
Fouille de Séquences

MARC PLANTEVIT

UCB LYON 1 - LIRIS - CNRS



LIRIS



`marc.plantevit@liris.cnrs.fr`

Outline

- 1 Introduction
- 2 Définitions
- 3 Algorithmes
- 4 Web Usage Mining
- 5 Conclusions

Règles d'association : rappel

Trans. ID	Items
1	A,D
2	A,C
3	A,B,C
4	A,B,E,F

Règles d'association : rappel

Trans. ID	Items
1	A,D
2	A,C
3	A,B,C
4	A,B,E,F

- Itemsets : A,B ou B,E,F

Règles d'association : rappel

Trans. ID	Items
1	A,D
2	A,C
3	A,B,C
4	A,B,E,F

- Itemsets : A,B ou B,E,F
- Support d'un itemset :
 - $support(AD) = 1$
 - $support(AC) = 2$

Règles d'association : rappel

Trans. ID	Items
1	A,D
2	A,C
3	A,B,C
4	A,B,E,F

- Itemsets : A,B ou B,E,F
- Support d'un itemset :
 - $support(AD) = 1$
 - $support(AC) = 2$
- Itemsets fréquents ($minsupp = 50\%$) :
 - $\{A, C\}$ et $\{A, B\}$ sont fréquents

Règles d'association : rappel

Trans. ID	Items
1	A,D
2	A,C
3	A,B,C
4	A,B,E,F

- Itemsets : A,B ou B,E,F
- Support d'un itemset :
 - $support(AD) = 1$
 - $support(AC) = 2$
- Itemsets fréquents ($minsupp = 50\%$) :
 - $\{A, C\}$ et $\{A, B\}$ sont fréquents
- Pour $minsupp = 50\%$ et $minconf = 50\%$, on a les règles suivantes :
 - $A \rightarrow C$ [50%, 50%]
 - $C \rightarrow A$ [50%, 100%]
 - $A \rightarrow B$ et $B \rightarrow A$?

Association vs. motif séquentiel

Règles d'association

- **ensemble d'items** qui apparaît fréquemment dans une base de données
- Corrélation entre ces items

Motifs séquentiels

- **Liste d'ensemble d'items** qui apparaît fréquemment dans une base de données
- L'idée principale est de comprendre les habitudes d'un consommateur, etc.

Analyse du panier de la ménagère

Analyse du panier de la ménagère

The image shows a receipt from Star Market with several annotations. The receipt text is as follows:

STOW STAR
(978) 897-5140

TERRA SALT & VING	1.99 F
TERRA SLT&PEPR C	1.99 F
HD PREN 12 CHOCO WP	.89 F
SUB-TOTAL	4.87
TAX	.00
BALANCE	4.87

STAR MARKET
GREAT ROAD
STOW, MA

VISA/MASTERCARD

CHANGE

9/04/00 6:24 PM 0152 04 0429 109

SIGN UP TODAY TO GET EXTRA SAVINGS
WITH THE STAR MARKET ADVANTAGE CARD
KEITH SHAW - STORE MANAGER
THANK YOU

Annotations on the left side of the receipt:

- Localization**: Points to the store name and phone number.
- Items bought**: Points to the list of items.
- Identification**: Points to the payment method.
- Date, time**: Points to the date and time stamp.

Encore

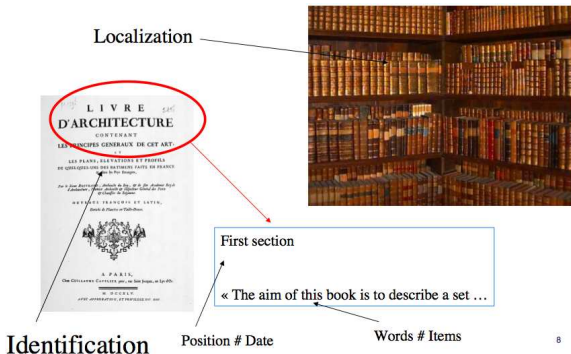
Textes

S_1 : ⟨ "s-Tomosyn", "as", "well", "as", "m-tomosyn", "was", "mainly", ... ⟩

S_2 :

⟨ "This", "region", "was", "necessary", "but", "not", "sufficient", "for", ... ⟩

...



Autres exemples de séquences de données

Biological Sequences : DNA, RNA and Proteins

*GAATTCTCTGTAACACTAAGCTCTCTTCCTCAAACCAGAGGTAGA
ATGTGTAATAATTTACAGAATTTCTAGACTTCAACGATCTGATTTTT
ATTTATTTTTATTTTTTTCAGGTTGAGACTGAGCTAAAGTTAATCTG*

Autres exemples de séquences de données

Biological Sequences : DNA, RNA and Proteins

```
GAATTCTCTGTAACACTAAGCTCTCTTCCTCAAACCAGAGGTAGA  
ATGTGTAATAATTTACAGAATTTCTAGACTTCAACGATCTGATTTTT  
ATTTATTTTTATTTTTTTCAGGTTGAGACTGAGCTAAAGTTAATCTG
```

Event Sequences : Weblogs, System Traces, Purchase Histories and Sales Histories

- weblog :
 $\langle 100, a \rangle, \langle 100, b \rangle, \langle 200, a \rangle, \langle 300, b \rangle, \langle 200, b \rangle, \langle 400, a \rangle, \langle 100, a \rangle, \langle 400, b \rangle$
- customer purchase history :
 $\langle 223100, 05/26/06, 10am, CentralStation, \{ WholeMealBread, AppleJuice \} \rangle,$
 $\langle 225101, 05/26/06, 11am, CentralStation, \{ Burger, Pepsi, Banana \} \rangle$
- Storewide sales history :
 $\langle 97100, 05/06, \{ \langle Apple : \$85K \rangle, \langle Bread : \$100K \rangle, \langle Cereal : \$150K \rangle, \dots \} \rangle,$
 $\langle 90089, 05/06, \{ \langle Apple : \$65K \rangle, \langle Bread : \$105K \rangle, \langle Diaper : \$20K \rangle, \dots \} \rangle$

Pourquoi l'extraction de motifs séquentiels ?

Un important problème en DM avec de nombreuses applications :

- Analyse des achats des consommateurs
- Analyse ADN,
- Conséquences des catastrophes naturelles,
- Web log mining,
- Détection de tendances

Outline

- 1 Introduction
- 2 Définitions**
- 3 Algorithmes
- 4 Web Usage Mining
- 5 Conclusions

Soit $\mathcal{I} = \{i_1, i_2, \dots, i_n\}$ un ensemble de littéraux appelés **items**

- $\mathcal{I} = \{\text{PC, TV, bière, couches, etc.}\}$

Soit $\mathcal{I} = \{i_1, i_2, \dots, i_n\}$ un ensemble de littéraux appelés **items**

- $\mathcal{I} = \{\text{PC, TV, bière, couches, etc.}\}$

Transaction

Une transaction est composée de :

- Un identifiant (*id_client*),
- Une date de transaction,
- Un ensemble d'items (**itemsets**)
- Exemple : $T = [Cid1, (Beer, Cake, Diaper)_5]$

Séquence et séquence de données

Séquence :

Une **liste ordonnée** d'itemsets

- Exemple : $\langle (TV), (DVD\text{-}Reader, Star\ War\ box) \rangle$

Séquence et séquence de données

Séquence :

Une **liste ordonnée** d'itemsets

- Exemple : $\langle (TV), (DVD\text{-}Reader, Star\ War\ box) \rangle$

Séquence de données

Soient T_1, T_2, \dots, T_m , les transactions d'un client C , la séquence de données de C est la paire :

$$(C, \langle itemset(T_1), itemset(T_2), \dots, itemset(T_m) \rangle)$$

Ex : Les achats d'un consommateur au cours du temps.

Séquence et séquence de données

Séquence :

Une **liste ordonnée** d'itemsets

- Exemple : $\langle (TV), (DVD\text{-}Reader, Star\ War\ box) \rangle$

Séquence de données

Soient T_1, T_2, \dots, T_m , les transactions d'un client C , la séquence de données de C est la paire :

$$(C, \langle itemset(T_1), itemset(T_2), \dots, itemset(T_m) \rangle)$$

Ex : Les achats d'un consommateur au cours du temps.

Base de séquences (de données)

$SDB =$

- $(C_1, \langle (Coffe, Mustard), (Coke), (Diaper) \rangle)$
- $(C_2, \langle (Beer)(Coke) \rangle)$

Inclusion

Soient deux séquences $S_1 = \langle a_1, a_2, \dots, a_n \rangle$ et $S_2 = \langle b_1, b_2, \dots, b_{n'} \rangle$, S_1 est une sous-séquence de S_2 (S_2 superséquence de S_1 , noté $S_1 \preceq S_2$) si il existe des entiers $1 \leq i_1 < i_2 < \dots < i_n \leq n'$ tels que :

$$a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$$

Exemple

- $S_1 = \langle (10)(20, 30), (40), (20) \rangle$
- $S_2 = \langle (20), (40) \rangle \preceq S_1$
- $S_3 = \langle (20), (30) \rangle \not\preceq S_1$

Support d'une séquence

L'occurrence d'une séquence n'est considérée qu'une seule fois dans une séquence de données

Ex : $\langle(20)\rangle$ et $\langle(10)(20)(20, 30), (40), (20)\rangle$

On dit aussi que la séquence de données **contient** la séquence $\langle(20)\rangle$.

Support d'une séquence

Le support d'une séquence s dans une base de données SDB correspond au nombre de séquences de données qui contiennent s .

Exemple

Customers	Date1	Date2	Date3	Date4
C1	10 20 30	20 40 50	10 20 60	10 40
C2	10 20 30	10 20 30		20 30 60
C3	20 30 50		10 40 60	10 20 30
C4	10 30 60	20 40	10 20 60	50

- $Absolute_Support(\langle(10, 30)(20)(20, 60)\rangle) = 3$
- $Relative_Support(\langle(10, 30)(20)(20, 60)\rangle) = \frac{3}{4} = 75\%$

Problématique

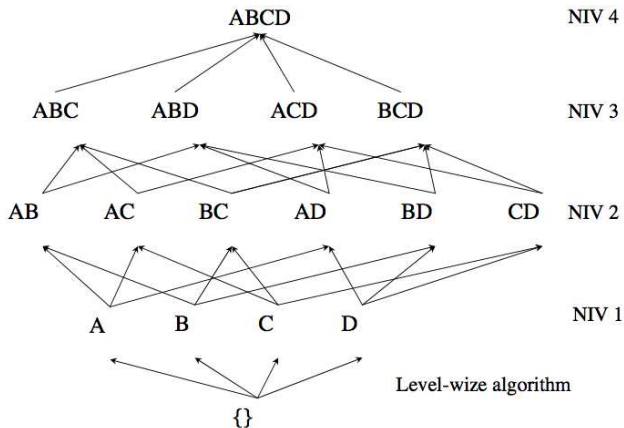
Etant donné un ensemble de séquences de données SDB et un seuil de support minimum $minsupp$, le but de l'extraction des motifs séquentiels est de trouver l'ensemble complet des séquences qui ont un support supérieur ou égal à $minsupp$.

- Le nombre de motifs séquentiels « cachés » dans les séquences de données peut être très importants
- Un algorithme d'extraction doit donc :
 - trouver **l'ensemble complet** des motifs séquentiels
 - être efficace (passage à l'échelle) en limitant le nombre de passes sur les données
 - permettre d'incorporer un ensemble varié de contraintes définies par l'utilisateur (syntaxiques, temporelles, etc.)

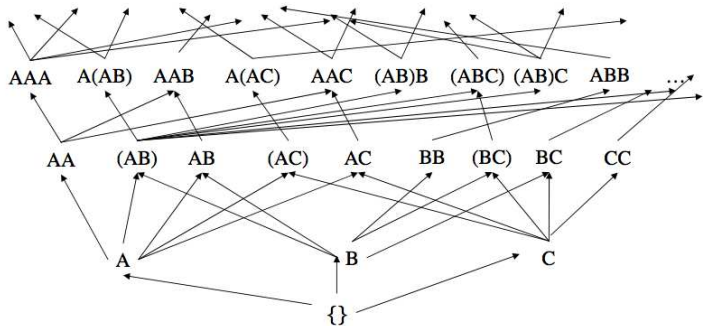
Outline

- 1 Introduction
- 2 Définitions
- 3 Algorithmes**
- 4 Web Usage Mining
- 5 Conclusions

Itemset : l'espace de recherche (un treillis)



Motifs séquentiels : l'espace de recherche



Propriété Apriori des motifs séquentiels

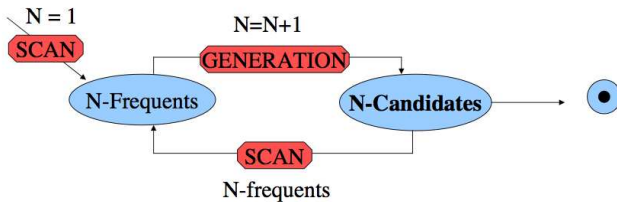
Une propriété fondamentale [Agrawal et Srikant, 94]

- Si une séquence S n'est pas fréquente
- Alors aucune de ses super-séquences est fréquente
- Ex : $\langle hb \rangle$ non fréquent alors $\langle hab \rangle$ $\langle (ah)b \rangle$ non plus

Seq. ID	Séquence
10	$\langle (bd)cb(ac) \rangle$
20	$\langle (bf)(ce)b(fg) \rangle$
30	$\langle (ah)(bf)abf \rangle$
40	$\langle (be)(ce)d \rangle$
50	$\langle a(bd)bcb(ade) \rangle$

minsupp = 2

Approche générale : générer/élaguer



Génération de candidats : deux types d'extension

S-Extension

On ajoute un nouvel itemset (une séquence) à la séquence considérée

- $\langle\langle a, b \rangle(c)\rangle \rightarrow \langle\langle a, b \rangle(c)(d)\rangle$

I-Extension

On ajoute un nouvel item à la séquence considérée dans un itemset déjà créé.

- $\langle\langle a, b \rangle(c)\rangle \rightarrow \langle\langle a, b \rangle(c, d)\rangle$

$L=1$

While ($\text{Result}_L \neq \text{NULL}$)

 Candidate Generate

 Prune

 Test

$L=L+1$

Une passe sur les données pour extraire les séquences composées d'un unique item (dans un unique itemset).

Seq. ID	Sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bcb(ade)>

Cand	Sup
<a>	3
	5
<c>	4
<d>	3
<e>	3
<f>	2
<g>	1
<h>	1

Les processus global

5th scan : 1 candidate
1 length-5 seq pattern

<(bd)cba>

4th scan : 8 candidates
6 length-4 seq pat

<abba> <(bd)bc> ...

3rd scan : 46 candidates
19 length-3 seq pat.

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

2nd scan : 51 candidates
19 length-2 seq pat.

<aa> <ab> ... <af> <ba> <bb> ... <ff> <(ab)> ... <(ef)>

1st scan : 8 candidates
6 length-1 seq pattern

<a> <c> <d> <e> <f> <g> <h>

Génération des candidats de longueur 2

S-Extension
51 2-Candidates

	<a>		<c>	<d>	<e>	<f>
<a>	<aa>	<ab>	<ac>	<ad>	<ae>	<af>
	<ba>	<bb>	<bc>	<bd>	<be>	<bf>
<c>	<ca>	<cb>	<cc>	<cd>	<ce>	<cf>
<d>	<da>	<db>	<dc>	<dd>	<de>	<df>
<e>	<ea>	<eb>	<ec>	<ed>	<ee>	<ef>
<f>	<fa>	<fb>	<fc>	<fd>	<fe>	<ff>

I-Extension

	<a>		<c>	<d>	<e>	<f>
<a>		<(ab)>	<(ac)>	<(ad)>	<(ae)>	<(af)>
			<(bc)>	<(bd)>	<(be)>	<(bf)>
<c>				<(cd)>	<(ce)>	<(cf)>
<d>					<(de)>	<(df)>
<e>						<(ef)>
<f>						

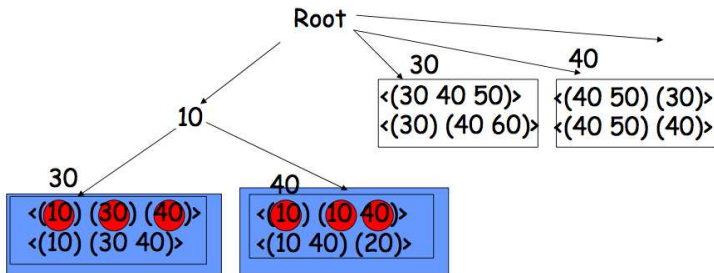
Without the
antimonotonic property
 $8*8+8*7/2=92$
candidates

L'étape la plus coûteuse

- Les candidats sont stockés en mémoire centrale.
- Il faut limiter les accès disques sur la base de données
- Possibilité de charger la base en mémoire quand c'est possible

Comment stocker efficacement les candidats ?

A Hash tree structure

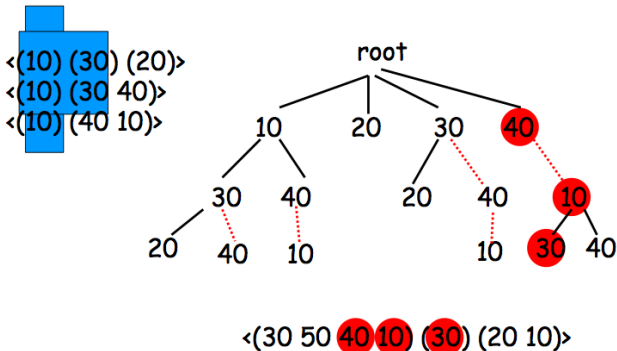


$$S = \langle (10)(30)(10, 40) \rangle$$

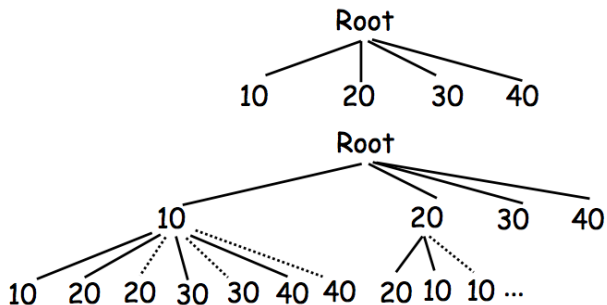
PSP (Prefix Tree for SP) [Masseglia et al. 98]

Une structure plus efficace basée sur un arbre des préfixes.

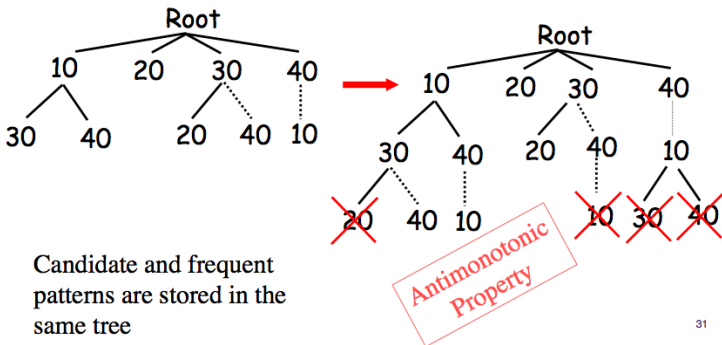
- 2 types d'arc



PSP : génération des 2-candidats



PSP : génération des k-candidats ($k > 2$)



31

SPADE (Sequential PAttern Discovery using Equivalent Class) [Zaki 2001]

- représentation verticale de la base de données
- base de séquences transformée en un ensemble de triplet (item, SID,EID)
- Extraction de motifs est effectuée avec une génération de candidats basée sur Apriori

SPADE : pour un calcul du support plus efficaces

SID	EID	Items
1	1	a
1	2	abc
1	3	ac
1	4	d
1	5	cf
2	1	ad
2	2	c
2	3	bc
2	4	ae
3	1	ef
3	2	ab
3	3	df
3	4	c
3	5	b
4	1	e
4	2	g
4	3	af
4	4	c
4	5	b
4	6	c

a		b		...
SID	EID	SID	EID	...
1	1	1	2	
1	2	2	3	
1	3	3	2	
2	1	3	5	
2	4	4	5	
3	2			
4	3			

ab			ba			...
SID	EID (a)	EID(b)	SID	EID (b)	EID(a)	...
1	1	2	1	2	3	
2	1	3	2	3	4	
3	2	5				
4	3	5				

aba				...
SID	EID (a)	EID(b)	EID(a)	...
1	1	2	3	
2	1	3	4	

- Une surgénération de candidats :
 - Pour 1000 1-séquences fréquentes, on génère $1000 \times 1000 \times \frac{1000 \times 999}{2} = 1,499,500$ 2 candidats
 - De multiples passages sur la bases de données
 - Approche en largeur sont coûteuse en gestion mémoire
- Pour extraire de longues séquences, ce type d'approche n'est pas adapté :
 - un nombre exponentiel de sous séquences candidates générées
 - pour une séquence de longueur 100 : $2^{100} - 1 \approx 10^{30}$

Approche « pattern growth »

- pas de génération de candidats
- extraction d'items fréquents sur des bases projetées
- approche gloutonne en profondeur

Préfixe et suffixe

- $\langle a \rangle$, $\langle aa \rangle$ et $\langle a(abc) \rangle$ sont des **préfixes** de la séquence $\langle a(abc)(ac)d(cf) \rangle$

Préfixe et suffixe

- $\langle a \rangle$, $\langle aa \rangle$ et $\langle a(abc) \rangle$ sont des **préfixes** de la séquence $\langle a(abc)(ac)d(cf) \rangle$
- Etant donnée la séquence $\langle a(abc)(ac)d(cf) \rangle$

Préfixe et suffixe

- $\langle a \rangle$, $\langle aa \rangle$ et $\langle a(abc) \rangle$ sont des **préfixes** de la séquence $\langle a(abc)(ac)d(cf) \rangle$
- Etant donnée la séquence $\langle a(abc)(ac)d(cf) \rangle$

Préfixe	Suffixe (Prefix-Based Projection)
$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle$
$\langle a\underline{a} \rangle$	$\langle (_bc)(ac)d(cf) \rangle$
$\langle a\underline{b} \rangle$	$\langle (_c)(ac)d(cf) \rangle$

- Des items particuliers : $(_b)$

Extraction de motifs séquentiels par projection des préfixes

- Etape 1 : extraction des 1-séquences fréquentes :

$\langle a \rangle, \langle b \rangle, \langle c \rangle, \langle d \rangle, \langle e \rangle, \langle f \rangle$

- Etape 2 : l'ensemble complet des motifs séquentiels peut être partitionné en 6 sous-ensembles :
 - ceux de préfixe $\langle a \rangle$,
 - ceux de préfixe $\langle b \rangle$,
 - ceux de préfixe $\langle c \rangle$,
 - ceux de préfixe $\langle d \rangle$,
 - ceux de préfixe $\langle e \rangle$,
 - ceux de préfixe $\langle f \rangle$.

Trouver les motifs séquentiels de préfixe $\langle a \rangle$

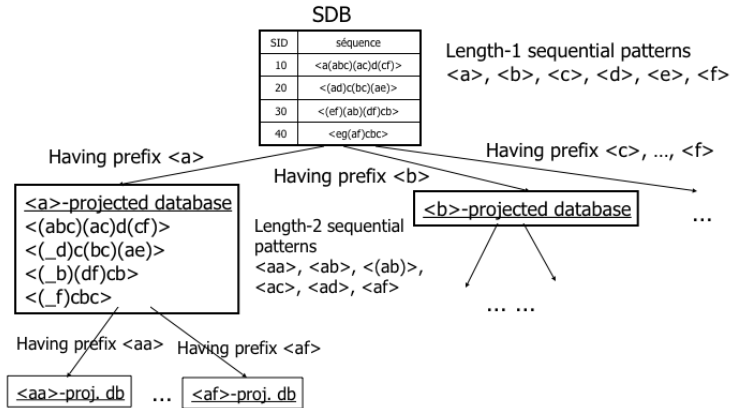
Simplement considérer les projections par rapport à $\langle a \rangle$:

- $\langle (abc)(ac)d(cf) \rangle$,
- $\langle (-d)c(bc)(ae) \rangle$,
- $\langle (-b)(df)cb \rangle$,
- $\langle (-f)cbc \rangle$

Motifs séquentiels de longueur 2 de préfixe $\langle a \rangle$:

- $\langle aa \rangle$,
- $\langle ab \rangle$,
- $\langle (ab) \rangle$,
- $\langle ac \rangle$,
- $\langle ad \rangle$,
- $\langle af \rangle$

Complétude de PrefixSpan



Efficacité de PrefixSpan

- Pas de génération de candidats
- La taille des bases projetées diminue sans cesse
- Le principal coût de PrefixSpan : Construire les bases projetées
 - Est amélioré à l'aide de **pseudo projections**

Pseudo projection

- Quand la base de données considérée peut être gérée en mémoire, **utilisation de pointeurs** pour les projections
- Pointeur sur les séquences
- Offset sur le suffixe

$s = \langle a(abc)(ac)d(cf) \rangle$

$s|_{\langle a \rangle} : (, 2) \langle (abc)(ac)d(cf) \rangle$

$s|_{\langle ab \rangle} : (, 4) \langle (_c)(ac)d(cf) \rangle$

Pseudo projection vs. projection physique

- La pseudo projection évite de stocker physiquement les suffixes
 - Efficacité (temps d'exécution et espace mémoire) quand la base peut être stocker en mémoire
- Quid des bases qui ne peuvent pas être stockées en mémoire ?
 - Intégration de projections physiques et de pseudo projections
 - Utilisation de pseudo projections quand la base tient en mémoire (la taille des bases projetées diminuant)

Motifs séquentiels fermés

Définition ?

Définition ?

Des motivations doubles

- Réduire le nombre de motifs (redondants) tout en maintenant la qualité de la connaissance extraite (pas de perte d'information)
- Améliorer l'extraction des motifs séquentiels en introduisant des propriétés d'élagages supplémentaires.

Motifs séquentiels fermés

Définition ?

Des motivations doubles

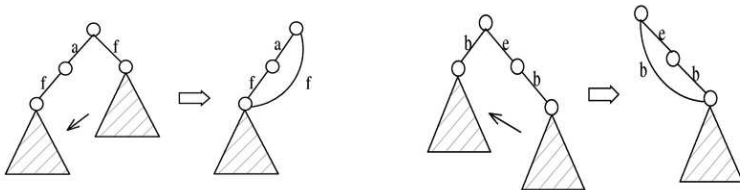
- Réduire le nombre de motifs (redondants) tout en maintenant la qualité de la connaissance extraite (pas de perte d'information)
- Améliorer l'extraction des motifs séquentiels en introduisant des propriétés d'élagages supplémentaires.

Deux approches :

- CloSpan (Yan et al. 2003)
- Bide (Wang et al. 2007)

Eviter de parcourir plusieurs fois les mêmes bases projetées

« Using Backward Subpattern and Backward Superpattern pruning to prune redundant search space »



Les limites de la gestion d'ensemble de candidats

- post-traitement : $O(n^2)$

Il faut éviter de gérer des candidats

BIDE : Idée Principale

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

BIDE : Idée Principale

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

- vers l'avant inter itemset $S' = \langle s_1, s_2, \dots, s_g, \{e'\} \rangle$

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

- 1 vers l'avant inter itemset $S' = \langle s_1, s_2, \dots, s_g, \{e'\} \rangle$
- 2 vers l'avant intra itemset $S' = \langle s_1, s_2, \dots, s_g \cup \{e'\} \rangle$

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

- 1 vers l'avant inter itemset $S' = \langle s_1, s_2, \dots, s_g, \{e'\} \rangle$
- 2 vers l'avant intra itemset $S' = \langle s_1, s_2, \dots, s_g \cup \{e'\} \rangle$
- 3 vers l'arrière inter itemset $S' = \langle s_1, s_2, \dots, s_i, \{e'\}, s_{i+1}, \dots, s_g \rangle$

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

- 1 vers l'avant inter itemset $S' = \langle s_1, s_2, \dots, s_g, \{e'\} \rangle$
- 2 vers l'avant intra itemset $S' = \langle s_1, s_2, \dots, s_g \cup \{e'\} \rangle$
- 3 vers l'arrière inter itemset $S' = \langle s_1, s_2, \dots, s_i, \{e'\}, s_{i+1}, \dots, s_g \rangle$
- 4 vers l'arrière intra itemset $S' = \langle s_1, s_2, \dots, s_i \cup \{e'\}, s_{i+1}, \dots, s_g \rangle$

BIDE : Idée Principale

Extension d'une g - k -séquence préfixe $\langle s_1, s_2, \dots, s_g \rangle$:

- 1 vers l'avant inter itemset $S' = \langle s_1, s_2, \dots, s_g, \{e'\} \rangle$
- 2 vers l'avant intra itemset $S' = \langle s_1, s_2, \dots, s_g \cup \{e'\} \rangle$
- 3 vers l'arrière inter itemset $S' = \langle s_1, s_2, \dots, s_i, \{e'\}, s_{i+1}, \dots, s_g \rangle$
- 4 vers l'arrière intra itemset $S' = \langle s_1, s_2, \dots, s_i \cup \{e'\}, s_{i+1}, \dots, s_g \rangle$

Séquence Close :

- Si pas d'extension qui conserve le support de la séquence

Intervalles de $\langle s_1, s_2, \dots, s_g \rangle$

Exhiber des items qui apparaissent tous les $i^{\text{èmes}}$ intervalles

$\underbrace{l_1}, s_1, \underbrace{l_2}, s_2, \underbrace{l_3}, s_3, \dots, s_{g-1}, \underbrace{l_g}, s_g$

Intervalles de $\langle s_1, s_2, \dots, s_g \rangle$

Exhiber des items qui apparaissent tous les $i^{\text{èmes}}$ intervalles

$\underbrace{l_1}, s_1, \underbrace{l_2}, s_2, \underbrace{l_3}, s_3, \dots, s_{g-1}, \underbrace{l_g}, s_g$

Une séquence peut apparaître plusieurs fois dans une séquence de données

- Séquence $\langle (a, b)(a, c) \rangle$
- Séquence de données
 $\langle (a, b)(a, c)(a, b)(a, c)(a, b)(a, c)(a, b)(a, c) \rangle$

Intervalles de $\langle s_1, s_2, \dots, s_g \rangle$

Exhiber des items qui apparaissent tous les $i^{\text{èmes}}$ intervalles

$\underbrace{l_1}_{s_1}, \underbrace{l_2}_{s_2}, \underbrace{l_3}_{s_3}, \dots, \underbrace{l_g}_{s_{g-1}}, s_g$

Une séquence peut apparaître plusieurs fois dans une séquence de données

- Séquence $\langle (a, b)(a, c) \rangle$
- Séquence de données $\langle (a, b)(a, c)(a, b)(a, c)(a, b)(a, c)(a, b)(a, c) \rangle$

Maximiser ces intervalles

Intervalles de $\langle s_1, s_2, \dots, s_g \rangle$

Exhiber des items qui apparaissent tous les $i^{\text{èmes}}$ intervalles

$\underbrace{l_1}_{s_1}, \underbrace{l_2}_{s_2}, \underbrace{l_3}_{s_3}, \dots, \underbrace{l_g}_{s_{g-1}}, s_g$

Une séquence peut apparaître plusieurs fois dans une séquence de données

- Séquence $\langle (a, b)(a, c) \rangle$
- Séquence de données $\langle (a, b)(a, c)(a, b)(a, c)(a, b)(a, c)(a, b)(a, c) \rangle$

Maximiser ces intervalles

$l_1 : \langle (a, b)(a, c) \rangle$

$\langle \underbrace{(a, b)(a, c)(a, b)(a, c)(a, b)(a, c)}_{(a, b)(a, c)} (a, b)(a, c) \rangle$

Intervalles de $\langle s_1, s_2, \dots, s_g \rangle$

Exhiber des items qui apparaissent tous les $i^{\text{èmes}}$ intervalles

$\underbrace{l_1}_{s_1}, \underbrace{l_2}_{s_2}, \underbrace{l_3}_{s_3}, \dots, \underbrace{l_g}_{s_{g-1}}, s_g$

Une séquence peut apparaître plusieurs fois dans une séquence de données

- Séquence $\langle (a, b)(a, c) \rangle$
- Séquence de données
 $\langle (a, b)(a, c)(a, b)(a, c)(a, b)(a, c)(a, b)(a, c) \rangle$

Maximiser ces intervalles

$l_2 : \langle (a, b)(a, c) \rangle$

$\langle (a, b) \underbrace{(a, c)(a, b)(a, c)(a, b)(a, c)(a, b)}_{(a, c)(a, b)(a, c)(a, b)} (a, c) \rangle$

Les motifs séquentiels sous contraintes

Extraction de motifs sous contraintes :

Permet de surmonter deux difficultés majeures de l'extraction de motifs :

- améliorer les connaissances extraites
- améliorer le processus d'extraction

Les contraintes : plus proche de l'utilisateur

- Les contraintes peuvent permettre de modéliser l'intérêt de l'utilisateur.

Une large gamme de contraintes

- item constraint
- length constraint
- aggregate constraint
- etc.

Des contraintes spécifiques aux séquences

Des contraintes spécifiques aux séquences

- min-gap
- max-gap
- durée (window size)

Exemple

Au tableau

Algorithmes d'extraction de motifs séquentiels sous contraintes

- Spirit [Garofalakis et al, 02]
- [Pei et al, 02]

Motivations

Difficulté de fixer le bon seuil de support :

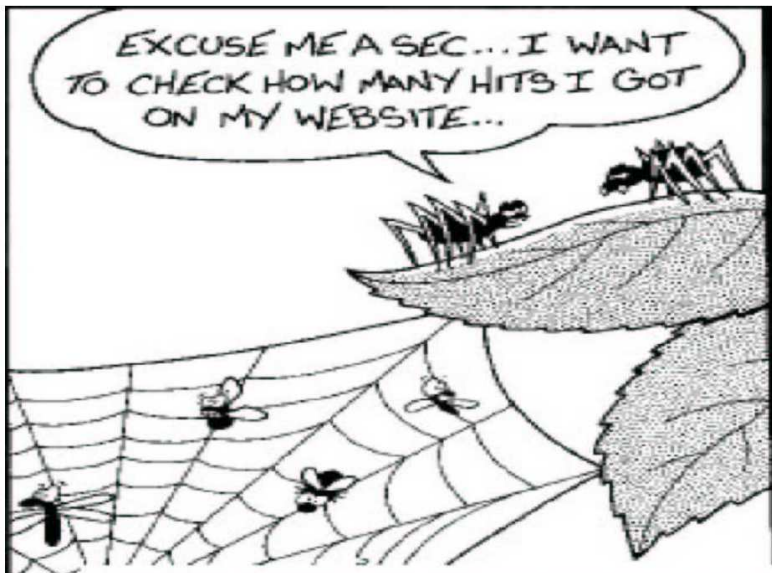
- coca-cola non fréquent
- pepsi non fréquent
- MAIS : soda fréquent

Peu d'approches permettant la prise en compte des hiérarchies dans l'extraction de motifs séquentiels

Outline

- 1 Introduction
- 2 Définitions
- 3 Algorithmes
- 4 Web Usage Mining**
- 5 Conclusions

Web Usage Mining



La découverte de motifs intéressants sur la navigation des visiteurs à partir de logs d'un site web

- Les pages contiennent des informations
- Les liens sont des « routes »
- Comment les gens naviguent sur Internet ?
- **WUM = clickstream analysis**

Hyper liens dynamiques

- générer des hyper liens dynamiques en regardant la navigation du visiteur

Applications

Hyper liens dynamiques

- générer des hyper liens dynamiques en regardant la navigation du visiteur

Personnalisation

Fournir au visiteur ce **qu'il veut** sans le lui demander explicitement

- générer des recommandations pour chaque utilisateur pendant sa navigation

Applications

Hyper liens dynamiques

- générer des hyper liens dynamiques en regardant la navigation du visiteur

Personnalisation

Fournir au visiteur ce **qu'il veut** sans le lui demander explicitement

- générer des recommandations pour chaque utilisateur pendant sa navigation

Site web adaptatif

Changer le site web en fonction des motifs extraits

Amélioration de site web

- comparer les caractéristiques des « clients » et des « non clients »
- changer le site afin que les « non clients » deviennent « clients »
- e.g. changer des hyper liens entre certaines pages

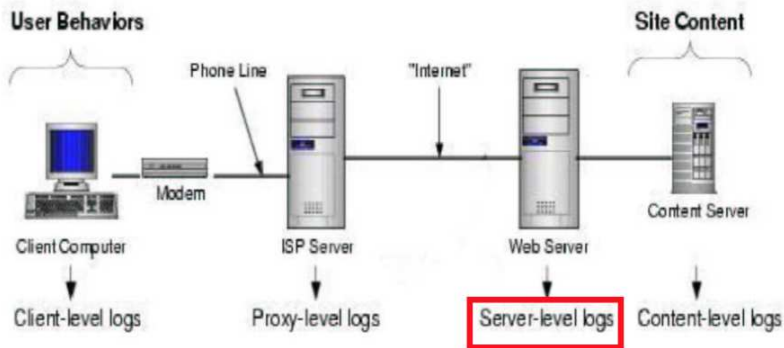
Web Usage Mining

- 75% des Français achètent une raquette de tennis et moins de trois après, des chaussures de sports.
- hyper lien dynamique



Log ou logs

Informations sur les chemins de navigations sont disponibles dans des fichiers de logs



Web logs

IP or domain name User Id Date and Time Request

123.456.78.9 - - [24/Oct/1999:19:13:44 -0400] "GET /Images/tagline.gif HTTP/1.0"

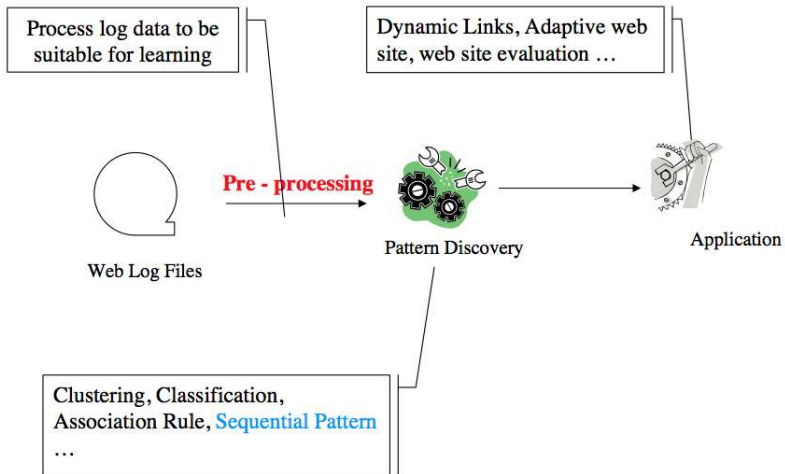
200 1449 <http://www.teced.com/> "Mozilla/4.51 [en] (Win98;l)"

Status File Size URL Browser Cookies

Web logs

IP Address	Time	Method/URL/Protocol	Sta	Size	Referred	Agent
123.456.78.9	[25/Apr/1998:03:04:41 -0500	GET A.html HTTP/1.0	20 0	3290	-	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:05:34 -0500	GET B.html HTTP/1.0	20 0	2050	A.html	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:05:39 -0500	GET L.html HTTP/1.0	20 0	4130	-	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:06:02 -0500	GET F.html HTTP/1.0	20 0	5096	B.html	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:06:58 -0500	GET A.html HTTP/1.0	20 0	3290	-	Mozilla/3.01 (X11, I, IRIX6.2, IP22)
123.456.78.9	[25/Apr/1998:03:07:42 -0500	GET B.html HTTP/1.0	20 0	2050	A.html	Mozilla/3.01 (X11, I, IRIX6.2, IP22)
123.456.78.9	[25/Apr/1998:03:07:55 -0500	GET R.html HTTP/1.0	20 0	8140	L.html	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:09:50 -0500	GET C.html HTTP/1.0	20 0	1820	A.html	Mozilla/3.01 (X11, I, IRIX6.2, IP22)
123.456.78.9	[25/Apr/1998:03:10:02 -0500	GET O.html HTTP/1.0	20 0	2270	F.html	Mozilla/3.01 (Win95, I)
123.456.78.9	[25/Apr/1998:03:10:45 -0500	GET J.html HTTP/1.0	20 0	9430	C.html	Mozilla/3.01 (X11, I, IRIX6.2, IP22)

KDD for WUM



Filtrage, nettoyage

- Status code (1xx : Informational, 2xx : Success, 3xx : Redirection, 4xx : Client Error, 5xx : Server Error)
- Requêtes automatiques (bots, performance monitoring systems)
- Supprimer toutes les entrées représentant des requêtes pour des objets graphiques, etc.
- Supprimer les entrées générées par des spiders/crawlers (google, etc.)

User Identification - Session Identification

- Comment identifier un utilisateur unique ? une transaction d'un utilisateur ?
- Pb :
 - Les identifiants des utilisateurs sont souvent supprimés par mesure de sécurité
 - Les adresses IP sont parfois cachées derrière des proxy
 - → web logs pas toujours fiables

User Identification - Session Identification

- Comment identifier un utilisateur unique ? une transaction d'un utilisateur ?
- Pb :
 - Les identifiants des utilisateurs sont souvent supprimés par mesure de sécurité
 - Les adresses IP sont parfois cachées derrière des proxy
 - → web logs pas toujours fiables
- Solutions : identification du client, cookies côté client, etc.

Identifier une session

Temps de navigation :

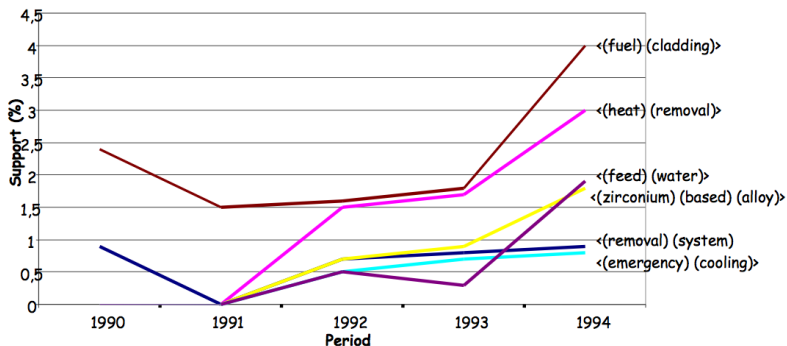
- Durée totale d'une session (≤ 30 minutes)
- Durée par page (≤ 10 minutes par page)

Sources de données :

- fichiers logs
- mais aussi cookies, bd d'utilisateurs (clients)

Analyse de tendance avec des motifs séquentiels

Patents « Induced Nuclear Reactions:
Processes, Systems and Elements »



Outline

- 1 Introduction
- 2 Définitions
- 3 Algorithmes
- 4 Web Usage Mining
- 5 Conclusions**

Nouveaux challenges

- privacy preserving mining (k anonymous sequential patterns)
- d'autres types de motifs (arbres, graphes)
- motifs séquentiels inattendus
- motifs séquentiels dans les flux de données.