

Règles d'association

Sommaire

Exemple : Panier de la ménagère

Définitions

A-Priori

Algorithmes génétiques

Résumé

L. Jourdan – Aide à la décision

Exemple : Analyse du panier de la ménagère

- Découverte d'**associations** et de **corrélations** entre les articles achetés par les clients en analysant les achats effectués (panier)

Lait, Oeufs, Céréale, Lait



Client 2

Lait, Oeufs, Sucre, Pain



Client 1

Oeufs, Sucre



Client 3

L. Jourdan – Aide à la décision

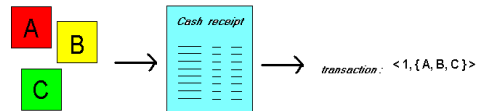
Exemple : Analyse du panier de la ménagère

Etant donnée :

- Une base de données de **transactions** de clients, où chaque transaction est représentée par un ensemble d'articles -**set of items**- (ex., produits)

Trouver :

- Groupes d'articles (itemset) achetés **fréquemment** (ensemble)



L. Jourdan – Aide à la décision

Exemple : Analyse du panier de la ménagère

Extraction d'informations sur le comportement de clients

- Si achat de riz + vin blanc ALORS achat de poisson (avec une grande probabilité)

Intérêt de l'information : peut suggérer ...

- Disposition des produits dans le magasin
- Quels produits mettre en promotion, gestion de stock, ...

Approche applicable dans d'autres domaines

- Cartes de crédit, e-commerce, ...
- Services des compagnies de télécommunication
- Services bancaires
- Traitements médicaux, ...

L. Jourdan – Aide à la décision

Règles d'associations

Recherche de règles d'association :

- Découvrir des patterns, corrélations, associations fréquentes, à partir d'ensembles d'items contenus dans des base de données.

Compréhensibles : Facile à comprendre

Utiles : Aide à la décision

Efficaces : Algorithmes de recherche

Applications :

- Analyse des achats de clients, Marketing, Accès Web, Design de catalogue, Génomique, etc.

L. Jourdan – Aide à la décision

Règles d'associations

Formats de représentation des règles d'association :

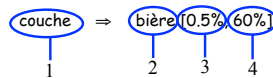
- couches \Rightarrow bière [0.5%, 60%]
- achète:couches \Rightarrow achète:bière [0.5%, 60%]
- "SI achète couches ALORS achète bière dans 60% de cas. Les couches et la bière sont tous deux achetés dans 0.5% des transactions de la base de données."

Autres représentations (utilisée dans l'ouvrage de Han) :

- achète(x, "couches") \Rightarrow achète(x, "bière") [0.5%, 60%]

L. Jourdan – Aide à la décision

Règles d'associations



"SI achète couche, ALORS achète bière, dans 60% de cas, dans 0.5% de la base"

Condition, partie gauche de la règle

Conséquence, partie droite de la règle

Support, fréquence ("partie gauche et droite sont présentes ensemble dans la base")

Confiance ("si partie gauche de la règle est vérifiée, probabilité que la partie droite de la règle soit vérifiée")

L. Jourdan – Aide à la décision

Règles d'associations

- Support** : % d'instances de la base vérifiant la règle.

$$\text{support}(A \Rightarrow B [s, c]) = p(A \text{ et } B) = \text{support}(\{A, B\})$$

- Confiance** : % d'instances de la base vérifiant l'implication

$$\text{confiance}(A \Rightarrow B [s, c]) = p(B|A) = p(A \text{ et } B) / p(A) = \frac{\text{support}(\{A, B\})}{\text{support}(\{A\})}$$

L. Jourdan – Aide à la décision

Exemple

TID	Items
1	Pain, Lait
2	Bière, Couches, Pain, Oeufs
3	Bière, Coca, Couches, Lait
4	Bière, Pain, Couches, Lait
5	Coca, Pain, Couches, Lait

$\{Couches, Lait\} \Rightarrow_{s, \alpha} Bière$

Règle : $X \Rightarrow_{s, \alpha} y$

Support : $s = \frac{\sigma(X \cup y)}{|T|} (s = P(X, y))$

Confiance : $\alpha = \frac{\sigma(X \cup y)}{\sigma(X)} (\alpha = P(y|X))$

$$s = \frac{\sigma(Couches, Lait, Bière)}{\text{Nombre total d'instances}} = \frac{2}{5} = 0.4$$

$$\alpha = \frac{\sigma(Couches, Lait, Bière)}{\sigma(Couches, Lait)} = 0.66$$

L. Jourdan – Aide à la décision

Règles d'associations

Support minimum σ :

- Elevé \Rightarrow peu d'itemsets fréquents
- \Rightarrow peu de règles valides qui ont été souvent vérifiées
- Réduit \Rightarrow plusieurs règles valides qui ont été rarement vérifiées

Confiance minimum γ :

- Elevée \Rightarrow peu de règles, mais toutes "pratiquement" correctes
- Réduite \Rightarrow plusieurs règles, plusieurs d'entre elles sont "incertaines"

Valeurs utilisées : $\sigma = 2 - 10\%$, $\gamma = 70 - 90\%$

L. Jourdan – Aide à la décision

Recherche de règles d'association

Données d'entrée : liste d'achats

Achat = liste d'articles (longueur variable)

	Produit A	Produit B	Produit C	Produit D	Produit E
Achat 1	*			*	
Achat 2	*	*	*		
Achat 3	*				*
Achat 4	*			*	*
Achat 5		*		*	

L. Jourdan – Aide à la décision

Recherche de règles d'association

Tableau de co-occurrence : combien de fois deux produits ont été achetés ensemble ?

	Produit A	Produit B	Produit C	Produit D	Produit E
Produit A	4	1	1	2	2
Produit B	1	2	1	1	0
Produit C	1	1	1	0	0
Produit D	2	1	0	3	1
Produit E	2	0	0	1	2

L. Jourdan – Aide à la décision

Illustration / Exemple

- **Règle d'association** :
 - Si A alors B (règle 1)
 - Si A alors D (règle 2)
 - Si D alors A (règle 3)
- **Supports** :
 - Support(1)=20% ; Support(2)=Support(3)=40%
- **Confiances** :
 - Confiance(2) = 50% ; Confiance(3) = 67%
- On préfère la règle 3 à la règle 2.

L. Jourdan – Aide à la décision

Description de la méthode

- Support et confiance ne sont pas toujours suffisants
- Ex : Soient les 3 articles A, B et C

article	A	B	C	A et B	A et C	B et C	A, B et C
support	45%	42,5%	40%	25%	20%	15%	5%

- Règles à 3 articles : même support 5%
- **Confiance**
 - Règle : Si A et B alors C = 0.20
 - Règle : Si A et C alors B = 0.25
 - Règle : Si B et C alors A = 0.33

L. Jourdan – Aide à la décision

Recherche de règles

- Soient une liste de n articles et de m achats.
- 1. Calculer le nombre d'occurrences de chaque article.
- 2. Calculer le tableau des co-occurrences pour les paires d'articles.
- 3. Déterminer les règles de niveau 2 en utilisant les valeurs de support, confiance et amélioration.
- 4. Calculer le tableau des co-occurrences pour les triplets d'articles.
- 5. Déterminer les règles de niveau 3 en utilisant les valeurs de support, confiance et amélioration
- ...

L. Jourdan – Aide à la décision

Complexité

- Soient :
 - n : nombre de transactions dans la BD
 - m : Nombre d'attributs (items) différents
- Complexité
 - Nombre de règles d'association : $O(m \cdot 2^{m-1})$
 - Complexité de calcul : $O(n \cdot m \cdot 2^m)$

L. Jourdan – Aide à la décision

Réduction de la complexité

- n de l'ordre du million (parcours de la liste nécessaire)
- Taille des tableaux en fonction de m et du nombre d'articles présents dans la règle

	2	3	4
n	$n(n-1)/2$	$n(n-1)(n-2)/6$	$n(n-1)(n-2)(n-3)/24$
100	4950	161 700	3 921 225
10000	$5 \cdot 10^7$	$1.7 \cdot 10^{11}$	$4.2 \cdot 10^{14}$

- Conclusion de la **règle restreinte** à un sous-ensemble de l'ensemble des articles vendus.
 - **Exemple** : articles nouvellement vendues.
- Création de **groupes** d'articles (différents niveaux d'abstraction).
- **Elagage** par support minimum.

L. Jourdan – Aide à la décision

Illustration sur une BD commerciale

Attribut	Compteur
Pain	4
Coca	2
Lait	4
Bière	3
Couches	4
Oeufs	1

Attributs (1-itemsets)

Itemset	Compteur
{Pain,Lait}	3
{Pain,Bière}	2
{Pain,Couches}	3
{Lait,Bière}	2
{Lait,Couches}	3
{Bière,Couches}	3

paires (2-itemsets)

Itemset	Compteur
{Pain,Lait,Couches}	3
{Lait,Couches,Bière}	2

Triplets (3-itemsets)

Support Minimum = 3

Si tout sous-ensemble est considéré,
 $C_1^6 + C_2^6 + C_3^6 = 41$
 En considérant un seuil support min,
 $6 + 6 + 2 = 14$

L. Jourdan – Aide à la décision

L'algorithme Apriori [Agrawal93]

- Deux étapes
 - Recherche des k-itemsets fréquents (support \geq MINSUP)
 - (Pain, Fromage, Vin) = 3-itemset
 - Principe : Les sous-itemsets d'un k-itemset fréquent sont obligatoirement fréquents
 - Construction des règles à partir des k-itemsets trouvés
 - Une règle fréquente est retenue si et seulement si sa confiance \geq MINCONF
 - Exemple : ABCD fréquent
 - AB \rightarrow CD est retenue si sa confiance \geq MINCONF

L. Jourdan – Aide à la décision

Recherche des k-itemsets fréquents (1)

Exemple

- $I = \{A, B, C, D, E, F\}$
- $T = \{AB, ABCD, ABD, ABDF, ACDE, BCDF\}$
- MINSUP = 1/2

Calcul de L1 (ensemble des 1-itemsets)

- $C_1 = I = \{A, B, C, D, E, F\}$ // C1 : ensemble de 1-itemsets candidats
- $s(A) = s(B) = 5/6$, $s(C) = 3/6$, $s(D) = 5/6$, $s(E) = 1/6$, $s(F) = 2/6$
- $L_1 = \{A, B, C, D\}$

Calcul de L2 (ensemble des 2-itemsets)

- $C_2 = L_1 \times L_1 = \{AB, AC, AD, BC, BD, CD\}$
- $s(AB) = 4/6$, $s(AC) = 2/6$, $s(AD) = 4/6$, $s(BC) = 2/6$, $s(BD) = 4/6$, $s(CD) = 3/6$
- $L_2 = \{AB, AD, BD, CD\}$

L. Jourdan – Aide à la décision

Recherche des k-itemsets fréquents (2)

- Calcul de L_3 (ensemble des 3-itemsets)
 - $C_3 = \{ABD\}$ ($ABC \notin C_3$ car $AC \notin L_2$)
 - $s(ABD) = 3/6$
 - $L_3 = \{ABD\}$
- Calcul de L_4 (ensemble des 4-itemsets)
 - $C_4 = \emptyset$
 - $L_4 = \emptyset$
- Calcul de L (ensembles des itemsets fréquents)
 - $L = \cup L_i = \{A, B, C, D, AB, AD, BD, CD, ABD\}$

L. Jourdan – Aide à la décision

L'algorithme Apriori

```

L1 = {1-itemsets fréquents};
for (k=2; Lk-1 ≠ ∅; k++) do
    Ck = apriori_gen(Lk-1);
    forall instances t ∈ T do
        Ct = subset(Ck, t);
        forall candidats c ∈ Ct do
            c.count++;
        Lk = { c ∈ Ck / c.count ≥ MINSUP }
    L = ∪ Li;
```

L. Jourdan – Aide à la décision

La procédure Apriori_gen

```

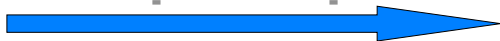
{ Jointure  $L_{k-1} * L_{k-1}$  ; k-2 éléments communs }
insert into  $C_k$ ;
select p.item1, p.item2, ..., p.itemk-1, q.itemk-1
from  $L_{k-1p}$ ,  $L_{k-1q}$ 
where p.item1=q.item1, ..., p.itemk-2=q.itemk-2
    , p.itemk-1 < q.itemk-1
forall itemsets c ∈  $C_k$  do
    forall (k-1)-itemsets s ⊂ c do
        if s ∉  $L_{k-1}$  then
            delete c from  $C_k$ ;
```

L. Jourdan – Aide à la décision

Apriori - Exemple

Base de données D		C_1		L_1	
TID	Items	itemset	sup.	itemset	sup.
100	1 3 4	{1}	2	{1}	2
200	2 3 5	{2}	3	{2}	3
300	1 2 3 5	{3}	3	{3}	3
400	2 5	{4}	1	{5}	3
		{5}	3		

Scan D

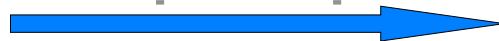


L. Jourdan – Aide à la décision

Apriori - Exemple

C_2		C_2		L_2	
itemset		itemset	sup.	itemset	sup.
{1 2}		{1 2}	1	{1 3}	2
{1 3}		{1 3}	2	{2 3}	2
{1 5}		{1 5}	1	{2 5}	3
{2 3}		{2 3}	2	{3 5}	2
{2 5}		{2 5}	3		
{3 5}		{3 5}	2		

Scan D



L. Jourdan – Aide à la décision

Apriori - Exemple

C_3		L_3	
itemset		itemset	sup.
{2 3 5}		{2 3 5}	2

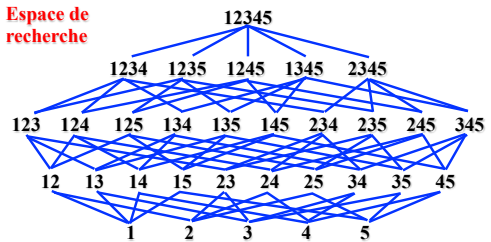
Scan D



L. Jourdan – Aide à la décision

Apriori - Exemple

Espace de recherche



L. Jourdan – Aide à la décision

Apriori - Exemple

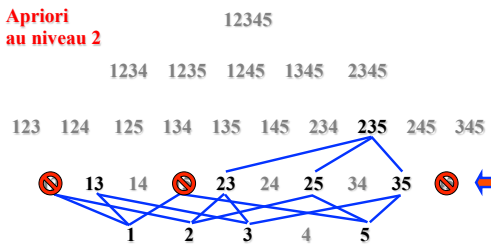
Apriori
au Niveau 1



L. Jourdan – Aide à la décision

Apriori - Exemple

Apriori
au niveau 2



L. Jourdan – Aide à la décision

Génération des règles à partir des itemsets

Pseudo-code :

- **pour** chaque itemset fréquent l
 - générer tous les sous-itemsets non vides s de l
- **pour** chaque sous-itemset non vide s de l
 - produire la règle " $s \Rightarrow (l-s)$ " si $\text{support}(l)/\text{support}(s) \geq \text{min_conf}$ ", où min_conf est la confiance minimale

- **Exemple** : itemset fréquent $l = \{abc\}$,
Sous-itemsets $s = \{a, b, c, ab, ac, bc\}$

- $a \Rightarrow bc, b \Rightarrow ac, c \Rightarrow ab$
– $ab \Rightarrow c, ac \Rightarrow b, bc \Rightarrow a$

L. Jourdan – Aide à la décision

Génération des règles à partir des itemsets

Règle 1 à mémoriser :

- La génération des itemsets fréquents est une opération **coûteuse**
- La génération des règles d'association à partir des itemsets fréquents est **rapide**

Règle 2 à mémoriser :

- Pour la génération des itemsets, le **seuil support** est utilisé.
- Pour la génération des règles d'association, le **seuil confiance** est utilisé.

Complexité en pratique ?

- A partir d'un exemple réel (petite taille) ...
- Expériences réalisées sur un serveur Alpha Citum 4/275 avec 512 MB de RAM & Red Hat Linux release 5.0 (kernel 2.0.30)

L. Jourdan – Aide à la décision

Apriori - Complexité

Phase coûteuse : Génération des candidats

- Ensemble des candidats de grande taille :
 - 10^4 1-itemset fréquents génèrent 10^7 candidats pour les 2-itemsets
 - Pour trouver un itemset de taille 100, e.x., $\{a_1, a_2, \dots, a_{100}\}$, on doit générer $2^{100} \approx 10^{30}$ candidats.
- Multiple scans de la base de données :
 - Besoin de $(n+1)$ scans, n est la longueur de l'itemset le plus long

L. Jourdan – Aide à la décision

Apriori - Complexité

En pratique :

- Pour l'algorithme Apriori basique, le nombre d'attributs est généralement plus critique que le nombre de transactions
- **Par exemple** :
 - 50 attributs chacun possédant 1-3 valeurs, 100.000 transactions (not very bad)
 - 50 attributs chacun possédant 10-100 valeurs, 100.000 transactions (quite bad)
 - 10.000 attributs chacun possédant 5-10 valeurs, 100 transactions (very bad...)
- **Notons** :
 - Un attribut peut avoir plusieurs valeurs différentes
 - Les algorithmes traitent chaque paire attribut-valeur comme un attribut (2 attributs avec 5 valeurs \rightarrow "10 attributs")

Quelques pistes pour résoudre le problème ...

L. Jourdan – Aide à la décision

Apriori – Réduction de la complexité

Suppression de transactions :

- Une transaction qui ne contient pas de k -itemsets fréquents est inutile à traiter dans les parcours (scan) suivants.

Partitionnement :

- Tout itemset qui est potentiellement fréquent dans une BD doit être potentiellement fréquent dans au moins une des partitions de la BD.

Echantillonnage :

- Extraction à partir d'un sous-ensemble de données, décroître le seuil support

L. Jourdan – Aide à la décision

Apriori - Avantages

- **Résultats clairs** : règles faciles à interpréter.
- **Simplicité de la méthode**
- **Aucune hypothèse préalable** (Apprentissage non supervisé)
- **Introduction du temps** : méthode facile à adapter aux séries temporelles. Ex : Un client ayant acheté le produit A est susceptible d'acheter le produit B dans deux ans.

L. Jourdan – Aide à la décision

Apriori - Inconvénients

- **Coût de la méthode** : méthode coûteuse en temps
- **Qualité des règles** : production d'un nombre important de règles triviales ou inutiles.
- **Articles rares** : méthode non efficace pour les articles rares.
- **Adapté aux règles binaires**
- **Apriori amélioré**
 - Variantes de Apriori : DHP, DIC, etc.
 - Partition [Savasere et al. 1995]
 - Eclat et Clique [Zaki et al. 1997]
 - ...

L. Jourdan – Aide à la décision

Typologie des règles

- **Règles d'association binaires**
 - Forme : *if C then P*. C, P : ensembles d'objets
- **Règles d'association quantitatives**
 - Forme : *if C then P*
 - C = terme1 & terme2 & ... & termen
 - P = termen+1
 - termei = <attributj, op, valeur> ou <attributj, op, valeur_de, valeur_a>
 - Classes : valeurs de P
 - Exemple : *if ((Age>30) & (situation=marié)) then prêt=prioritaire*
- **Règles de classification généralisée**
 - Forme : *if C then P*, P=p1, p2, ..., pm P: attribut but
- etc.

L. Jourdan – Aide à la décision

Règles d'association – Résumé

- Probablement la contribution la plus significative de la communauté KDD
- Méthodes de recherche de règles :
 - A-priori
 - Algorithmes génétiques
- Plusieurs travaux ont été publiés dans ce domaine

L. Jourdan – Aide à la décision

Règles d'association – Résumé

Plusieurs issues ont été explorées : intérêt d'une règle, optimisation des algorithmes, parallélisme et distribution, ...

Directions de recherche :

- Règles d'associations pour d'autres types de données : données spatiales, multimedia, séries temporelles, ...

L. Jourdan – Aide à la décision

Critères pour les règles

Mesure	Formule	Effet
Support S	$\frac{C \text{ et } P}{N}$	% transactions qui contiennent C et P
Confiabilité C	$\frac{C \text{ et } P}{C}$	Probabilité conditionnelle
Intérêt I	$\frac{C \text{ et } P}{C \times P}$	Privilège les motifs rares (ayant un support faible)
Conviction V	$\frac{C \times \bar{P}}{C \text{ et } \bar{P}}$	Mesure la faiblesse de (C, not P) V >> 1 : P se passe avec C
Piatetsky-Shapiro's	$C \text{ et } P - C \times P$	Mesure la dépendance
Surprise R	$\frac{(C \text{ et } P - C \text{ et } \bar{P})}{P}$	Cherche des règles étonnantes Mesure l'infirmité(C, NOT P)

L. Jourdan – Aide à la décision