

# Algorithms that Search for Association Rules

## What are Association Rules?

- Association Rules have some similarity with Classification Rules
  - ... similar idea: `left-hand-side`  $\Rightarrow$  `right-hand-side`
  - the difference between Association and Classification rules
  - ... concerns the information in the: `right-hand-side`
- ... `right-hand-side` of an **Association Rule**
  - predicts the **value of one or more attributes** (any attribute)
  - ... algorithm searches for those attributes; *without user intervention*
- ... `right-hand-side` of a **Classification Rule**
  - predicts the **value of a single attribute** (*defined by the user*)
- A classification algorithm can search for association rules
  - given that it executes once for every combination of attributes, with every possible combination of values, on the right-hand side!
  - ... an enormous number of association rules to be pruned...

## When to Search for Association Rules?


- Whenever our problem can be formulated within the scope of
  - **market-basket** analysis
- The **market-basket analysis** goal is to
  - “**find groups of items that tend to occur together in transactions**”
  - e.g., supermarket checkout records are analyzed to detect associations among items that people purchase
- ... in e-commerce it may be used for Web page personalization
  - find that 40% of users visit pages A, B and C, and that 75% (of those 40%) have similar behavior pattern of always visiting C after A and B
  - ... based on that rule, a dynamic link could be created for users who are likely to be interested in page C
- ... the previous association rule could be expressed as,
  - $(A \text{ and } B) \Rightarrow C \mid \text{support}=40\% \mid \text{confidence}=75\% \mid$

## ... search for Association Rules is Unsupervised

- The goal is not predict a target value
  - because there is no predefined target (or class attribute)
  - ... no class values are provided
- The goal is to find the intrinsic relations among data
  - algorithms explore the interconnectedness of the data
  - ... as if all attributes were considered as “possible classes”
- Each rule associates two sets (LHS and RHS) of attribute values
  - i.e., left-hand-side  $\Rightarrow$  right-hand-side
  - both the LHS and the RHS are extracted without user intervention
- The association rules also provide classification support?
  - in fact the set of rules having the same attribute (e.g., attrA) in the RHS can be regarded as classification rules (and attrA as a class)
  - ... but classification methods should be used for classification problems!

“a classical story”

## Stories – Beer and Diapers



- ♦ Diapers and Beer. Most famous example of market basket analysis for the last few years. If you buy diapers, you tend to buy beer.
- T. Blischok headed Terradata's Industry Consulting group.
- K. Heath ran self joins in SQL (1990), trying to find two itemsets that have baby items, which are particularly profitable.
- Found this pattern in their data of 50 stores/90 day period.
- Unlikely to be significant, but it's a nice example that explains associations well.

*RF Ronny Kohavi ICML 1997*

Whether a real story or a “fabrication” it is as effective illustrative example

## Problem Formulation

The **market-basket analysis** goal is to,  
“find groups of **items** that tend to occur together in **transactions**”

- $I = \{ i_1, i_2, \dots, i_M \}$ , where
  - the  $I$  is a set of  $M$  **items**
  - each  $i_j$  is an item
- $T = \{ t_1, t_2, \dots, t_N \}$ 
  - the  $T$  is a set of  $N$  **transactions**
  - each  $t_j$  is a transaction
  - and  $t_j \subseteq I$  (i.e.,  $t_j$  is a set of items)

## An example – transactions in supermarket

the recorded  
supermarket transactions

TID	ITEM
t1	beer
t1	milk
t2	bread
t2	butter
t3	bread
t3	butter
t3	jelly
t4	bread
t4	butter
t4	bread
t5	beer
t5	bread

- $I = \{ \text{beer, bread, butter, jelly, milk} \}$ 
  - the set of **items** sold in the store
- $T = \{ t1, t2, t3, t4, t5 \}$ 
  - the set of **baskets recorded** at checkout
- each **basket and the items purchased**
  - $t1 = \{ \text{beer, milk} \}$
  - $t2 = \{ \text{bread, butter} \}$
  - $t3 = \{ \text{bread, butter, jelly} \}$
  - $t4 = \{ \text{bread=2, butter} \}$  // **describe quantity**
  - $t5 = \{ \text{beer, bread} \}$

## ... transactions in supermarket – different representation

the recorded  
supermarket transactions

TID	ITEM
t1	beer
t1	milk
t2	bread
t2	butter
t3	bread
t3	butter
t3	jelly
t4	bread
t4	butter
t4	bread
t5	beer
t5	bread

a tabular **sparse** representation of the same information

	beer	bread	butter	jelly	milk
t1	1				1
t2		1	1		
t3		1	1	1	
t4		2	1		
t5	1	1			

.tab (Orange file)

a **compact** representation  
of the same information

t1	beer milk
t2	bread butter
t3	bread butter jelly
t4	bread=2 butter
t5	beer bread

.basket (Orange file)



## ... how to generate the different representations?

the recorded  
supermarket transactions

TID	ITEM
t1	beer
t1	milk
t2	bread
t2	butter
t3	bread
t3	butter
t3	jelly
t4	bread
t4	butter
t4	bread
t5	beer
t5	bread

apply the  
**pivot table**  
operator

a tabular **sparse** representation  
of the same information

Count		Labels				
Row		beer	bread	butter	jelly	milk
t1		1				1
t2			1	1		
t3			1	1	1	
t4			2	1		
t5		1	1			

a **compact** representation  
of the same information

implement (e.g., in Python) a script  
to **transform table** into this format

```
beer milk  
bread butter  
bread butter jelly  
bread=2 butter  
beer bread
```

## ... another example – documents in a collection

the collection  
of documents

*doc1*: student, teach, school

*doc2*: student, school

*doc3*: teach, soccer, city, game

*doc4*: tennis, soccer

*doc5*: game, school, tennis

*doc6*: student, tennis

- $I = \{ \text{student, teach, school, city, game, tennis, soccer} \}$ 
  - the set of **keywords** in the documents' collection
- $T = \{ t1, t2, t3, t4, t5, t6 \}$ 
  - the set of **documents recorded** in the collection
- each document is treated as a bag (**basket**) of words
  - $t1 = \{ \text{student, teach, school} \}$
  - $t2 = \{ \text{student, school} \}$
  - $t3 = \{ \text{tennis, soccer, city, game} \}$
  - $t4 = \{ \text{tennis, soccer} \}$
  - $t5 = \{ \text{game, school, tennis} \}$
  - $t6 = \{ \text{student, tennis} \}$

## ... recall – What is Market Basket Analysis?

- Understanding behavioral patterns (e.g., of customers)
  - “what items were bought together?”
  - ... “what is in each shopping car/basket?”
- The basket data is a collection of items bought in a transaction
  - an “itemset”
- How does this data differ from a transaction database?
  - apply a “pivot table” operator to the transaction (operational) table
- The overall goal is to generate qualified decisions and strategies
  - what to put on sale?
  - how to place merchandise on shelves for maximizing profit?
  - how to segment customers based on their buying patterns?

## Association Rule and Itemset – Definitions

A transaction ***t*** ***contains*** ***X***, a set of items (itemset) in ***I***  
if  $X \subseteq t$

An ***association rule*** is an implication of the form

$$LHS \Rightarrow RHS$$

where  $LHS, RHS \subset I$  and  $LHS \cap RHS = \emptyset$

An ***itemset*** is a set of items

e.g.,  $X = \{ \text{milk, bread, cereal} \}$  is an itemset

A ***k-itemset*** is an itemset with ***k*** items

e.g.,  $X = \{ \text{milk, bread, cereal} \}$  is a 3-itemset

## Two (Main) Measures for “Rule Evaluation”

- **Support** of rule  $LHS \Rightarrow RHS$

- measures how often the collection of items in an association, i.e.,  $LHS$  and  $RHS$ , occur together as a percentage of all the transactions

$$\begin{aligned}\text{support}( LHS \Rightarrow RHS ) &= \\ &= P( LHS, RHS ) = \\ &= \# \text{tuples-with-both-}LHS\text{-and-}RHS / \# \text{total-of-tuples}\end{aligned}$$

- **Confidence** of rule  $LHS \Rightarrow RHS$

- measures how likely it is that  $RHS$  occurs when  $LHS$  has occurred

$$\begin{aligned}\text{confidence}( LHS \Rightarrow RHS ) &= \\ &= P( RHS \mid LHS ) = \\ &= \# \text{tuples-with-both-}LHS\text{-and-}RHS / \# \text{tuples-with-}LHS\end{aligned}$$

## The Goal of the Association Rules' Mining

Generate **all association rules**, from the dataset, that **satisfy**, the user defined, **minimum support** (*min\_sup*), and **minimum confidence** (*min\_conf*)

- An itemset satisfies minimum support if the occurrence frequency
  - of the itemset is greater or equal to ***min\_sup***
- The “**frequent itemset**” concept
  - designates an itemset that **satisfies minimum support**
- A “**strong rule**” is one that satisfies, both
  - a minimum **support threshold** and a minimum **confidence threshold**
- Rules originating from the same itemset have the same support
  - but can have different confidence

... example – transaction, itemset, support, confidence

$$T = \begin{array}{l|l} t_1 & ABCD \\ t_2 & BC \\ t_3 & AC \\ t_4 & AC \\ t_5 & ABCD \\ t_6 & ABC \end{array}$$

Itemset	Support	Confidence
A	?	?
B	?	?
C	?	?
D	?	?
AB	?	?
AC	?	?
AD	?	?
CD	?	?
ABC	?	?
ACD	?	?
ABCD	?	?
...		

... support and confidence

$$T = \begin{array}{l|l} t_1 & ABCD \\ t_2 & BC \\ t_3 & AC \\ t_4 & AC \\ t_5 & ABCD \\ t_6 & ABC \end{array}$$

Itemset	Support	Confidence
A	5 / 6	–
B	4 / 6	–
C	6 / 6	–
D	2 / 6	–
AB	3 / 6	A ⇒ B: 3 / 5 B ⇒ A: 3 / 4
AC	5 / 6	...
AD	2 / 6	...
CD	2 / 6	...
ABC	3 / 6	<div> <div>AB ⇒ C: 3 / 3</div> <div>CB ⇒ A: 3 / 3</div> <div>AC ⇒ B: 3 / 5</div> </div> <div> <div>A ⇒ BC: 3 / 5</div> <div>B ⇒ AC: 3 / 4</div> <div>C ⇒ AB: 3 / 6</div> </div>
ACD	2 / 6	...
ABCD	2 / 6	...
...		



## Are the “support” and “confidence” enough?

- Support for  $A \Rightarrow B = P(A, B)$ 
  - “how frequently the items in the rule occur together”
- Confidence for  $A \Rightarrow B = P(B | A) = P(A, B) / P(A)$ 
  - “probability of both the antecedent and the consequent appearing in the same transaction”
- ... support and confidence are used to determine if a rule is valid
  - however, there are times when both of these measures may be high, and yet still produce a rule that is not useful...
- e.g., orange-juice  $\Rightarrow$  milk | support 30% | confidence 75%
  - sounds like an excellent rule, and in most cases, it would be!
  - but, *what if customers in general buy milk 90% of the time?*
  - ... *then, orange juice customers are actually less likely to buy milk than customers in general!*

## Additional measure – Lift (Improvement)

- The “lift” (improvement) indicates the strength of a rule
  - over the random co-occurrence of the antecedent and the consequent, given their individual support
- ... “lift” provides information about the improvement,
  - i.e., the increase in probability of the consequent given the antecedent
- the “lift” is defined as

$$\begin{aligned}\text{lift}( LHS \Rightarrow RHS ) &= \\ &= \text{support}(LHS \Rightarrow RHS) / ( \text{support}(LHS) * \text{support}(RHS) ) \\ &= \text{confidence}(LHS \Rightarrow RHS) / \text{support}(RHS) \\ &= P(RHS \mid LHS) / P(RHS)\end{aligned}$$

if Lift  $\geq 1$ , then the *LHS* and the *RHS* are positively correlated,  
otherwise, the *LHS* and the *RHS* are negatively correlated

## ... the “Lift” (Improvement) – an example

- ... there are times when both of the “support” and “confidence” measures are high, and yet still produce a rule that is not useful
- e.g., orange-juice  $\Rightarrow$  milk | support 30% | confidence 75%
  - sounds like an excellent rule, and in most cases, it would be!
  - but, *what if customers in general buy milk 90% of the time?*
  - ... *then, orange juice customers are actually less likely to buy milk than customers in general!*
- The “Lift” for the previous rule is given by
$$\begin{aligned}\text{lift}(\text{orange-juice} \Rightarrow \text{milk}) &= \\ &= \text{confidence}(\text{orange-juice} \Rightarrow \text{milk}) / \text{support}(\text{milk}) \\ &= 75\% / 90\% = \mathbf{0.83} < \mathbf{1}\end{aligned}$$
- ... the increase in probability of the consequent given antecedent is ***lower than 1***, so those ***items are negatively correlated***

... example – support, confidence and lift

$$T = \begin{array}{l|l} t_1 & ABCD \\ t_2 & BC \\ t_3 & AC \\ t_4 & AC \\ t_5 & ABCD \\ t_6 & ABC \end{array}$$

Itemset	Support	Confidence	Lift
A	5/6	–	–
B	4/6	–	–
C	6/6	–	–
D	2/6	–	–
AB	3/6	A ⇒ B: 3/5 B ⇒ A: 3/4	?
AC	5/6	...	...
AD	2/6	...	...
CD	2/6	...	...
ABC	3/6	AB ⇒ C: 3/3   A ⇒ BC: 3/5 CB ⇒ A: 3/3   B ⇒ AC: 3/4 AC ⇒ B: 3/5   C ⇒ AB: 3/6	?
ACD	2/6	...	...
ABCD	2/6	...	...
...			...

... example – support, confidence and lift

$T =$   
 $t_1$  | ABCD  
 $t_2$  | BC  
 $t_3$  | AC  
 $t_4$  | AC  
 $t_5$  | ABCD  
 $t_6$  | ABC

Itemset	Support	Confidence	Lift
A	5/6	–	–
B	4/6	–	–
C	6/6	–	–
D	2/6	–	–
AB	3/6	$A \Rightarrow B: 3/5$ $B \Rightarrow A: 3/4$	$A \Rightarrow B: (3/5)/(4/6)=9/10$ $B \Rightarrow A: (3/4)/(5/6)=9/10$
AC	5/6	...	...
AD	2/6	...	...
CD	2/6	...	...
ABC	3/6	$AB \Rightarrow C: 3/3$   $A \Rightarrow BC: 3/5$ $CB \Rightarrow A: 3/3$   $B \Rightarrow AC: 3/4$ $AC \Rightarrow B: 3/5$   $C \Rightarrow AB: 3/6$	$AB \Rightarrow C: (3/3)/(6/6)=1$ ...
ACD	2/6	...	...
ABCD	2/6	...	...
...			...

## The “Lift” measure – characteristics

- If  $\text{Lift} > 1$  then
  - the rule is better at predicting the result than “just guessing”
- If  $\text{Lift} < 1$  then
  - the rule is doing worse than “just guessing”
- If  $\text{Lift} = 1$  then
  - the rule is “similar to guessing”

### *Synthesis*

Any rule with an improvement (lift) of less than 1 does not indicate a real cross-selling opportunity, no matter how high its support and confidence, because it actually offers less ability to predict a purchase than does random chance.

... “be careful” with the “Lift” measure

*Rules that hold nearly 100% of the time may not get the highest possible lift.*

An example:

5% of the people (in the database) are retired.

90% of the people (in the database) are more than 2 years old.

now, consider the rule:

`retired  $\Rightarrow$  more-than-2-years-old`

then, the lift for that rule is:  $0.05 / (0.05 \cdot 0.9) = 1.11$

which is only slightly above 1 for that rule (that holds nearly 100% of the time)!

## Some other measures

- **Leverage** of rule  $LHS \Rightarrow RHS$

- measures the proportion of additional elements covered by both the premise (left side) and consequence (right side) above the expected

$$\begin{aligned}\text{leverage}( LHS \Rightarrow RHS ) &= \\ &= \text{support}( LHS, RHS ) - \text{support}( LHS ) * \text{support}( RHS ) \\ &= P( LHS, RHS ) - P( LHS ) * P( RHS )\end{aligned}$$

- **Coverage** of rule  $LHS \Rightarrow RHS$

- measures how likely it is that the  $LHS$  occurs

$$\begin{aligned}\text{coverage}( LHS \Rightarrow RHS ) &= \\ &= \text{support}( LHS ) = \\ &= P( LHS )\end{aligned}$$

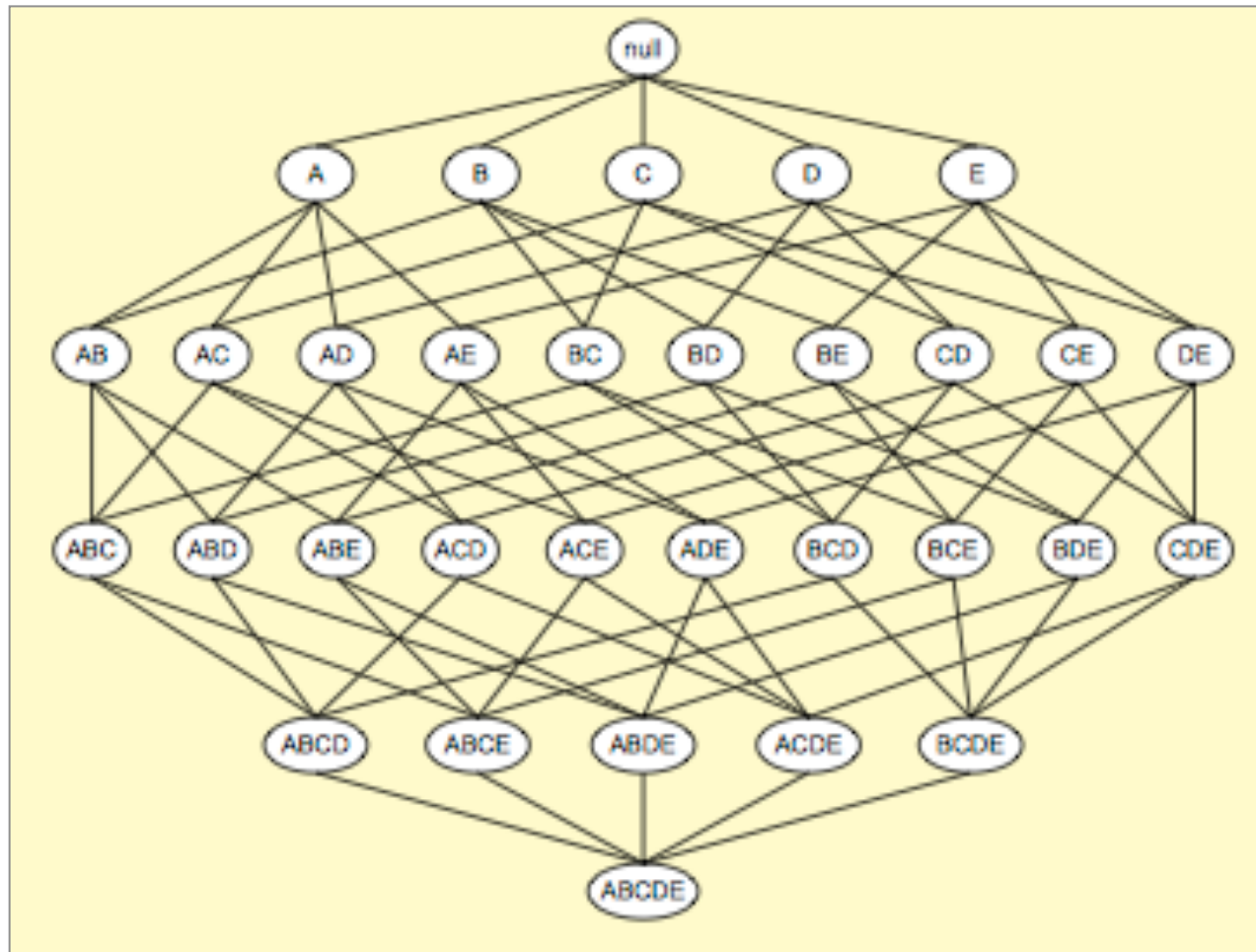


## How to search for (i.e., mining of) association rules?

- Step 1: “frequent itemset generation”
  - generate all itemsets with support  $\geq \textit{min\_sup}$  (e.g., provided by user)
- Step 2: “rule generation”
  - generate rules, from each frequent itemset, with confidence  $\geq \textit{min\_conf}$
  - ... each rule is a binary partition (into LHS, RHS) of a frequent itemset

The “frequent itemset generation is computationally expensive.

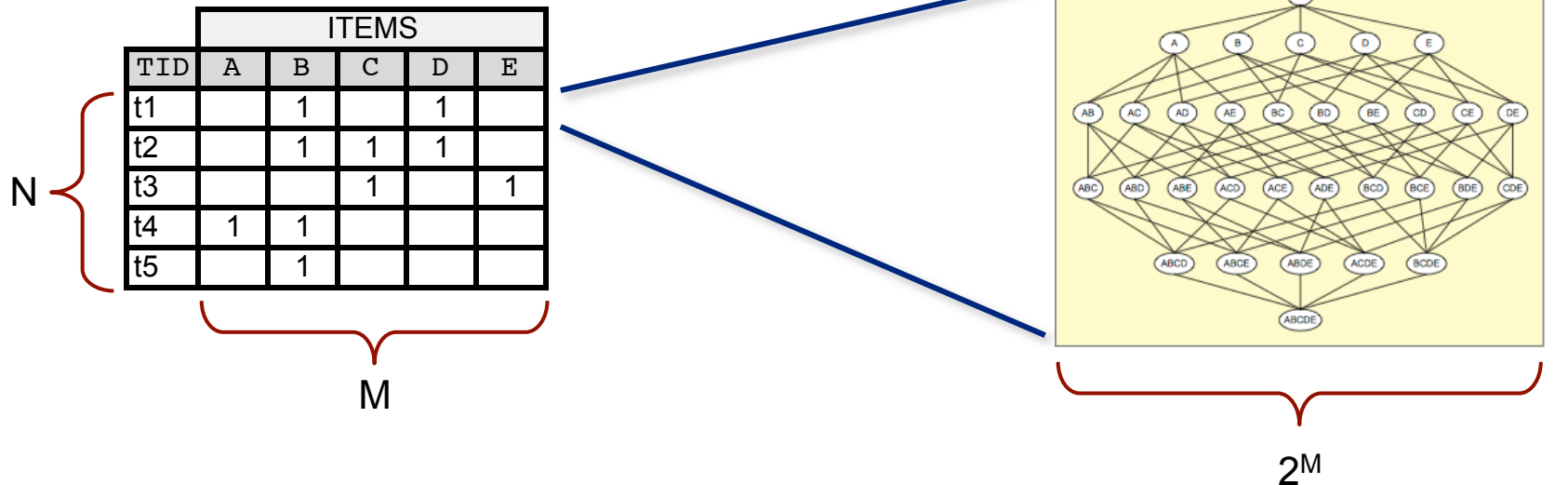
## The itemset space (an example with 5 items)



Given  $M$  items there are  $2^M$  possible candidate itemsets.

e.g., for 5 items we have a total of  $2^5=32$  itemsets

## The “frequent itemset generation”



a “brute force” approach:

Each itemset is a candidate frequent itemset

Count the support of each candidate by scanning the database  
i.e., match each (of the  $2^M$ ) candidate itemset against all the ( $N$ ) transactions!

complexity is  $O(N * 2^M)$ ; exponentially expensive with  $M$

## Strategies for the “frequent itemset generation”

- Reduce the number (  $2^M$  ) of candidate itemsets to analyze
  - use pruning techniques
  - ... e.g., as in the APRIORI algorithm
- Reduce the number (  $N$  ) of transactions to scan
  - filter the number of transactions as the size of itemset increases
- Reduce the number (  $N * 2^M$  ) of comparisons
  - use efficient data structures to store candidates or transactions
  - no need to match every candidate against every transaction
  - ... e.g., as in the APRIORI algorithm

## Algorithms to search for association rules

- The APRIORI algorithm
  - by Agrawal, R., T. Imielinski, and A. Swami
  - “Mining Association Rules between Sets of Items in Large Databases”; ACM SIGMOD International Conference on Management of Data“
  - ... available in Orange DataMining
  - “(...) *two algorithms for induction of association rules, a standard APRIORI algorithm for sparse (basket) data analysis and a variant for attribute-value data sets*”; Cf., <http://orange.biolab.si/doc/reference/Orange.associate/>
- Other algorithms
  - DHP: Dynamic Hash and Pruning
  - FP-Growth: Frequent Pattern Growth
  - H-Mine: Hyper-structure mining of frequent-patterns

## The APRIORI principle

*If an **itemset** is **frequent**, then all of its **subsets** must also be **frequent***

or

*If an **itemset** is **infrequent**, then all of its **supersets** must also be **infrequent***

- The APRIORI principle holds due to this property of “support”:

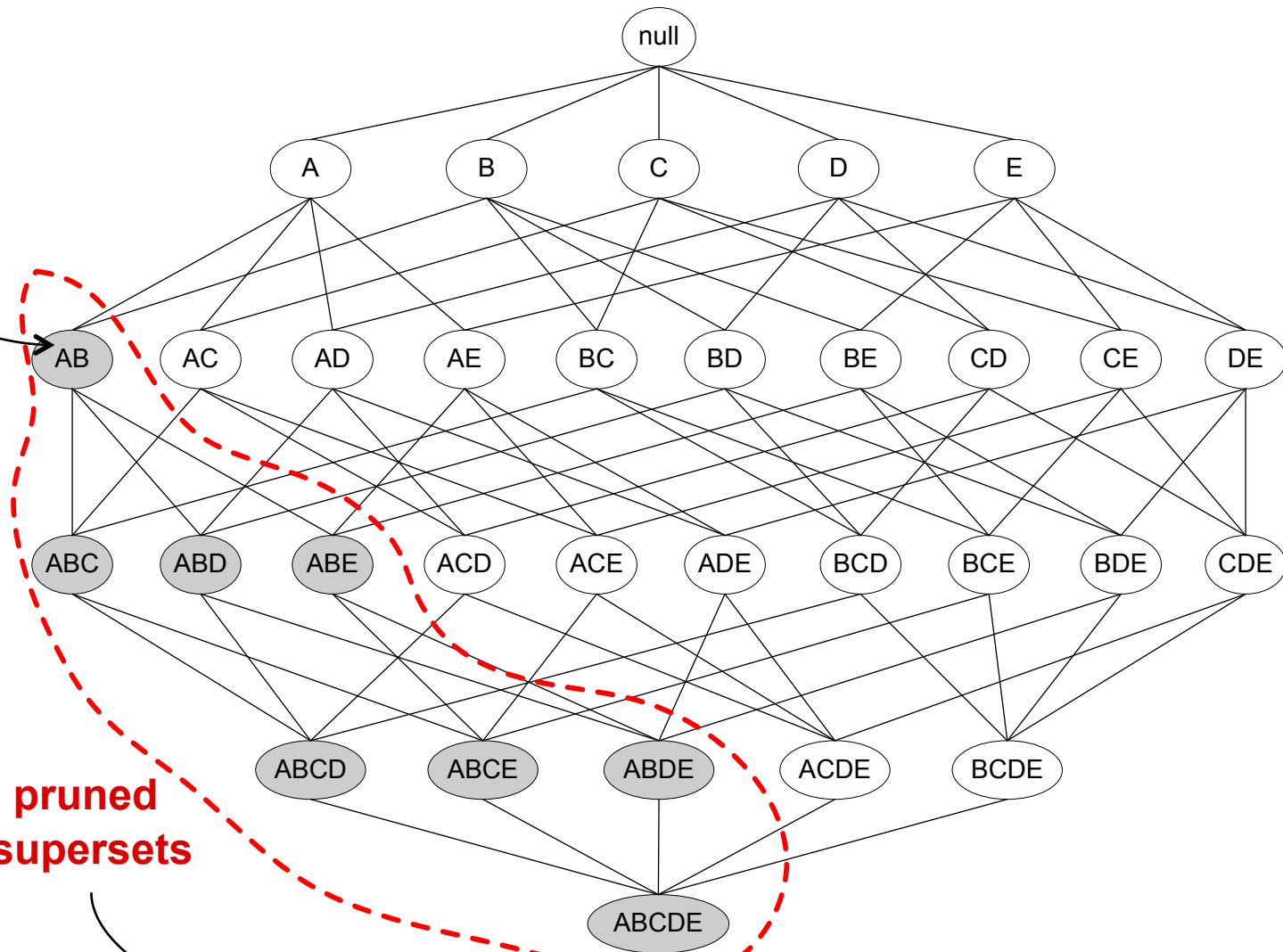
$$\forall X, Y: (X \subseteq Y) \Rightarrow (\text{support}(X) \geq \text{support}(Y))$$

- ... support of an itemset never exceeds the support of its subsets
  - this is known as the **anti-monotone** property of support
  - ... increase itemset dimension decreases, or maintains, its support

# Illustration of the APRIORI principle

found to be infrequent

pruned  
supersets



## The APRIORI algorithm – example

Consider the following “recorded database transactions”.

Given that “**support = 2**”, generate the “**frequent itemset**”

*(to simplify, “support” is presented before dividing it by the total number of transactions)*

*Hint: apply the “APRIORI principle”*

the recorded  
database transactions

TID	ITEM
t1	A
t1	C
t1	D
t2	B
t2	C
t2	E
t3	A
t3	B
t3	C
t3	E
t4	B
t4	E



## ... example – a representation step

build a sparse representation of the transaction database  
makes it easier to analyze the data

the recorded  
database transactions

TID	ITEM
t1	A
t1	C
t1	D
t2	B
t2	C
t2	E
t3	A
t3	B
t3	C
t3	E
t4	B
t4	E

apply the  
**pivot table**  
operator

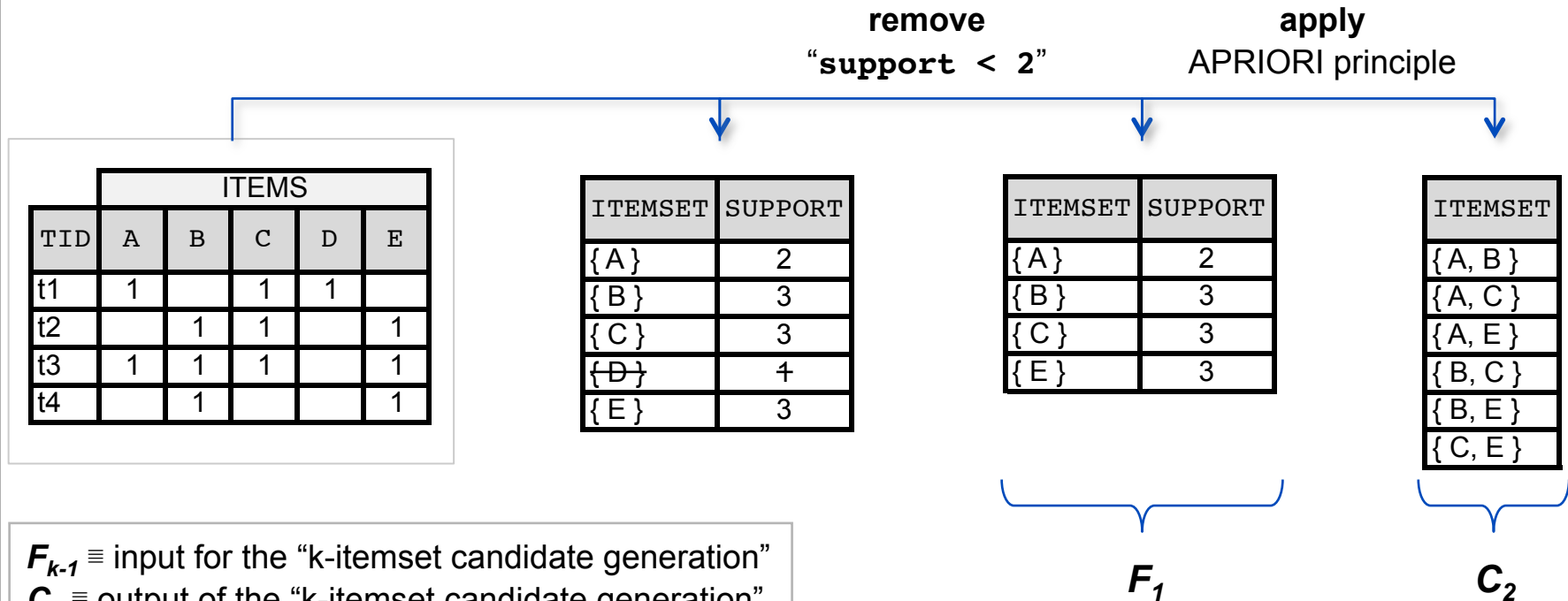
a tabular **sparse** representation  
of the same information

TID	ITEMS				
	A	B	C	D	E
t1	1		1	1	
t2		1	1		1
t3	1	1	1		1
t4		1			1

... example – apply the “APRIORI principle”

Given that “**support** = 2”, generate the “**frequent itemset**”

*Hint: apply the “APRIORI principle”*



$F_{k-1} \equiv$  input for the “k-itemset candidate generation”  
 $C_k \equiv$  output of the “k-itemset candidate generation”

see next slide...

## ... example – apply the “APRIORI principle” (cont.)

TID	ITEMS				
	A	B	C	D	E
t1	1		1	1	
t2		1	1		1
t3	1	1	1		1
t4		1			1

Given that “**support = 2**”, generate the “**frequent itemset**”

*Hint: apply the “APRIORI principle”*

ITEMSET	ITEMSET	SUPPORT	ITEMSET	SUPPORT	ITEMSET	ITEMSET	SUPPORT	ITEMSET	SUPPORT
{A, B}	<del>{A, B}</del>	4	{A, C}	2	{A, B, C}	<del>{A, B, C}</del>	4	<del>{B, C, E}</del>	2
{A, C}	{A, C}	2	{B, C}	2	{A, B, E}	<del>{A, B, E}</del>	4		
{A, E}	<del>{A, E}</del>	4	{B, E}	3	{A, C, E}	<del>{A, C, E}</del>	4		
{B, C}	{B, C}	2	{C, E}	2	{B, C, E}	{B, C, E}	2		
{B, E}	{B, E}	3							
{C, E}	{C, E}	2							

$F_{k-1}$   $\equiv$  input for the “k-itemset candidate generation”  
 $C_k$   $\equiv$  output of the “k-itemset candidate generation”

$F_2$

$C_3$

Notice that {A, B, C} and {A, B, E} could never belong to “frequent itemset” because {A, B} was already excluded.  
 The same for {A, C, E} and {A, E}.

## Details on how to generate “frequent itemset” candidates

- An important **assumption** – the “**ordering of items**”
  - meaning that **items** are stored in **lexicographical order**
  - ... which is a **total order relation**
- The items’ ordering is used by the algorithm in each item set  
 $w = ( w[1], w[2], \dots, w[k] )$  is a tuple representing a k-itemset  $w$ ,  
where  $w[1] < w[2] < \dots < w[k]$  according to the total order
- ... for example, in the 3-itemset  $w = ( B, C, E )$  we have
  - we have  $w[1]=B < w[2]=C < w[3]=E$
  - according to the lexicographical total ordering of items

... how to generate “frequent itemset” candidates?

The “**generate-candidate**” function:

**input:**  $F_{k-1}$ ; a  $(k-1)$ -itemset

**output:**  $C_k$ ; a superset, called “candidates”, of the set of all frequent  $k$ -itemsets.

The “**generate-candidate**” function has 2 steps.

- **step1: self-joining**
  - generate all possible candidate itemsets  $C_k$  of length  $k$ .
- **step2: pruning**
  - remove all candidates  $C_k$  that can't be frequent (i.e., are non-frequent)
  - ... a superset of a non-frequent itemset is also non-frequent!

... the “generate-candidate” function

```
function generate-candidate(  $F_{k-1}$  )  
   $C_k \leftarrow \emptyset$ ;  
  forall  $f_1, f_2 \in F_{k-1}$   
    with  $f_1 = \{i_1, \dots, i_{k-2}, i_{k-1}\}$   
    and  $f_2 = \{i_1, \dots, i_{k-2}, i'_{k-1}\}$   
    and  $i_{k-1} < i'_{k-1}$  do  
       $c \leftarrow \{i_1, \dots, i_{k-1}, i'_{k-1}\}$ ;           // join the  $f_1$  and  $f_2$  itemsets  
       $C_k \leftarrow C_k \cup \{c\}$ ;  
      for each  $(k-1)$ -subset  $s$  of  $c$  do  
        if (  $s \notin F_{k-1}$  ) then  
          delete  $c$  from  $C_k$ ;           // prune superset,  $c$ , of non-frequent  $s$   
        end  
      end  
  end  
  return  $C_k$ ;
```

## ... example of “generate-candidate” function

The “**generate-candidate**” function:

**input:**  $F_{k-1}$ ; a  $(k-1)$ -itemset

**output:**  $C_k$ ; a superset, called “candidates”, of the set of all frequent  $k$ -itemsets.

Apply the “**generate-candidate**” function to  $F_3$  and generate  $C_4$ .

**input:**  $F_{4-1} = \{ \{1, 2, 3\}, \{1, 2, 4\}, \{1, 3, 4\}, \{1, 3, 5\}, \{2, 3, 4\} \}$

**output:**  $C_4 = ?$

## ... example (“generate-candidate”)

**input:**  $F_{4-1} = \{ \{1, 2, 3\}, \{1, 2, 4\}, \{1, 3, 4\}, \{1, 3, 5\}, \{2, 3, 5\} \}$   
**output:**  $C_4 = ?$

```

function generate-candidate(  $F_{k-1}$  )
     $C_k \leftarrow \emptyset$ ;
    forall  $f_1, f_2 \in F_{k-1}$ 
        with  $f_1 = \{i_1, \dots, i_{k-2}, i_{k-1}\}$ 
        and  $f_2 = \{i_1, \dots, i_{k-2}, i'_{k-1}\}$ 
        and  $i_{k-1} < i'_{k-1}$  do
             $c \leftarrow \{i_1, \dots, i_{k-1}, i'_{k-1}\}$ ;    // join  $f_1$  and  $f_2$ 
             $C_k \leftarrow C_k \cup \{c\}$ ;
            for each  $(k-1)$ -subset  $s$  of  $c$  do
                if (  $s \notin F_{k-1}$  ) then
                    delete  $c$  from  $C_k$ ;    // prune
            end
        end
    return  $C_k$ ;
    
```

**input**

$F_{4-1} = \{$   
 $\{1, 2, 3\},$   
 $\{1, 2, 4\},$   
 $\{1, 3, 4\},$   
 $\{1, 3, 5\},$   
 $\{2, 3, 4\} \}$

**after join**

$C_4 = \{$   
 $\{1, 2, 3, 4\},$   
 $\{1, 3, 4, 5\} \}$

**after prune**

$C_4 = \{$   
 $\{1, 2, 3, 4\} \}$

$\{1, 3, 4, 5\}$  is removed  
 because  $\{1, 4, 5\}$  is not in  $F_3$ .



... how to “generate-rules” from the “frequent itemsets”?

**“frequent itemsets”  $\neq$  “association rules”**

we need **one more step** to generate **association rules**

**for each** “frequent itemset”  $X$

**for each** proper nonempty subset  $A$  of  $X$  (i.e.,  $\forall A \subset X$  and  $X \neq \emptyset$ )

let  $B = X - A$

**if** confidence(  $A \Rightarrow B$  )  $\geq$  **min\_conf**

**then**  $A \Rightarrow B$  is an association rule

*recall that*

**support(  $A \Rightarrow B$  ) = support(  $A \cup B$  ) = support(  $X$  )**

**confidence(  $A \Rightarrow B$  ) = support(  $A \cup B$  ) / support(  $A$  )**

... example of “generate-rules”

**suppose {2, 3, 4} is a “frequent itemset”  
with support = 50%**

**Which rules can be generated from that itemset  
and what is the confidence values of each rule?**

*for the purpose of this example assume that:*

- each 1-itemset has confidence=80%
- each 2-itemset has confidence=50%

... example (“generate-rules”)

**suppose {2, 3, 4} is a “frequent itemset”  
with support = 50%**

Proper non-empty subsets and  
*assumed* confidence values:

{2}	confidence=80%
{3}	confidence=80%
{4}	confidence=80%
{2, 3}	confidence=50%
{2, 4}	confidence=50%
{3, 4}	confidence=50%

**The list of the corresponding  
association rules:**

$2 \Rightarrow 3,4$	confidence=50%/80%=62.5%
$3 \Rightarrow 2,4$	confidence=50%/80%=62.5%
$4 \Rightarrow 2,3$	confidence=50%/80%=62.5%
$2,3 \Rightarrow 4$	confidence=50%/50%=100%
$2,4 \Rightarrow 3$	confidence=50%/50%=100%
$3,4 \Rightarrow 2$	confidence=50%/50%=100%

## The “generate-rules” process – a synthesis

- In order to generate  $A \Rightarrow B$  we need to have
  - the value of  $\text{confidence}(A \Rightarrow B)$ , that can be computed
  - ... from the values of  $\text{support}(A \cup B)$  and  $\text{support}(A)$
- So, all the required information for confidence computation
  - was already recorded during itemset “generate-candidate” process
  - ... no need to scan the data anymore for “generate-rule” process
- ... the “generate-rules” process
  - is less time consuming than the itemset “generate-candidates” process

... example – the “generate-rule” process

*the “generate-rule” process:*

**for each** “frequent itemset”  $X$

**for each** proper nonempty subset  $A$  of  $X$  (i.e.,  $\forall A \subset X$  and  $X \neq \emptyset$ )

let  $B = X - A$

**if** confidence(  $A \Rightarrow B$  )  $\geq$  **min\_conf**

**then**  $A \Rightarrow B$  is an association rule

**If {A, B, C, D} is a “frequent itemset” then  
what is the entire set of candidate rules?**

**Apply the above “generate-rule” process**

... example – the result of the “generate-rule” process

**for each** “frequent itemset”  $X$

**for each** proper nonempty subset  $A$  of  $X$  (i.e.,  $\forall A \subset X$  and  $X \neq \emptyset$ )

let  $B = X - A$

**if** confidence( $A \Rightarrow B$ )  $\geq$  **min\_conf**

**then**  $A \Rightarrow B$  is an association rule

If  $\{A, B, C, D\}$  is a “frequent itemset” then the candidate rules are:

$A, B, C \Rightarrow D$        $A, B, D \Rightarrow C$        $A, C, D \Rightarrow B$        $B, C, D \Rightarrow A$

$A, B \Rightarrow C, D$        $A, C \Rightarrow B, D$        $A, D \Rightarrow B, C$

$B, C \Rightarrow A, D$        $B, D \Rightarrow A, C$

$C, D \Rightarrow A, B$

$A \Rightarrow B, C, D$        $B \Rightarrow A, C, D$        $C \Rightarrow A, B, D$        $D \Rightarrow A, B, C$

If  $|X| = k$ , then there are  $2^k - 2$  candidate association rules  
(ignoring  $\emptyset \Rightarrow X$  and  $X \Rightarrow \emptyset$ )

... “confidence” and the anti-monotone property

In general “confidence” does not have an anti-monotone property:  
e.g.,  $\text{confidence}(ABC \Rightarrow D)$  can be higher or lower than  $\text{confidence}(AB \Rightarrow DE)$

But, “confidence” of rules generated from the same itemset,  $X$ ,  
have an anti-monotone property:

e.g.,  $X = \{A, B, C, D\}$

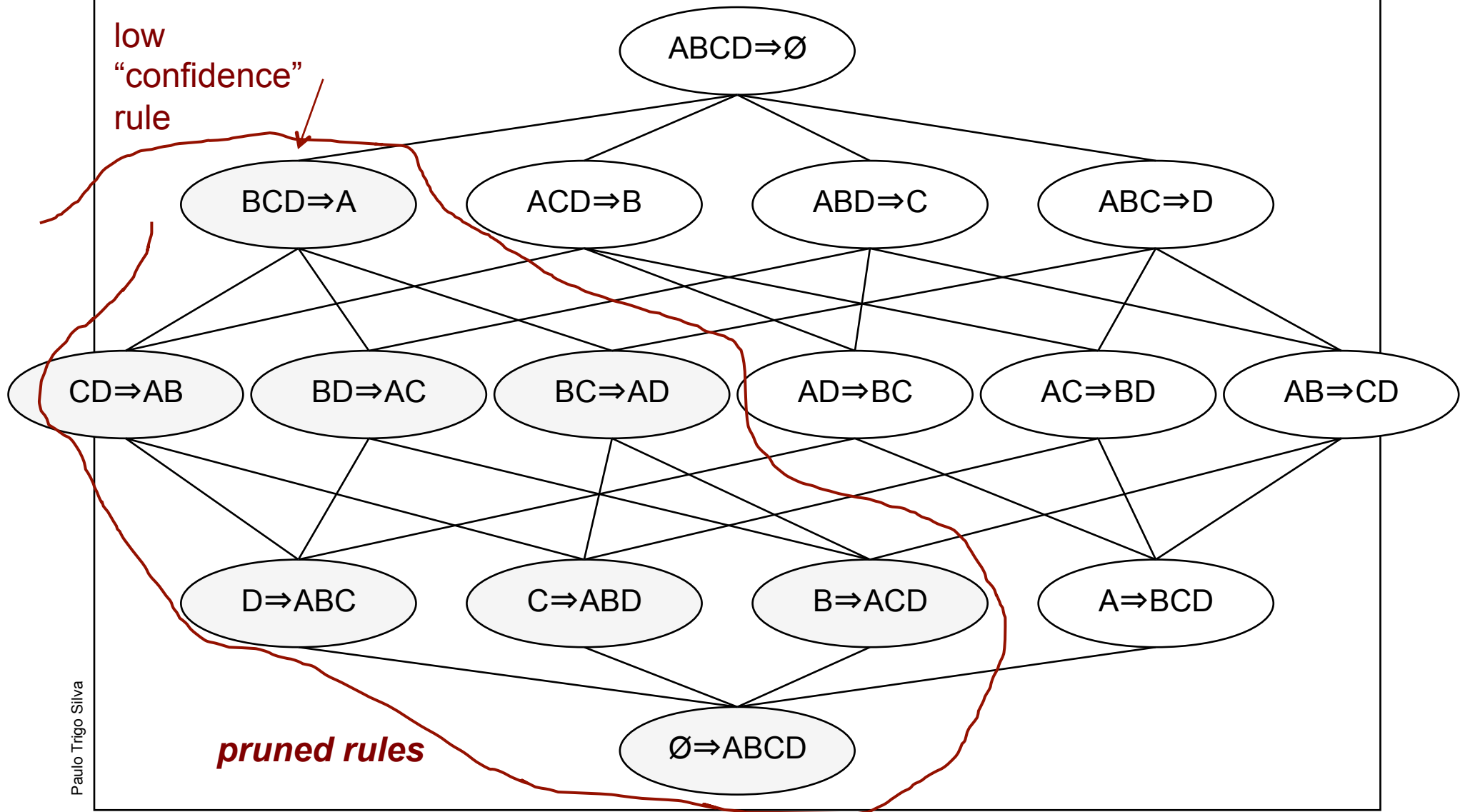
$\text{confidence}(ABC \Rightarrow \mathbf{D}) \geq \text{confidence}(AB \Rightarrow \mathbf{CD}) \geq \text{confidence}(A \Rightarrow \mathbf{BCD})$

***The “confidence” anti-monotone property:***

**“confidence” is anti-monotone with respect to the number of items  
on the RHS of the rule.**

**i.e., increase RHS dimension decreases, or maintains, its confidence**

... the “generate-rule” space





## A more efficient way to “generate-rule”

- A candidate rule is generated by merging two rules
  - that share the same prefix in the rule consequent
  - and have a non disjoint antecedent
  - ... i.e., **same prefix in RHS** of both rules
  - ... and **intersection of both LHS is not empty**
- ...  $\text{join}( CD \Rightarrow AB, BD \Rightarrow AC )$  would produce the candidate rule:  
 $D \Rightarrow ABC$   
(LHS with common items of both rules; RHS includes both last items)
- ... if rule  $D \Rightarrow ABC$  does not have the required confidence
  - then prune this rule
- recall that support counts have already been computed
  - during the “frequent itemset” generation step

## Characteristics of the APRIORI method and an extension

- Uses the same minimum support constraint for all items
  - i.e., a single `min_sup` threshold value
  - in many scenarios some items appear more frequently than others
  - e.g., people buy “cooking pan” and “coffee machine” much less frequently than they buy bread and milk!
- If frequency of items vary a lot, then
  - if `min_sup` is set too high, rules involving rare items will not be found
  - if `min_sup` is set too low, may cause combinatorial explosion as those frequent items will be associated with each other in all possible ways
- The original algorithm can be extended to
  - deal with “multiple `min_sup` thresholds”
  - ... where each item can have a “minimum item support” (MIS)
  - ... the APRIORI principle must consider the total order relation of MIS