

#### **DATA MINING**

# LES RÈGLES ASSOCIATIVES

**Mohamed Heny SELMI** 

medheny.selmi@esprit.tn

#### **OBJECTIFS**



- ✓ Rechercher les associations consiste à rechercher les règles de type :
   « Si pour un individu, la variable A = Xa, la variable B = Xb, etc, alors, dans 80% des cas, la variable Z = Xz, cette configuration se rencontrant pour 30 % des individus »
- ✓ Repérer des règles liant les données avec un bon niveau de probabilité
  - découverte de relations fines entre attributs (ou variables)
  - généralisation des dépendances fonctionnelles
- ✓ Mettre en évidence les produits / des articles achetés ensemble
- ✓ Transcrire la connaissance sous forme de règles d'association
- ✓ Règles du style : < si [P(tid,X) := prémisse] alors [P(tid,Y) := conséquence] >
- ✓ Différents types de règles
  - origine « panier de la ménagère »
  - étendues aux tables multiples et aux attributs continus



#### **ANALYSE DES TICKETS DE CAISSE**

N° Transaction (Caddie)		Contenu du	caddie	
1	Poulet	Moutarde	Œufs	Pates
2	Moutarde	Œufs		
3	Pain	Beurre	Poulet	
4	Pates			
5	Pain	Lait	Beurre	
6	Œufs	Pain		
7	Confiture			



Une observation = un caddie

Ne tenir compte que de la présence des produits : peu importe leur quantité

Dans un caddie : le nombre de produits est variables

La liste des produits est immense et variable



### **TABLEAU DES TRANSACTIONS**

- ✓ Mettre en évidence les produits / des articles achetés ensemble.
- ✓ Transcrire la connaissance sous forme de règles d'association

 $\underline{si}[P(tid,X) := prémisse] \underline{alors}[P(tid,Y) := conséquence]$ 

N° Transaction (Caddie)		Contenu du	ı caddie	
1	Poulet	Moutarde	Œufs	Pates

#### <u>si</u> Poulet <u>et</u> Moutarde <u>alors</u> Œufs <u>et</u> Pates

N° Transaction (Caddie)		Contenu du caddie
6	Œufs	Pain

si Œufs alors Pain

## **TABLEAU BINAIRE**



N° Transaction (Caddie)		Contenu du	ı caddie	
1	Poulet	Moutarde	Œufs	Pates
2	Moutarde	Œufs		
3	Pain	Beurre	Poulet	
4	Pates			
5	Pain	Lait	Beurre	
6	Œufs	Pain		
7	Confiture			

	P1	P2	Р3	P4	P5	P6	P7	P8
1	1	1	1	1	0	0	0	0
2	0	1	1	0	0	0	0	0
3	1	0	0	0	1	1	0	0
4	0	0	0	1	0	0	0	0
5	0	0	0	0	1	1	1	0
6	0	0	1	0	1	0	0	0
7	0	0	0	0	0	0	0	1

désignation
P1 = Poulet
P2 = Moutarde
P3 = Œufs
P4 = Pates
P5 = Pain
P6 = Beurre
P7 = Lait
P8 = Confiture



#### **CODAGE DISJONCTIF COMPLET**

Observation	Taille	Corpulence
1	Petit	Mince
2	Grand	Enveloppé
3	Grand	Mince



Observation	Taille = Petit	Taille = Grand	Corpulence = Mince	Corpulence = Enveloppé
1	1	0	1	0
2	0	1	0	1
3	0	1	1	0



Dès que l'on peut se ramener à des données o/1:

Il est possible de construire des règles d'association

### PASSAGE EN FORME DISJONCTIVE COMPLÈTE



Catégoriel, qualitatif, discret: type marché, entreprises, taux, appartenance, ...

	Marché		Р	Ε	Α
C <sub>1</sub>	Part.	C <sub>1</sub>	1		
C <sub>2</sub>	Autre	C <sub>2</sub>			1
C <sub>3</sub>	Part.	C <sub>3</sub>	1		
<b>C</b> <sub>4</sub>	Part.	C <sub>4</sub>	1		
<b>C</b> <sub>5</sub>	Entr.	<b>C</b> <sub>5</sub>		1	

• Continu, quantitatif: virement, âge, température, consommation, pourcentage, ...

	Dom		d0	d1	d2	d3	d4
<b>c</b> <sub>1</sub>	1100	<b>C</b> <sub>1</sub>			1		
C <sub>2</sub>	0	C <sub>2</sub>	1				
<b>C</b> <sub>3</sub>	2200	<b>C</b> <sub>3</sub>				1	
C <sub>4</sub>	800	<b>C</b> <sub>4</sub>		1			
<b>C</b> <sub>5</sub>	3800	<b>C</b> <sub>5</sub>					1

## CRITÈRES D'ÉVALUATION DES RÈGLES D'ASSOCIATION



#### **SUPPORT**

#### CONFIANCE

- √ indicateur de « fiabilité »
- ✓ probabilité absolue :

 $P(X \cup Y)$ 

Règle d'association :  $p1 \rightarrow p2$ 

✓ ||X U Y||/ ||BD|| = % de transactions vérifiant la règle

- ✓ Indicateur de « précision »
- ✓ probabilité conditionnelle : P(Y/X)
- ✓ ||X U Y||/||X|| = % de transactions vérifiant l'implication

$$sup(R1) = 2 : en termes absolus$$

ou sup(R1) = 2 / 6 = 33%: en termes relatifs

Conf(R1) = 
$$sup(R1) / sup(antécédant R1)$$
  
=  $sup(p1 \rightarrow p2) / sup(p1) = 2 / 4 = 50 %$ 

Caddie	р1	p2	р3	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1



« Bonne » règle = règle avec un support et une confiance élevée



#### **ANALYSE DES TICKETS**

{ "crème" } → { "pain" }



ID	PRODUITS
1	pain, crème, eau
2	crème
3	pain, crème, vin
4	eau
5	crème, eau

Support = Prob. (crème et pain):

$$Sup = \frac{\text{nom(tran.contenant crème et pain)}}{\text{nom\_total(tran.)}} = \frac{2}{5} = 0.4$$

Confiance = Prob(crème et pain / crème):

Conf = 
$$\frac{\text{nom(tran. contenant crème et pain)}}{\text{nom(tran. contenant crème)}} = \frac{2}{4} = 0.5 = \frac{\text{sup(crème et pain)}}{\text{sup(crème)}}$$

## DÉMARCHE D'EXTRACTION DES RÈGLES D'ASSOCIATION



Paramètres: Fixer un degré d'exigence sur les règles à extraire

- ✓ Support min. (exp. 2 transactions)
- ✓ Confiance min. (exp. 75%)



L'idée est surtout de contrôler (limiter) le nombre de règles produites

Démarche: Construction en deux temps

- ✓ recherche des itemsets fréquents (support >= support min.)
- √ à partir des itemsets fréquents, produire les règles (conf. >= conf. min.)

Caddie	p1	р2	p3	р4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1

#### **Quelques définitions:**

- item = produit
- itemset = ensemble de produits (ex. {p1,p3})
- sup(itemset) = nombre de transactions d'apparition simultanée des produits (ex. sup{p1,p3} = 4)
- card(itemset) = nombre de produits dans l'ensemble (ex. card{p1,p3} = 2)



#### **ALGORITHME APRIORI [AGRAWAL94]**

#### Première passe:

- ✓ recherche des 1-ensembles fréquents
- ✓ un compteur par produits

# L'algorithme génère un candidat de taille k à partir de deux candidats de taille k-1 différents par le dernier élément

✓ procédure apriori-gen

#### Passe k:

- ✓ comptage des k-ensemble fréquents candidats
- ✓ sélection des bons candidats



## **APRIORI - FRÉQUENTS ITEMSETS**

```
L1 = { frequent 1-ensemble };
for (k = 2 ; Lk-1 \neq \emptyset ; k++) do
       Ck = apriori-gen(Lk-1); // Generate new candidates
       foreach transactions t ∈ DB do
               {// Counting
               Ct = { subset(Ck, t) }; // get subsets of t candidates
               foreach c \in Ct do c.count++;
       Lk = \{c \in Ck \mid c.count >= minsup \}; // Filter candidates
F = \{Lk\};
```



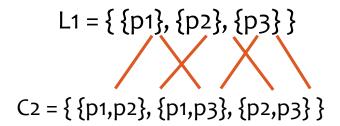


On va prendre la valeur du Support minimal = 3

Caddie	p1 p2		рЗ	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1



1-Itemsets	Support
{p1}	4
{p2}	3
{p3}	5
{p4}	1







On va prendre le Support minimal = 3

Caddie	р1	p2	рЗ	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1



2-Itemsets	Support
{p1,p2}	2
{p1,p3}	4
{p2,p3}	3





On va prendre le Support minimal = 3

Caddie	p1 p2		p3	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1



3-Itemsets	Support
{p1,p2, p3}	2

$$L_3 = \emptyset$$

$$F = \{ \{p1\}, \{p2\}, \{p3\}, \{p1,p3\}, \{p2,p3\} \}$$





On va prendre le pourcentage de la confiance minimale = 65%

Caddie	р1	p2	рЗ	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	0	1	0
5	0	1	1	0
6	0	0	0	1

p1
$$\rightarrow$$
p3: confiance = 4/4 = **100** % p3 $\rightarrow$ p1: confiance = 4/5 = **80** %

p3
$$\rightarrow$$
p1: confiance = 4/5 = **80** %



$$p2 \rightarrow p3$$
: confiance =  $3/3 = 100 \%$ 

p2
$$\to$$
p3: confiance = 3/3 = **100** % p3 $\to$ p2: confiance = 3/5 = **60** %



## INDICATEUR DE PERTINENCE DES RÈGLES MESURE D'INTÉRÊT : LIFT D'UNE RÈGLE

L'amélioration apportée par une règle, par rapport à une réponse au hasard est appelée « lift » et vaut :

- Le lift est une bonne mesure de performance de la règle d'association.
- Le lift est la confiance de la règle divisée par la valeur espérée de la confiance.

#### Interprétation du lift :

- Un lift supérieur à 1 : Indique une corrélation positive
- Un lift de 1 indique une corrélation nulle
- Un lift inférieur à 1 : Indique une corrélation négative

T26	A	В	С	D	E
T1245	В	С	E	F	
T156	В	E			
T2356	A	В	D		
T145	С	D			

lift 
$$(C \rightarrow B) = 5/6 < 1$$

lift 
$$(B \rightarrow E) = 6/5 > 1$$



#### ETUDE DE CAS DE RECHERCHE D'ASSOCIATIONS INTÉRESSANTES



- ✓ Le principe de l'algorithme est de rechercher l'ensemble L₁ de tous les items
  - apparaissant dans au moins **S**<sub>min</sub> **x m** transactions.
- ✓ Puis, parmi C₂ qui est le produit cartésien de L₁ avec lui-même, on construit
  - l'ensemble L<sub>2</sub> de tous les couples d'items apparaissant dans au moins
  - **S**<sub>min</sub> **x m** transactions.
- ✓ L'algorithme s'arrête quand L<sub>k</sub> est vide.

#### UTILITÉ DES RÈGLES D'ASSOCIATION









La course à la fidélisation des clients



Réductions personnalisées à la caisse



**Profil-client** 



Le test des nouveaux produits



Le panier moyen



Le parcours magasin



Cartes de fidélité