

POLITECNICO DI MILANO

Association Rules: Data Mining and Text Mining

Prof. Pierluca Lanzi

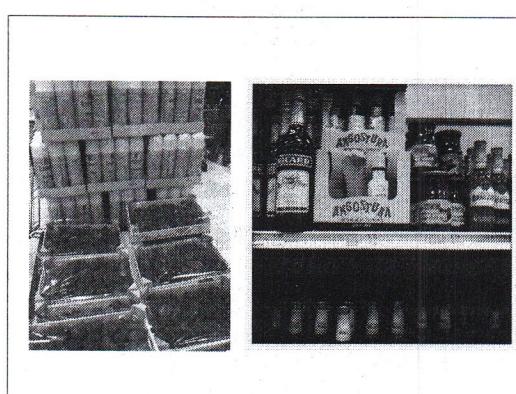
Study Material

- "Data Mining and Analysis" by Zaki & Meira
 - Chapter 8
 - Section 9.1
 - Chapter 10 up to Section 10.2.1 included
 - Section 12.1
- <http://www.dataminingbook.info>

Data Mining
and Analysis
Techniques for Advanced Data Mining Applications

Prof. Pierluca Lanzi

POLITECNICO DI MILANO



Market-Basket Transactions

Bread	Jam	Steak	Jam
Peanuts	Soda	Jam	Soda
Milk	Chips	Soda	Peanuts
Fruit	Milk	Chips	Milk
Jam	Fruit	Bread	Fruit
Jam	Fruit		
Soda	Soda		
Chips	Chips		
Milk	Milk		
Bread			

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

What is Association Rule Mining?

- Finding frequent patterns and associations among sets of items in transaction databases, relational databases, or other information repositories
- Applications
 - Basket data analysis
 - Cross-marketing
 - Catalog design
 - ...

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

What is Association Rule Mining?

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in a transaction
- Consider the following data set with 8 transactions
- Examples of itemsets
 - {Bread, Milk}
 - {Soda, Chips}
- Examples of association rules
 - {Bread} \Rightarrow {Milk}
 - {Soda} \Rightarrow {Chips}
 - {Bread} \Rightarrow {Jam}

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

We see that bread and milk are often together. Maybe there is an association rule that says "if you have one you also have (probably) the other"

most frequent items
 "Itemsets" are the fundamental pattern in association rule mining

What is An Association Rule?

- Implication of the form $X \Rightarrow Y$, where X and Y are itemsets
- Example, {Bread} \Rightarrow {Milk}
- Two rule evaluation metrics
 - Support
 - Confidence

evaluation metrics for association rules
 $(Bread) \Rightarrow \{Milk\}$

Support
 Fraction of transactions that contain both X and Y
 $s = \frac{\sigma(\{Bread, Milk\})}{\# \text{ of transactions}} = 0.38$

Confidence
 Measures how often items in Y appear in transactions that contain X
 $c = \frac{\sigma(\{Bread, Milk\})}{\sigma(\{Bread\})} = 0.75$

= how often bread and milk happen together
 $"P(Y, X)"$

= how often milk appears in transaction containing bread
 (probability of Y given that we already have X)
 $"P(Y|X)"$

What is the Goal?

- Given a set of transactions T, the goal of association rule mining is to find all rules having
 - support \geq minsup threshold
 - confidence \geq minconf threshold
- Brute-force approach
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the minsup and minconf thresholds
- Brute-force approach is computationally prohibitive!

- e.g. $p\% = 99\%$.
- the support must guarantee the items that go together at least 99% of times
 - the confidence must guarantee the item Y that is taken at least 99% of the times if it's already taken the item X
- \rightarrow the support must guarantee for example: "find me items that go together at least p% of times"
- \rightarrow same thing for the confidence

In order to obtain association rules we must first focus on itemsets (frequent itemsets)

Mining Association Rules

the support is the same
the confidence is different

- $\{\text{Bread, Jam}\} \Rightarrow \{\text{Milk}\}$ s=0.4 c=0.75
- $\{\text{Milk, Jam}\} \Rightarrow \{\text{Bread}\}$ s=0.4 c=0.75
- $\{\text{Bread}\} \Rightarrow \{\text{Milk, Jam}\}$ s=0.4 c=0.75
- $\{\text{Jam}\} \Rightarrow \{\text{Bread, Milk}\}$ s=0.4 c=0.6
- $\{\text{Milk}\} \Rightarrow \{\text{Bread, Jam}\}$ s=0.4 c=0.5
- They are all binary partitions of the itemset {Milk, Bread, Jam}
- Thus, rules originating from the same itemset have,
 - Same support but
 - Can have different confidence

Prof. Pierluca Lanzi POLITECNICO DI MILANO

We can decouple the support and confidence requirements!

First, find the itemsets with support $\geq \text{minsup}$

Next, generate the rules with confidence $\geq \text{minconf}$

What are Frequent Itemsets?

- Frequent Itemset
 - An itemset whose support is greater than or equal to the minsup threshold
- Support
 - Fraction of transactions that contain an itemset
 - $s(\{\text{Milk, Bread}\}) = 3/8$
 $s(\{\text{Soda, Chips}\}) = 4/8$
- Support count (σ)
 - Raw count of occurrence of an itemset
 - For example, $\sigma(\{\text{Milk, Bread}\}) = 3$ $\sigma(\{\text{Soda, Chips}\}) = 4$

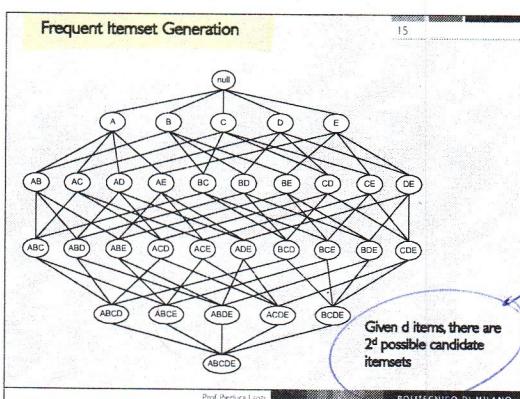
TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Mining Association Rules: a Two Step Approach

- Frequent Itemset Generation
 - Generate all itemsets whose support $\geq \text{minsup}$
- Rule Generation
 - Generate high confidence rules from frequent itemset
 - Each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is computationally expensive

generation of all possible itemsets starting from: A,B,C,D,E



Brute Force

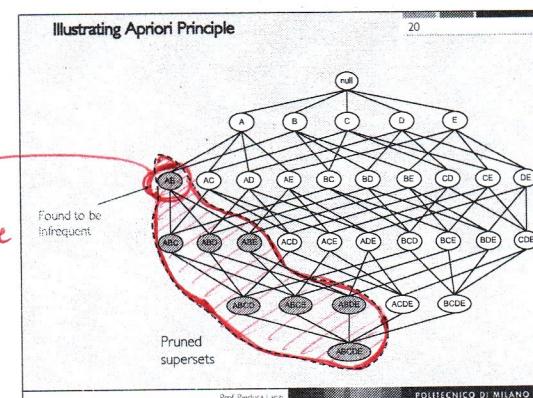
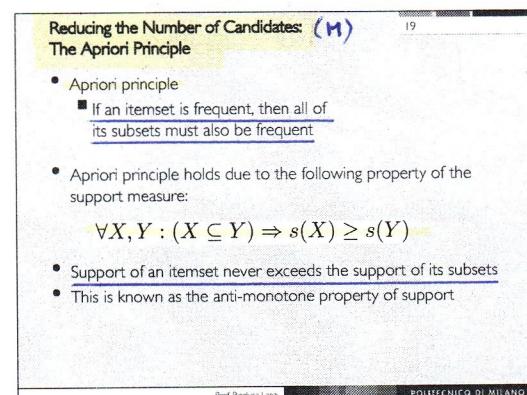
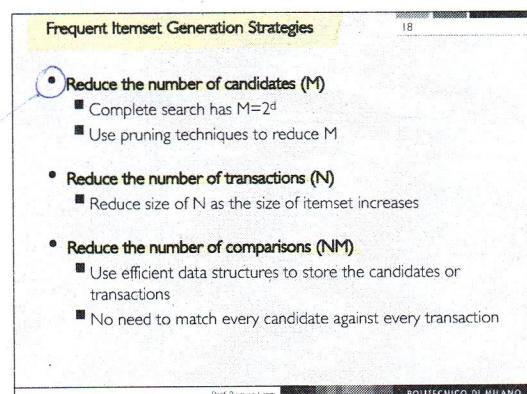
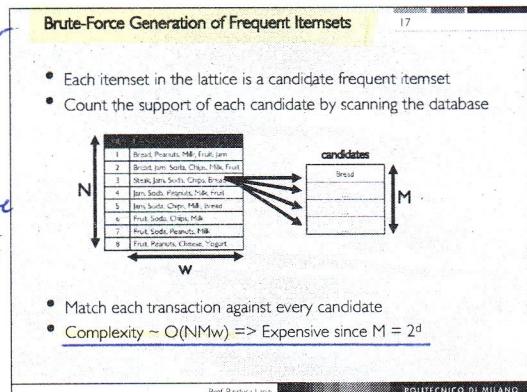
```

1  $\mathcal{F} \leftarrow \emptyset$  // set of frequent itemsets
2 foreach  $X \subseteq I$  do
3    $sup(X) \leftarrow \text{COMPUTESUPPORT}(X, D)$ 
4   if  $sup(X) \geq minsup$  then
5      $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X, sup(X))\}$ 
6 return  $\mathcal{F}$ 

COMPUTESUPPORT ( $X, D$ ):
1  $sup(X) \leftarrow 0$ 
2 foreach  $(t, l(t)) \in D$  do
3   if  $X \subseteq l(t)$  then
4      $sup(X) \leftarrow sup(X) + 1$ 
5 return  $sup(X)$ 

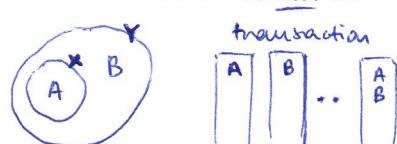
```

Prof. Pierluca Lanzi POLITECNICO DI MILANO



As soon as we find out that AB is not frequent we can prune and eliminate everything that was AB

This means that if, for example, we have apples (A) and bananas (B) then the number of times that we find A is for sure greater or equal to the number of times for which we find A and B.



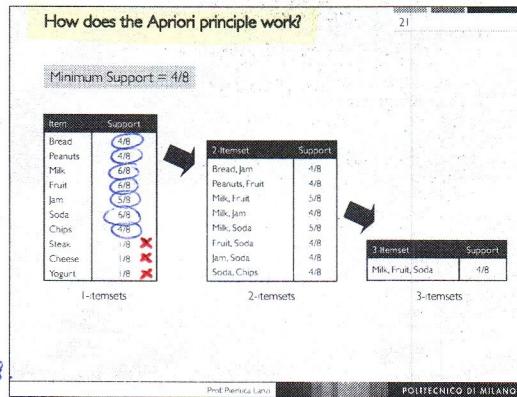
$(\# A) \geq (\#\{A, B\})$
 $\# =$ number of times we find ... in transactions

In this way, if already A is not frequent we know that A and B cannot be frequent!

This helps in the way that: if we find an itemset which is not frequent then EVERY itemset which contains the non-frequent itemset will be non frequent!

good for pruning candidates

1. We establish a minimum support, in this case: 4/8
2. We start with itemsets of 1 item and we evaluate their support
3. We consider for the second step only the itemsets of 2 that has support $\geq 4/8$. We create itemsets of 2 items and we evaluate their support
4. We consider from 3. only the itemsets with support $\geq 4/8$. Then we proceed to form itemsets of 3 elements (starting from the elements of 3. (the "survived") [and so on]



Example from the Textbook

Given the following database and a min support of 3, generate all the frequent itemsets

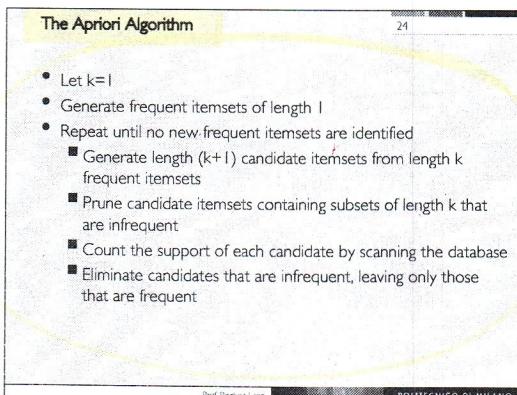
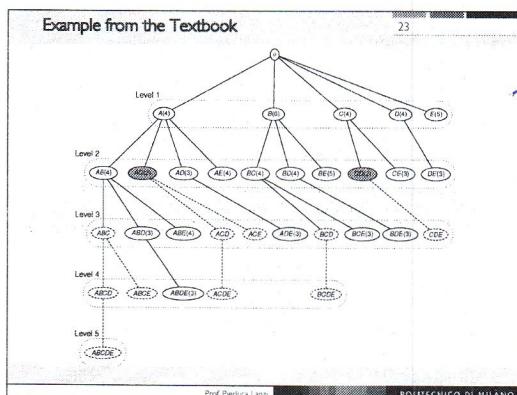
D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

Binary Database

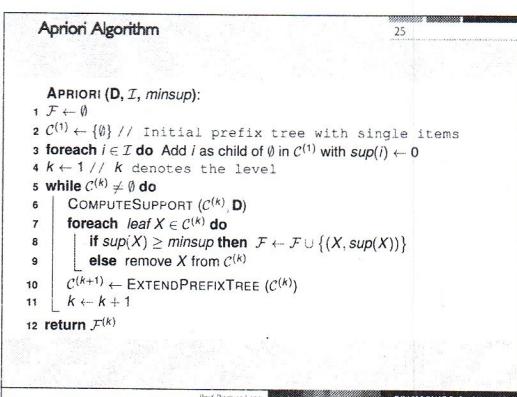
t	I(t)
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

Transaction Database

Prof. Pierluca Lanzi POLITECNICO DI MILANO



0. set the minimum support
1. $k=1$
2. Generate all k -itemsets (if $k>1$ generate all k -itemsets containing only the preceding survived $(k-1)$ -itemsets)
3. Compute the support for all the new itemsets
4. Prune the itemsets with support < minimum support
5. If there are no new frequent itemset then end, otherwise $k++$ and go to 2.



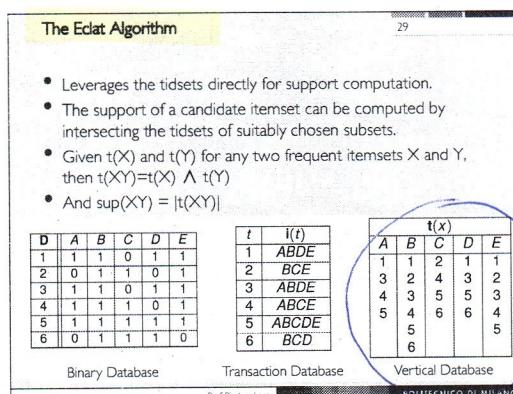
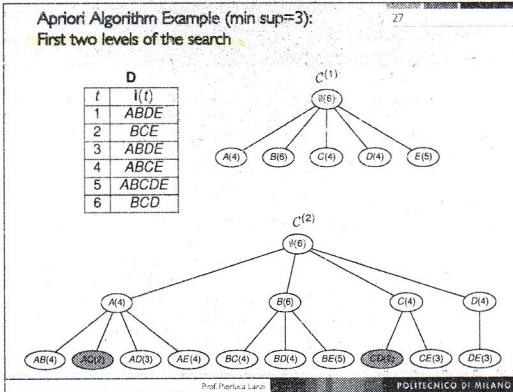
Apriori Algorithm

```

    COMPUTESUPPORT ( $C^{(k)}$ , D):
    1 foreach  $\langle t, I(t) \rangle \in D$  do
    2   foreach k-subset  $X \subseteq I(t)$  do
    3     if  $X \in C^{(k)}$  then  $sup(X) \leftarrow sup(X) + 1$ 

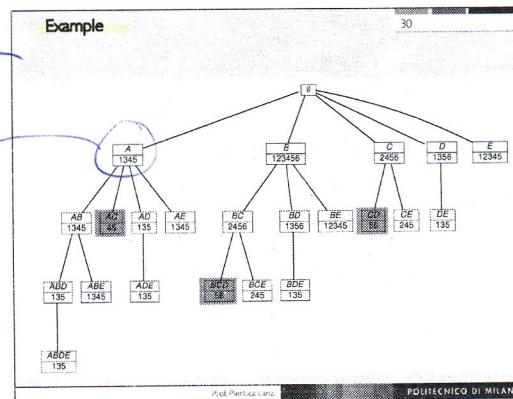
    EXTENDPREFIXTREE ( $C^{(k)}$ ):
    1 foreach leaf  $X_a \in C^{(k)}$  do
    2   foreach leaf  $X_b \in SIBLING(X_a)$ , such that  $b > a$  do
    3      $X_{ab} \leftarrow X_a \cup X_b$ 
     // prune candidate if there are any infrequent subsets
    4     if  $X_b \in C^{(k)}$ , for all  $X_j \subset X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
    5       Add  $X_{ab}$  as child of  $X_a$  with  $sup(X_{ab}) \leftarrow 0$ 
    6     if no extensions from  $X_a$  then
    7       remove  $X_a$ , and all ancestors of  $X_a$  with no extensions, from  $C^{(k)}$ 
    8 return  $C^{(k)}$ 
  
```

Prof. Perluca Lanza POLITECNICO DI MILANO



VERTICAL DATABASE
for every item we list the number of the transactions that contain it

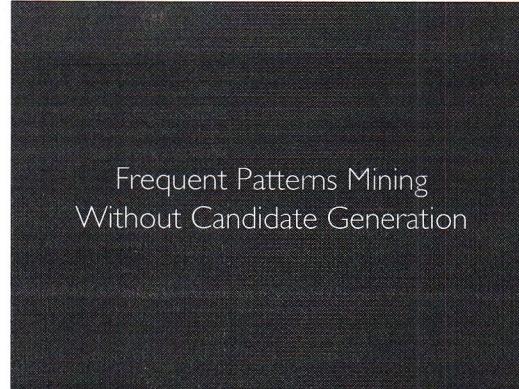
Now: the support of an item, for example A, is just the number of transactions that the item has, in this case, for A, 4.
if we want to find the support of AB we just do the intersection of the two lists. (with binary representation it's really fast)



```

// Initial Call:  $\mathcal{F} \leftarrow \emptyset, P \leftarrow \{(i, t(i)) \mid i \in \mathcal{I}, |t(i)| \geq \text{minsup}\}$ 
ECLAT( $P, \text{minsup}, \mathcal{F}$ ):
1 foreach  $(X_a, t(X_a)) \in P$  do
2    $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X_a, \text{sup}(X_a))\}$ 
3    $P_a \leftarrow \emptyset$ 
4   foreach  $(X_b, t(X_b)) \in P$ , with  $X_b > X_a$  do
5      $X_{ab} = X_a \cup X_b$ 
6      $t(X_{ab}) = t(X_a) \cap t(X_b)$ 
7     if  $\text{sup}(X_{ab}) \geq \text{minsup}$  then
8        $P_a \leftarrow P_a \cup \{(X_{ab}, t(X_{ab}))\}$ 
9   if  $P_a \neq \emptyset$  then ECLAT( $P_a, \text{minsup}, \mathcal{F}$ )

```



Apriori's Performance Bottlenecks

- The core of the Apriori algorithm
 - Use frequent $(k-1)$ -itemsets to generate candidate frequent k -itemsets
 - Use database scan and pattern matching to collect counts for the candidate itemsets
- It may need to generate huge candidate sets
 - 10^4 frequent 1-itemset will generate 10^7 candidate 2-itemsets
 - To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, needs to generate $2^{100} \sim 10^{30}$ candidates.
 - Multiple scans of database, it needs $(n+1)$ scans, n is the length of the longest pattern

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Mining Frequent Patterns Without Candidate Generation

- Compress a large database into a compact, Frequent-Pattern tree (FP-tree) structure
 - Highly condensed, but complete for frequent pattern mining
 - Avoid costly database scans
- Use an efficient, FP-tree-based frequent pattern mining method
- A divide-and-conquer methodology: decompose mining tasks into smaller ones
- Avoid candidate generation: sub-database test only

Prof. Pierluca Lanzi POLITECNICO DI MILANO

The FP-Growth Algorithm

- Leave the generate-and-test paradigm of Apriori
- Data sets are encoded using a compact structure, the FP-tree
- Frequent itemsets are extracted directly from the FP-tree
- Major Steps to mine FP-tree
 - Construct the frequent pattern tree
 - For each frequent item i compute the projected FP-tree
 - Recursively mine conditional FP-trees and grow frequent patterns obtained so far
 - If the conditional FP-tree contains a single path, simply enumerate all the patterns

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Part I

(p. 231)

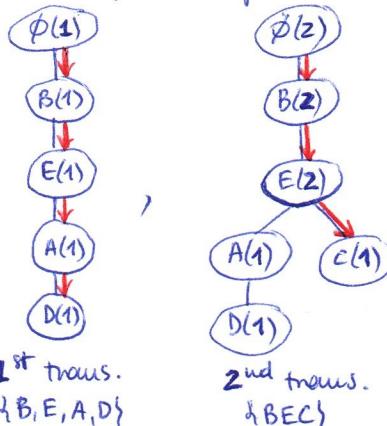
We take the table and we obtain an order based on how many times it appears (+ alphabetically):

{B(6), E(5), A(4), C(4), D(4)}

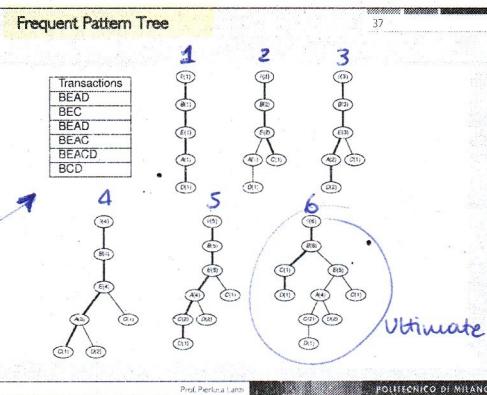
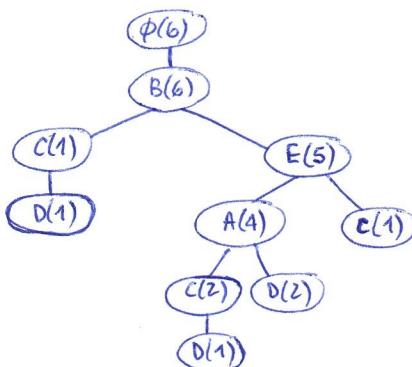
→ New order: BEACD

→ We obtain a new table with this order (new table of transactions)

Then we start building the tree with the 1st transaction. The leaves will count how many times we pass through them:



Ultimate:



Part II

From here we focus on the **CONDITIONAL FP-tree** (for D). Which are the trees (subtrees) from the ultimate ending with D?

BCD

BEACD

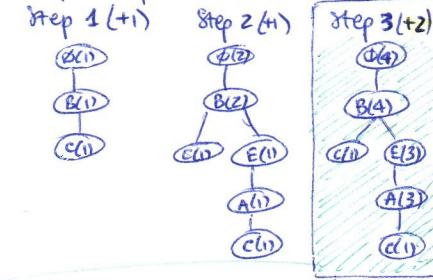
BEAD

$\text{cnt}(D) = 1$

$\text{cnt}(D) = 1$

$\text{cnt}(D) = 2$

And we now build the projected frequent pattern tree for D:



This is the conditional FP-tree for D. This can generate all the possible combinations of the items that generate frequent itemsets for D.

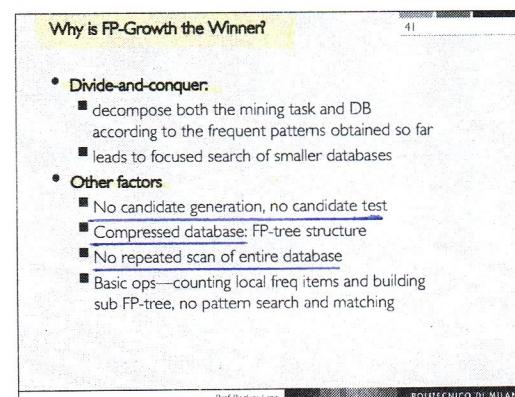
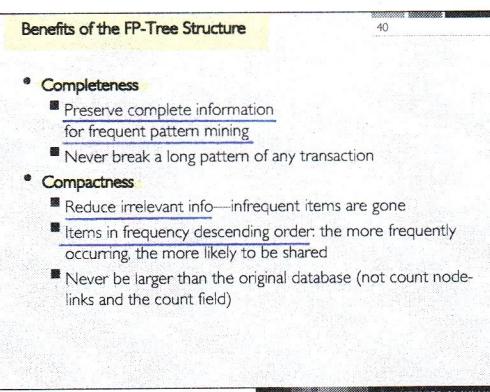
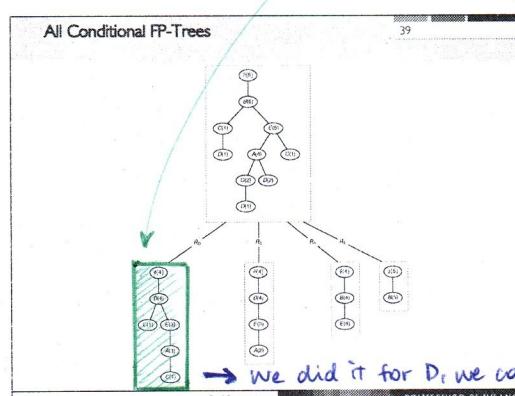
! So, given D we can have:
 $\text{BEA, EA, BE, E, A, B}$

Because the most frequent (single) items given D are A, B, E

→ These are all the patterns that are prefixes of frequent itemsets that end with D.

→ Notice that we didn't consider C because the minimum support (we fixed as) 3.

→ we did it for D, we can proceed with all the others.



Rule Generation

Rule Generation

43

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

{ we have all the possible combinations (which are a lot)

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

How to efficiently generate rules from frequent itemsets?

44

- Confidence does not have an anti-monotone property
- $c(ABC \Rightarrow D)$ can be larger or smaller than $c(AB \Rightarrow D)$
- However, confidence of rules generated from the same itemset has an anti-monotone property
- $L = \{A, B, C, D\}$: $c(ABC \Rightarrow D) \geq c(AB \Rightarrow CD) \geq c(A \Rightarrow BCD)$
- Confidence is anti-monotone with respect to the number of items on the right-hand side of the rule

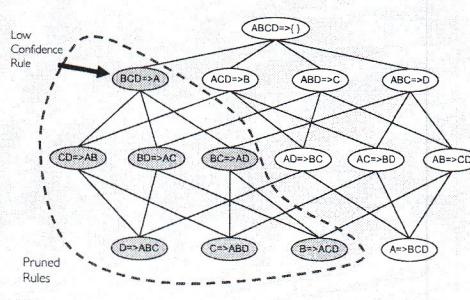
→ However, since it's more difficult to predict the presence of two product instead of the presence of one, we still have a rule:
 $c(ABC \Rightarrow D) \geq c(A \Rightarrow BCD)$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Rule Generation for Apriori Algorithm

45



Notation: "head" → "tail"

When we start, we start with only one element on the tail and then, going on with levels, we add elements to the tail

confidence = how many times the tail appears when we have the head (conditional probability)

Association Rules Generation

46

```

ASSOCIATIONRULES ( $F$ , minconf):
1 foreach  $Z \in F$ , such that  $|Z| \geq 2$  do
2    $A \leftarrow \{X \mid X \subseteq Z, X \neq \emptyset\}$ 
3   while  $A \neq \emptyset$  do
4      $X \leftarrow$  maximal element in  $A$ 
5      $A \leftarrow A \setminus X$  // remove  $X$  from  $A$ 
6      $c \leftarrow sup(Z)/sup(X)$ 
7     if  $c \geq minconf$  then
8       print  $X \rightarrow Y, sup(Z), c$ 
9     else
10       $A \leftarrow A \setminus \{W \mid W \subseteq X\}$ 
           // remove all subsets of  $X$  from  $A$ 

```

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

analogue to the idea with supports

Example of Rule Generation (1)

- Suppose that we computed the following frequent itemsets and we have a minconf of 0.9

sup	itemsets
6	B
5	E, BE
4	A, C, D, AB, AE, BC, BD, ABE
3	AD, CE, DE, ABD, ADE, BCE, BDE, ABDE

- We start from the largest itemset (ABDE) and consider all the possible subsets $\mathcal{A} = \{\text{ABDE}(3), \text{ABD}(3), \text{ABE}(4), \text{ADE}(3), \text{BDE}(3), \text{AB}(4), \text{AD}(3), \text{AE}(4), \text{BD}(4), \text{BE}(5), \text{DE}(3), \text{A}(4), \text{B}(6), \text{D}(4), \text{E}(5)\}$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Example of Rule Generation (2)

- The first subset is $X = \text{ABD}$, and the confidence of $\text{ABD} \Rightarrow \text{E}$ is $3/3 = 1.0$, so we output the rule.
- The next subset is $X = \text{ABE}$, but the corresponding rule $\text{ABE} \Rightarrow \text{D}$ is not strong since $c(\text{ABE} \Rightarrow \text{D}) = 3/4 = 0.75$
- We can thus remove from \mathcal{A} all subsets of ABE 
- The updated set of antecedents is therefore $\mathcal{A} = \{\text{ADE}(3), \text{BDE}(3), \text{AD}(3), \text{BD}(4), \text{DE}(3), \text{D}(4)\}$
- When the algorithm ends, it will output the rules:
 $\text{ABD} \Rightarrow \text{E}$ conf=1.0; $\text{ADE} \Rightarrow \text{B}$ conf=1.0; $\text{BDE} \Rightarrow \text{A}$, conf=1.0;
 $\text{AD} \Rightarrow \text{BE}$ conf=1.0; $\text{DE} \Rightarrow \text{AB}$ conf=1.0

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Rule Assessment Measures

How to Set an Appropriate Minsup?

- If minsup is set too high, we could miss itemsets involving interesting rare items (e.g., expensive products)
- If minsup is set too low, it is computationally expensive and the number of itemsets is very large
- A single minimum support threshold may not be effective

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

If we set min sup = 1 then we're looking for itemsets that are always there, it's not a great discover

If we set min sup = 0.01 then items are very easily frequent \rightarrow we obtain a huge amount of frequent itemsets. How can we select among them?

Lift

- Lift is the ratio of the observed joint probability of X and Y to the expected joint probability if they were statistically independent.
 $\text{lift}(X \Rightarrow Y) = \frac{\text{sup}(X, Y)}{\text{sup}(X)\text{sup}(Y)} = \text{conf}(X \Rightarrow Y)/\text{sup}(Y)$
- Lift is a measure of the deviation from stochastic independence (if it is 1 then X and Y are independent)
- Lift also measures the surprise of the rule. A lift close to 1 means that the support of a rule is expected considering the supports of its components.
- We typically look for values of lift that are much larger (i.e., above expectation) or much smaller (i.e., below expectation) than one.

$$\approx \frac{\text{P}(Y|X)}{\text{P}(Y)}$$

this is not a proper union/inters...
it's an "AND"

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Summarizing Itemsets

we want to find something that summarizes informations. For example: if we say that "ABDE" is a frequent itemset we also say that {A, B, D, E} and all their combinations are frequent itemsets.

Maximal frequent itemsets are interesting because all their subsets are frequent too

Maximal Frequent Itemsets

- A frequent itemset X is called maximal if it has no frequent supersets
- The set of all maximal frequent itemsets, given as

$$M = \{X \mid X \in F \text{ and } \exists Y \supset X, \text{ such that } Y \in F\}$$

- M is a condensed representation of the set of all frequent itemset F , because we can determine whether any itemset is frequent or not using M
- If there is a maximal itemset Z such that $X \subseteq Z$, then X must be frequent, otherwise Z cannot be frequent
- However, M alone cannot be used to determine $\text{sup}(X)$, we can only use to have a lower-bound, that is, $\text{sup}(X) \geq \text{sup}(Z)$ if $X \subseteq Z \in M$.

Prof. Perluca Lanz | POLITECNICO DI MILANO

Once we have a maximal frequent itemset we're interested in knowing when/where the support changes.

We know that all the supersets have smaller supports than their subsets.

If a subset is such that ALL its supersets have STRICTLY smaller supports then it's called CLOSED FREQ. ITEMSET.

→ If we have a closed frequent itemset X then every itemset containing X has a STRICTLY smaller support than X .

Closed Frequent Itemsets

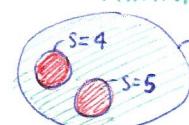
- An itemset X is closed if all supersets of X have strictly less support, that is,

$$\text{sup}(X) > \text{sup}(Y), \text{ for all } Y \supset X$$

- The set of all closed frequent itemsets C is a condensed representation, as we can determine whether an itemset X is frequent, as well as the exact support of X using C alone.

Prof. Perluca Lanz | POLITECNICO DI MILANO

CLOSED FREQ. ITEMSETS
MINIMAL GENERATOR



minimal generator = all the subsets have a strictly larger support than the minimal

Minimal Generators

- A frequent itemset X is a minimal generator if it has no subsets with the same support

$$G = \{X \mid X \in F \text{ and } \exists Y \subset X, \text{ such that } \text{sup}(X) = \text{sup}(Y)\}$$

- Thus, all subsets of X have strictly higher support, that is,

$$\text{sup}(X) < \text{sup}(Y)$$

Prof. Perluca Lanz | POLITECNICO DI MILANO

we start from:

Transaction database

Tid	Itemset
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

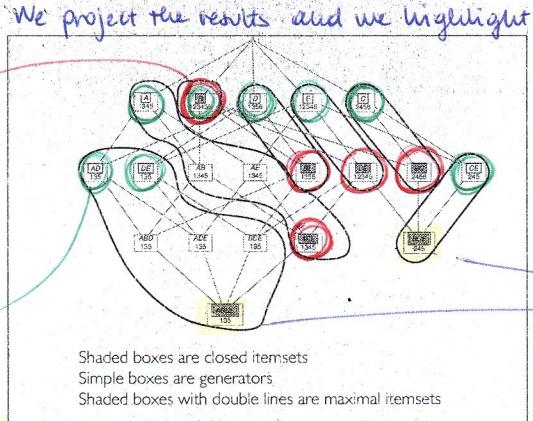
and we obtain:

Frequent itemsets (minsup = 3)

SUP	Itemsets
6	B
5	E, BE
4	A, C, D, AB, AE, BC, BD, ABE
3	AD, CE, DE, ABD, ADE, BCE, BDE, ABDE

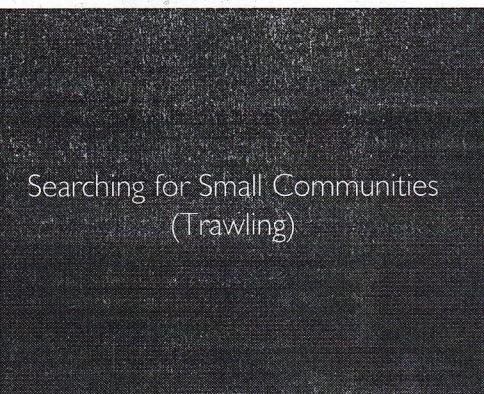
for example: B is a dotted itemset because every itemset that contains B has a lower support

for example: AD is a generator because it has no subsets with the same support ($s(A)=4, s(D)=4$)

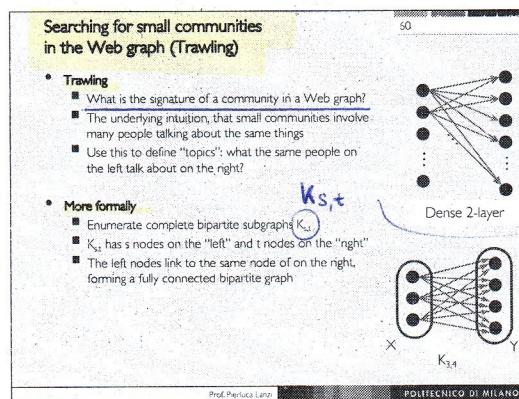


To check that an itemset is a MAXIMAL ITEMSET we look if it is at the end of the tree, to check if it's a GENERATOR we have to look at the upper levels of the itemsets (since we look at the subsets), to check if it's a CLOSED ITEMSET we look at the lower level (since we look at the bigger itemsets containing the one we're analyzing)

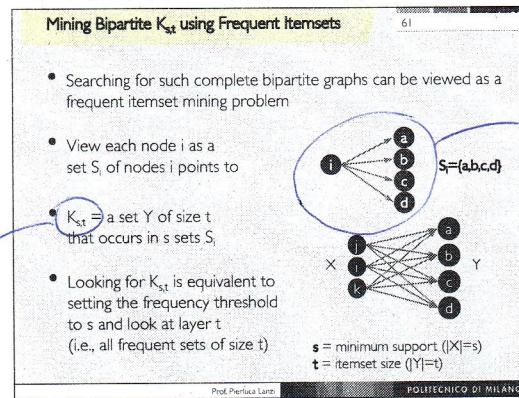
these are maximal because there's nothing below. if there would have been then they wouldn't be maximal
(These 2 itemsets somehow already summarize the whole tree)



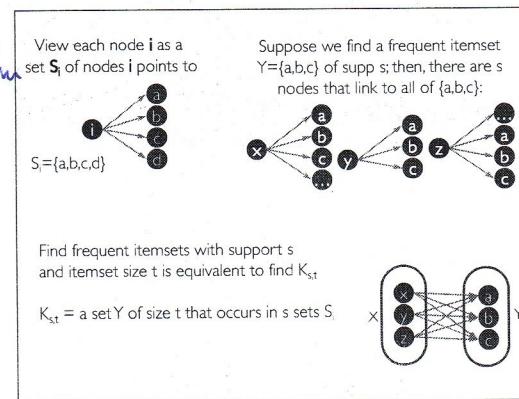
We can trawl it as an ASSOCIATION MINING problem



community of S talking about t topics

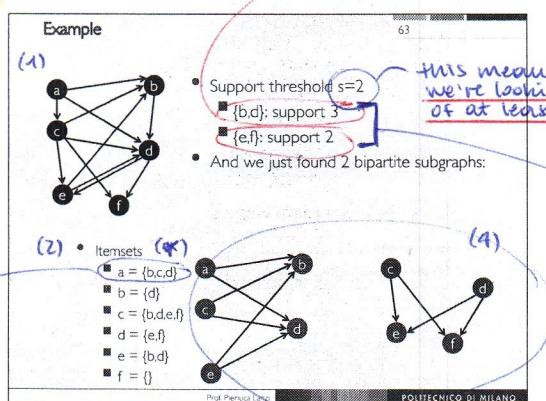


for the node i we build a set of all the nodes to which i is connected: $\{a, b, c, d\}$ *



- From the graph (1) generate the database (2)
- From the database (2) generate the frequent itemsets (3)
- From the frequent itemsets (3) generate the subgraphs (4)

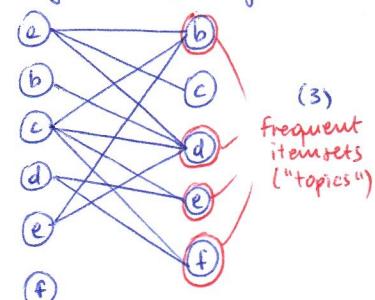
This because a is connected to b, c, d



$\{b,d\}$ supports 3 means that are at least 3 elements that contains $\{b,d\}$ (in fact: a,c,e)

this means that basically we're looking for communities of at least 2 people discussing stuff

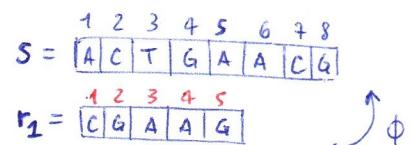
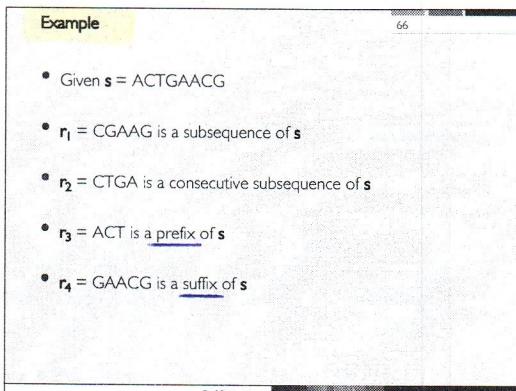
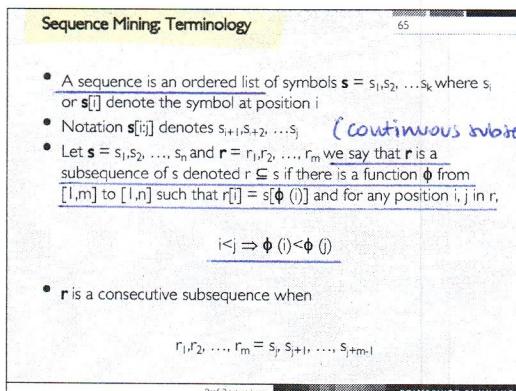
Starting from (*) we get:



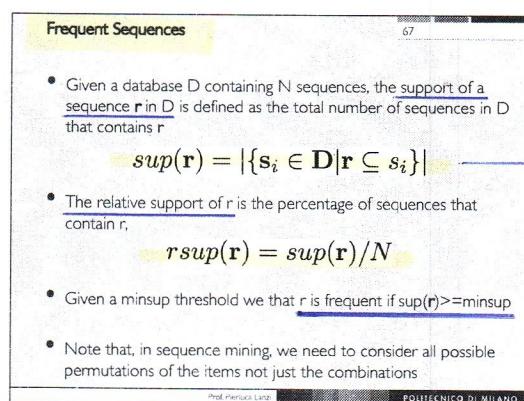
We can divide the community

Mining Frequent Sequences

This is based on the idea that an itemset it's not a set but a sequence (the order inside the transaction is important)



In this case ϕ is:
 $\phi(1) = 2$ $\phi(4) = 6$
 $\phi(2) = 4$ $\phi(5) = 8$
 $\phi(3) = 5$
 $\Rightarrow \exists \phi \Rightarrow r_2$ is a subsequence



this check how many times we see a sequence (subsequence) in the database

Example

- Given the following sequence database

Id	Sequence
s ₁	CAGAAGT
s ₂	TGACAG
s ₃	GAAGT

- With a minsup of 3, the set of frequent subsequences is
 - (A(3), G(3), T(3))
 - AA(3), AG(3), GA(3), GG(3)
 - AAG(3), GAA(3), GAG(3)
 - GAAG(3)

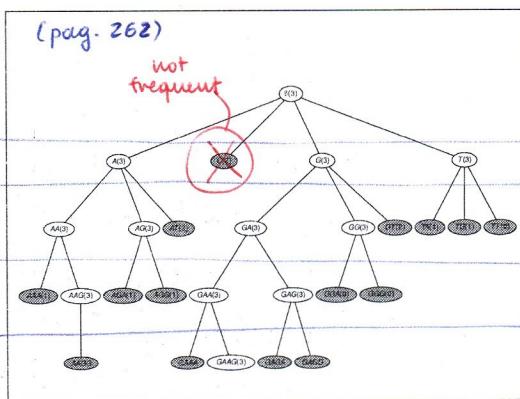
Support? 3 : S₁ : CAGAAGT
S₂ : TGACAG
S₃ : GAAGT

for example: what is the support of A? 3. Even if A appears more than once in every sequence we don't care (we don't consider repetitions)

for example: CAG has a support of 2 \Rightarrow it's not frequent (minsup = 3)

Level-Wise Sequence Mining: GSP Algorithm

- Searches the sequence prefix tree using a level-wise (breadth-first search).
- Given the set of frequent sequences at level k, the algorithm generates the candidate for level k+1 and compute the support of each candidate and prune the not frequent ones.
- For each sequence s_i in D, we check if a candidate r is a subsequence of s_i, if it is the support of r is incremented. Once the frequent sequences at level k are computed the k+1 candidates are generated.
- For each leaf r_a, the sequence is extended with the last symbol of any other leaf r_b that shares the same prefix (it has the same parent) so r_{ab} = r_a + r_b[k]. If r_{ab} is infrequent, we prune it.



Notice: this time the order (in the sequence) matters, so from A we'll start AA, AG, AT but for AG we won't connect to G! from G we'll start GA! (In apriori we wouldn't do it!)

→ same procedure as with apriori but the number of candidates is even larger (this is an issue)

1st level: only 1-sequences
2nd level: 2-sequences
3rd
4th

The GSP Algorithm

```

GSP (D, Σ, minsup):
1  $\mathcal{F} \leftarrow \emptyset$ 
2  $\mathcal{C}^{(1)} \leftarrow \{\}$  // Initial prefix tree with single symbols
3 foreach  $s \in \Sigma$  do Add s as child of  $\emptyset$  in  $\mathcal{C}^{(1)}$  with  $sup(s) \leftarrow 0$ 
4  $k \leftarrow 1$  // k denotes the level
5 while  $\mathcal{C}^{(k)} \neq \emptyset$  do
6   COMPUTESUPPORT ( $\mathcal{C}^{(k)}, D$ )
7   foreach leaf  $r \in \mathcal{C}^{(k)}$  do
8     if  $sup(r) \geq minsup$  then  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(r, sup(r))\}$ 
9     else remove  $r$  from  $\mathcal{C}^{(k)}$ 
10   $\mathcal{C}^{(k+1)} \leftarrow EXTENDPREFIXTREE (\mathcal{C}^{(k)})$ 
11   $k \leftarrow k + 1$ 
12 return  $\mathcal{F}^{(k)}$ 

```

The GSP Algorithm

```

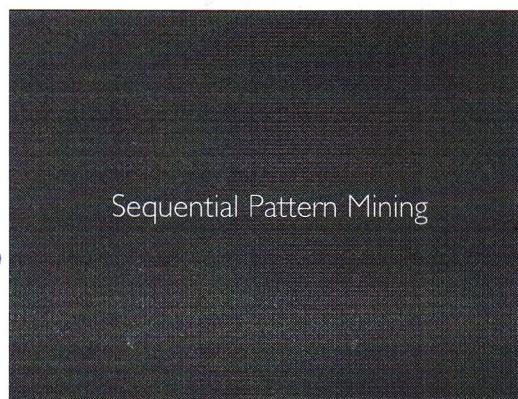
COMPUTESUPPORT ( $\mathcal{C}^{(k)}, D$ ):
1 foreach  $s_i \in D$  do
2   foreach  $r \in \mathcal{C}^{(k)}$  do
3     if  $r \subseteq s_i$  then  $sup(r) \leftarrow sup(r) + 1$ 

EXTENDPREFIXTREE ( $\mathcal{C}^{(k)}$ ):
1 foreach leaf  $r_a \in \mathcal{C}^{(k)}$  do
2   foreach leaf  $r_b \in CHILDREN(PARENT(r_a))$  do
3      $r_{ab} \leftarrow r_a + r_b[k]$  // extend  $r_a$  with last item of  $r_b$ 
     // prune if there are any infrequent subsequences
4     if  $r_{ab} \in \mathcal{C}^{(k)}$ , for all  $r_c \subset r_{ab}$ , such that  $|r_c| = |r_{ab}| - 1$  then
5       Add  $r_{ab}$  as child of  $r_a$  with  $sup(r_{ab}) \leftarrow 0$ 
6     if no extensions from  $r_a$  then
7       remove  $r_a$ , and all ancestors of  $r_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
8 return  $\mathcal{C}^{(k)}$ 

```

*

- frequent itemset: we go to Esselunga and buy stuff, what do we buy together?
- frequent sequences: we have a DNA sequence and we look for frequent patterns inside of sequence (it's different from before because a sequence is not just an itemset)
- sequential pattern mining: we look for sequences of itemsets and we're going to look at relationships among items



→ association among transactions

we had **FREQUENT ITEMSETS** where we had transactions of itemsets and we were looking for frequent subsets, then we had **FREQUENT SEQUENCES** where every itemset is a sequence, now we have **SEQUENTIAL PATTERN MINING** where we don't consider just a sequence but sequences. *

Sequential Pattern Mining

74

- Association rules do not consider the order of transactions, however, in many applications ordering is significant
- In market basket analysis, it is interesting to know whether people buy some items in sequence
- Example
 - Buying bed first and then bed sheets later
 - Navigational patterns of users in a Web site from sequences of page visits of users

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Examples of Sequence

75

- **Web sequence**
 - < {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera} {Shopping Cart} {Order Confirmation} {Return to Shopping} >
- **Sequence of initiating events causing the accident at 3-mile Island:**
 - <{clogged resin} {outlet valve closure} {loss of feedwater} {condenser polisher outlet valve shut} {booster pumps trip} {main waterpump trips} {main turbine trips} {reactor pressure increases}>
- **Sequence of books checked out at a library:**
 - <{Fellowship of the Ring} {The Two Towers} {Return of the King}>

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Basic Concepts

76

- Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of items.
- A sequence is an ordered list of itemsets.
- We denote a sequence s by $\langle a_1, a_2, \dots, a_n \rangle$, where a_i is an itemset, also called an element of s .
- An element (or an itemset) of a sequence is denoted by $\langle x_1, x_2, \dots, x_n \rangle$, where $x_i \in I$ is an item.
- We assume without loss of generality that items in an element of a sequence are in lexicographic order

Prof. Pierluca Lanzi POLITECNICO DI MILANO

For example:
sequence of 3 elements:

bread, milk	$a_1 = \{x_1, x_2\}$
water	a_2
cookies	a_3

When we perform sequential pattern mining we look for relations between elements of one itemsets and elements of another one

Sequential Pattern Mining

77

- Sequence databases consist of sequences of ordered elements or events (with or without time)
- Sequential Pattern Mining is the mining of frequently occurring ordered events or subsequences as patterns
- Sequences consists of nominal (categorical) data. When sequence are based on numerical data, then we apply time series analysis

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Basic concepts

- The size of a sequence is the number of elements (or itemsets) in the sequence
- The length of a sequence is the number of items in the sequence (a sequence of length k is called k -sequence)
- Given two sequences, $s_1 = \langle a_1, a_2, \dots, a_v \rangle$ and $s_2 = \langle b_1, b_2, \dots, b_w \rangle$
- s_1 is a subsequence of s_2 if there exist integers $1 \leq j_1 < j_2 < \dots < j_r \leq v$ such that $a_1 \subseteq b_{j_1} \wedge a_2 \subseteq b_{j_2} \wedge \dots \wedge a_r \subseteq b_{j_r}$

a_j, b_j are itemsets!

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Example

- Let $I = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$.
- The sequence $\langle \{3\} \{4, 5\} \{8\} \rangle$ is contained in (or is a subsequence of) $\langle \{6\} \{3, 7\} \{9\} \{4, 5, 8\} \{3, 8\} \rangle$
- In fact, $\{3\} \subseteq \{3, 7\}$, $\{4, 5\} \subseteq \{4, 5, 8\}$, and $\{8\} \subseteq \{3, 8\}$
- However, $\langle \{3\} \{8\} \rangle$ is not contained in $\langle \{3, 8\} \rangle$ or vice versa
- The size of the sequence $\langle \{3\} \{4, 5\} \{8\} \rangle$ is 3, and the length of the sequence is 4.

$$\text{size} = \# \text{ itemsets}, \text{length} = \# \text{ items}$$

Prof. Pierluca Lanzi POLITECNICO DI MILANO

This means: we once bought 3 and in the transaction after we bought 8 $\Rightarrow \langle \{3\}, \{8\} \rangle$ (3 happened at some point and later happened 8)
On the other hand: $\langle \{3, 8\} \rangle$ means that 3 and 8 happened together, which means that there is no sequential aspect of this.

Goal of Sequential Pattern Mining

- The input is a set S of input data sequences (or sequence database)
- The problem of mining sequential patterns is to find all the sequences that have a user-specified minimum support
- Each such sequence is called a frequent sequence, or a sequential pattern
- The support for a sequence is the fraction of total data sequences in S that contains this sequence

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Terminology

- Itemset**
 - Non-empty set of items
 - Each itemset is mapped to an integer
- Sequence**
 - Ordered list of itemsets
- Support for a Sequence**
 - Fraction of total customers that support a sequence.
- Maximal Sequence**
 - A sequence that is not contained in any other sequence.
- Large Sequence**
 - Sequence that meets minisup

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Frequent Sequences

- Given a database $D = \{s_1, \dots, s_N\}$ of N sequences and a sequence r , we define the support count of r as $\text{sup}(r) = |\{s_i \mid r \text{ is a subsequence of } s_i\}|$
- And the support (or relative support), $\text{rsup}(r) = \text{sup}(r)/N$

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Example 1.

Customer ID	Transaction Time	Items Bought
1	June 25 '93	30
1	June 30 '93	90
2	June 10 '93	10,20
2	June 15 '93	30
2	June 20 '93	40,60,70
3	June 25 '93	30,50,70
4	June 25 '93	30
4	June 30 '93	40,70
4	July 25 '93	90
5	June 12 '93	90

Customer ID	Customer Sequence
1	<(30)>
2	<(10 20)(30)(40 60 70)>
3	<(30, 50, 70)>
4	<(30)(40 70)(90)>
5	<(90)>

Maximal seq with support > 40%
<(30)(90)>
<(30)(40 70)>

Note: Use Minisup of 40%, no less than two customers must support the sequence
 <(10 20)(30)> Does not have enough support (Only by Customer #2)
 <(30)>, <(70)>, <(30)(40)> ... are not maximal.

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Example 2.

- Given the following sequence database and a minsup of 2

SID	Sequence
10	<(a)(ab)(ac)(d)(cf)>
20	<(ad)(c)(bc)(ae)>
30	<(e)(ab)(df)(g)(b)>
40	<(e)(g)(af)(c)(b)(c)>

- <(a)(bc)(d)(c)> is a subsequence of <(a)(ab)(ac)(d)(cf)>
- <(ab)(c)> is a frequent sequence or a sequential pattern

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Example 3.

- Given the sequence database

SID	Sequence
10	<(bd)(c)(a)(ac)>
20	<(b)(c)(ce)(b)(f(g))>
30	<(ah)(b)(b)(a)(b)(f)>
40	<(ba)(ce)(d)>
50	<(a)(bd)(b)(c)(b)(ade)>

- Given a minsup of 2, <(bd)cb> is a sequential pattern (or frequent sequence)

actually:
 <(b d)(c)(b)>

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Apriori principle still applies

If a sequence S is not frequent,
 then every super-sequence of S is not
frequent

In the previous example,

How do we proceed?
 Level-wise search!

Mining Sequential Patterns (Outline of the basic principle)

- Level-by-level do
 - Generate candidate sequences
 - Scan database to collect support counts
 - Use Apriori property to prune candidates
- Only generate candidates satisfying Apriori property
- Limitations
 - It can generate a massive amount of candidates (1000 frequent 1-sequences can generate around 1.5 million candidates)
 - It requires many scans of the database
- Accordingly, as for association rules, there are several prefix-based extensions (e.g. PrefixSpan) that deal with all these limitations

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Example	
Seq ID	Sequence
10	<(a)(b)(b)(a)>
20	<(b)(a)(c)(b)>
30	<(a)(b)(a)(b)(b)>
40	<(a)(c)(b)>
50	<(a)(b)(c)(c)(b)(a)d>

all possible letters appearing in the database

1st scan - 8 candidates
6 sequential patterns (length 1)

2nd scan - 51 candidates
19 sequential patterns (length 2)
10 candidates were not in the database

3rd scan - 46 candidates
19 sequential patterns (length 3)
20 candidates were not in the database

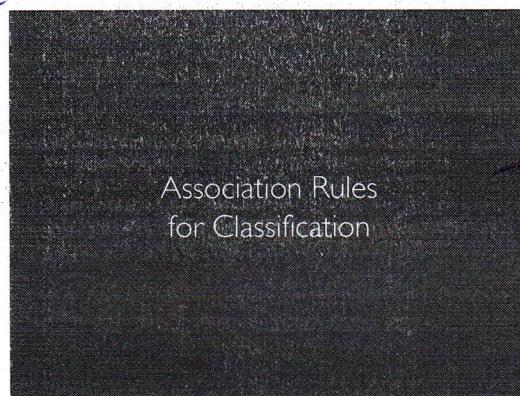
<a> <c> <d> <a> <c> <d>

<a> ... <ff> <a> ... <ff>

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

since we consider all the combinations also considering the order we may end with a lot of candidates which are not even in the database (\Rightarrow we waste a lot of time)

\Rightarrow That's why we prefer the FP-tree approach: there is no candidates generation



Association rules means: $X \Rightarrow Y$
(if we buy X then we buy Y).
We can translate it (use it) for classification if we consider X as a set of features and Y a label

Mining Association Rules for Classification (the CBA algorithm)	
•	Association rule mining assumes that the data consist of a set of transactions. Thus, the typical tabular representation of data used in classification must be mapped into such a format.
•	Association rule mining is then applied to the new dataset and the search is focused on association rules in which the tail identifies a class label
•	$X \Rightarrow c_i$ (where c_i is a class label)
•	The association rules are pruned using the pessimistic error-based method used in C4.5
•	Finally, rules are sorted to build the final classifier.

The Weather Dataset					
Outlook	Temp	Humidity	Windy	Play	
Sunny	Hot	High	False	No	
Sunny	Hot	High	True	No	
Overcast	Hot	High	False	Yes	
Rainy	Mild	High	False	Yes	
Rainy	Cool	Normal	False	Yes	
Rainy	Cool	Normal	True	No	
Overcast	Cool	Normal	True	Yes	
Sunny	Mild	High	False	No	
Sunny	Cool	Normal	False	Yes	
Rainy	Mild	Normal	False	Yes	
Sunny	Mild	Normal	True	Yes	
Overcast	Mild	High	True	Yes	
Overcast	Hot	Normal	False	Yes	
Rainy	Mild	High	True	No	

one row is one transaction.
We translate the dataset in one-hot-encoding approach
 \Rightarrow we obtain a binary dataset,
which was our representation of the dataset for association rule

CLASS (LABELS)

CBA Model	
outlook=overcast ==> play=yes	
humidity=normal windy=False ==> play=yes	
outlook=rainy windy=False ==> play=yes	
outlook=sunny humidity=high ==> play=no	
outlook=rainy windy=True ==> play=no	
(default class is the majority class)	

This is better than a decision tree or a logistic regression (in terms of READABILITY)

this transform leads to this: