

## CS522 Advanced Database Systems Mining Frequent Patterns

Chengyu Sun  
California State University, Los Angeles

## Sales Transactions

TID	Transactions
1	Beef, Chicken, Milk
2	Beef, Cheese
3	Cheese, Boots
4	Beef, Chicken, Cheese
5	Beef, Chicken, Clothes, Cheese, Milk
6	Chicken, Clothes, Milk
7	Chicken, Clothes, Milk
8	Beef, Milk

## Support Count

- ◆ The support count, or frequency, of an itemset is the number of the transactions that contain the itemset
  - Item, Itemset, and Transaction
- ◆ Examples:
  - $\text{support\_count}(\{\text{beef}\})=5$
  - $\text{support\_count}(\{\text{beef, chicken, milk}\})=??$

## Frequent Itemset

- ◆ An itemset is frequent if its support count is greater than or equals to a minimum support count threshold
  - $\text{support\_count}(X) \geq \text{min\_sup}$

## The Need for Closed Frequent Itemsets

- ◆ Two transactions
  - $\langle a_1, a_2, \dots, a_{100} \rangle$  and  $\langle a_1, a_2, \dots, a_{50} \rangle$
- ◆  $\text{min\_sup}=1$
- ◆ # of frequent itemsets??

## Closed Frequent Itemset

- ◆ An itemset  $X$  is closed if there exists no *proper superset* of  $X$  that has the same support count
- ◆ A closed frequent itemset is an itemset that is both *closed* and *frequent*

## Closed Frequent Itemset Example

- ◆ Two transactions
  - $\langle a_1, a_2, \dots, a_{100} \rangle$  and  $\langle a_1, a_2, \dots, a_{50} \rangle$
- ◆  $\text{min\_sup}=1$
- ◆ Closed frequent itemset(s)??

## Maximal Frequent Itemset

- ◆ An itemset  $X$  is a maximal frequent itemset if  $X$  is frequent and there exists no *proper superset* of  $X$  that is also frequent
- ◆ Example: if  $\{a, b, c\}$  is a maximal frequent itemset, which one of these *cannot* be a MFI
  - $\{a, b, c, d\}$ ,  $\{a, c\}$ ,  $\{b, d\}$

## Maximal Frequent Itemset Example

- ◆ Two transactions
  - $\langle a_1, a_2, \dots, a_{100} \rangle$  and  $\langle a_1, a_2, \dots, a_{50} \rangle$
- ◆  $\text{min\_sup}=1$
- ◆ Maximal frequent itemset(s)??
- ◆ Maximal frequent itemset vs. closed frequent itemset??

## From Frequent Itemsets to Association Rules

- ◆  $\{\text{chicken}, \text{cheese}\}$  is a frequent set
- ◆  $\{\text{chicken}\} \Rightarrow \{\text{cheese}\}??$
- ◆ Or is it  $\{\text{cheese}\} \Rightarrow \{\text{chicken}\}??$

## Association Rules

- ◆  $A \Rightarrow B$ 
  - $A$  and  $B$  are itemsets
  - $A \cap B = \emptyset$

## Support

- ◆ The support of  $A \Rightarrow B$  is the percentage of the transactions that contain  $A \cup B$

$$\text{support}(A \Rightarrow B) = P(A \cup B) = \frac{\text{support\_count}(A \cup B)}{|D|}$$

$P(A \cup B)$  is the probability that a transaction contains  $A \cup B$   
 $D$  is the set of the transactions

## Confidence

- ◆ The confidence of  $A \Rightarrow B$  is the percentage of the transactions containing **A** that also contains **B**

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support\_count}(A \cup B)}{\text{support\_count}(A)}$$

## Support and Confidence Example

- ◆  $\{ \text{chicken} \} \Rightarrow \{ \text{cheese} \} ??$
- ◆  $\{ \text{cheese} \} \Rightarrow \{ \text{chicken} \} ??$

## Strong Association Rule

- ◆ An association rule is strong if it satisfies both a minimum support threshold ( $\text{min\_sup}$ ) and a minimum confidence threshold ( $\text{min\_conf}$ )
- ◆ Why do we need both *support* and *confidence*??

## Association Rule Mining

- ◆ Find strong association rules
  - Find all frequent itemsets
  - Generate strong association rules from the frequent itemsets

## The Apriori Property

- ◆ All nonempty subsets of a frequent itemset must also be frequent
- ◆ Or, if an itemset is not frequent, its supersets cannot be frequent either

## Finding Frequent Itemsets – The Apriori Algorithm

- ◆ Given  $\text{min\_sup}$
- ◆ Find the frequent 1-itemsets  $L_1$
- ◆ Find the frequent k-itemsets  $L_k$  by joining the itemsets in  $L_{k-1}$
- ◆ Stop when  $L_k$  is empty

## Apriori Algorithm Example

beef	1
chicken	2
milk	3
cheese	4
boots	5
clothes	6

◆ Support 25%

TID	Transactions
1	1, 2, 3
2	1, 4
3	4, 5
4	1, 2, 4
5	1, 2, 6, 4, 3
6	2, 6, 3
7	2, 6, 3
8	1, 3

## L<sub>1</sub>

- ◆ Scan the data once to get the count of each item
- ◆ Remove the items that do not meet min\_sup

C <sub>1</sub>	support_count	L <sub>1</sub>
{1}	5	{1}
{2}	5	{2}
{3}	5	{3}
{4}	4	{4}
{5}	1	
{6}	3	{6}

## L<sub>2</sub>

- ◆ C<sub>2</sub> = L<sub>1</sub> × L<sub>1</sub>
- ◆ Scan the dataset again for the support\_count of C<sub>2</sub>, then remove non-frequent itemsets from C<sub>2</sub>, i.e. C<sub>2</sub> → L<sub>2</sub>

C <sub>2</sub>	support_count	L <sub>2</sub>
{1,2}	3	{1,2}
{1,3}	3	{1,3}
{1,4}	3	{1,4}
{1,6}	1	
{2,3}	4	{2,3}
{2,4}	2	{2,4}
{2,6}	3	{2,6}
{3,4}	1	
{3,6}	3	{3,6}
{4,6}	1	

## L<sub>3</sub>

- ◆ ??

## From L<sub>k-1</sub> to C<sub>k</sub>

- ◆ Let l<sub>i</sub> be an itemset in L<sub>k-1</sub>, and l<sub>i</sub>[j] be the jth item in l<sub>i</sub>
- ◆ Items in an itemset are sorted, i.e. l<sub>i</sub>[1] < l<sub>i</sub>[2] < ... < l<sub>i</sub>[k-1]
- ◆ l<sub>1</sub> and l<sub>2</sub> are joinable if
  - Their first k-2 items are the same, and
  - l<sub>1</sub>[k-1] < l<sub>2</sub>[k-1]

## From C<sub>k</sub> to L<sub>k</sub>

- ◆ Reduce the size of C<sub>k</sub> using the Apriori property
  - any (k-1)-subset of a candidate must be frequent, i.e. in L<sub>k-1</sub>
- ◆ Scan the dataset to get the support counts

## Generate Association Rules from Frequent Itemsets

- ◆ For each frequent itemset  $l$ , generate all nonempty subset of  $l$
- ◆ For every nonempty subset  $s$  of  $l$ , output rule  $s \Rightarrow (l-s)$  if  $\text{conf}(s \Rightarrow (l-s)) \geq \text{min\_conf}$

## Confidence-based Pruning ...

- ◆  $\text{conf}(\{a, b\} \Rightarrow \{c, d\}) < \text{min\_conf}$ 
  - $\text{conf}(\{a\} \Rightarrow \{c, d\}) ??$
  - $\text{conf}(\{a, b, e\} \Rightarrow \{c, d\}) ??$
  - $\text{conf}(\{a\} \Rightarrow \{b, c, d\}) ??$

## ... Confidence-based Pruning

- ◆ If  $\text{conf}(s \Rightarrow (l-s)) < \text{min\_conf}$ , then  $\text{conf}(s' \Rightarrow (l-s')) < \text{min\_conf}$  where  $s' \subseteq s$ .
- ◆ Example:
  - $\text{conf}(\{a, b\} \Rightarrow \{c, d\}) < \text{min\_conf}$
  - ??

## Limitations of the Apriori Algorithm

- ◆ Multiple scans of the datasets
  - How many??
- ◆ Need to generate a large number of candidate sets

## FP-Growth Algorithm

- ◆ Frequent-pattern Growth
- ◆ Mine frequent itemsets *without candidate generation*

## FP-Growth Example

TID	Transactions	
1	I1, I2, I5	
2	I2, I4	
3	I2, I3, I6	
4	I1, I2, I4	
5	I1, I3	
6	I2, I3	
7	I1, I3	
8	I1, I2, I3, I5	
9	I1, I2, I3	

min\_sup=2

## L

- ◆ Scan the dataset and find the frequent 1-itemsets
- ◆ Sort the 1-itemsets by support count in descending order

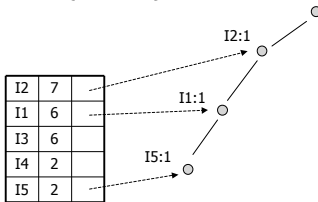
L
I2: 7
I1: 6
I3: 6
I4: 2
I5: 2

## FP-tree

- ◆ Each transaction is processed in  $\mathbb{L}$  order (why??) and becomes a branch in the FP tree
- ◆ Each node is linked from  $\mathbb{L}$

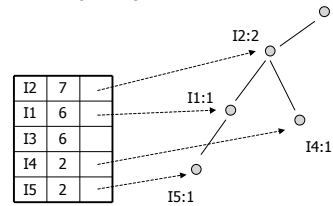
## FP-tree Construction ...

- ◆ T1: {I2,I1,I5}



## ... FP-tree Construction ...

- ◆ T2: {I2,I4}



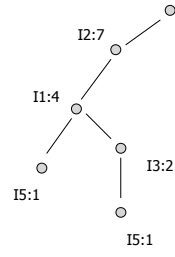
## ... FP-tree Construction

- ◆??

## Mining the FP-tree

- ◆ For each item  $i$  in  $\mathbb{L}$  (in ascending order), find the branch(s) in the FP tree that ends in  $i$  – Prefix Paths
  - If there is only one path, generate all the frequent patterns ended in  $i$
  - Else create the Conditional FP-tree for  $i$  and recursively run the mining algorithm on the conditional FP-tree

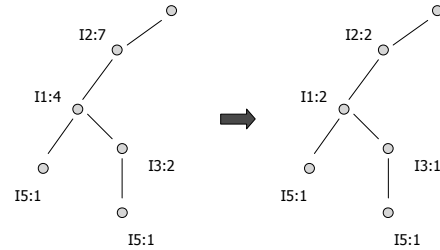
## Prefix Paths of I5



I5 is frequent  $\rightarrow$  {I5:2}

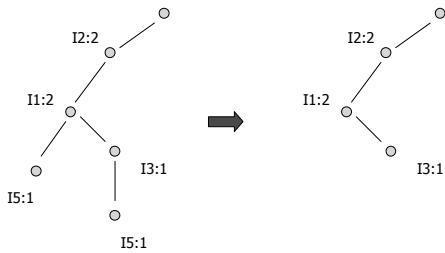
## From Prefix Paths to Conditional FP-tree

◆ Adjust the support counts



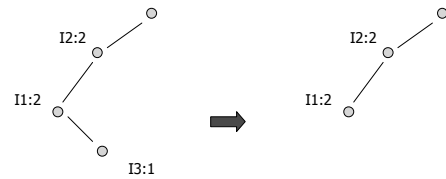
## From Prefix Paths to Conditional FP-tree

◆ Remove the suffix

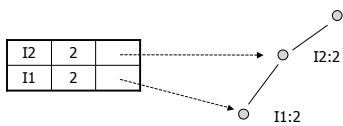


## From Prefix Paths to Conditional FP-tree

◆ Remove the infrequent items

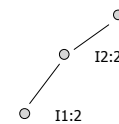


## Conditional FP-tree for I5



◆ A FP-tree with *suffix pattern* {I5}

## Prefix Paths of {I1,I5}



{I1,I5} is frequent and there is a single path  
 $\rightarrow$  {I1,I5:2}, {I2,I1,I5:2}

### Prefix Paths of {I2,I5}

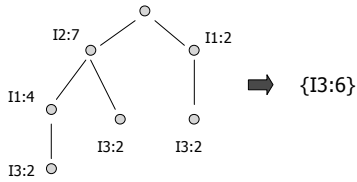


{I2,I5} is frequent and there is a single path  
 → {I2,I5:2}

### All Frequent Itemsets with Suffix I5

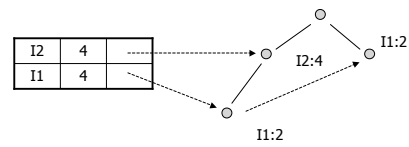
- ◆ {I5:2}
- ◆ {I1,I5:2}, {I2,I1,I5:2}
- ◆ {I2,I5:2}

### Mining The FP-tree – I3 ...



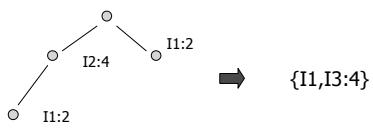
Prefix Paths of I3

### ... Mining The FP-tree – I3 ...



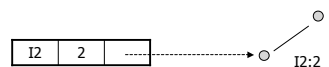
Conditional FP-tree with suffix pattern {I3}

### ... Mining The FP-tree – I3 ...



Prefix Paths of {I1,I3}

### ... Mining The FP-tree – I3 ...



Conditional FP-tree for {I2,I1,I3}



## ... Mining The FP-tree – I3 ...



Prefix Paths of  $\{I2,I1,I3\}$

## ... Mining The FP-tree – I3



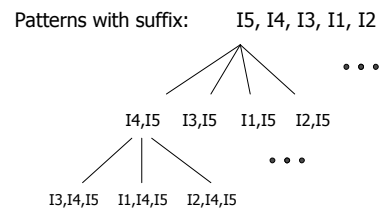
Prefix Paths of  $\{I2,I3\}$

## All Frequent Itemsets with Suffix I3

- ◆  $\{I3:6\}$
- ◆  $\{I1,I3:4\}$
- ◆  $\{I2,I1,I3:2\}$
- ◆  $\{I2,I3:4\}$

## About FP-tree Mining

- ◆ A divide-and-conquer approach



## Data Partitioning

- ◆ Divide dataset into  $n$  non-overlapping partitions such that *each partition fits into main memory*
- ◆ Find local frequent itemsets in each partition (1 scan)
  - Local  $min\_sup$ ??
- ◆ All local frequent itemsets form a candidate set
  - Will it include all the global frequent itemsets??
- ◆ Find global frequent itemsets from candidates (1 scan)

## Vertical Data Format

Item	TID_set
I1	T1,T4,T5,T7,T8,T9
I2	T1,T2,T3,T4,T6,T8,T9
I3	T3,T5,T6,T7,T8,T9
I4	T2,T4
I5	T1,T8

- ◆ And how does it help??

## Strong Association Rules Could Be Misleading ...

### ◆ Example:

- 10,000 transactions
- 6,000 transactions included games
- 7,500 transactions included videos
- 4,000 transactions included both

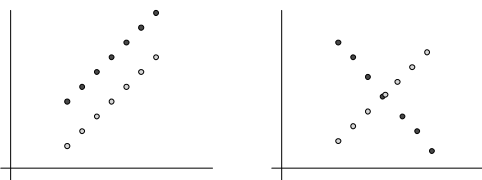
### ◆ {game} ⇒ {video}

- Support?? Confidence??

## ... Strong Association Rules Could Be Misleading

- ◆ Does buying game really imply buying video as well??

## Correlation



## Multiplication Rule

- ◆ If two events A and B are independent of each other

$$P(AB) = P(A)P(B)$$

## From Multiplication Rule to Lift

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

- ◆ lift({game}, {video}) = ??

## Problem of Lift

datasets	mc	m'c	mc'	m'c'	total	lift
A <sub>1</sub>	100	100	100	100	400	??
A <sub>2</sub>	100	100	100	1,000	1,300	??
A <sub>3</sub>	100	100	100	10,000	10,300	??
A <sub>4</sub>	100	100	100	100,000	100,300	??

mc: # of transactions that contain both milk and coffee  
 m'c: # of transactions that contain milk but not coffee  
 m'c': # of transactions that contain coffee but not milk  
 m'c'': # of transactions that contain neither milk nor coffee

## Null-invariant

- ◆ A *null-transaction* is a transaction that does not contain any of the itemsets being examined
- ◆ A correlation measure is null-invariant if its value is not affected by the number of null-transactions.

## Some Null-invariant Measures

Cosine Measure	$\frac{P(AB)}{\sqrt{P(A) \times P(B)}}$
All_confidence	$\min\{P(A B), P(B A)\}$
Max_confidence	$\max\{P(A B), P(B A)\}$
Kulczynski Measure	$\frac{1}{2}(P(A B) + P(B A))$

## Cosine vs. Lift

$$lift(A, B) = \frac{P(AB)}{P(A)P(B)} = \frac{\frac{sup(A \cup B)}{N}}{\frac{sup(A)}{N} \frac{sup(B)}{N}} = \frac{N \cdot sup(A \cup B)}{sup(A) \cdot sup(B)}$$

$$cosine(A, B) = \frac{P(AB)}{\sqrt{P(A) \times P(B)}} = \frac{\frac{sup(A \cup B)}{N}}{\sqrt{\frac{sup(A)}{N} \frac{sup(B)}{N}}} = \frac{sup(A \cup B)}{\sqrt{sup(A) \cdot sup(B)}}$$

cosine( {milk}, {coffee} ) = ??

## Ranges of Correlation Measures

	Perfectly positively correlated	Perfectly negatively correlated
Cosine		
All_conf		
Max_conf		
Kulc		

## Choosing Correlation Measures ...

datasets	mc	m'c	mc'	m'c'	cosine	all conf	max conf	Kulc
D <sub>1</sub>	10,000	1000	1000	100,000	0.91	0.91	0.91	0.91
D <sub>3</sub>	100	1000	1000	100,000	0.09	0.09	0.09	0.91
D <sub>4</sub>	1,000	1000	1000	100,000	0.5	0.5	0.5	0.5
D <sub>5</sub>	1,000	100	10,000	100,000	0.29	0.09	0.91	0.5
D <sub>6</sub>	1,000	10	100,000	100,000	0.10	0.01	0.99	0.5

## ... Choosing Correlation Measures

- ◆ Recommended:
  - Kulczynski + Imbalance Ratio (IR)

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

## Mining Sequential Patterns

- ◆  $\langle \{\text{computer}\}, \{\text{printer}\}, \{\text{printer cartridge}\} \rangle$
- ◆  $\langle \{\text{bread, milk}\}, \{\text{bread, milk}\}, \{\text{bread, milk}\} \dots \rangle$
- ◆  $\langle \{\text{home.jsp}\}, \{\text{search.jsp}\}, \{\text{product.jsp}\}, \{\text{product.jsp}\}, \{\text{search.jsp}\} \dots \rangle$

## Terminology and Notations

- ◆ Item, itemset
- ◆ Event = itemset
- ◆ A sequence is an ordered list of events
  - $\langle e_1 e_2 e_3 \dots e_i \rangle$
  - E.g.  $\langle (a)(abc)(bc)(d)(ac)(f) \rangle$
- ◆ The length of a sequence is the number of items in the sequence, i.e. *not the number of events*

## Sequences vs. Itemsets

- ◆  $\{a, b, c\}$ 
  - # of 3-itemset(s)??
  - # of 3-sequence(s)??

## Subsequence

- ◆  $A = \langle a_1 a_2 a_3 \dots a_n \rangle$
- ◆  $B = \langle b_1 b_2 b_3 \dots b_m \rangle$
- ◆ A is a *subsequence* of B if there exists  $1 \leq j_1 < j_2 < \dots < j_n \leq m$  such that  $a_1 \subseteq b_{j_1}, a_2 \subseteq b_{j_2}, \dots, a_n \subseteq b_{j_n}$

## Subsequence Example

- ◆  $s = \langle (abc) (de) (f) \rangle$
- ◆ Which of these are subsequences of  $s$ ??
  - $s_1 = \langle (ab) (d) \rangle$
  - $s_2 = \langle (ab) (f) \rangle$
  - $s_3 = \langle (ac) (f) \rangle$
  - $s_4 = \langle (abcde) \rangle$
  - $s_5 = \langle (a) (de) \rangle$
  - $s_6 = \langle (de) (a) (f) \rangle$

## Sequential Pattern

- ◆ If A is a subsequence of B, we say B *contains* A
- ◆ The support count of A is the number of sequences that contain A
- ◆ A is *frequent* if  $\text{support\_count}(A) \geq \text{min\_sup}$
- ◆ A frequent sequence is called a sequential pattern

## Apriori Property Again

- ◆ Every nonempty subsequence of a frequent sequence is frequent

## GSP Algorithm

- ◆ *Generalized Sequential Patterns*
- ◆ An extension of the Apriori algorithm for mining sequential patterns

## GSP Example

SID	Sequence	
1	<(a)(ab)(a)>	min_sup=2
2	<(a)(c)(bc)>	
3	<(ab)(c)(b)>	
4	<(a)(c)(c)>	

## L<sub>1</sub>

C <sub>1</sub>	support_count	L <sub>1</sub>
<(a)>	4	<(a)>
<(b)>	3	<(b)>
<(c)>	3	<(c)>

## L<sub>2</sub>

C <sub>2</sub>	support_count	L <sub>2</sub>
<(a)(a)>	1	
<(a)(b)>	3	<(a)(b)>
<(a)(c)>	3	<(a)(c)>
<(b)(a)>	1	
<(b)(b)>	1	
<(b)(c)>	1	
<(c)(a)>	0	
<(c)(b)>	2	<(c)(b)>
<(c)(c)>	2	<(c)(c)>
<(ab)>	2	<(ab)>
<(ac)>	0	
<(bc)>	1	

## From L<sub>k-1</sub> to C<sub>k</sub>

- ◆ Two sequences  $s_1$  and  $s_2$  are joinable if the subsequence obtained by dropping the first item in  $s_1$  is the same as the subsequence obtained by dropping the last item in  $s_2$
- ◆ The joined sequence is  $s_1$  concatenated with the last item  $i$  of  $s_2$ 
  - If the last two items in  $s_2$  are in the same event,  $i$  is merged into the last event of  $s_1$ ;
  - Otherwise  $i$  becomes a separate event

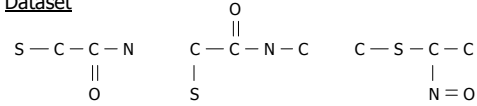


## Subgraph

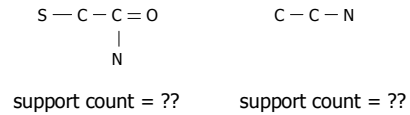
- ◆ A graph  $G' = (V', E')$  is a subgraph of another graph  $G = (V, E)$  if its vertex set  $V'$  is a subset of  $V$  and its edge set  $E'$  is a subset of  $E$

## Support

### Dataset



### Subgraph Patterns

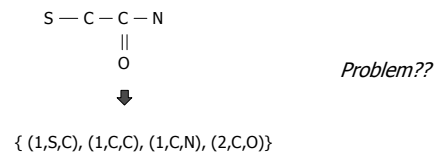


## Frequent Subgraph Mining

- ◆ Given a set of graphs and  $\text{min\_sup}$ , find all subgraphs  $g$  with  $\text{support}(g) \geq \text{min\_sup}$ 
  - Typically we only consider graphs that *undirected and connected*

## Transform Graph to Itemset

- ◆ Each combination of an edge label with its corresponding vertex labels is mapped to an "item"



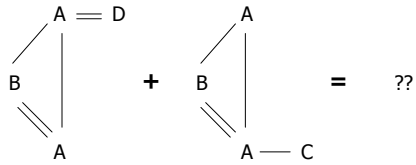
## Apriori-based Approach

- ◆ Candidate generation
- ◆ Candidate pruning
- ◆ Support counting
- ◆ Candidate elimination

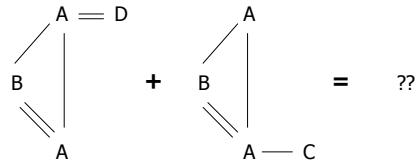
## Candidate Generation

- ◆ Merge two frequent  $(k-1)$ -subgraphs to form a candidate  $k$ -subgraph
  - What is  $k$ ??
- ◆ The two  $(k-1)$ -subgraphs must share a common  $(k-2)$ -subgraph, referred to as their *core*

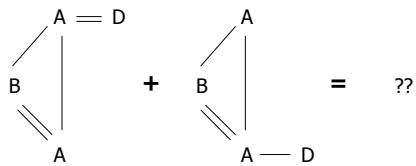
## Vertex Growing



## Edge Growing ...



## ... Edge Growing

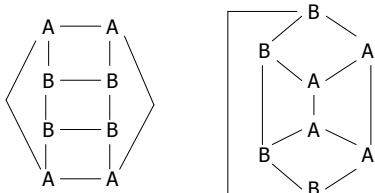


## Candidate Pruning

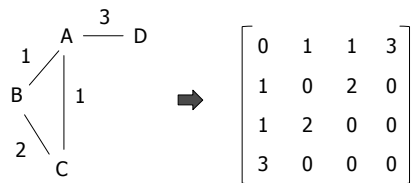
- ◆ Remove an edge from a candidate  $k$ -subgraph and check if the resulting  $(k-1)$ -subgraph is connected and frequent

## Graph Isomorphism Problem

- ◆ Determine whether two graphs are topologically equivalent, i.e. *isomorphic*

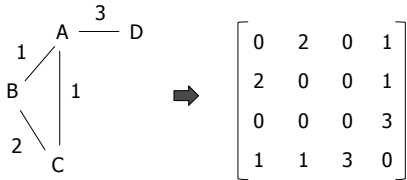


## Adjacency Matrix ...





## ... Adjacency Matrix



How many adjacency matrices can a graph with k vertices have??

## String Representation of an Adjacency Matrix

$$\begin{bmatrix} 0 & 1 & 1 & 3 \\ 1 & 0 & 2 & 0 \\ 1 & 2 & 0 & 0 \\ 3 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 2 & 0 & 1 \\ 2 & 0 & 0 & 1 \\ 0 & 0 & 0 & 3 \\ 1 & 1 & 3 & 0 \end{bmatrix}$$



112300



200113

## Graph Code

- ◆ A.K.A. Canonical label
- ◆ The string representation of the adjacency matrix that has the lowest (or highest) lexicographic value

## Support Counting

- ◆ Isomorphism test a candidate k-subgraph against the k-subgraphs of each graph

## Summary

- ◆ Frequent itemsets, association rules, sequential patterns, subgraph patterns
  - Measures: support, confidence, correlation
  - Algorithms: Apriori, FP-Growth, association rule generation, GPS
  - Optimizations: partitioning, vertical data format, various pruning techniques

## Readings

- ◆ Textbook Chapter 6