

Basic Concepts: Frequent Itemsets (Patterns)

- Itemset: A set of one or more items
- k-itemset: $X = \{x_1, \dots, x_k\}$
- Support (count) of X: frequency or the number of occurrences of an itemset X
- Relative support: s: The fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An itemset X is frequent if the support of X is no less than a minsup threshold (from as of)

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

- Let minsup = 50%
- Freq. 1-itemsets:
 - Beer: 3 (60%); Nuts: 3 (60%)
 - Diaper: 4 (80%); Eggs: 3 (60%)
- Freq. 2-itemsets:
 - Beer, Diaper: 3 (60%)

From Frequent Itemsets to 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

- Association rules: $X \rightarrow Y, c$
- Support, s: The probability that a transaction contains $X \cup Y$
- Confidence, c: The conditional probability that a transaction containing X also contains Y
- $c = \text{sup}(X \cup Y) / \text{sup}(X)$
- Association rule mining: Find all of the rules, $X \rightarrow Y$, with minimum support and confidence
- Frequent itemsets: Let minsup = 50%
- Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
- Freq. 2-itemsets: Beer, Diaper: 3
- Association rules: Let minconf = 50%
- Beer \rightarrow Diaper (60%, 100%)
- Diaper \rightarrow Beer (60%, 75%)

The Apriori Algorithm—An Example

TID	Items
10	A, C, D
20	A, B, C, E
30	A, B, C, E
40	A, B, C, E

Itemset	sup
(A, C)	2
(B, C)	2
(C, E)	2

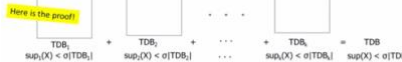
Itemset	sup
(A, C)	2
(B, C)	2
(C, E)	2

Itemset	sup
(A, C)	2
(B, C)	2
(C, E)	2

Itemset	sup
(A, C)	2
(B, C)	2
(C, E)	2

Partitioning: Scan Database Only Twice

- Theorem: Any itemset that is potentially frequent in TDB must be frequent in at least one of the partitions of TDB



- Method: (A. Savasere, E. Omiecinski and S. Navathe, VLDB'95)
- Scan 1: Partition database (how?) and find local frequent patterns
- Scan 2: Consolidate global frequent patterns (how?)
- Why does this method guarantee to scan TDB only twice?

Mining Frequent Patterns, Associations, and Correlations
 $\text{support}(X \Rightarrow Y) = P(X \cup Y)$
 $\text{confidence}(X \Rightarrow Y) = P(Y|X)$

$$\text{confidence}(A \Rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support}(A)} = \frac{\text{support.count}(A \cup B)}{\text{support.count}(A)}$$

- An itemset X is closed in a data set D if there exists no proper super-itemset Y such that Y has the same support count as X in D.
- An itemset X is a maximal frequent itemset (or max-itemset) in a data set D if X is frequent and there exists no frequent super pattern Y of X.

Mining Quantitative Associations

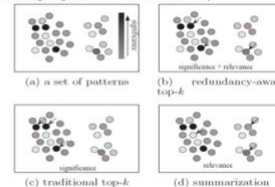
- Mining associations with numerical attributes
- Ex.: Numerical attributes: age and salary
- Methods
- Static discretization based on predefined concept hierarchies
- Data cube-based aggregation
- Dynamic discretization based on data distribution
- Clustering: Distance-based association
- First one-dimensional clustering, then association
- Deviation analysis:
- Gender = female \Rightarrow Wage: mean=\$7/hr (overall mean = \$9)

Mining Extraordinary Phenomena in Quantitative Association Mining

- Mining extraordinary (i.e., interesting) phenomena
- Ex.: Gender = female \Rightarrow Wage: mean=\$7/hr (overall mean = \$9)
- LHS: a subset of the population
- RHS: an extraordinary behavior of this subset
- The rule is accepted only if a statistical test (e.g., Z-test) confirms the inference with high confidence
- Subrule: Highlights the extraordinary behavior of a subset of the population of the super rule
- Ex.: (Gender = female) \wedge (South = yes) \Rightarrow mean wage = \$6.3/hr
- Rule condition can be categorical or numerical (quantitative rules)
- Ex.: Education in [14-18] (yrs) \Rightarrow mean wage = \$11.64/hr
- Efficient methods have been developed for mining such extraordinary rules (e.g., Aumann and Lindell@KDD'99)

Redundancy-Aware Top-k Patterns

- Desired patterns: high significance & low redundancy



- Method: Use MMS (Maximal Marginal Significance) for measuring the combined significance of a pattern set
- Xin et al., Extracting Redundancy-Aware Top-K Patterns, KDD'06

Exploring Vertical Data Format: ECLAT

- ECLAT (Equivalence Class Transformation): A depth-first search algorithm using set intersection [Zaki et al. @KDD'97]
- Tid-List: List of transaction-ids containing an itemset
- Vertical format: $t(e) = \{T_{10}, T_{20}, T_{30}\}$; $t(a) = \{T_{10}, T_{20}\}$; $t(ac) = \{T_{10}, T_{20}\}$
- Properties of Tid-Lists
- $t(X) = t(Y)$: X and Y always happen together (e.g., $t(ac) = t(cd)$)
- $t(X) \subseteq t(Y)$: transaction having X always has Y (e.g., $t(ac) \subseteq t(ce)$)
- Deriving frequent patterns based on vertical intersections
- Using diffset to accelerate mining
- Only keep track of differences of tids
- $t(e) = \{T_{10}, T_{20}, T_{30}\}$; $t(ce) = \{T_{10}, T_{20}\} \rightarrow \text{Diffset}(ce, e) = \{T_{30}\}$

TID	Itemset
10	a, c, d, e
20	a, b, e
30	b, c, e

Item	TidList
a	10, 20
b	20, 30
c	10, 30
d	10
e	10, 20, 30

FPGrowth: Mining Frequent Patterns by Pattern Growth

- Idea: Frequent pattern growth (FPGrowth)
- Find frequent single items and partition the database based on each such item
- Recursively grow frequent patterns by doing the above for each partitioned database (also called conditional database)
- To facilitate efficient processing, an efficient data structure, FP-tree, can be constructed
- Mining becomes
- Recursively construct and mine (conditional) FP-trees
- Until the resulting FP-tree is empty, or until it contains only one path—single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

Example: Construct FP-tree from a Transactional DB

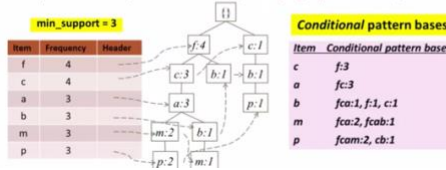
TID	Items in the Transaction	Ordered, frequent items
100	{f, a, c, d, g, h, m, p}	{f, c, a, m, p}
200	{a, b, c, f, h, m, p}	{f, c, a, h, m}
300	{b, f, h, g, c, m, p}	{f, h, m, p}
400	{b, c, a, h, m, p}	{f, h, m, p}
500	{a, f, c, g, d, h, p, m, p}	{f, c, a, m, p}

- Scan DB once, find single item frequent pattern: Let min_support = 3
- Sort frequent items in frequency descending order, f-list
- Scan DB again, construct FP-tree

Item	Frequency	Header
f	4	
c	4	
a	3	
b	3	
m	3	
p	3	

Divide and Conquer Based on Patterns and Data

- Pattern mining can be partitioned according to current patterns
- Patterns containing p: p's conditional database: $fca:m2, cb:1$
- Patterns having m but no p: m's conditional database: $fca:2, fcb:1$
-
- p's conditional pattern base: transformed prefix paths of item p



Defining Negative Correlation: Need Null-Invariance in Definition

- A good definition on negative correlation should take care of the null-invariance problem
- Whether two itemsets A and B are negatively correlated should not be influenced by the number of null-transactions
- A Kulczynski measure-based definition
- If itemsets A and B are frequent but $(P(A|B) + P(B|A))/2 < \epsilon$, where ϵ is a negative pattern threshold, then A and B are negatively correlated
- For the same needle package problem:
- No matter there are in total 200 or 10⁵ transactions
- If $\epsilon = 0.01$, we have $(P(A|B) + P(B|A))/2 = (0.01 + 0.01)/2 < \epsilon$

Mining Compressed Patterns

- Why mining compressed patterns?
- Too many scattered patterns but not so meaningful
- Pattern distance measure
- $\text{Dist}(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}$
- 6-clustering: For each pattern P, find all patterns which can be expressed by P and whose distance to P is within δ (δ -cover)
- All patterns in the cluster can be represented by P
- Method for efficient, direct mining of compressed frequent patterns (e.g., D. Xin, X. Han, X. Yan, H. Cheng, "On Compressing Frequent Patterns", Knowledge and Data Engineering, 60:5-29, 2007)

CLOSET+: Mining Closed Itemsets by Pattern-Growth

- Efficient, direct mining of closed itemsets
- Ex. Itemset merging: If Y appears in every occurrence of X, then Y is merged with X
- d-proj. db: $\langle aef, acf \rangle \rightarrow$ dcf-d-proj. db: {e}, thus we get: acfd:2
- Many other tricks (but not detailed here), such as
- Hybrid tree projection
- Bottom-up physical tree-projection

Interestingness Measure: Lift

- Measure of dependent/correlated events: lift
- $\text{lift}(B, C) = \frac{c(B \cup C)}{s(C)} = \frac{s(B \cup C)}{s(B) \times s(C)}$

	B	-B	Σ
C	400	350	750
-C	200	50	250
Σ	600	400	1000

- Lift(B, C) may tell how B and C are correlated
- Lift(B, C) = 1: B and C are independent
- > 1: positively correlated
- < 1: negatively correlated

Interestingness Measure: χ^2

- Another measure to test correlated events: χ^2
- $\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$

	B	-B	Σ
C	400 (450)	350 (300)	750
-C	200 (150)	50 (100)	250
Σ	600	400	1000

- General rules
- $\chi^2 = 0$: independent
- $\chi^2 > 0$: correlated, either positive or negative, so it needs additional test

Measure	Definition	Range	Null-Invariant
$\chi^2(A, B)$	$\sum_{i,j=0,1} \frac{(e(a_i b_j) - o(a_i b_j))^2}{e(a_i b_j)}$	$[0, \infty]$	No
$\text{Lift}(A, B)$	$\frac{s(A \cup B)}{s(A) \times s(B)}$	$[0, \infty]$	No
$\text{AllConf}(A, B)$	$\frac{s(A \cup B)}{\max\{s(A), s(B)\}}$	$[0, 1]$	Yes
$\text{Jaccard}(A, B)$	$\frac{s(A \cup B)}{s(A) + s(B) - s(A \cup B)}$	$[0, 1]$	Yes
$\text{Cosine}(A, B)$	$\frac{s(A \cup B)}{\sqrt{s(A) \times s(B)}}$	$[0, 1]$	Yes
$\text{Kulczynski}(A, B)$	$\frac{1}{2} \left(\frac{s(A \cup B)}{s(A)} + \frac{s(A \cup B)}{s(B)} \right)$	$[0, 1]$	Yes
$\text{MaxConf}(A, B)$	$\max\left\{ \frac{s(A)}{s(A \cup B)}, \frac{s(B)}{s(A \cup B)} \right\}$	$[0, 1]$	Yes

Many Null Transactions throws off Lift and Chi-square.

Imbalance Ratio with Kulczynski Measure

- IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications:

$$\text{IR}(A, B) = \frac{|s(A) - s(B)|}{s(A) + s(B) - s(A \cup B)}$$

Data set	m	-m	m-c	-m-c	Jaccard	Cosine	Kulc	IR
D_1	10,000	1,000	1,000	100,000	0.83	0.91	0.91	0
D_2	10,000	1,000	1,000	100	0.83	0.91	0.91	0
D_3	100	1,000	1,000	100,000	0.05	0.09	0.09	0
D_4	1,000	1,000	1,000	100,000	0.33	0.75	0.75	0
D_5	1,000	100	100,000	100,000	0.09	0.29	0.5	0.50
D_6	1,000	10	100,000	100,000	0.01	0.10	0.5	0.99

Multi-Level Freq Patterns - Efficient mining: Shared multi-level mining - Use the lowest min-support to pass down the set of candidates

Compressed Patterns

Pat. ID	Item-Set	Support
P1	{f, a, c, d, g, h, m, p}	205227
P2	{a, b, c, f, h, m, p}	205211
P3	{b, f, h, g, c, m, p}	101758
P4	{a, f, c, g, d, h, p, m, p}	161563
P5	{a, c, f, g, d, h, p, m, p}	161576

Compressed patterns: P1 and P2 are not compressed by P3, P4, P5, P6, P7, P8, P9, P10, P11, P12, P13, P14, P15, P16, P17, P18, P19, P20, P21, P22, P23, P24, P25, P26, P27, P28, P29, P30, P31, P32, P33, P34, P35, P36, P37, P38, P39, P40, P41, P42, P43, P44, P45, P46, P47, P48, P49, P50, P51, P52, P53, P54, P55, P56, P57, P58, P59, P60, P61, P62, P63, P64, P65, P66, P67, P68, P69, P70, P71, P72, P73, P74, P75, P76, P77, P78, P79, P80, P81, P82, P83, P84, P85, P86, P87, P88, P89, P90, P91, P92, P93, P94, P95, P96, P97, P98, P99, P100.

Distance of patterns: $\text{Dist}(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|} = 1 - \frac{101758}{205227 + 205211 - 101758} = 1 - \frac{101758}{408688} = 0.75$

Distance of patterns: $\text{Dist}(P_1, P_3) = 1 - \frac{|T(P_1) \cap T(P_3)|}{|T(P_1) \cup T(P_3)|} = 1 - \frac{101758}{205227 + 101758 - 101758} = 1 - \frac{101758}{408688} = 0.75$

Distance of patterns: $\text{Dist}(P_1, P_4) = 1 - \frac{|T(P_1) \cap T(P_4)|}{|T(P_1) \cup T(P_4)|} = 1 - \frac{101758}{205227 + 161563 - 101758} = 1 - \frac{101758}{465032} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_5) = 1 - \frac{|T(P_1) \cap T(P_5)|}{|T(P_1) \cup T(P_5)|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_6) = 1 - \frac{|T(P_1) \cap T(P_6)|}{|T(P_1) \cup T(P_6)|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_7) = 1 - \frac{|T(P_1) \cap T(P_7)|}{|T(P_1) \cup T(P_7)|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_8) = 1 - \frac{|T(P_1) \cap T(P_8)|}{|T(P_1) \cup T(P_8)|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_9) = 1 - \frac{|T(P_1) \cap T(P_9)|}{|T(P_1) \cup T(P_9)|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{10}) = 1 - \frac{|T(P_1) \cap T(P_{10})|}{|T(P_1) \cup T(P_{10})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{11}) = 1 - \frac{|T(P_1) \cap T(P_{11})|}{|T(P_1) \cup T(P_{11})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{12}) = 1 - \frac{|T(P_1) \cap T(P_{12})|}{|T(P_1) \cup T(P_{12})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{13}) = 1 - \frac{|T(P_1) \cap T(P_{13})|}{|T(P_1) \cup T(P_{13})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{14}) = 1 - \frac{|T(P_1) \cap T(P_{14})|}{|T(P_1) \cup T(P_{14})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{15}) = 1 - \frac{|T(P_1) \cap T(P_{15})|}{|T(P_1) \cup T(P_{15})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{16}) = 1 - \frac{|T(P_1) \cap T(P_{16})|}{|T(P_1) \cup T(P_{16})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{17}) = 1 - \frac{|T(P_1) \cap T(P_{17})|}{|T(P_1) \cup T(P_{17})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{18}) = 1 - \frac{|T(P_1) \cap T(P_{18})|}{|T(P_1) \cup T(P_{18})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{19}) = 1 - \frac{|T(P_1) \cap T(P_{19})|}{|T(P_1) \cup T(P_{19})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{20}) = 1 - \frac{|T(P_1) \cap T(P_{20})|}{|T(P_1) \cup T(P_{20})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{21}) = 1 - \frac{|T(P_1) \cap T(P_{21})|}{|T(P_1) \cup T(P_{21})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{22}) = 1 - \frac{|T(P_1) \cap T(P_{22})|}{|T(P_1) \cup T(P_{22})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{23}) = 1 - \frac{|T(P_1) \cap T(P_{23})|}{|T(P_1) \cup T(P_{23})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{24}) = 1 - \frac{|T(P_1) \cap T(P_{24})|}{|T(P_1) \cup T(P_{24})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{25}) = 1 - \frac{|T(P_1) \cap T(P_{25})|}{|T(P_1) \cup T(P_{25})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{26}) = 1 - \frac{|T(P_1) \cap T(P_{26})|}{|T(P_1) \cup T(P_{26})|} = 1 - \frac{101758}{205227 + 161576 - 101758} = 1 - \frac{101758}{465045} = 0.78$

Distance of patterns: $\text{Dist}(P_1, P_{27}) = 1 - \frac{|T(P_1) \cap T(P_{27})|}{|T(P_1) \cup T(P_{27})|} = 1 - \frac{101758$

GSP Mining and Pruning

- 5th scan: 1 cand. 1 length-5 seq. pat. $\langle \{b\}d \rangle$
- 4th scan: 8 cand. 7 length-4 seq. pat. $\langle \{a\}b \rangle$, $\langle \{b\}c \rangle$
- 3rd scan: 46 cand. 20 length-3 seq. pat. 20 cand. not in DB at all
- 2nd scan: 51 cand. 19 length-2 seq. pat. 10 cand. not in DB at all
- 1st scan: 8 cand. 6 length-1 seq. pat. $\langle a \rangle$, $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$
- Repeat (for each level (i.e., length-k))
 - Scan DB to find length-k frequent sequences
 - Generate length-(k+1) candidate sequences from length-k frequent sequences using Apriori
 - set $k = k+1$
 - Until no frequent sequence or no candidate can be found

Implementation Consideration: Pseudo-Projection vs. Physical Projection

- Major cost of PrefixSpan: Constructing projected DBs
- Suffixes largely repeating in recursive projected DBs
- When DB can be held in main memory, use pseudo projection
- No physically copying suffixes
- Pointer to the sequence
- Offset of the suffix
- But if it does not fit in memory
- Physical projection
- Suggested approach:
 - Integration of physical and pseudo-projection
 - Swapping to pseudo-projection when the data fits in memory

Pattern-Growth Approach

- Depth-first growth of subgraphs from k -edge to $(k+1)$ -edge, then $(k+2)$ -edge subgraphs
- Major challenge
 - Generating many duplicate subgraphs
- Major idea to solve the problem
 - Define an order to generate subgraphs
 - DFS spanning tree: Flatten a graph into a sequence using depth-first search
- gSpan (Yan & Han, ICDM'02)

gSPAN: Graph Pattern Growth in Order

- Right-most path extension in subgraph pattern growth
- Right-most path: The path from root to the right-most leaf (choose the vertex with the smallest index at each step)
- Reduce generation of duplicate subgraphs
- Completeness: The enumeration of graphs using right-most path extension is complete
- DFS code: Flatten a graph into a sequence using depth-first search

Constraints in General Data Mining

- A data mining query can be in the form of a meta-rule or with the following language primitives
- Knowledge type constraint
 - Ex.: Classification, association, clustering, outlier finding, ...
 - Data constraint — using SQL-like queries
 - Ex.: Find products sold together in NY stores this year
 - Dimension/level constraint
 - Ex.: In relevance to region, price, brand, customer category
 - Rule (or pattern) constraint
 - Ex.: Small sales (price < \$10) triggers big sales (sum > \$200)
 - Interestingness constraint
 - Ex.: Strong rules: $\min_sup \geq 0.02$, $\min_conf \geq 0.6$, $\min_correlation \geq 0.7$

Pattern Monotonicity and Its Roles

- A constraint c is **monotone**: If an itemset S satisfies the constraint c , so does any of its superset
- That is, we do not need to check c in subsequent mining
- Ex. 1: $c_1: \text{sum}(S.Price) \geq v$ is **monotone**
- Ex. 2: $c_2: \text{min}(S.Price) \leq v$ is **monotone**
- Ex. 3: $c_3: \text{range}(S.profit) \geq 15$ is **monotone**
- Itemset ab satisfies c_3
- So does every superset of ab

Sequential Pattern Mining in Vertical Data Format: The SPADE Algorithm

- A sequence database is mapped to: $\langle \text{SID}, \text{EID} \rangle$
 - Grow the subsequences (patterns) one item at a time by Apriori candidate generation
- | SID | Sequence | SID | EID | Itemset |
|-----|---------------------|-----|-----|----------------|
| 1 | $\langle a \rangle$ | 1 | 1 | $\{a\}$ |
| 2 | $\langle a \rangle$ | 1 | 2 | $\{a, b\}$ |
| 3 | $\langle a \rangle$ | 1 | 3 | $\{a, c\}$ |
| 4 | $\langle a \rangle$ | 1 | 4 | $\{a, d\}$ |
| 5 | $\langle a \rangle$ | 1 | 5 | $\{a, e\}$ |
| 6 | $\langle a \rangle$ | 1 | 6 | $\{a, f\}$ |
| 7 | $\langle a \rangle$ | 1 | 7 | $\{a, g\}$ |
| 8 | $\langle a \rangle$ | 1 | 8 | $\{a, h\}$ |
| 9 | $\langle a \rangle$ | 1 | 9 | $\{a, i\}$ |
| 10 | $\langle a \rangle$ | 1 | 10 | $\{a, j\}$ |
| 11 | $\langle a \rangle$ | 1 | 11 | $\{a, k\}$ |
| 12 | $\langle a \rangle$ | 1 | 12 | $\{a, l\}$ |
| 13 | $\langle a \rangle$ | 1 | 13 | $\{a, m\}$ |
| 14 | $\langle a \rangle$ | 1 | 14 | $\{a, n\}$ |
| 15 | $\langle a \rangle$ | 1 | 15 | $\{a, o\}$ |
| 16 | $\langle a \rangle$ | 1 | 16 | $\{a, p\}$ |
| 17 | $\langle a \rangle$ | 1 | 17 | $\{a, q\}$ |
| 18 | $\langle a \rangle$ | 1 | 18 | $\{a, r\}$ |
| 19 | $\langle a \rangle$ | 1 | 19 | $\{a, s\}$ |
| 20 | $\langle a \rangle$ | 1 | 20 | $\{a, t\}$ |
| 21 | $\langle a \rangle$ | 1 | 21 | $\{a, u\}$ |
| 22 | $\langle a \rangle$ | 1 | 22 | $\{a, v\}$ |
| 23 | $\langle a \rangle$ | 1 | 23 | $\{a, w\}$ |
| 24 | $\langle a \rangle$ | 1 | 24 | $\{a, x\}$ |
| 25 | $\langle a \rangle$ | 1 | 25 | $\{a, y\}$ |
| 26 | $\langle a \rangle$ | 1 | 26 | $\{a, z\}$ |
| 27 | $\langle a \rangle$ | 1 | 27 | $\{a, \dots\}$ |
| 28 | $\langle a \rangle$ | 1 | 28 | $\{a, \dots\}$ |
| 29 | $\langle a \rangle$ | 1 | 29 | $\{a, \dots\}$ |
| 30 | $\langle a \rangle$ | 1 | 30 | $\{a, \dots\}$ |
| 31 | $\langle a \rangle$ | 1 | 31 | $\{a, \dots\}$ |
| 32 | $\langle a \rangle$ | 1 | 32 | $\{a, \dots\}$ |
| 33 | $\langle a \rangle$ | 1 | 33 | $\{a, \dots\}$ |
| 34 | $\langle a \rangle$ | 1 | 34 | $\{a, \dots\}$ |
| 35 | $\langle a \rangle$ | 1 | 35 | $\{a, \dots\}$ |
| 36 | $\langle a \rangle$ | 1 | 36 | $\{a, \dots\}$ |
| 37 | $\langle a \rangle$ | 1 | 37 | $\{a, \dots\}$ |
| 38 | $\langle a \rangle$ | 1 | 38 | $\{a, \dots\}$ |
| 39 | $\langle a \rangle$ | 1 | 39 | $\{a, \dots\}$ |
| 40 | $\langle a \rangle$ | 1 | 40 | $\{a, \dots\}$ |
| 41 | $\langle a \rangle$ | 1 | 41 | $\{a, \dots\}$ |
| 42 | $\langle a \rangle$ | 1 | 42 | $\{a, \dots\}$ |
| 43 | $\langle a \rangle$ | 1 | 43 | $\{a, \dots\}$ |
| 44 | $\langle a \rangle$ | 1 | 44 | $\{a, \dots\}$ |
| 45 | $\langle a \rangle$ | 1 | 45 | $\{a, \dots\}$ |
| 46 | $\langle a \rangle$ | 1 | 46 | $\{a, \dots\}$ |
| 47 | $\langle a \rangle$ | 1 | 47 | $\{a, \dots\}$ |
| 48 | $\langle a \rangle$ | 1 | 48 | $\{a, \dots\}$ |
| 49 | $\langle a \rangle$ | 1 | 49 | $\{a, \dots\}$ |
| 50 | $\langle a \rangle$ | 1 | 50 | $\{a, \dots\}$ |

Ref: SPADE (Sequential Pattern Discovery using Equivalent Class)

[M. Zaki 2001]

CloSpan: Mining Closed Sequential Patterns

- A closed sequential pattern s : There exists no superpattern s' such that $s' \supset s$, and s' and s have the same support
- Which ones are closed? $\langle abc \rangle$: 20, $\langle abcd \rangle$: 20, $\langle abcde \rangle$: 15
- Why directly mine closed sequential patterns?
 - Reduce # of (redundant) patterns
 - Attain the same expressive power
- Property P_1 : If $s \supset s'$, s is closed iff two project DBs have the same size
- Explore **Backward Subpattern** and **Backward Superpattern** pruning to prune redundant search space
- Greatly enhances efficiency (Yan, et al., SDM'03)

Apriori-Based Approach

- The Apriori property (anti-monotonicity): A size- k subgraph is frequent if and only if all of its subgraphs are frequent
- A candidate size- $(k+1)$ edge/vertex subgraph is generated if its corresponding two k -edge/vertex subgraphs are frequent
- Iterative mining process:
 - Candidate-generation \rightarrow candidate pruning \rightarrow support counting \rightarrow candidate elimination

SpiderMine: Mining Top-K Large Structural Patterns in a Massive Network

- Large patterns are informative to characterize a large network (e.g., social network, web, or bio-network)
- Similar to pattern fusion, mining large patterns should not aim for completeness but for representativeness of the target results
- SpiderMine (Zhu et al., VLDB'11): Mine top- K largest frequent substructure patterns whose diameter is bounded by D_{\max} with a probability at least $1-\epsilon$
- General idea: Large patterns are composed of a number of small components ("spiders"), which will eventually connect together after some rounds of pattern growth
- r-Spider**: An r -spider is a frequent graph pattern P such that there exists a vertex u of P , and all other vertices of P are within distance r from u

Pattern Space Pruning with Pattern Anti-Monotonicity

- Constraint c is **anti-monotone**
 - If an itemset S violates constraint c , so does any of its superset
 - That is, mining on itemset S can be terminated
- Ex. 1: $c_1: \text{sum}(S.Price) \leq v$ is **anti-monotone**
- Ex. 2: $c_2: \text{range}(S.profit) \leq 15$ is **anti-monotone**
- Itemset ab violates c_2 ($\text{range}(ab) = 40$)
- So does every superset of ab
- Ex. 3: $c_3: \text{sum}(S.Price) \geq v$ is **not anti-monotone**
- Ex. 4: Is $c_4: \text{support}(S) \geq \sigma$ anti-monotone?
- Yes! Apriori pruning is essentially pruning with an anti-monotonic constraint!

TID	Transaction	Item	Profit
10	a, b, c, d, f, h	a	40
20	b, c, d, f, g, h	b	0
30	b, c, d, f, g	c	-20
40	a, c, e, f, g	d	-15
		e	-30
		f	-10
		g	20
		h	5

PrefixSpan: A Pattern-Growth Approach

- | SID | Sequence | Prefix | Suffix (Projection) |
|-----|---------------------|---------------------|--------------------------|
| 10 | $\langle a \rangle$ | $\langle a \rangle$ | $\langle \{a\}b \rangle$ |
| 20 | $\langle a \rangle$ | $\langle a \rangle$ | $\langle \{a\}c \rangle$ |
| 30 | $\langle a \rangle$ | $\langle a \rangle$ | $\langle \{a\}d \rangle$ |
| 40 | $\langle a \rangle$ | $\langle a \rangle$ | $\langle \{a\}e \rangle$ |
- Prefix and suffix
 - Given $\langle a \rangle$, $\langle \{a\}b \rangle$, $\langle \{a\}c \rangle$, $\langle \{a\}d \rangle$, $\langle \{a\}e \rangle$, ...
 - Prefixes**: $\langle a \rangle$, $\langle a \rangle$, $\langle a \rangle$, $\langle a \rangle$, $\langle a \rangle$, ...
 - Suffixes**: Prefixes-based projection
 - PrefixSpan Mining: Prefix Projections
 - Step 1: Find length-1 sequential patterns
 - $\langle a \rangle$, $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$, $\langle e \rangle$, $\langle f \rangle$
 - Step 2: Divide search space and mine each projected DB
 - $\langle a \rangle$ -projected DB
 - $\langle b \rangle$ -projected DB
 - ...
 - $\langle f \rangle$ -projected DB, ...

PrefixSpan: Mining Prefix-Projected DBs

- | SID | Sequence |
|-----|---------------------|
| 10 | $\langle a \rangle$ |
| 20 | $\langle a \rangle$ |
| 30 | $\langle a \rangle$ |
| 40 | $\langle a \rangle$ |
- Length-1 sequential patterns: $\langle a \rangle$, $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$, $\langle e \rangle$, $\langle f \rangle$
 - Length-2 sequential patterns: $\langle a \rangle$, $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$, $\langle e \rangle$, $\langle f \rangle$
 - Major strength of PrefixSpan:
 - No candidate subsets to be generated
 - Projected DBs keep shrinking

Frequent (Sub)Graph Patterns

- Given a labeled graph dataset $D = \{G_1, G_2, \dots, G_n\}$, the supporting graph set of a subgraph g is $D_g = \{G_i \mid g \subseteq G_i, G_i \in D\}$
- support(g) = $|D_g| / |D|$
- A (sub)graph g is **frequent** if support(g) $\geq \min_sup$
- Ex.: Chemical structures
- Alternative:
 - Mining frequent subgraph patterns from a single large graph or network

Feature-Based Similarity Search

- Decompose a query graph into a set of features
- Feature-based similarity measure
 - Each graph is represented as a feature vector $X = [x_1, x_2, \dots, x_n]$
 - Similarity is defined by the distance of their corresponding vectors
- If graph G contains the major part of a query graph q , G should share a number of common features with q
- Given a relaxation ratio, one can calculate the maximal number of features that can be missed!

Different Kinds of Constraints Lead to Different Pruning Strategies

- Constraints can be categorized as
 - Pattern space pruning constraints** vs. **data space pruning constraints**
 - Pattern space pruning constraints**
 - Anti-monotonic**: If constraint c is violated, its further mining can be terminated
 - Monotonic**: If c is satisfied, no need to check c again
 - Succinct**: If the constraint c can be enforced by directly manipulating the data
 - Convertible**: c can be converted to monotonic or anti-monotonic if items can be properly ordered in processing
 - Data space pruning constraints**
 - Data succinct**: Data space can be pruned at the initial pattern mining process
 - Data anti-monotonic**: If a transaction t does not satisfy c , then t can be pruned to reduce data processing effort

Data Space Pruning with Data Anti-Monotonicity

- A constraint c is **data anti-monotone**: In the mining process, if a data entry t cannot satisfy a pattern p under c , t cannot satisfy p 's superset either
- Data space pruning: Data entry t can be pruned
- Ex. 1: $c_1: \text{sum}(S.Profit) \geq v$ is **data anti-monotone**
- Let constraint c_1 be: $\text{sum}(S.Profit) \geq 25$
- $T_{30}: \{b, c, d, f, g\}$ can be removed since none of their combinations can make an S whose sum of the profit is ≥ 25
- Ex. 2: $c_2: \text{min}(S.Price) \leq v$ is **data anti-monotone**
- Consider $v = 5$ but every item in a transaction, say T_{30} , has a price higher than 10
- Ex. 3: $c_3: \text{range}(S.Profit) > 25$ is **data anti-monotone**

TID	Transaction	Item	Profit
10	a, b, c, d, f, h	a	40
20	b, c, d, f, g, h	b	0
30	b, c, d, f, g	c	-20
40	a, c, e, f, g	d	-15
		e	-30
		f	-10
		g	20
		h	5