dots$
- Scan 2: 51 candidates, 19 length-2 seq. patterns, 10 candidates not in DB at all
 $\langle aa \rangle, \langle ab \rangle, \dots, \langle af \rangle, \langle ba \rangle, \langle bb \rangle, \dots, \langle ff \rangle, \langle (ab) \rangle, \dots, \langle (ef) \rangle$
- Scan 1: 8 candidates, 6 length-1 seq. patterns
 $\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle, \langle g \rangle, \langle h \rangle$

Seq. ID	Sequence
10	$\langle (bd)cb(ac) \rangle$
20	$\langle (bf)(ce)b(fg) \rangle$
30	$\langle (ah)(bf)abf \rangle$
40	$\langle (be)(ce)d \rangle$
50	$\langle a(bd)bcb(ade) \rangle$

$min_sup = 2$

83

Candidate Generate-and-Test Drawbacks

- Huge set of candidate sequences generated
- Multiple Scans of entire database needed
 - Length of each candidate grows by one at each database scan

84

Prefix and Suffix (Projection)

- $\langle a \rangle, \langle aa \rangle, \langle a(ab) \rangle$ and $\langle a(abc) \rangle$ are prefixes of sequence $\langle a(abc)(ac)d(cf) \rangle$
- Given sequence $\langle a(abc)(ac)d(cf) \rangle$, we have:

Prefix	Suffix (Prefix-Based Projection)
$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle$
$\langle aa \rangle$	$\langle (_bc)(ac)d(cf) \rangle$
$\langle ab \rangle$	$\langle (_c)(ac)d(cf) \rangle$
$\langle bc \rangle$	$\langle d(cf) \rangle$
$\langle (bc) \rangle$	$\langle (ac)d(cf) \rangle$

85

Mining Sequential Patterns by Prefix Projections

- Step 1: find length-1 frequent sequential patterns
 - <a>, , <c>, <d>, <e>, <f>
- Step 2: divide search space. The complete set of sequential patterns can be partitioned into six subsets:
 - The ones having prefix <a>;
 - The ones having prefix ;
 - ...
 - The ones having prefix <f>

SID	sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

86

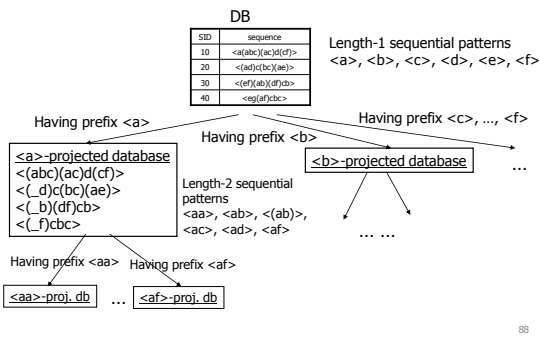
Finding Seq. Patterns with Prefix <a>

- Only need to consider projections w.r.t. <a>
 - <a>-projected database: <(abc)(ac)d(cf)>, <(_d)c(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
- Find all length-2 frequent seq. patterns having prefix <a>: <aa>, <ab>, <(ab)>, <ac>, <ad>, <af>
 - Further partition into those 6 subsets
 - Having prefix <aa>;
 - Having prefix <ab>;
 - Having prefix <(ab)>;
 - ...
 - Having prefix <af>

SID	sequence
10	<a(abc)(ac)d(cf)>
20	<(ad)c(bc)(ae)>
30	<(ef)(ab)(df)cb>
40	<eg(af)cbc>

87

Completeness of PrefixSpan



88

Efficiency of PrefixSpan

- No candidate sequence needs to be generated
- Projected databases keep shrinking
- Major cost of PrefixSpan: constructing projected databases
 - Can be improved by pseudo-projections

89

Pseudo-Projection

- Major cost of PrefixSpan: projection
 - Postfixes of sequences often appear repeatedly in recursive projected databases
- When (projected) database can be held in memory, use pointers
 - Pointer to the sequence, offset of the postfix
- Why is this a bad idea when the (projected) database does not fit in memory?
 - $$s = \text{a(abc)(ac)d(cf)} >$$

$$s|<a>: (, 2) \text{ <(abc)(ac)d(cf)>}$$

$$s|<ab>: (, 4) \text{ <(_c)(ac)d(cf)>}$$

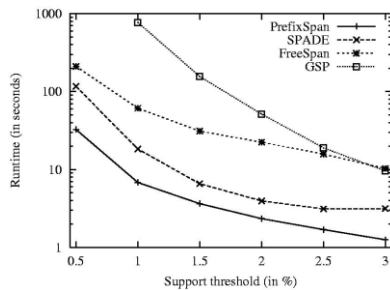
90

Pseudo-Projection vs. Physical Projection

- Pseudo-projection avoids physically copying postfixes
 - Efficient in running time and space when database can be held in main memory
- Not efficient when database cannot fit in main memory
 - Disk-based random access
- Suggested Approach:
 - Integration of physical and pseudo-projection
 - Swapping to pseudo-projection when the data set fits in memory

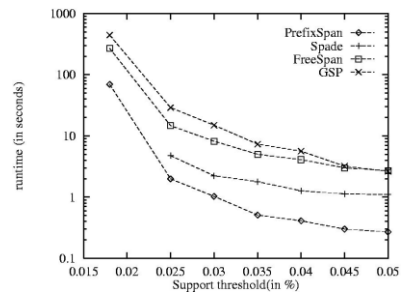
91

Performance on Data Set C10T8S8I8



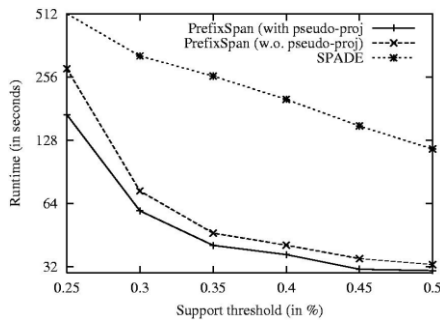
92

Performance on Data Set Gazelle



93

Effect of Pseudo-Projection



94

Sequence Mining Variations

- Multidimensional and multilevel patterns
- Constraint-based mining of sequential patterns
- Periodicity analysis
- Mining biological sequences
 - Hot research area, major topic by itself
- All these not discussed in class; see book
- Some of my own research: finding relevant sequences in bursty data; see paper

95

Frequent-Pattern Mining: Summary

- Important task in data mining
- Scalable frequent pattern mining methods
 - Apriori (itemsets, candidate generation & test)
 - GSP (sequences, candidate generation & test)
 - Projection-based (FP-growth for itemsets, PrefixSpan for sequences)
- Mining a variety of rules and interesting patterns

132

Frequent-Pattern Mining: Research Problems

- Mining fault-tolerant frequent, sequential and structured patterns
 - Patterns allows limited faults (insertion, deletion, mutation)
- Mining truly interesting patterns
 - Surprising, novel, concise,...
- Application exploration
 - E.g., DNA sequence analysis and bio-pattern classification

133