

# **CS 412 Intro. to Data Mining Chapter 7 : Advanced Frequent Pattern Mining** Qi Li, Computer Science, Univ. Illinois at Urbana-Champaign, 2018



### **Chapter 7 : Advanced Frequent Pattern Mining**

Mining Diverse Patterns



- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining
- Graph Pattern Mining
- Pattern Mining Application: Mining Software Copy-and-Paste Bugs

#### **Summary**

# **Mining Diverse Patterns**

- Mining Multiple-Level Associations
- Mining Multi-Dimensional Associations
- Mining Quantitative Associations
- Mining Negative Correlations
- Mining Compressed and Redundancy-Aware Patterns

## **Mining Multiple-Level Frequent Patterns**

- Items often form hierarchies
  - E.g.: Dairyland 2% milk;Wonder wheat bread
- How to set min-support thresholds?



- Uniform min-support across multiple levels (reasonable?)
- Level-reduced min-support: Items at the lower level are expected to have lower support
- **G** Efficient mining: *Shared* multi-level mining
  - Use the lowest min-support to pass down the set of candidates

### Redundancy Filtering at Mining Multi-Level Associations

- Multi-level association mining may generate many redundant rules
- Redundancy filtering: Some rules may be redundant due to "ancestor" relationships between items
  - $\square \quad milk \Rightarrow wheat bread [support = 8\%, confidence = 70\%] (1)$
  - $\square 2\% \text{ milk} \Rightarrow \text{wheat bread [support = 2\%, confidence = 72\%] (2)}$ 
    - □ Suppose the 2% milk sold is about ¼ of milk sold in gallons
      - (2) should be able to be "derived" from (1)

### Redundancy Filtering at Mining Multi-Level Associations

- $\square milk \Rightarrow wheat bread [support = 8\%, confidence = 70\%] (1)$
- 2% milk ⇒ wheat bread [support = 2%, confidence = 72%] (2) A rule is redundant if its support is close to the "expected" value, according to its "ancestor" rule, and it has a similar confidence as its "ancestor"
- Rule (1) is an ancestor of rule (2), which one to prune?

# **Customized Min-Supports for Different Kinds of Items**

- We have used the same min-support threshold for all the items or item sets to be mined in each association mining
- In reality, some items (e.g., diamond, watch, ...) are valuable but less frequent
- It is necessary to have customized min-support settings for different kinds of items
- One Method: Use group-based "individualized" min-support
  - □ E.g., {diamond, watch}: 0.05%; {bread, milk}: 5%; ...
  - How to mine such rules efficiently?
    - Existing scalable mining algorithms can be easily extended to cover such cases

# **Mining Multi-Dimensional Associations**

- Single-dimensional rules (e.g., items are all in "product" dimension)
  - $\square \quad buys(X, "milk") \Rightarrow buys(X, "bread")$
- □ Multi-dimensional rules (i.e., items in ≥ 2 dimensions or predicates)
  - Inter-dimension association rules (*no repeated predicates*)
    - □ age(X, "18-25")  $\land$  occupation(X, "student")  $\Rightarrow$  buys(X, "coke")
  - Hybrid-dimension association rules (repeated predicates)
    - □ age(X, "18-25")  $\land$  buys(X, "popcorn")  $\Rightarrow$  buys(X, "coke")
- Attributes can be categorical or numerical
  - Categorical Attributes (e.g., *profession, product*: no ordering among values): Data cube for inter-dimension association
  - Quantitative Attributes: Numeric, implicit ordering among values discretization, clustering, and gradient approaches

# **Mining Quantitative Associations**

- Mining associations with numerical attributes
  - **E.g.:** Numerical attributes: age and salary
- Methods
  - Static discretization based on predefined concept hierarchies
    - Discretization on each dimension with hierarchy
      - □ age:  $\{0-10, 10-20, ..., 90-100\} \rightarrow \{young, mid-aged, old\}$
  - 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
  - **E.g.:** 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
  - □ E.g.: (Gender = female) ^ (South = yes)  $\Rightarrow$  mean wage = \$6.3/hr
- □ Rule condition can be categorical or numerical (quantitative rules)
  - □ E.g.: Education in [14-18] (yrs)  $\Rightarrow$  mean wage = \$11.64/hr

### **Rare Patterns**

- Rare patterns
  - □ Very low support but interesting (e.g., buying Rolex watches)
  - How to mine them? Setting individualized, group-based min-support thresholds for different groups of items

# **Negative Patterns**

- Negative patterns
- Negatively correlated: Unlikely to happen together
- Ex.: Since it is unlikely that the same customer buys both a Ford Expedition (an SUV car) and a Ford Fusion (a hybrid car), buying a Ford Expedition and buying a Ford Fusion are likely negatively correlated patterns
- □ How to define negative patterns?
- A support-based definition of negative correlated patterns
- If itemsets A and B are both frequent but rarely occur together, i.e., sup(A U B) << sup(A) × sup(B)</p>

Does this remind you the definition of *lift*?

### **Defining Negative Correlated Patterns**

- □ Is this a good definition for large transaction datasets?
- Ex.: Suppose a store sold two needle packages A and B 100 times each, but only one transaction contained both A and B
  - □ When there are in total 200 transactions, we have
    - □  $s(A \cup B) = 0.005, s(A) \times s(B) = 0.25, s(A \cup B) << s(A) \times s(B)$
  - □ But when there are 10<sup>5</sup> transactions, we have
    - □  $s(A \cup B) = 1/10^5$ ,  $s(A) \times s(B) = 1/10^3 \times 1/10^3$ ,  $s(A \cup B) > s(A) \times s(B)$
  - What is the problem?—Null transactions: The support-based definition is not null-invariant!

#### Defining Negative Correlation: Need Null-Invariance in Definition

- A good definition on negative correlation should take care of the nullinvariance 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

 $(s(A \cup B)/s(A) + s(A \cup B)/s(B))/2 < \epsilon$ ,

where  $\varepsilon$  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 10<sup>5</sup> transactions
    - □ If  $\epsilon$  = 0.01, we have

 $(s(A \cup B)/s(A) + s(A \cup B)/s(B))/2 = (0.01 + 0.01)/2 < \epsilon$ 

# **Mining Compressed Patterns**

Pat-ID	Item-Sets	Support
P1	{38,16,18,12}	205227
P2	{38,16,18,12,17}	205211
P3	{39,38,16,18,12,17}	101758
P4	{39,16,18,12,17}	161563
Р5	{39,16,18,12}	161576

- Why mining compressed patterns? Too many scattered patterns but not so meaningful
- Pattern distance measure

 $Dist(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}$ 

- δ-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

- Closed patterns
  - P1, P2, P3, P4, P5
  - Emphasizes too much on support
- Max-patterns
  - □ P3: information loss
- Desired output (a good balance):
  - □ P2, P3, P4

### **Redundancy-Aware Top-k Patterns**

Desired patterns: high significance & low redundancy



(c) traditional top-k

(d) summarization

- 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

# **Chapter 7 : Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining
- Graph Pattern Mining
- Pattern Mining Application: Mining Software Copy-and-Paste Bugs

#### **Summary**





# **Constraint-Based Pattern Mining**

- Why Constraint-Based Mining?
- Different Kinds of Constraints: Different Pruning Strategies
- Constrained Mining with Pattern Anti-Monotonicity
- Constrained Mining with Pattern Monotonicity
- Constrained Mining with Convertible Constraints
- Constrained Mining with Data Anti-Monotonicity
- Constrained Mining with Succinct Constraints
- Handling Multiple Constraints

# Why Constraint-Based Mining?

- Pattern mining in practice: Often a user-guided, interactive process
  - User directs what to be mined using a data mining query language (or a graphical user interface), specifying various kinds of constraints
- □ What is constraint-based mining?
  - Mine together with user-provided constraints
- Why constraint-based mining?
  - User flexibility: User provides constraints on what to be mined
  - Optimization: System explores such constraints for mining efficiency
    - E.g., Push constraints deeply into the mining process

#### Various Kinds of User-Specified Constraints in Data Mining

- □ Knowledge type constraint—Specifying what kinds of knowledge to mine
  - **E.g.**: Classification, association, clustering, outlier finding, ...
- **Data constraint**—using SQL-like queries
  - **E.g.:** Find products sold together in NY stores this year
- Dimension/level constraint—similar to projection in relational database
  - **E.g.:** In relevance to region, price, brand, customer category
- Interestingness constraint—various kinds of thresholds
  - **E**.g.: Strong rules: min\_sup  $\geq$  0.02, min\_conf  $\geq$  0.6, min\_correlation  $\geq$  0.7
- Rule (or pattern) constraint

The focus of this study

E.g.: Small sales (price < \$10) triggers big sales (sum > \$200)

### **Pattern Space Pruning with Pattern Anti-Monotonicity**

Item	Price	Profit
а	100	40
b	40	0
С	150	-20
d	35	-15
е	55	-30
f	45	-10
g	80	20
h	10	5

Note: item.price > 0 Profit can be negative

#### A 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
- E.g. 1:  $c_1$ : sum(S.price)  $\leq v$  is anti-monotone

- E.g. 2:  $c_2$ : range(S.profit)  $\leq$  15 is anti-monotone
- Itemset *ab* violates c<sub>2</sub> (range(ab) = 40)
- So does every superset of *ab*

### **Pattern Space Pruning with Pattern Anti-Monotonicity**

TID	Transa	action	E.g. 3. $c_3$ : sum(S.Price) $\ge v$ is not anti-monotone
10	a, b, c, o	d, f, h	
20	b, c, d, t	f, g, h	
30	b, c, d, t	f, g	
40	a, c, e, f	;, g	
min_	sup = 2		E.g. 4. Is $c_4$ : support(S) $\geq \sigma$ anti-monotone?
Item	Price	Profit	Yes! Apriori pruning is essentially pruning with an ant
а	100	40	monotonic constraint!
b	40	0	
С	150	-20	
d	35	-15	
е	55	-30	
f	45	-10	
g	80	20	
h	10	5	

### **Pattern Monotonicity and Its Roles**

TID	Transaction			
10	a, b, c, c	l, f, h		
20	b, c, d, f	; g, h		
30	b, c, d, f	, g		
40	a, c, e, f	, g		
min_sup = 2				
Item	Price Profit			
а	100	40		
b	40	0		
С	150	-20		
d	35	-15		
е	55	-30		
f	45	-10		
g	80	20		
	•••			

- 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
- E.g. 1:  $c_1$ : sum(S.Price)  $\geq v$  is monotone

• E.g. 2:  $c_2$ : min(S.Price)  $\leq v$  is monotone

- E.g. 3:  $c_3$ : range(S.profit)  $\geq$  15 is monotone
  - Itemset *ab* satisfies  $c_3$
  - So does every superset of *ab*

### **Apriori for Pattern Anti-Monotone Constraint**



### **Convertible Constraints: Ordering Data in Transactions**

TID	Transaction			
10	a, b, c, o	d, f, h		
20	a, b, c, o	d, f, g, h		
30	b, c, d, t	f, g		
40	a, c, e, f	, g		
min_s	sup = 2			
Item	Price Profit			
а	100	40		
b	40	0		
С	150	-20		
d	35	-15		
е	55	-30		
f	45	-5		
g	80	30		
h	10 5			

- Convert tough constraints into (anti-)monotone by proper ordering of items in transactions
- Examine c<sub>1</sub>: avg(S.profit) > 20
- Order items in (profit) value-descending order
  - <a, g, f, b, h, d, c, e>
- An itemset *ab* violates c<sub>1</sub> (avg(ab) = 20)
  - So does ab\* (i.e., ab-projected DB)
  - C<sub>1</sub>: anti-monotone if patterns grow in the right order!

# **Can item-reordering work for Apriori?**

TID	Transaction			
10	a, b, c, (	d <i>,</i> f, h		
20	a, b, c, (	d, f, g, h		
30	b, c, d, <sup>-</sup>	f, g		
40	a, c, e, f	f, g		
min_s	sup = 2			
Item	Price Profit			
а	100	40		
b	40	40 0		
С	150 -20			
d	35 -15			
е	55 -30			
f	45 -5			
g	80 30			
h	10 5			

#### constraint: avg(S.profit) > 20



avg(gf) = 12.5 < 20, avg(af) = 17.5 < 20, avg(ag) = 35 > 20

- But avg(agf) = 21.7 > 20
- Apriori will not generate "agf" as a candidate

# **Data Space Pruning with Data Anti-Monotonicity**

TID	Transaction
10	a, b, c, d, f, h
20	b, c, d, f, g, h
30	b, c, d, f, g
40	a, c, e, f, g

- A constraint c is *data anti-monotone*: In the mining process, if a data entry *t* cannot contribute to a pattern *p* satisfying *c*, *t* cannot contribute to *p*'s superset either
- Data space pruning: Data entry *t* can be pruned

min_sup = 2				
Item	Price	Profit		
а	100	40		
b	40	0		
С	150	-20		
d	35	-15		
е	55	-30		
f	45	-10		
g	80	20		
h	10	5		

- E.g. 1:  $c_1$ : sum(S.Profit)  $\ge v$  is data anti-monotone
- □ Let constraint  $c_1$  be:  $sum(S.Profit) \ge 25$ 
  - T<sub>30</sub>: {b, c, d, f, g} can be removed since none of their combinations can make an S whose sum of the profit is ≥ 25

# **Data Space Pruning with Data Anti-Monotonicity**

TID	Transaction
10	a, b, c, d, f, h
20	b, c, d, f, g, h
30	b, c, d, f, g
40	a, c, e, f, g

- A constraint c is *data anti-monotone*: In the mining process, if a data entry *t* cannot contribute to a pattern *p* satisfying *c*, *t* cannot contribute to *p*'s superset either
- Data space pruning: Data entry *t* can be pruned

min_sı		
Item	Price	Profit
а	100	40
b	40	0
С	150	-20
d	35	-15
е	55	-30
f	45	-10
g	80	20
h	10	5

- **E.g.** 2:  $c_2$ : *min*(*S*.*Price*)  $\leq v$  is data anti-monotone
  - Consider v = 5 but every item in a transaction, say T<sub>50</sub>, has a price higher than 10

**E.g.** 3:  $c_3$ : range(S.Profit) > 25 is data anti-monotone

# **Data Space Pruning Should Be Explored Recursively**

TID	Transaction	Item	Profit	
10	a, b, c, d, f, h	а	40	
20	b, c, d, f, g, h	b	0	
30	b, c, d, f, g	С	-20	
40	a, c, e, f, g	d	-15	
mir	1 sup = 2	е	-30	
		f	-10	
		g	20	
		h	5	

- Example. c<sub>3</sub>: *range(S.Profit) > 25*
- We check b's projected database
- But item "a" is infrequent (sup = 1)



- After removing "a (40)" from  $T_{10}$ 
  - **T**<sub>10</sub> cannot satisfy  $c_3$  any more
    - □ Since "b (0)" and "c (−20), d (−15), f (−10), h (5)"

By removing  $T_{10}$ , we can also prune "h" in  $T_{20}$  b's-proj. DB



### **Data Space Pruning Should Be Explored Recursively**



□ Note:  $c_3$  prunes  $T_{10}$  effectively only after "a" is pruned (by min-sup) in b's projected DB

# Succinctness: Pruning Both Data and Pattern Spaces

- Succinctness: If the constraint *c* can be enforced by directly manipulating the data
- E.g. 1: To find those patterns without item *i* 
  - Remove *i* from DB and then mine (pattern space pruning)
- E.g. 2: To find those patterns containing item *i*
- Mine only *i*-projected DB (data space pruning)
- E.g. 3:  $c_3$ : min(S.Price)  $\leq v$  is succinct
  - Start with only items whose price ≤ v and remove transactions with high-price items only (pattern + data space pruning)
  - E.g. 4:  $c_4$ : sum(S.Price)  $\ge v$  is not succinct
  - It cannot be determined beforehand since sum of the price of itemset S keeps increasing

#### **Apriori + Succinct Constraint**



#### **Constrained FP-Growth: Push a Succinct Constraint Deep**



### Different Kinds of Constraints Lead to Different Pruning Strategies

In summary, constraints can be categorized as Pattern space pruning constraints vs. data space pruning constraints

<ul> <li>Anti-monotonic: If constraint c is violated, its further mining can be terminated</li> <li>Monotonic: If c is satisfied, no need to check c again</li> <li>Convertible: c can be converted to monotonic or anti-monotonic if items can be properly ordered in processing</li> <li>Succinct: If the constraint c can be enforced by directly manipulating the data</li> </ul>	Patter	rn space pruning constraints	Data space pruning constraints		
		<ul> <li>Anti-monotonic: If constraint c is violated, its further mining can be terminated</li> <li>Monotonic: If c is satisfied, no need to check c again</li> <li>Convertible: c can be converted to monotonic or anti-monotonic if items can be properly ordered in processing</li> <li>Succinct: If the constraint c can be enforced by directly manipulating the data</li> </ul>	•	<ul> <li>Data succinct: Data</li> <li>space can be pruned at</li> <li>the initial pattern</li> <li>mining process</li> <li>Data anti-monotonic: If</li> <li>a transaction t does not</li> <li>satisfy c, then t can be</li> <li>pruned to reduce data</li> <li>processing effort</li> </ul>	

## **How to Handle Multiple Constraints?**

- It is beneficial to use multiple constraints in pattern mining
- But different constraints may require potentially conflicting item-ordering
- If there exists conflict ordering between  $c_1$  and  $c_2$ 
  - Try to sort data and enforce *one constraint* first (which one?)
  - Then enforce the other constraint when mining the projected databases
- E.g.  $c_1$ : avg(S.profit) > 20, and  $c_2$ : avg(S.price) < 50
  - Assum c<sub>1</sub> has more pruning power
    - Sort in profit descending order and use c<sub>1</sub> first
  - For each project DB, sort trans. in price ascending order and use c<sub>2</sub> at mining

# **Summary: Constraint-Based Pattern Mining**

- Why Constraint-Based Mining?
- Different Kinds of Constraints: Different Pruning Strategies
- Constrained Mining with Pattern Anti-Monotonicity
- Constrained Mining with Pattern Monotonicity
- Constrained Mining with Convertible Constraints
- Constrained Mining with Data Anti-Monotonicity
- Constrained Mining with Succinct Constraints
- Handling Multiple Constraints
### **Chapter 7 : Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining



- Graph Pattern Mining
- Pattern Mining Application: Mining Software Copy-and-Paste Bugs

#### **Summary**

## **Sequential Pattern Mining**

Sequential Pattern and Sequential Pattern Mining

- GSP: Apriori-Based Sequential Pattern Mining
- SPADE: Sequential Pattern Mining in Vertical Data Format
- PrefixSpan: Sequential Pattern Mining by Pattern-Growth
- CloSpan: Mining Closed Sequential Patterns
- Constraint-Based Sequential-Pattern Mining

### **Sequence Databases & Sequential Patterns**

- Sequential pattern mining has broad applications
  - Customer shopping sequences
    - Purchase a laptop first, then a digital camera, and then a smartphone, within 6 months
  - Medical treatments, natural disasters (e.g., earthquakes), science & engineering processes, stocks and markets, ...
  - □ Weblog click streams, calling patterns, ...
  - Software engineering: Program execution sequences, ...
  - Biological sequences: DNA, protein, ...
- Transaction DB, sequence DB vs. time-series DB
- Gapped vs. non-gapped sequential patterns
  - Shopping sequences, clicking streams vs. biological sequences

## **Sequential Pattern and Sequential Pattern Mining**

Sequential pattern mining: Given a set of sequences, find the complete set of frequent subsequences (i.e., satisfying the min\_sup threshold)

#### A sequence: < (ef) (ab) (df) c b >

- An <u>element</u> may contain a set of <u>items</u> (also called <u>events</u>)
- \* Items within an element are unordered and we list them alphabetically

#### A sequence database

SID	Sequence
10	<a(<u>abc)(a<u>c</u>)d(cf)&gt;</a(<u>
20	<(ad)c(bc)(ae)>
30	<(ef)( <u>ab</u> )(df) <u>c</u> b>
40	<eg(af)cbc></eg(af)cbc>

#### **Sequential Pattern and Sequential Pattern Mining**

Sequential pattern mining: Given a set of sequences, find the complete set of frequent subsequences (i.e., satisfying the min\_sup threshold)

<a(bc)dc> is a **subsequence** of <<u>a</u>(a<u>bc</u>)(ac)<u>d(c</u>f)>

Given support threshold min\_sup = 2, <(ab)c> is a <u>sequential pattern</u>

#### A sequence database

SID	Sequence
10	<a(<u>abc)(a<u>c</u>)d(cf)&gt;</a(<u>
20	<(ad)c(bc)(ae)>
30	<(ef)( <u>ab</u> )(df) <u>c</u> b>
40	<eg(af)cbc></eg(af)cbc>

#### **Sequential Pattern Mining Algorithms**

- Algorithm requirement: Efficient, scalable, finding complete set, incorporating various kinds of user-specific constraints
- The Apriori property still holds: If a subsequence  $s_1$  is infrequent, none of  $s_1$ 's super-sequences can be frequent
- Representative algorithms
  - **GSP** (Generalized Sequential Patterns): Srikant & Agrawal @ EDBT'96)
  - □ Vertical format-based mining: **SPADE** (Zaki@Machine Leanining'00)
  - Pattern-growth methods: PrefixSpan (Pei, et al. @TKDE'04)
- □ Mining closed sequential patterns: CloSpan (Yan, et al. @SDM'03)
- Constraint-based sequential pattern mining (to be covered in the constraint mining section)

#### **GSP: Apriori-Based Sequential Pattern Mining**



#### 44

#### **GSP: Apriori-Based Sequential Pattern Mining**

#### Generate length-2 candidate sequences

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

cotc

singleton x singleton

			sets			
	<a></a>	<b></b>	<c></c>	<d></d>	<e></e>	<f></f>
<a></a>		<(ab)>	<(ac)>	<(ad)>	<(ae)>	<(af)>
<b></b>			<(bc)>	<(bd)>	<(be)>	<(bf)>
<c></c>				<(cd)>	<(ce)>	<(cf)>
<d></d>					<(de)>	<(df)>
<e></e>						<(ef)>
<f></f>						
					-	

#### Apriori Pruning

 w/o pruning: 8\*8 + 8\*7/2 = 92
 length-2 candidates
 w/ pruning: 6\*6 + 6\*5/2 = 51
 length-2 candidates

min_	_ <i>sup</i> = 2	

sup

Cand.



#### **GSP Mining and Pruning**

<(bd)cba>

5<sup>th</sup> scan: 1 cand. 1 length-5 seq. pat.

4<sup>th</sup> scan: 8 cand. 7 length-4 seq. pat.

3<sup>rd</sup> scan: 46 cand. 20 length-3 seq. pat. 20 cand. not in DB at all

2<sup>nd</sup> scan: 51 cand. 19 length-2 seq. pat. 10 cand. not in DB at all

1<sup>st</sup> scan: 8 cand. 6 length-1 seq. pat.

<abba> <(bd)bc> ...

<abb> <aab> <aba> <bab> ...

<aa> <ab> ... <af> <ba> <bb> ... <ff> <(ab)> ... <(ef)>

<a> <b> <c> <d> <e> <f> <g> <h>

Remove			
Candidates not in DB			
Candidates < min_sup			

	$min_sup = 2$
SID	Sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bcb(ade)></a(bd)bcb(ade)>

#### **GSP Mining and Pruning**

- 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

GSP (Generalized Sequential Patterns): Srikant & Agrawal @ EDBT'96)

#### Sequential Pattern Mining in Vertical Data Format: The SPADE Algorithm

- A sequence database is mapped to: <SID, EID>
- Grow the subsequences (patterns) one item at a time by Apriori candidate generation

SID	Sequence	
1	<a(<u>abc)(a<u>c</u>)d(cf)&gt;</a(<u>	
2	<(ad)c(bc)(ae)>	
3	<(ef)( <u>ab</u> )(df) <u>c</u> b>	
4 <eg(af)cbc></eg(af)cbc>		
min_sup = 2		
f: SPADE ( <u>S</u> equenti		

Ref: SPADE (<u>S</u>equential <u>PA</u>ttern <u>D</u>iscovery using <u>E</u>quivalent Class) [M. Zaki 2001]

SID	EID	Items
1	1	a
1	2	$^{\rm abc}$
1	3	ac
1	4	d
1	5	$\mathbf{cf}$
2	1	$\operatorname{ad}$
2	2	с
2	3	$\mathbf{bc}$
2	4	ae
3	1	ef
3	2	$^{\rm ab}$
3	3	df
3	4	$\mathbf{c}$
3	5	b
4	1	e
4	2	ģ
4	3	$\mathbf{af}$
4	4	с
4	5	b
4	6	с

TATA

	a	k	)				
SID	$\operatorname{EID}$	$\mathbf{SID}$	EID				
1	1	1	2				
1	2	2	3				
1	3	3	2				
2	1	3	<b>5</b>				
2	4	4	<b>5</b>				
3	2						
4	3						
	$^{\rm ab}$				ba		
SID	EID (a)	EID(b	)	SID	EID (b)	EID(a)	
1	1	2		1	2	3	
2	1	3		2	3	4	
3	2	5					
4	3	5					
		_					
		aba	ւ				
SID	EID (	a) E	ID(b	)	EID(a)		
1	1		<b>2</b>		3		
2	1		3		4		

### **PrefixSpan: A Pattern-Growth Approach**

SID	Sequence	min_sup =	= 2
10	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>	Prefix	Suffix (Projection)
20	<(ad)c(bc)(ae)>	<a></a>	<(abc)(ac)d(cf)>
30	<(ef)(ab)(df)cb>	<aa></aa>	<(_bc)(ac)d(cf)>
40	$<(c_1)(a_2)(a_1)(a_3)$	<ab></ab>	<(_c)(ac)d(cf)>
40			

- PrefixSpan Mining: Prefix Projections
  - □ Step 1: Find length-1 sequential patterns
    - <a>, <b>, <c>, <d>, <e>, <f>
  - Step 2: Divide search space and mine each projected DB
    - <a>-projected DB,
    - <b>-projected DB,

<f>-projected DB, ...

- Prefix and suffix
  - Given <a(abc)(ac)d(cf)>
  - Prefixes: <a>, <aa>,<a(ab)>, <a(abc)>, ...
  - Suffix: Prefixes-based projection

PrefixSpan (Prefix-projected Sequential pattern mining) Pei, et al. @TKDE'04

<sup>...</sup> 

# **PrefixSpan: Mining Prefix-Projected DBs**



#### 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

 $s = \langle a(abc)(ac)d(cf) \rangle$   $s = \langle abc)(ac)d(cf) \rangle$   $\langle abc)(ac)d(cf) \rangle$   $\langle abb \rangle$   $s | \langle ab \rangle \rangle$   $\langle (abc)(ac)d(cf) \rangle$ 

#### **CloSpan: Mining Closed Sequential Patterns**

- □ A closed sequential pattern s: There exists no superpattern s' such that s' ⊃ s, and s' and s have the same support
- □ Which ones are closed? <abc>: 20, <abcd>:20, <abcd>: 15

### **CloSpan: Mining Closed Sequential Patterns**

- □ Why directly mine closed sequential patterns?
  - Reduce # of (redundant) patterns
  - Attain the same expressive power
- Property P<sub>1</sub>: If s > s<sub>1</sub>, 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)



## **CloSpan: When Two Projected DBs Have the Same Size**



#### **Constraint-Based Sequential-Pattern Mining**

- Share many similarities with constraint-based itemset mining
- Anti-monotonic: If S violates *c*, the super-sequences of S also violate *c*
- sum(S.price) < 150; min(S.value) > 10
- Monotonic: If S satisfies c, the super-sequences of S also do so
  - element\_count (S) > 5; S  $\supseteq$  {PC, digital\_camera}
- □ Data anti-monotonic: If a sequence  $s_1$  with respect to S violates  $c_3$ ,  $s_1$  can be removed
  - □  $c_3$ : sum(S.price) ≥ v
- **Succinct:** Enforce constraint c by explicitly manipulating data
  - □ S  $\supseteq$  {i-phone, MacAir}
- **Convertible:** Projection based on the sorted value not sequence order
  - value\_avg(S) < 25; profit\_sum (S) > 160
  - max(S)/avg(S) < 2; median(S) min(S) > 5

#### **Timing-Based Constraints in Seq.-Pattern Mining**

- Order constraint: Some items must happen before the other
  - □ {algebra, geometry}  $\rightarrow$  {calculus} (where " $\rightarrow$ " indicates ordering)
  - Anti-monotonic: Constraint-violating sub-patterns pruned
- Min-gap/max-gap constraint: Confines two elements in a pattern
  - $\Box \quad E.g., mingap = 1, maxgap = 4$
  - Succinct: Enforced directly during pattern growth
- Max-span constraint: Maximum allowed time difference between the 1<sup>st</sup> and the last elements in the pattern
  - □ E.g., maxspan (S) = 60 (days)
  - Succinct: Enforced directly when the 1<sup>st</sup> element is determined
- Window size constraint: Events in an element do not have to occur at the same time: Enforce max allowed time difference
  - E.g., window-size = 2: Various ways to merge events into elements

#### **Episodes and Episode Pattern Mining**

- Episodes and regular expressions: Alternative to seq. patterns
  - Serial episodes: AB a total order relationship: first A then B
- □ E.g. Given a large shopping sequence database, one may like to find
  - □ Suppose the pattern order follows the template (A|B)C\*(D E), and
  - Sum of the prices of A, B, C\*, D, and E is greater than \$100, where C\* means C appears \*-times
  - How to efficiently mine such episode patterns?

## **Chapter 7 : Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining
- Graph Pattern Mining



Pattern Mining Application: Mining Software Copy-and-Paste Bugs

#### **Summary**

#### What Is Graph Pattern Mining?

Chem-informatics:

Mining frequent chemical compound structures

ŇН



**M.Renz P.Kroumlger** Social networks, web communities, tweets, ... J.Aszligfalg\A.Zuumlfle **P.Kunath C.Boumlhm** Finding frequent research collaboration subgraphs A.Pryakhin **R.Schneider H.Kriegel B.BraunmuumIller M.Schubert R.Ramakrishnan H.G.Molina S.Brecheisen B.Seeger** M.Poumltke M.Pfeifle .Silberschatz Ioannidis J.D.Ullman **C.Heinz** J.Widom J.Sander M.Livnv M.N.Garofalakis Y.Papakonstantinou **Y.Sagiv** S.Sudarshan C.Beeri

#### Frequent (Sub)Graph Patterns

- □ Given a labeled graph dataset D = {G<sub>1</sub>, G<sub>2</sub>, ..., G<sub>n</sub>), the supporting graph set of a subgraph g is D<sub>g</sub> = {G<sub>i</sub> |  $g \subseteq G_i$ , G<sub>i</sub> ∈ D}
  - $\square \quad \text{support}(g) = |\mathsf{D}_g| / |\mathsf{D}|$
- □ A (sub)graph g is **frequent** if  $support(g) \ge min_sup$
- Ex.: Chemical structures
- Alternative:
  - Mining frequent subgraph patterns from a single large graph or network



## **Applications of Graph Pattern Mining**

- Bioinformatics
  - Gene networks, protein interactions, metabolic pathways
- Chem-informatics: Mining chemical compound structures
- □ Social networks, web communities, tweets, ...
- Cell phone networks, computer networks, ...
- U Web graphs, XML structures, Semantic Web, information networks
- □ Software engineering: Program execution flow analysis
- Building blocks for graph classification, clustering, compression, comparison, and correlation analysis
- Graph indexing and graph similarity search

#### Graph Pattern Mining Algorithms: Different Methodologies

- Generation of candidate subgraphs
  - Apriori vs. pattern growth (e.g., FSG vs. gSpan)
- Search order
  - Breadth vs. depth
- Elimination of duplicate subgraphs
  - Passive vs. active (e.g., gSpan [Yan & Han, 2002])
- Support calculation
  - Store embeddings (e.g., GASTON [Nijssen & Kok, 2004], FFSM [Huan, Wang, & Prins, 2003], MoFa [Borgelt & Berthold, ICDM'02])
- Order of pattern discovery
  - □ Path  $\rightarrow$  tree  $\rightarrow$  graph (e.g., GASTON [Nijssen & Kok, 2004])

## **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



#### Candidate Generation: Vertex Growing vs. Edge Growing

Methodology: Breadth-search, Apriori joining two size-k graphs

□ Many possibilities at generating size-(*k*+1) candidate graphs



Generating new graphs with one more vertex

- AGM (Inokuchi, Washio, & Motoda, PKDD'00)
- Generating new graphs with one more edge
  - □ FSG (Kuramochi & Karypis, ICDM'01)
- Performance shows via edge growing is more efficient

# **Pattern-Growth Approach**

- Depth-first growth of subgraphs from k-edge to (k+1)-edge, then (k+2)-edge subgraphs
  (k+2)-edge
- 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 <u>complete</u>
- DFS code: Flatten a graph into a sequence using depth-first search



## Why Mine Closed Graph Patterns?

- Challenge: An **n**-edge frequent graph may have 2<sup>n</sup> subgraphs
- Motivation: Explore *closed frequent subgraphs* to handle graph pattern explosion problem
- A frequent graph G is *closed* if there exists no supergraph of G that carries the same support as G



If this subgraph is *closed* in the graph dataset, it implies that none of its frequent super-graphs carries the same support

- Lossless compression: Does not contain non-closed graphs, but still ensures that the mining result is complete
- Algorithm CloseGraph: Mines closed graph patterns directly

## **CloseGraph: Directly Mining Closed Graph Patterns**

CloseGraph: Mining closed graph patterns by extending gSpan (Yan & Han, KDD'03)



At what condition can we stop searching their children, i.e., early termination?

- Suppose G and G<sub>1</sub> are frequent, and G is a subgraph of G<sub>1</sub>
- If in any part of the graph in the dataset where G occurs, G<sub>1</sub> also occurs, then we need not grow G (except some special, subtle cases), since none of G's children will be closed except those of G<sub>1</sub>

#### **Experiment and Performance Comparison**

- The AIDS antiviral screen compound dataset from NCI/NIH
- The dataset contains 43,905 chemical compounds
- Discovered patterns: The smaller minimum support, the bigger and more interesting subgraph patterns discovered



## **Chapter 7 : Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining
- Graph Pattern Mining
- Pattern Mining Application: Mining Software Copy-and-Paste Bugs



#### **Summary**

### **Pattern Mining Application: Software Bug Detection**

#### Mining rules from source code

- Bugs as deviant behavior (e.g., by statistical analysis)
- Mining programming rules (e.g., by frequent itemset mining)
- Mining function precedence protocols (e.g., by frequent subsequence mining)
- Revealing neglected conditions (e.g., by frequent itemset/subgraph mining)
- Mining rules from revision histories
  - By frequent itemset mining
- Mining copy-paste patterns from source code
  - □ Find copy-paste bugs (e.g., CP-Miner [Li et al., OSDI'04]) (to be discussed here)
    - Reference: Z. Li, S. Lu, S. Myagmar, Y. Zhou, "<u>CP-Miner</u>: A Tool for Finding Copy-paste and Related Bugs in Operating System Code", OSDI'04

# **Application Example: Mining Copy-and-Paste Bugs**

- Copy-pasting is common
  - 12% in Linux file system
  - 19% in X Window system
- Copy-pasted code is error-prone
- □ Mine *"forget-to-change"* bugs by sequential pattern mining
  - Build a sequence database from source code
  - Mining sequential patterns
  - Finding mismatched identifier names & bugs

```
Courtesy of Yuanyuan Zhou@UCSD
```

```
void init prom meminit(void)
  for (i=0; i<n; i++) {
    total[i].adr = list[i].addr;
     total[i].bytes = list[i].size;
     total[i].more = &total[i+1];
                                     Code copy-and-
  . . . . . .
                                     pasted but forget
for (i=0; i<n; i++) {
                                     to change "id"!
     taken[i].adr = list[i].addr;
     taken[i].bytes = list[i].size,
     taken[i].more = &total[i+1];
      (Simplified example from linux-
```

2.6.6/arch/sparc/prom/memory.c)

## **Building Sequence Database from Source Code**

#### (mapped to) Imapped to Statement

- Tokenize each component
  - □ Different operators, constants, key words
     → different tokens
  - □ Same type of identifiers  $\rightarrow$  same token
- $\square$  Program  $\rightarrow$  A long sequence
  - Cut the long sequence by blocks



#### Hash values

65	for (i=0; i <n; i++)="" th="" {<=""></n;>
16	total[i].adr = list[i].addr;
16	total[i].bytes = list[i].size;
71	total[i].more = &total[i+1];
	}
65	for (i=0; i <n; i++)="" td="" {<=""></n;>
16	taken[i].adr = list[i].addr;
16	taken[i].bytes = list[i].size;
71	taken[i].more = &total[i+1];
	}

Final sequence DB: (65) (16, 16, 71) ... (65) (16, 16, 71)
#### Sequential Pattern Mining & Detecting "Forget-to-Change" Bugs

- Modification to the sequence pattern mining algorithm
  - Constrain the max gap



f (a1)

f (a2);

f (a3);

conflict

f1 (b1);

f1 (b2);

**£**2 (b3);

- Composing Larger Copy-Pasted Segments
  - Combine the neighboring copy-pasted segments repeatedly
- Find conflicts: Identify names that cannot be mapped to the corresponding ones
  - E.g., 1 out of 4 "total" is unchanged, unchanged ratio = 0.25
  - □ If 0 < *unchanged ratio* < *threshold*, then report it as a bug
- CP-Miner reported many C-P bugs in Linux, Apache, ... out of millions of LOC (lines of code)

### **Chapter 7 : Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
- Constraint-Based Frequent Pattern Mining
- Sequential Pattern Mining
- Graph Pattern Mining
- Pattern Mining Application: Mining Software Copy-and-Paste Bugs



## **Summary: Advanced Frequent Pattern Mining**

- Mining Diverse Patterns
  - Mining Multiple-Level Associations
  - Mining Multi-Dimensional Associations
  - Mining Quantitative Associations
  - Mining Negative Correlations
  - Mining Compressed and Redundancy-Aware Patterns
- Sequential Pattern Mining
  - Sequential Pattern and Sequential Pattern Mining
  - GSP: Apriori-Based Sequential Pattern Mining
  - SPADE: Sequential Pattern Mining in Vertical Data Format
  - PrefixSpan: Sequential Pattern Mining by Pattern-Growth
  - CloSpan: Mining Closed Sequential Patterns

- Constraint-Based Frequent Pattern Mining
  - Why Constraint-Based Mining?
  - Constrained Mining with Pattern Anti-Monotonicity
  - Constrained Mining with Pattern Monotonicity
  - Constrained Mining with Data Anti-Monotonicity
  - Constrained Mining with Succinct Constraints
  - Constrained Mining with Convertible Constraints
  - Handling Multiple Constraints
  - Constraint-Based Sequential-Pattern Mining
- Graph Pattern Mining
  - Graph Pattern and Graph Pattern Mining
  - Apriori-Based Graph Pattern Mining Methods
  - gSpan: A Pattern-Growth-Based Method
  - CloseGraph: Mining Closed Graph Patterns
- Pattern Mining Application: Mining Software Copyand-Paste Bugs

## **References: Mining Diverse Patterns**

- R. Srikant and R. Agrawal, "Mining generalized association rules", VLDB'95
- Y. Aumann and Y. Lindell, "A Statistical Theory for Quantitative Association Rules", KDD'99
- K. Wang, Y. He, J. Han, "Pushing Support Constraints Into Association Rules Mining", IEEE Trans. Knowledge and Data Eng. 15(3): 642-658, 2003
- D. Xin, J. Han, X. Yan and H. Cheng, "On Compressing Frequent Patterns", Knowledge and Data Engineering, 60(1): 5-29, 2007

76

- D. Xin, H. Cheng, X. Yan, and J. Han, "Extracting Redundancy-Aware Top-K Patterns", KDD'06
- □ J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent Pattern Mining: Current Status and Future Directions", Data Mining and Knowledge Discovery, 15(1): 55-86, 2007
- F. Zhu, X. Yan, J. Han, P. S. Yu, and H. Cheng, "Mining Colossal Frequent Patterns by Core Pattern Fusion", ICDE'07

#### **References: Constraint-Based Frequent Pattern Mining**

- R. Srikant, Q. Vu, and R. Agrawal, "Mining association rules with item constraints", KDD'97
- R. Ng, L.V.S. Lakshmanan, J. Han & A. Pang, "Exploratory mining and pruning optimizations of constrained association rules", SIGMOD'98
- G. Grahne, L. Lakshmanan, and X. Wang, "Efficient mining of constrained correlated sets", ICDE'00
- J. Pei, J. Han, and L. V. S. Lakshmanan, "Mining Frequent Itemsets with Convertible Constraints", ICDE'01
- J. Pei, J. Han, and W. Wang, "Mining Sequential Patterns with Constraints in Large Databases", CIKM'02
- F. Bonchi, F. Giannotti, A. Mazzanti, and D. Pedreschi, "ExAnte: Anticipated Data Reduction in Constrained Pattern Mining", PKDD'03
- F. Zhu, X. Yan, J. Han, and P. S. Yu, "gPrune: A Constraint Pushing Framework for Graph Pattern Mining", PAKDD'07

## **References: Sequential Pattern Mining**

- R. Srikant and R. Agrawal, "Mining sequential patterns: Generalizations and performance improvements", EDBT'96
- M. Zaki, "SPADE: An Efficient Algorithm for Mining Frequent Sequences", Machine Learning, 2001
- J. Pei, J. Han, B. Mortazavi-Asl, J. Wang, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, "Mining Sequential Patterns by Pattern-Growth: The PrefixSpan Approach", IEEE TKDE, 16(10), 2004
- X. Yan, J. Han, and R. Afshar, "CloSpan: Mining Closed Sequential Patterns in Large Datasets", SDM'03
- □ J. Pei, J. Han, and W. Wang, "Constraint-based sequential pattern mining: the patterngrowth methods", J. Int. Inf. Sys., 28(2), 2007
- M. N. Garofalakis, R. Rastogi, K. Shim: Mining Sequential Patterns with Regular Expression Constraints. IEEE Trans. Knowl. Data Eng. 14(3), 2002
- H. Mannila, H. Toivonen, and A. I. Verkamo, "Discovery of frequent episodes in event sequences", Data Mining and Knowledge Discovery, 1997

# **References: Graph Pattern Mining**

- C. Borgelt and M. R. Berthold, Mining molecular fragments: Finding relevant substructures of molecules, ICDM'02
- □ J. Huan, W. Wang, and J. Prins. Efficient mining of frequent subgraph in the presence of isomorphism, ICDM'03
- A. Inokuchi, T. Washio, and H. Motoda. An apriori-based algorithm for mining frequent substructures from graph data, PKDD'00
- □ M. Kuramochi and G. Karypis. Frequent subgraph discovery, ICDM'01
- S. Nijssen and J. Kok. A Quickstart in Frequent Structure Mining can Make a Difference. KDD'04
- N. Vanetik, E. Gudes, and S. E. Shimony. Computing frequent graph patterns from semistructured data, ICDM'02
- X. Yan and J. Han, gSpan: Graph-Based Substructure Pattern Mining, ICDM'02
- X. Yan and J. Han, CloseGraph: Mining Closed Frequent Graph Patterns, KDD'03
- X. Yan, P. S. Yu, J. Han, Graph Indexing: A Frequent Structure-based Approach, SIGMOD'04
- X. Yan, P. S. Yu, and J. Han, Substructure Similarity Search in Graph Databases, SIGMOD'05

