



		Agenda
•	Introduction	
•	Concept	
•	Examples	
•	Algorithms	
•	Temporal Association Rules	
•	Software	
	• WEKA	
	ARTool	► <i>m</i> ~2
•	Conclusions	
Taller o	le Sistemas Multiagentes - 2002	Mg. Silvia Schiaffino

#### Introduction

- Data Mining is the process of finding interesting **trends** or **patterns** in large datasets in order to guide future decisions.
- Data Mining is the discovery of knowledge and useful information from the large amounts of data stored in databases.
- Related to exploratory data analysis (area of statistics) and knowledge discovery (area in artificial intelligence, machine learning).
- Data Mining is characterized by having VERY LARGE
   datasets.

Taller de Sistemas Multiagentes - 2002

#### Some real problems

- Associations between products bought in a store
  - milk

•



- beers and diapersDeciding on product discounts
- Figuring out how to keep customer at lowest cost to company
- Personal recommendation page at Amazon
- Stock market trends

Association Rules: describing association relationships among the attributes in the set of relevant data.

Taller de Sistemas Multiagentes - 2002







Goal Find association rules with high support and high confidence. Support and confidence values are specified by the user. Remember: Finding such a rule does not mean that there must be a relationship between the left and right sides. A person must evaluate such rules by hand. Example: {milk} → {bread}

 Formal statement of the problem

  $I = \{i_1, i_2, \dots, i_m\}$  is a set of items

 D is a set of transactions T 

 Each transaction T is a set of items (subset of I)

 TID is a unique identifier that is associated with each transaction

 The problem is to generate all association rules that have support and confidence greater than the user-specified minimum support and minimum confidence

Taller de Sistemas Multiagentes - 2002

Taller de Sistemas Multiagentes - 2002

Mg. Silvia Schiaffino

#### Problem decomposition

Mg. Silvia Schiaffino

The problem can be decomposed into two subproblems:

- Find all sets of items (*itemsets*) that have support (number of transactions) greater than the minimum support (*large itemsets*).
- 2. Use the *large itemsets* to generate the desired rules. For each *large itemset i*, find all non-empty subsets, and for each subset *a* generate a rule a ==> (I-a) if its confidence is greater than the minimum confidence.

Taller de Sistemas Multiagentes - 2002

	Notatior
k-itemset: itemset containining k items	
L <sub>k</sub> : set of frequent (large) k-itemsets. Eac itemset; ii) support count	h member has to fields: i)
$C_k$ : set of candidate k-itemsets (potencia has to fields: i) itemset; ii) support co	Illy frequent. Each member unt
Taller de Sistemas Multiagentes - 2002	Mg. Silvia Schiaffino

#### General Algorithm

- In the first pass, the support of each individual item is counted, and the *large* ones are determined
- In each subsequent pass, the *large* itemsets determined in the previous pass is used to generate new itemsets called *candidate* itemsets.
- The support of each *candidate* itemset is counted, and the *large* ones are determined.
- 4. This process continues until no new *large* itemsets are found.

Taller de Sistemas Multiagentes - 2002



Example Database  $C_2$  $L_1$ TID Items Itemset Support Itemset Support 100 134 {1} 2 {1 3}\* 2 200 235 {2} 3 {1 4} 1 300 1235 {3} 3 {3 4} 1 400 25 {5} {2 3}\* 2 3 C. {2 5}\* 3 Itemset Support {3 5}\* 2 {1 3 4} 1 {1 2} 1 {2 3 5}\* 2 {1 5} 1 {135} 1 Taller de Sistemas Multiagentes - 2002 Mg. Silvia Schiaffino

#### SETM Algorithm

Mg. Silvia Schiaffino

Candidate itemsets are generated on-the-fly as the database is scanned, but counted at the end of the pass.

- New candidate itemsets are generated the same way as in AIS 1. algorithm, but the TID of the generating transaction is saved with the candidate itemset in a sequential structure.
- At the end of the pass, the support count of *candidate* itemsets is 2. determined by aggregating this sequential structure

It has the same disadvantage of the AIS algorithm.

Another disadvantage is that for each candidate itemset, there are as many entries as its support value

Taller de Sistemas Multiagentes - 2002

								Ex	ample
Dat	abase			L	-1		C <sub>2</sub>		
TID	Items	;	Item	set	Support	1	Itemset	TID	
100	134		{1	}	2		{1 3}	100	
200	235		12	3	3		{1 4}	100	
200	4.0.0	_		, ,	-		{3 4}	100	
300	123	<u>^</u>	{3	3	3	-	{2 3}	200	
400	2 5	പറ	{5	}	3		{2 5}	200	
		<b>C</b> <sub>3</sub>					{3 5}	200	
	lte	emset	TID				{1 2}	300	
	{1 3 4}		100				{1 3}	300	
	(104)						{1 5}	300	
	{2 3 5}		200				{2 3}	300	
	{	{1 3 5}					{2 5}	300	
	{:	2 3 5}	300				{3 5}	300	
er de Sistem	as Multiage	ntes - 2002					{2 5} Mg	400 Silvia S	chiaffino











Example **C**<sub>2</sub> Database  $L_1$ Support Itemset Itemset Support TID Items {1 2} 100 {1} 2 134 {1 3}\* 2 200 235 {2} 3 {1 5} 1 300 1235 {3} 3 {2 3}\* 2 400 25 {5} 3 {2 5}\* 3 C' C'2 {3 5}\* 100 200 {2 3 5} {1 3} 300 {2 3 5} 200  $\{2\;3\},\,\{2\;5\},\,\{3\;5\}$ C<sub>2</sub> 300 {1 2}, {1 3}, {1 5}, {2 3}, Itemset Support {2 5}, {3 5} {2 3 5}\* 2 400 {2.5} Taller de Sistemas Multiagentes - 2002 Mg. Silvia Schiaffino

#### Apriori Hybrid

Performance Analysis demonstrated that:

- · Apriori does better than AprioriTid in the earlier passes.
- AprioriTid does better than Apriori in the later passes.

Hence, a hybrid algorithm can be designed that uses Apriori in the initial passes and switches to AprioriTid when it expects that the set C' will fit in memory.

Taller de Sistemas Multiagentes - 2002





Mg. Silvia Schiaffino



Taller de Sistemas Multiagentes - 2002

















#### Several Refinements...

- Improved algorithms for the discovery of large itemsets
- Improved algorithms for the discovery of generalized
   and quantitative association rules
- New measures for other types of data (different from basket data)
- Temporal considerations

Taller de Sistemas Multiagentes - 2002

# In the calculus of support the denominator represents the total number of transactions in a time period when the involved items may have not existed. For example, a supermarket may start selling new products. Some itemsets may become obsolote, and thus,

Some itemsets may become obsolote, and thu some rules may become outdated.



Mg. Silvia Schiaffino

Taller de Sistemas Multiagentes - 2002

#### Solution

Consider the lifetime of an item, this being the period between the first and the last time the items appears in transactions in the database.

The basic idea is to limit the search for frequent itemsets to the lifetime of the itemset's members.

Each rule has an associated time frame, corresponding to the lifetime of the items participating in the rule.

Taller de Sistemas Multiagentes - 2002





#### **Temporal Association Rule**

A temporal association rule is an expression of the form  $X \rightarrow Y$  [t1, t2], where [t1, t2] is a time frame corresponding to the lifespan of  $X \cup Y$  expressed in a granularity determined by the user.

A temporal association rule has three factors associated with it: support, temporal support and confidence

(support and confidence have to be redefined in this model!!)

Taller de Sistemas Multiagentes - 2002

### Cyclic Association Rules

Association rules that exist in certain time intervals and thus display cyclic variations over time.

Example: If we compute association rules over monthly sales data, we may observe seasonal variation where certain rules are true at approximately the same month each year. Similarly, association rules can also display regular hourly, daily, weekly, etc. Variation that is cyclical in nature.



Taller de Sistemas Multiagentes - 2002

	[Ozden98] Approach
	It is assumed that transactional data is timestamped and that time intervals are specified by te user to divide the data into disjoint segments.
	Cyclical Association Rule: if the ruel has the minimum confidence and support at regular time intervals
	Mining Association Rules Problem: the generation of all the cycles of an association rule
alle	er de Sistemas Multiagentes - 2002 Mg. Silvia Schiaffino





Construction into addising framework     Construction into addising framework     Construction into addising framework     Construction	A Construction for a finite structure of the second st		WE
No.         No. <th>An and a second second</th> <th>An An America (a La Seconda Se</th> <th>Aller Martinester</th>	An and a second	An An America (a La Seconda Se	Aller Martinester
			Annu Annu Annu Annu Annu Annu Annu Annu















	WEKA: Generate	ed Rule
Minimum support: 0.05		
Minimum metric <confidence>: 0.9</confidence>		
Number of cycles performed: 17		
Generated sets of large itemsets:		
Size of set of large itemsets L(1): 1	2	
Size of set of large itemsets L(2): 4	7	
Size of set of large itemsets L(3): 3	9	
Size of set of large itemsets L(4): 6		
Best rules found:		
1. humidity=normal windy=FALSE	4 ==> play=yes 4 conf:(1)	
<ol><li>temperature=cool 4 ==&gt; humidit</li></ol>	ty=normal 4 conf:(1)	
<ol><li>outlook=overcast 4 ==&gt; play=ye</li></ol>	es 4 conf:(1)	
<ol><li>temperature=cool play=yes 3 ==</li></ol>	=> humidity=normal 3 conf:(1)	
<ol><li>outlook=rainy windy=FALSE 3 =</li></ol>	=> play=yes 3 conf:(1)	
<ol><li>outlook=rainy play=yes 3 ==&gt; w</li></ol>	indy=FALSE 3 conf:(1)	
<ol><li>outlook=sunny humidity=high 3</li></ol>	==> play=no 3 conf:(1)	
8. outlook=sunny play=no 3 ==> hi	umidity=high 3 conf:(1)	
<ol><li>temperature=cool windy=FALSE</li></ol>	2 ==> humidity=normal play=yes 2	cont:(1)
10. temperature=cool humidity=nor	mal windy=FALSE 2 ==> play=yes 2	conf:(1)
aller de Sistemas Multiagentes - 2002	Mg. S	ilvia Schiaffino

	Software: ARTool
ARtool represents a collection of	of algorithms and tools
for the mining of association databases.	rules in binary
http://www.cs.umb.edu/~laur/AF	RTool

Mg. Silvia Schiaffino

		ARTool
Differe	ent algorithms for generating frequent it	emsets:
•	Apriori	
•	Closure	
•	ClosureOpt	
•	FPGrowth	
Differe	ent algorithms for rules generation:	
•	Apriori Rules	
•	Cover Rules	
•	Cover Rules Opt.	
uller de Sistem	as Multiagentes - 2002	Mg. Silvia Schiaffino











Note:         Note: <th< th=""><th></th><th></th><th></th><th>AR</th><th>Tool</th></th<>				AR	Tool
Note         Note <th< th=""><th>Property lines</th><th></th><th></th><th></th><th></th></th<>	Property lines				
Name         Name         Dependence         Dependence           Image: Space of the state of the	States I want have \$	A CONTRACTOR OF A CONTRACTOR OFTA CONT			
Provide a set of the set of	territory land -	Brief Brief	152.530.5205	hep:	
Production Conference of the Conference of	Rener report int	ni kanonini kali kanonini hydroxechini shina adal gi doversitini shina adal gi doversitini shina adal doversiti shina adal kali shina atali aji shina adal doversitini			
And the set of table and the description of the set of					
Land triblen i 1996 on 50% for the problem in the second s	whether is the Bernstein Souther				
ize disputer to a close exclusion or classes (Classes)/Classes/Allowerski konstituers appellist, somi ter etamologica (AD) energi energ	and a construction of the second s	a di Cangda languin con alt Na panala d'Angela dalta			
in the sector with the date of the sector of	lan dina dipertita dan din propis Tana dipensitiwa (1963)	el en la banne d'Alabamet d'Analise génération de ser de la companya de la companya de la companya de la compa	inge oppeliji, ind		
Navalla sana kangi Najadagi yana da sana	landing succession to the set				
	and a second				
	and were the				







	ARTool
tera tren ten	
BALANCE PERSONNELLA INCOMENTATIONS	PROVING OCCUPANT Name Officers
Resentances of Deformation Section of Controller Trademating Respective Into No. Spectral Trademating Respective Into No. Sciences Trademating Respective Into No. Sciences Into No. 2017 Section 2017 Section 2017 Section 2017 Section 2017 Section 2017 Section 2017 Section 2017 Respective Respective Into No. Respective	i forminger gant fill, store
Taller de Sistemas Multiagentes - 2002	Mg. Silvia Schiaffino



	ARTool
tana tan	
Table Freed East Annual Lances and Second Annu	Minute         Operating         Supple         Operating
nanovala Associa Asso International Statu Manganata Markati Ang Manganata Associa Statu Manganata Associa Statu	nakaanah kendinan sagarti Laat nidu asambara 10. door







## Bibliography

- R. Agrawal, T. Imielinski, A. Swami: "Mining Associations between Sets of Items in Massive Databases", Proc. of the ACM-SIGMOD 1993 Int'l Conference on Management of Data, Washington D.C., May 1993, 207-216.
- R. Agrawal, R. Srikant: "Fast Algorithms for Mining Association Rules", Proc. of the 20th Int'l Conference on Very Large Databases, Santiago, Chile, Sept. 1994. Expanded version available as IBM Research Report RJ9839, June 1994
- Yingjiu Li, Peng Ning, X. Sean Wang, Sushil Jajodia, "Discovering Calendarbased Temporal Association Rules," (Full version) in Proceedings of the 8th International Symposium on Temporal Representation and Reasoning (TIME 01), pages 111-118, Italy, June 2001.

Taller de Sistemas Multiagentes - 2002