

# Les Motifs séquentiels multidimensionnels étoilés

$M^2SP$  : une nouvelle approche

**Marc Plantevit**<sup>\*</sup>,<sup>1</sup> Anne Laurent<sup>\*</sup>,<sup>1</sup> Maguelonne Teisseire<sup>\*</sup>,<sup>1</sup> Dominique Laurent<sup>\*\*</sup>,<sup>1</sup> Y.W Choong<sup>\*\*</sup>,<sup>\*\*\*</sup>

\*LIRMM, Montpellier

\*\* LICP, Cergy Pontoise

\*\*\* HELP, Kuala Lumpur

20 mai 2005

---

<sup>1</sup>la virgule est un *opérateur commutatif n-aire*

# Plan

- 1 Introduction
- 2 Approches existantes et leurs limites
- 3 M<sup>2</sup>SP
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations
- 5 Conclusion

# Plan

- 1 **Introduction**
- 2 Approches existantes et leurs limites
- 3 M<sup>2</sup>SP
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations
- 5 Conclusion

# Introduction

## Motivations

- Vers une extraction d'information plus riche
- Développement d'OLAP
- limites des motifs séquentiels "classiques"

## exemple

- Motif séquentiel :  $\langle \{ \textit{parapluie}, \textit{écharpe} \}, \{ \textit{crème solaire} \} \rangle$
- Motif séquentiel multidimensionnel : Ce que l'on veut extraire  
 $\langle \{ ( \textit{Paris}, \textit{Professeur}, \textit{parapluie} ), ( \textit{Paris}, \textit{Professeur}, \textit{écharpe} ) \}$   
 $\{ ( \textit{Palavas}, \textit{Professeur}, \textit{crème solaire} ) \} \rangle$

# Plan

- 1 Introduction
- 2 Approches existantes et leurs limites**
- 3 M<sup>2</sup>SP
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations
- 5 Conclusion

# Les approches existantes abordant le problème des motifs séquentiels multidimensionnels

## *J Han, H Pinto & co*

- l'approche fondatrice
- recherche de motifs séquentiels *intra pattern*

## *Yu & Chen*

- extraction d'information sur des données spécifiques (weblog) organisées en différents niveaux d'agrégation

## *de Amo, Furtado & co*

- approche basée sur la logique temporelle

# J Han, H Pinto & co

## Motif séquentiel multidimensionnel $\beta$ :

$$\beta = \underbrace{\text{PATTERN}}_{\text{déf. sur 1 ou pls. dim.}} + \underbrace{\text{séquence d'items}}_{\text{déf. sur une seule dim.}}$$

## Exemple : Problème du "panier de la ménagère"

- Pattern :
  - informations relatives aux clients
  - Age, ville, etc
- Séquence d'items : les produits achetés
- (*Paris, Sportif*,  $\langle$ tee-shirt JO-2012, drapeau JO-2012 $\rangle$ )

# Les Patterns

- Pattern = n-uplet défini sur n dimensions  $\mathcal{D}_i$

$$t = (d_1, \dots, d_n) \quad d_i \in \text{Dom}(D_i) \cup \{*\}$$

- \* : un joker

## Principe

- Essai d'instancier le plus possible de dimensions afin d'obtenir des pattern maximalelement spécifiques (moins \* possible)



## 2 méthodes équivalentes d'extraction

### patterns fréquents puis séquences

- 1 Pattern  $p$   $\xrightarrow{\text{identifie}}$  un ensemble de transactions  $T$
- 2 Dans  $T$  on va extraire des motifs séquentiels classiques  $s$
- 3  $(p, s)$  est un motif séquentiel multidimensionnel

### méthode réciproque : séquences fréquentes puis patterns

- 1 Sur la dimension *Produits* extraire les motifs séquentiels classiques.
- 2 Matching des séquences fréquentes avec les patterns qui les identifient

# exemples

## exemples de motifs séquentiels multidimensionnels

$minsupp = 2$

- $\langle (business, middle), a \rangle$
- $\langle (chicago), bf \rangle$

cid	Cust-Grp	City	Age-grp	product-sequence
10	<b>business</b>	Boston	<b>middle</b>	$\langle (bd)cba \rangle$
20	professional	<b>Chicago</b>	young	$\langle bf(ce)(fg) \rangle$
30	<b>business</b>	<b>Chicago</b>	<b>middle</b>	$\langle (ah)abf \rangle$
40	education	New York	retired	$\langle (be)(ce) \rangle$

# Critique



- approche fondatrice
- permet une analyse plus fine



- pas de prise en compte de valeurs de mesure
  - que des information *intra pattern* (pas de corrélations possibles entre patterns différents )
- On ne pas extraire des informations du type :

(*N.Y., business, age\_middle, <Clou>*), (*N.Y., Retired, age\_old, <pneu>*)

# Plan

- 1 Introduction
- 2 Approches existantes et leurs limites
- 3 M<sup>2</sup>SP**
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations
- 5 Conclusion

# Base de données multidimensionnelles (ROLAP)

**Cellule** =  $\langle (d_1, \dots, d_n), \mu \rangle$

- $\forall i \in [1 \dots n], d_i \in \text{Dom}(D_i)$
- $\mu \in \text{Dom}(M), \mu$  est la mesure de la cellule.

## (Hyper-)Cube

Un hypercube  $C$  de données est un ensemble de cellules qui sont définies sur les mêmes dimensions.

$$C : \text{Dom}(D_1) \times \text{Dom}(D_2) \times \dots \times \text{Dom}(D_n) \longrightarrow M$$

# exemple

## Représentation d'un hypercube de données sous forme bidimensionnelle

Date	Cust-Grp	City	Age	Product	Measure
1	Educ	Chigago	A	clou	50
1	Educ	Chigago	B	pneu	2
1	Educ	Los Angeles	A	clou	30
1	Reti.	Miami	C	clou	20
1	Reti.	Miami	C	marteau	2
2	Educ	Chigago	B	rustine	10
2	Educ	Chigago	B	pneu	3
2	Educ	Los Angeles	A	clou	20
3	Educ	Los Angeles	B	rustine	15



## Sous-cube de données

- un cube de données peut être partitionné en différents **sous-cubes** en fonction de certaines dimensions

## Partition par rapport à *Cust-Grp* et *City*

Date	Cust-Grp	City	Age	Product	Measure
1	Educ	Chigago	A	clou	50
1	Educ	Chigago	B	pneu	2
2	Educ	Chigago	B	rustine	10
2	Educ	Chigago	B	pneu	3
1	Educ	Los Angeles	A	clou	30
2	Educ	Los Angeles	A	clou	20
3	Educ	Los Angeles	B	rustine	15
1	Reti.	Miami	C	clou	20
1	Reti.	Miami	C	marteau	2

# Données manipulées

## Partition des dimensions

$$D = \mathcal{D}_{\mathcal{F}} \oplus \mathcal{D}_{\mathcal{R}} \oplus \mathcal{D}_{\mathcal{A}} \oplus \{D_t\}$$

- $D_t$  pour la dimension temporelle
- $\mathcal{D}_{\mathcal{A}}$  pour les dimensions dites d'**analyse**
- $\mathcal{D}_{\mathcal{R}}$  pour les dimensions dites de **référence**
- $\mathcal{D}_{\mathcal{F}}$  pour les dimensions **ignorées**.



## Notation

Chaque cellule  $cell = \langle (d_1, \dots, d_n), \mu \rangle$  d'un cube peut être notée :

$$cell = \langle (f, r, a, t), \mu \rangle$$

avec :

- $f$  la restriction sur  $\mathcal{D}_{\mathcal{F}}$  de  $cell$
- $r$  la restriction sur  $\mathcal{D}_{\mathcal{R}}$  de  $cell$
- $a$  la restriction sur  $\mathcal{D}_{\mathcal{A}}$  de  $cell$
- $t$  la restriction sur  $D_t$  de  $cell$

# M<sup>2</sup>SP : item, itemset, séquence multidimensionnels

## Définitions

- item multidimensionnel :

$$e = \langle a, \mu \rangle$$

- itemset multidimensionnel :

$$i = \{e_1, \dots, e_p\}$$

- séquence multidimensionnelle :

$$s = \langle i_1, \dots, i_l \rangle$$

# Support d'une séquence

Un cube de données  $\mathcal{C}_c$  supporte une séquence  $s = \langle i_1, \dots, i_l \rangle$  si et seulement si :

$$\forall j = 1 \dots l, \exists d_j \in D_{i_j}, \quad \forall e = \langle a, \mu_e \rangle \in i_j, \\ \exists \text{cell} = \langle (f, a, r, d_j), \mu \rangle \in \mathcal{C}_c \quad \text{avec } d_1 \prec_t d_2 \prec_t \dots \prec_t d_l \\ \text{et } \mu_e = \mu$$

- **Support d'une séquence** Soient  $s$  un séquence, un cube DC composé d'un ensemble de sous-cubes  $\mathcal{C}_c$  construits à partir de  $\mathcal{D}_{\mathcal{R}}$  est :

$$\text{support}(s) = \frac{|\{\mathcal{C}_c \text{ supportant } s\}|}{|\{\mathcal{C}_c\}|}$$

- séquence fréquente
- k-fréquent

# Exemple récapitulatif

$$\mathcal{D}_{\mathcal{R}} = \{Cust - Grp, City\}$$

- Soient  $\mathcal{D}_{\mathcal{A}} = \{Age, Product\}$      $suppmin = 0.3$

- Sous-cubes définis sur  $\mathcal{D}_{\mathcal{R}}$  :

Cust-Grp	City
Educ	Chigago
Educ	Los Angeles
Reti.	Miami

- Soit la séquence

$$s = \langle \{(A, clou, 50)(B, pneu, 2)\}, \{(B, rustine, 10)\} \rangle$$

- cherchons le support

# Sous-cube 1

Date	Cust-Grp	City	Age	Product	Measure
1	Educ	Chigago	<b>A</b>	<b>clou</b>	<b>50</b>
1	Educ	Chigago	<b>B</b>	<b>pneu</b>	<b>2</b>
2	Educ	Chigago	B	pneu	3
2	Educ	Chigago	<b>B</b>	<b>rustine</b>	<b>10</b>

- 2 dates différentes

→ le sous-cube 1 supporte la séquence s

## exemple (suite)

### sous-cube 2 :

Date	Cust-Grp	City	Age	Product	Measure
1	Educ	Los Angeles	<b>A</b>	<b>clou</b>	30
2	Educ	Los Angeles	A	clou	20
3	Educ	Los Angeles	B	rustine	15

pas possible : trop peu de nuplets par date pour pouvoir supporter l'itemset de la séquence

### sous-cube 3 :

Date	Cust-Grp	City	Age	Product	Measure
1	Reti.	Miami	C	clou	20
1	Reti.	Miami	C	marteau	2

pas assez de date pour pouvoir supporter la séquence

## Support

$$\text{Support}(s) = \frac{1}{3} \longrightarrow \text{fréquente}$$

généralisation : les valeurs "jokers"

## vers de nouveaux horizons . . .

### Problème :

- $\langle \textit{clou}, \textit{Chicago}, \mu \rangle, \langle \textit{clou}, \textit{Miami}, \mu \rangle, \langle \textit{clou}, \textit{L.A.}, \mu \rangle$  non fréquents
- $\langle \textit{clou}, (\textit{Chicago} \vee \textit{Miami} \vee \textit{L.A.}), \mu \rangle$  fréquent
- Que faire ?

### Solution :

- Réduire ponctuellement les axes de recherche : \* un méta-symbole

généralisation : les valeurs "jokers"

# M<sup>2</sup>SP- $\alpha$ : Item multidimensionnel $\alpha$ -étoilé (\*)

Un item multidimensionnel  $\alpha$ -étoilé est de la forme :

$$e = \langle a, \mu \rangle \text{ où } \forall d_i \in a, d_i \in \text{Dom}(D_i) \cup \{*\}$$

- 1  $\forall e = \langle a, \mu \rangle, \exists d_{i_j} \in a \text{ tq } d_{i_j} \neq *$   
(au moins une dimension fixée)
- 2  $\forall d_{i_j} = *, \nexists x \in \text{Dom}(D_{i_j}) \text{ tq } e' = \langle a_{[d_{i_j}/x]}, \mu \rangle, \text{ support}(e') \geq \text{suppmin}$   
(rôle de seconde chance)



## \* et support

Un cube de données  $\mathcal{C}_c$  supporte une séquence

$s = \langle i_{s1}, \dots, i_{sI} \rangle$  si et seulement si :

$\forall j = 1 \dots I, \exists d_j \in D_t, \forall e = \langle (a_{i_1}, \dots, a_{i_m}), \mu_e \rangle \in i_j$

$\exists cell = \langle (f, (x_{i_1}, \dots, x_{i_m}), r, d_j), \mu \rangle \in \mathcal{C}_c,$

avec  $a_i = x_i$  ou  $a_i = *$  et  $d_1 \prec_t d_2 \prec_t \dots \prec_t d_I$

généralisation : les valeurs "jokers"

# M<sup>2</sup>SP- $\mu$ : Item multidimensionnel « $\mu$ -étoilé » ( $\ast$ )

## définition

Un item multidimensionnel  $\mu$ -étoilé est de la forme :

$$e = \langle a, \mu \rangle$$

où

- $\mu \in \text{Dom}(M) \cup \ast$

## $\ast$ et support

Un cube de données  $\mathcal{C}_c$  supporte une séquence  $s = \langle i_{s_1}, \dots, i_{s_l} \rangle$  si et seulement si :

$$\forall j = 1 \dots l, \exists d_j \in D_t, \forall e = \langle a, \mu_e \rangle \in i_j, \exists \text{cell} = \langle (f, a, \mu) \in \mathcal{C}_c,$$

$$\text{avec } \mu_e = \ast \text{ ou } \mu_e = \mu$$

# Algorithmes

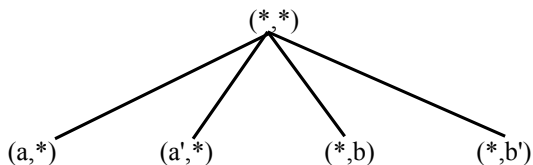
## Génération des items fréquents

- obtention des **items maximale**ment spécifiques
- génération par niveau

## Génération des séquences fréquentes

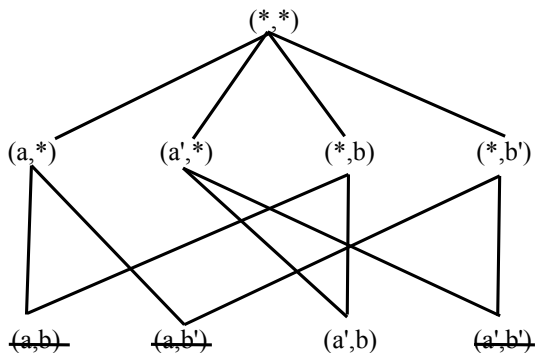
- méthode générer/élaguer
- utilisation de l'algorithme psp

# Exemple : génération des items fréquents

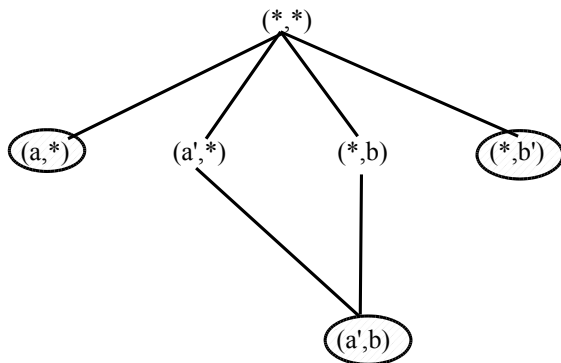




# Exemple : génération des items fréquents



# Exemple : la bordure



# Calcul du support

- prérequis : prétraitement des données (cube by *date*,  $D_1, \dots, D_n$ )
- calcul du support d'une séquence
  - pour chaque sous-cube :  
**supportCount**( $s, DC, \mathcal{D}_{\mathcal{R}}, comptage$ )
  - dans un sous-cube donné : **supportCube**( $s, C, comptage$ )
- ancrage ( $\sigma_{condition}(C) \mapsto C'$  avec  $C' \subseteq C$ )

## complexité

- $n_C$  est le nombre de cellules du cube  $C$
- $m = |\mathcal{D}_{\mathcal{A}}|$  est le nombre de dimensions des items multidimensionnels.
- supportCount :  $O(n_C \times m \times \log n_C)$ .
- supportCube :  $O(l \times n_{max} \times m \times \log n_{max})$



# Plan

- 1 Introduction
- 2 Approches existantes et leurs limites
- 3 M<sup>2</sup>SP
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations**
- 5 Conclusion



## généralisation des motifs séquentiels "classiques"

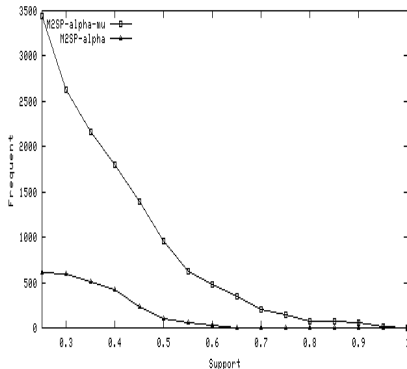
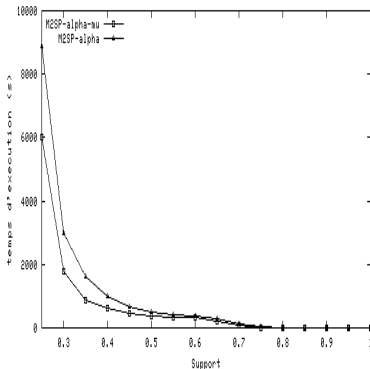
- $\mathcal{D}_{\mathcal{R}}$  = dimension relative aux clients
- $\mathcal{D}_{\mathcal{A}}$  = dimension relative aux produits
- $D_t$  = la date
- $\mathcal{D}_{\mathcal{F}}$  = le reste

## généralisation

- tous les fréquents extraits par J.Han, peuvent être extraits par  $M^2SP_{\alpha-\mu}$
- on en trouve plus (*inter pattern*)

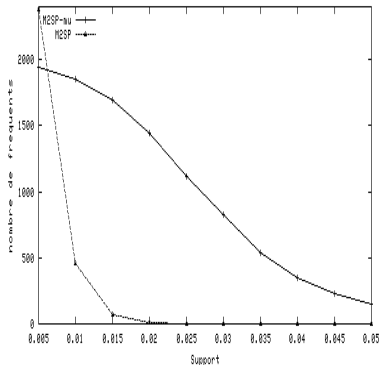
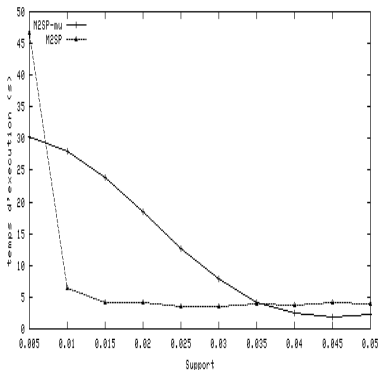


# Support M<sup>2</sup>SP- $\alpha(-\mu)$



- La mesure : une contrainte

# Support M<sup>2</sup>SP(- $\mu$ )

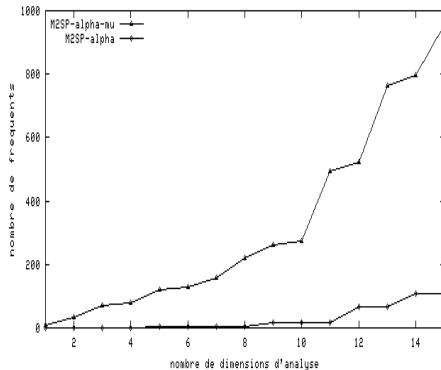


- Une contrainte supplémentaire : pour des extractions sur des données spécifiques



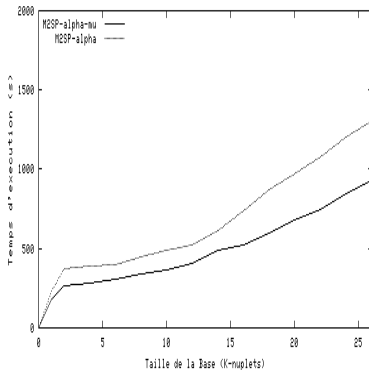
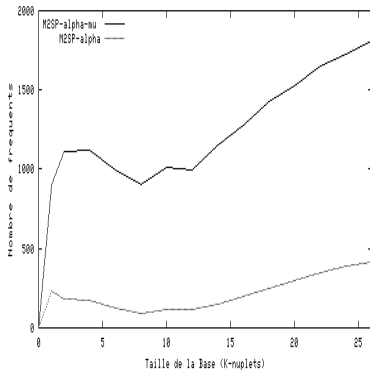
## Propriété

- le nombre de fréquents croît en fonction du nombre de dimensions





# Passage à l'échelle



# Plan

- 1 Introduction
- 2 Approches existantes et leurs limites
- 3 M<sup>2</sup>SP
  - Notre modèle de données
  - définitions et formalismes
  - généralisation : les valeurs "jokers"
  - Algorithmes et exemples
- 4 Résultats théoriques et expérimentations
- 5 **Conclusion**

# Conclusion et perspectives

## une approche plus multidimensionnelle et plus générale

- *inter cube*
- prise en compte des quantités possibles ou non(⊗)
- plusieurs types de motifs séquentiels
- une plus grande liberté dans le choix des axes

## perspectives de travail

- utilisation de la théorie des sous-ensembles flous
- choix automatique des axes de références et d'analyse
- relation d'ordre (spatiale, spatio-temporelle, ...)
- classifieur multidimensionnel
- motifs séquentiels multidimensionnels généralisés (contraintes de temps, ...)