



INSTITUT NATIONAL
DE RECHERCHE
EN INFORMATIQUE
ET EN AUTOMATIQUE



centre de recherche PARIS - ROCQUENCOURT

Analyse des données évolutives : application aux données d'usage du Web

Alzenny GOMES DA SILVA

Directeur de thèse : Pr Edwin DIDAY
Co-directeur de thèse : Dr Yves LECHEVALLIER

24 septembre 2009

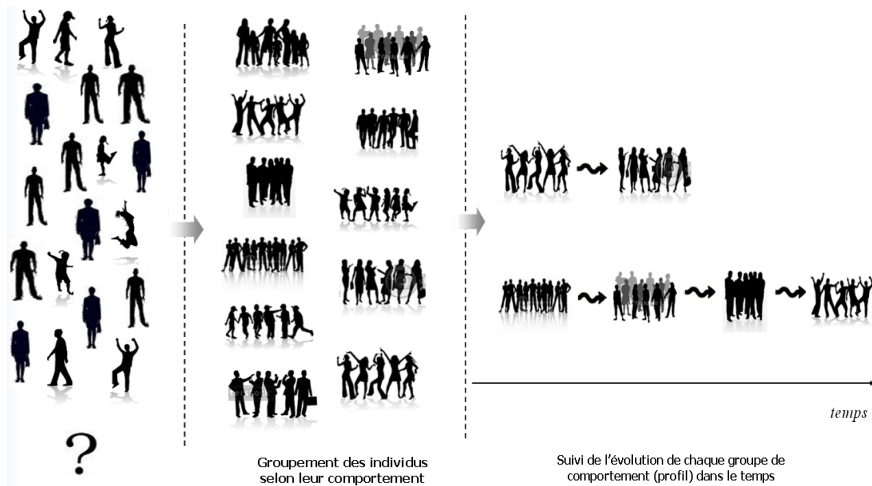
- 1 Introduction
- 2 Analyses exploratoires des stratégies de classification
- 3 Approche de classification pour la détection et le suivi des changements sur des données évolutives
- 4 Expérimentation
- 5 Conclusion

Contexte/Motivation

- Le **Web** : source de données volumineuses et dynamiques
- L'accès aux pages Web a une nature dynamique
 - Changement du contenu et de la structure du site Web
 - Changement **d'intérêt/comportement** des utilisateurs
- Le comportement d'accès d'un internaute peut dépendre :
 - de l'heure et du jour de la semaine
 - des événements saisonniers
 - des événements ponctuels (crises, catastrophes, compétitions sportives)
 - etc.



Notre problématique



La fouille de données d'usage du Web (Web Usage Mining, WUM)

Ensemble de techniques basées sur la fouille de données pour analyser l'usage d'un site Web (Cooley *et al.*, 1999).

Le log Web

Fichier enregistré par le serveur Web contenant les traces d'usage laissées par les internautes lors des connections

- l'adresse IP du client
- la date et l'heure de la requête
- la page consultée
- le code de statut (valeur 200 en cas de réussite)
- le nombre d'octets transmis
- la page précédemment visitée
- le navigateur Web et le système d'exploitation de la machine cliente

L'unité base de notre analyse

Navigation : ensemble de requêtes (clics) proches en provenance d'un même utilisateur

Positionnement de la thèse



- Soit :

E : ensemble de données à classifier

$H^{(b)}$: sous-ensemble de E tel que $H^{(b)} \neq \emptyset$,

$$\forall i \neq j : H^{(i)} \cap H^{(j)} = \emptyset \text{ et } \bigcup_{b=1}^B H^{(b)} = E$$

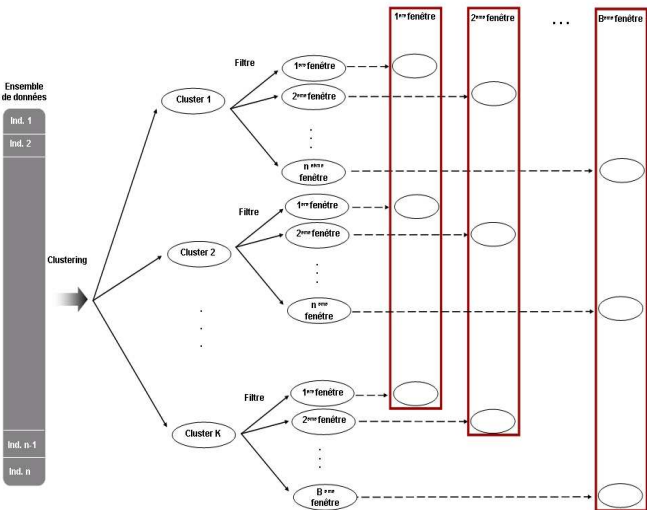
Partitionnement des données par paquets (fenêtres)

- 1 Partitionnement par nombre d'effectifs constant** (fenêtre logique)
on fixe le nombre d'individus qui doivent être contenus dans chaque fenêtre.
- 2 Partitionnement par intervalle de temps constant** (fenêtre temporelle)
on fixe un intervalle de temps durant lequel les données analysées seront enregistrées dans une fenêtre, par exemple 30 minutes, 2 heures, 1 jour, etc.

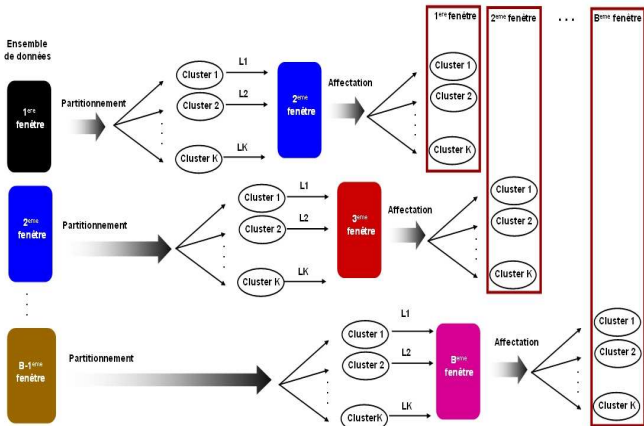
Stratégies de classification analysées

- 1 Classification globale (CG)
- 2 Classification locale précédente (CGL_1)
- 3 Classification locale dépendante (CGL_2)
- 4 Classification locale indépendante (CLI)

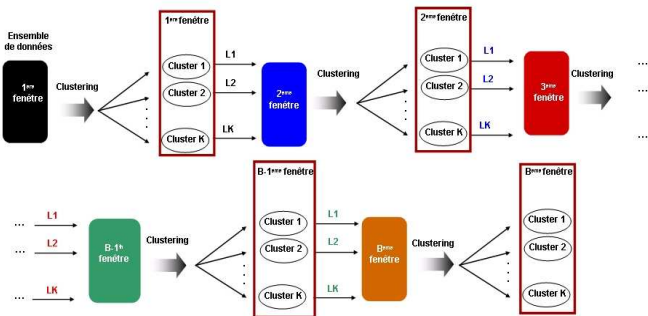
Classification globale (CG)



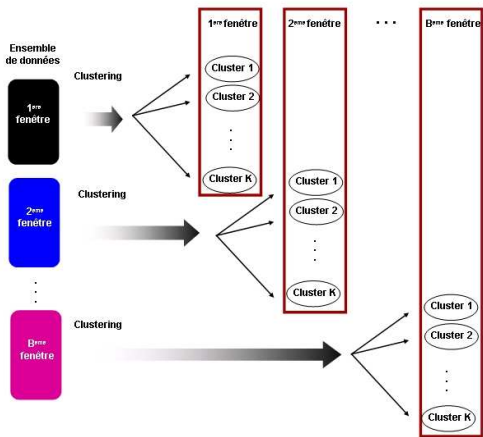
Classification locale précédente (CGL_1)



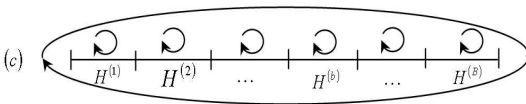
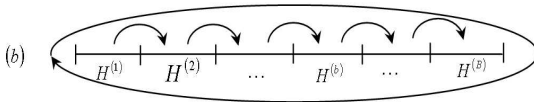
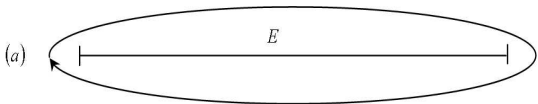
Classification locale dépendante (CGL_2)



Classification locale indépendante (CLI)



Comparatif des stratégies de classification



(a) Classification **globale** CG sur l'ensemble E [MND (Diday, 1971)]

(b) Classification **locale précédente** CGL_1 sur les paquets $H^{(b)}$ de E [MNDS (Diday, 1975)]

(c) Classification **locale dépendante** CGL_2 sur les paquets $H^{(b)}$ de E [MNSO]

Analyse d'un site Web académique¹

Aperçu du site Web du Centre d'Informatique (CIn) de l'UFPE :

URL analysée

Menu de options

Home

- The Informatics' center
- Mission
- Objectives
 - In the informatics' center
 - In the university
 - In the region
 - In the country
 - In the world
- Organization chart
- Departments
- Courses
- Resolutions
- Localization
- Tunnel of the time

People

- Direction
- Coordinators
- Professors
- Support
- Management
- Administrative

Graduation

- Information of the graduation
- How to be a student
- How to ingress in the pos-graduation
- Coordinator/Secretariat

Pos-graduation

- Registrations and Information
 - Subjects for Mester 2005
 - Address
 - Documents
 - How Internal regimen
 - Internal regimen (11/1/1996)
 - Internal regimen (02/1992)
 - Norms of qualification and theses proposition
 - School registration 2/2000
 - School registration 1/2000
 - School registration 1/2001
 - Disciplines summaries
 - Masters results
 - School registration and pre-registration
 - Orientation of School registration
 - PhD results.

- Description
- Professors group
- Lines of research
- Specialization
 - Information
 - Concentration areas
 - Course's Structure
 - Professors group
 - Computational Resources
- Masters
- PhD
- APC
- Administrative
- Financial resources
- Events

Extension

- Partners
- Events
- Services
- UFPE for all

Research

- Groups of research
- Publications
 - Theses and dissertations
 - Technician reports
 - Tutorials, manuals and class notes
 - Book chapters
 - Complete books
 - Complete papers in periodics
 - Complete Works in proceedings
 - Summary in proceedings
- Projects of research
- Events
 - Events related to the Informatics' center
 - Events of interest
- Honors and prizes

Infrastructure

Services

News

- Releases
- Clipping
- Contacts

Phones

- Professors
- Sectors
- Employees

1. <http://www.cin.ufpe.br/>

Site Web du Centre d'Informatique (CIn) de l'UFPE

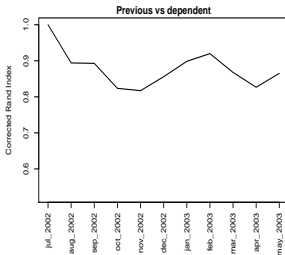
- Composition :
 - pages statiques (pages personnelles, pages de support de cours, etc.)
 - 90 pages pages dynamiques gérées par des *servlets* en Java
- Période analysée :
 - du 01 juillet 2002 à 00 :10 :25 jusqu'au 15 mai 2003 à 13 :22 :27

Paramètres de l'analyse

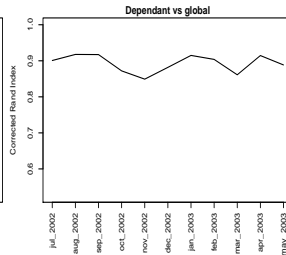
distance :	Euclidienne
nombre de clusters K :	10
nombre d'initialisations :	50
type de fenêtre :	temporelle de taille égale à 1 mois
données d'entrée :	tableau <i>navigations</i> × <i>variables statistiques</i>
total d'individus :	138 536 navigations

Valeurs de RC entre les résultats des stratégies de classification deux-à-deux :

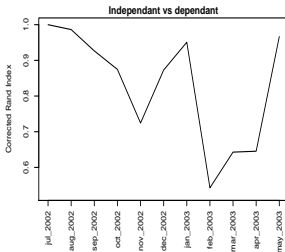
MNDS sur $H^{(b)}$
versus
MNDSO sur $H^{(b)}$



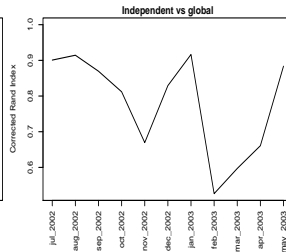
MND sur E
versus
MNDSO sur $H^{(b)}$



MND sur $H^{(b)}$ (CLI)
versus
MNDSO sur $H^{(b)}$

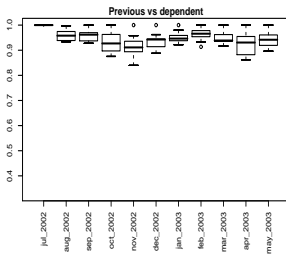


MND sur $H^{(b)}$ (CLI)
versus
MND sur E

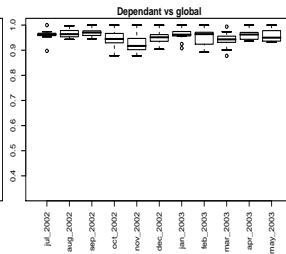


Valeurs de F-mesure entre les résultats des stratégies de classification deux-à-deux :

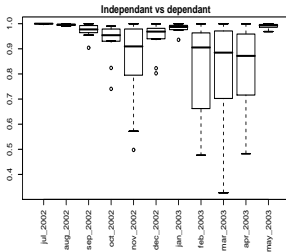
MNDS sur $H^{(b)}$
versus
MNDSO sur $H^{(b)}$



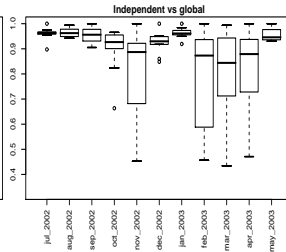
MND sur E
versus
MNDSO sur $H^{(b)}$



MND sur $H^{(b)}$ (CLI)
versus
MNDSO sur $H^{(b)}$

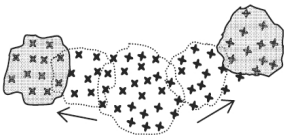


MND sur $H^{(b)}$ (CLI)
versus
MND sur E

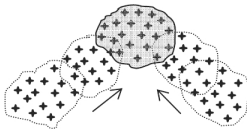


Interprétation des changements

- **Scission ou division d'un cluster** : un certain nombre d'individus appartenant à un cluster de la partition P_1^t migre vers d'autres clusters dans la partition P_2^t



- **Fusion ou absorption de deux ou plusieurs clusters** : les individus de différents clusters de la partition P_1^t migrent vers un même cluster de la partition P_2^t



Interprétation des changements

- **Disparition d'un cluster** : un cluster de la partition P_1^t n'est plus présent dans la partition P_2^t

cluster obsolète



- **Apparition d'un cluster** : un cluster inexistant dans la partition P_1^t émerge dans la partition P_2^t



nouveau cluster



Détermination du nombre de classes

- Question largement exploitée dans la littérature d'apprentissage automatique (Milligan & Cooper, 1985) (Dubes, 1987) (Halkidi & Vazirgiannis, 2001)
- *"Une classification ne peut pas être variée ou fautive, ni probable ou improbable, mais seulement profitable ou non profitable"* (Williams & Lance, 1965)
- Découverte du **vrai** nombre de classes ?
 - Nombre de clusters **acceptable/utile** dans un but précis
- Segmenter les individus de manière à faire ressortir des préférences d'usage (pages Web, produits de commerce, etc.)

Les indices de détermination du nombre de classes comparés dans notre étude :

- 1 la pseudo-statistique de Calinski et Harabasz (CH) [ou $G1^*$ sur R]
- 2 l'indice de Baker et Hubert (BH) [ou $G2^*$ sur R]
- 3 l'indice de Hubert et Levin (HL) [ou $G3^*$ sur R]
- 4 l'indice Silhouette (S)
- 5 l'indice de Davies et Bouldin (DB)
- 6 l'indice de Krzanowski et Lai (KL)
- 7 l'indice de Hartigan (H)
- 8 la statistique Gap (G)
- 9 stratégie de détermination du nombre de clusters basée sur les dérivées première et seconde (D)

*Package **clusterSim** de **R** :

<http://cran.r-project.org/web/packages/clusterSim>

Stratégie de détermination du nombre de classes

- 1 Appliquer le SOM (Kohonen, 1995) avec une grille contenant 200 neurones initialisés par une ACP (Elemento, 1999)
- 2 Appliquer une CAH avec le critère de Ward sur les prototypes (neurones) finaux de la SOM (Murtagh, 1995)
- 3 Tracer le graphique des gains d'inertie intra-classe obtenus à chaque itération de l'algorithme CAH
- 4 Le nombre de classes à retenir sera obtenu par :
 - L'application des différents indices de détermination du nombre de classes [1-8]
 - Le **coude** dans la décroissance des valeurs du gain d'inertie intra-classe, repéré à l'aide des dérivées (différences) premières et secondes [D] (Lebart *et al.*, 1995)

- 1 Introduction
- 2 Analyses exploratoires des stratégies de classification
- 3 Approche de classification pour la détection et le suivi des changements sur des données évolutives
 - Méthodologie de génération de données
 - Application de l'approche sur des données artificielles
- 4 Expérimentation
- 5 Conclusion

Expérimentation sur des données artificielles

Méthodologie de génération de données

Un algorithme de base de génération des données artificielles d'usage

```

1 Pour c dans l'intervalle[1,K] {
2   Pour chaque individu i de la classe c {
3      $x_i \leftarrow \{0, \dots, 0\}$ 
4      $nbClic \leftarrow \text{random}(nbClicMIN, nbClicMAX)$ 
5     Pour y dans l'intervalle[1, nbClic] {
6        $z \leftarrow \text{random}(0, 1)$ 
7       Pour j dans l'intervalle[1, p] {
8         Si ( $z \leq \text{cumul}(p_c^j)$ ) {
9            $x_i^j \leftarrow x_i^j + 1$ 
10          Sortir de la boucle la plus interne
11        }
12      }
13    }
14  }
15 }

```

Exemple :

$$c_1 = (p_1^1, p_1^2, p_1^3) = (0.3, 0.4, 0.3)$$

$$\text{cumul}(c_1) = \begin{array}{|c|c|c|} \hline p_1^1 & (p_1^1 + p_1^2) & (p_1^1 + p_1^2 + p_1^3) \\ \hline 0.3 & 0.7 & 1.0 \\ \hline \end{array}$$

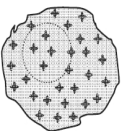
$$\left. \begin{array}{l} nbClic = 10 \\ \rightarrow \end{array} \right\} \begin{cases} z = \text{random}(0, 1) \\ P(z \leq \text{cumul}(p_c^j)) \end{cases}$$

Expérimentation sur des données artificielles

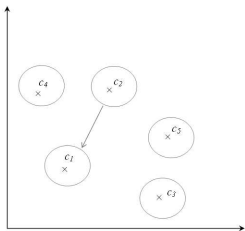
Méthodologie de génération de données

Trois algorithmes spécialisés dans la simulation de changements sur des données artificielles

- 1 Changement lié à l'effectif des classes
 - Rétrécissement d'une classe (facteur β_1)
 - Grossissement d'une classe (facteur β_2)



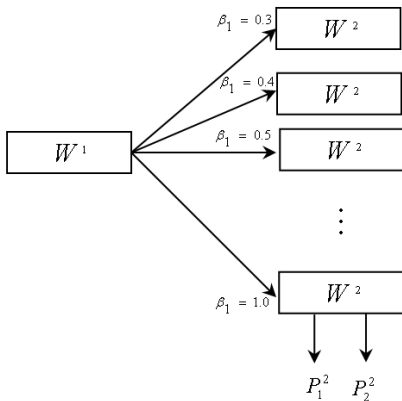
- 2 Changement lié à la position des classes (facteur β_3) dans l'espace des données



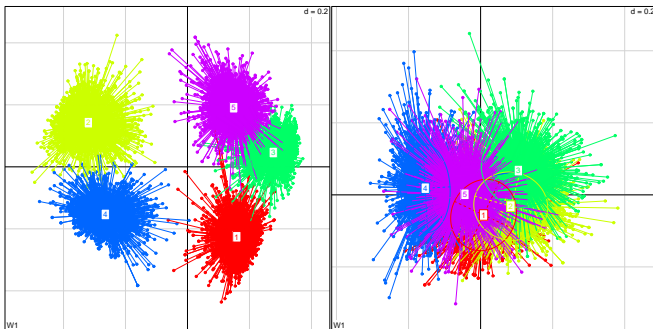
Expérimentation sur des données artificielles

Méthodologie de génération de données

Génération de la fenêtre W^2 à partir de W^1 (cas de rétrécissement) :



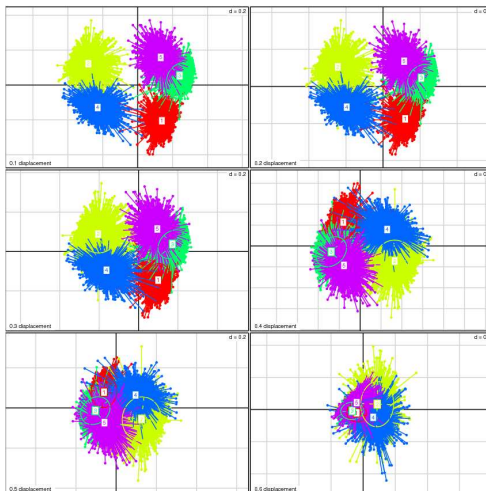
Projection des classes artificielles contenues dans W^1



classes bien séparées

classes recouvrantes

Simulation de l'effet de rapprochement des classes



Expérimentation sur des données artificielles

Le tableau suivant discrimine le nombre de clusters suggéré par chaque indice pour chaque fenêtre de données analysée.

Fenêtre	CH	BH	HL	KL	DB	H	G	S	D
1	5	5	10	4	5	5	5	5	5
2	5	30	5	4	5	5	8	5	5
3	5	30	6	4	5	5	6	5	5
4	2	30	14	4	30	5	6	4	5
5	2	30	9	4	30	5	5	5	5
6	2	30	9	9	29	5	8	7	5
7	10	30	9	9	30	2	2	10	4
%correct	42.85	14.28	14.28	0.00	42.85	85.71	28.57	57.14	85.71

Meilleure performance présentée par :

- l'indice de Hartigan
- la stratégie basée sur des dérivées

Expérimentation sur des données réelles

- Centre d'informatique de l'UFPE (Recife, Brésil)²
- Structure dynamique contenant 90 pages organisées sur 10 thèmes : *The Informatics' center, People, Graduation, Pos-graduation, Extension, Research, Infrastructure, Services, News et Phones.*

Paramètres

méthode de classification :	SOM (Kohonen, 1995) initialisée par une ACP suivie d'une CAH avec le critère de Ward
distance :	Euclidienne
type de fenêtre :	temporelle de taille égale à 1 mois
données d'entrée :	tableau <i>navigations</i> × <i>thèmes de pages</i>
nombre de thèmes de pages :	10
filtre appliqué :	navigations ayant un minimum de 10 clics
total d'individus :	28 220 navigations

2. <http://www.cin.ufpe.br/>

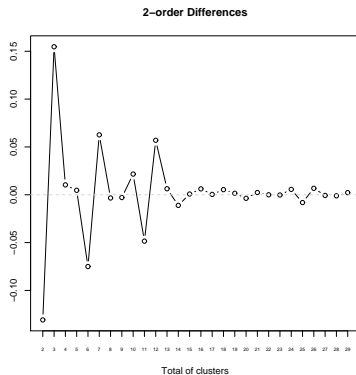
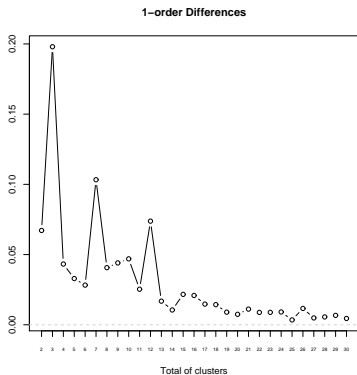
Expérimentation sur des données réelles

Nombre de classes suggéré par les différents indices pour le jeu de données d'usage du CIn-UFPE :

Fenêtre	CH	BH	HL	KL	DB	H	G	S	D
1	30	30	6	2	29	4	3	3	3
2	3	30	2	2	30	5	5	21	3
3	2	2	6	7	2	3	2	2	5
4	2	2	11	3	2	4	2	2	4
5	30	30	12	9	26	2	2	29	3
6	29	30	8	8	30	3	2	30	3
7	2	2	13	24	2	10	3	2	2
8	30	30	8	8	30	2	2	30	4
9	30	25	11	2	30	3	3	25	3
10	30	30	2	4	25	5	3	30	3
11	30	30	7	4	28	2	2	30	5

Résultat obtenu par le critère des dérivées pour la première fenêtre analysée

Les différences premières et secondes des inerties intra-classes relatives des partitions contenant de 2 à 30 classes :



Conclusion de l'étude de cas sur les indices de détermination du nombre de classes

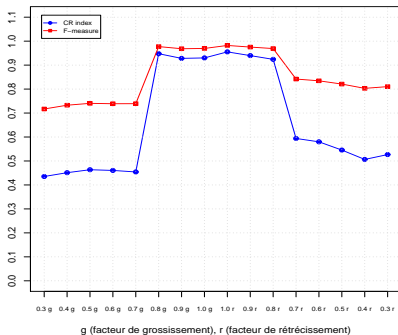
- Les différents indices ne sont pas toujours d'accord sur le nombre de clusters à retenir sur un même jeu de données
 - le choix du nombre de clusters s'avère être une tâche délicate
- Après plusieurs analyses sur différents jeux de données :
 - La **stratégie basée sur les dérivées** semble être très efficace

- 1 Introduction
- 2 Analyses exploratoires des stratégies de classification
- 3 Approche de classification pour la détection et le suivi des changements sur des données évolutives
 - Méthodologie de génération de données
 - Application de l'approche sur des données artificielles
- 4 Expérimentation
- 5 Conclusion

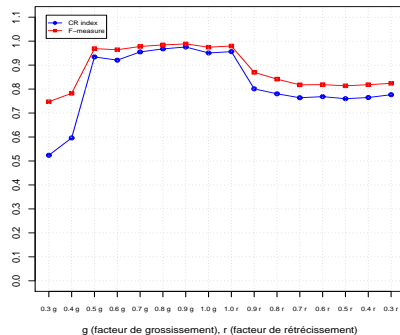
Application de l'approche sur des données artificielles

Grossissement et rétrécissement de la classe cible

Valeurs de l'indice corrigé de Rand et de la F-mesure :



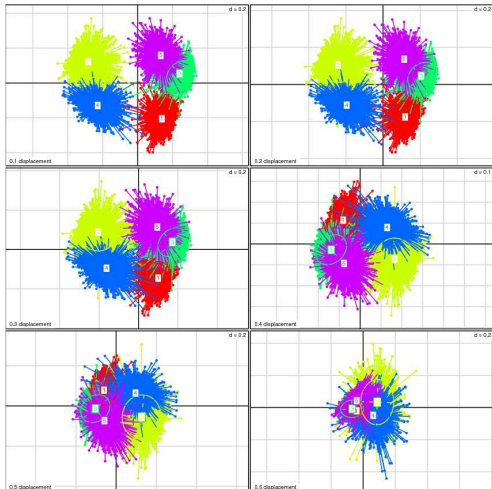
classes bien séparées



classes recouvrantes

Application de l'approche sur des données artificielles

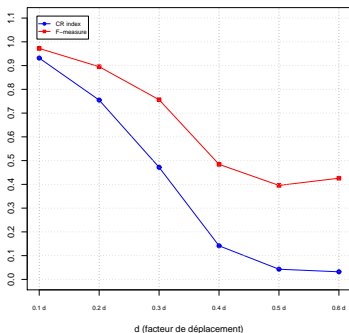
Simulation de l'effet de déplacement des classes



Application de l'approche sur des données artificielles

Simulation de l'effet de déplacement des classes

Valeurs de l'indice corrigé de Rand et de la F-mesure :

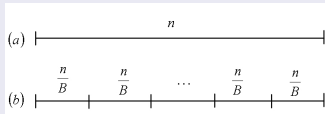


Conclusion de l'étude de cas sur les données artificielles

- L'approche proposée se montre très efficace pour la détection des changements
- Les deux indices utilisés pour la comparaison de partitions (Rand corrigé et F-measure) se montrent capables de repérer les changements dans le temps
 - L'indice corrigé de Rand étant plus sensible aux changements que la F-mesure

Caractérisation et avantages de l'approche de classification

- Stratégie du type *diviser pour régner*
- Réduction de la complexité de la méthode de classification. Ex. : $O(n^2) \rightarrow BO\left(\left(\frac{n}{B}\right)^2\right)$



- Indépendance de la méthode de classification
- Processus de détection des changements basé sur l'extension
- Application des algorithmes non-incrémentaux dans un processus incrémental
- Utilisation d'un passé récent (mémoire)
- Résumé des données au cours du temps

Schéma de la base de données

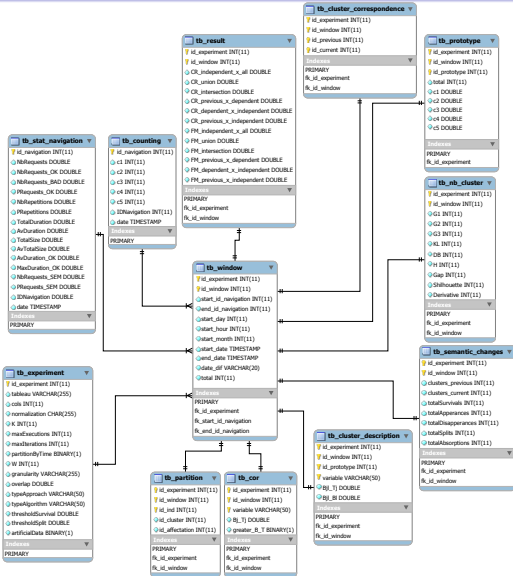


Schéma de la base de données

N°	Variable	Signification
1	NbRequests	Nombre de clics effectués durant la navigation
2	NbRequests_OK	Nombre de requêtes réussies (statut = 200) dans la navigation
3	NbRequests_Bad	Nombre de requêtes échouées (statut ≠ 200) dans la navigation
4	PRequests_OK	Pourcentage de requêtes réussies (= $NbRequests_OK / NbRequests$)
5	NbRepetitions	Nombre de requêtes répétées dans la navigation
6	PRepetitions	Pourcentage de répétitions (= $NbRepetitions / NbRequests$)
7	DureeTotale	Durée totale de la navigation (en secondes)
8	MDuree	Moyenne de la durée des requêtes (= $DureeTotale / NbRequests$)
9	MDuree_OK	Moyenne de la durée des requêtes réussies (= $DureeTotale_OK / NbRequests_OK$)
10	NbRequests_Sem	Nombre de requêtes liées aux pages de la structure sémantique du site
11	PRequests_Sem	Pourcentage de requêtes liées aux pages de la structure sémantique du site (= $NbRequests_Sem / NbRequests$)
12	TotalSize	Total d'octets transférés durant la navigation
13	MSize	Moyenne d'octets transférés durant la navigation (= $TotalSize / NbRequests_OK$)
14	DureeMax_OK	Durée maximale parmi les durées des requêtes réussies

- tb_result
- tb_partition
- tb_cluster_correspondance
- tb_prototype
- tb_nb_cluster
- tb_semantic_changes
- tb_cluster_description
- tb_cor

tb_stat_navigation	
id_navigation	INT(11)
NbRequests	DOUBLE
NbRequests_OK	DOUBLE
NbRequests_BAD	DOUBLE
PRequests_OK	DOUBLE
NbRepetitions	DOUBLE
PRepetitions	DOUBLE
TotalDuration	DOUBLE
AvDuration	DOUBLE
TotalSize	DOUBLE
AvTotalSize	DOUBLE
AvDuration_OK	DOUBLE
MaxDuration_OK	DOUBLE
NbRequests_SEM	DOUBLE
PRequests_SEM	DOUBLE
IDNavigation	DOUBLE
date	TIMESTAMP
Indexes	
PRIMARY	

tb_experiment	
id_experiment	INT(11)
tableau	VARCHAR(255)
cols	INT(11)
normalization	CHAR(255)
K	INT(11)
maxExecutions	INT(11)
maxIterations	INT(11)
partitionByTime	BINARY(1)
W	INT(11)
granularity	VARCHAR(255)
overlap	DOUBLE
typeApproach	VARCHAR(50)
typeAlgorithm	VARCHAR(50)
thresholdSurvival	DOUBLE
thresholdSplit	DOUBLE
artificialData	BINARY(1)
Indexes	
PRIMARY	

tb_counting	
id_navigation	INT(11)
c1	INT(11)
c2	INT(11)
c3	INT(11)
c4	INT(11)
c5	INT(11)
IDNavigation	INT(11)
date	TIMESTAMP
Indexes	
PRIMARY	

ID	theme1	theme2	...	themej	...	themej-1	themej	date-heure
1	c1 ¹	c2 ¹	c ^{j-1} 1	c ^j 1	11/MM/AAAA hh:mm:ss
2	c2 ²	c2 ²	c ^{j-1} 2	c ^j 2	11/MM/AAAA hh:mm:ss
...
i	c ⁱ	c ⁱ	c ^{j-1} i	c ^j i	11/MM/AAAA hh:mm:ss
...

tb_window	
id_experiment	INT(11)
id_window	INT(11)
start_id_navigation	INT(11)
end_id_navigation	INT(11)
start_day	INT(11)
start_hour	INT(11)
start_month	INT(11)
start_date	TIMESTAMP
end_date	TIMESTAMP
date_diff	VARCHAR(20)
total	INT(11)
Indexes	
PRIMARY	
fk_id_experiment	
fk_start_id_navigation	
fk_end_id_navigation	

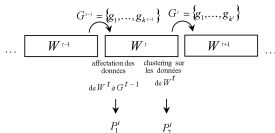
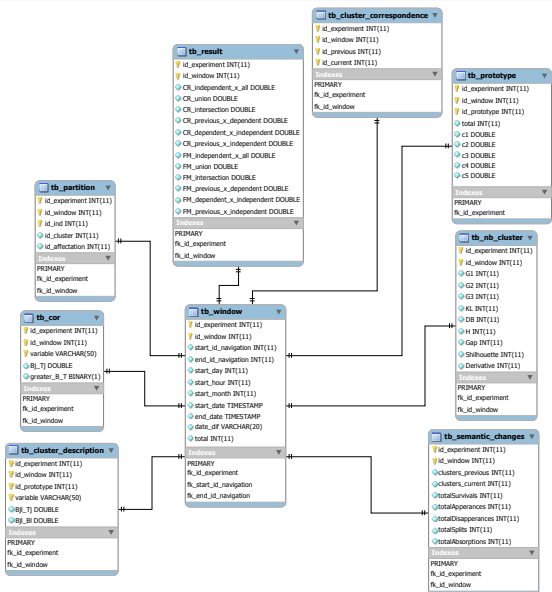


Schéma de la base de données



Analyse d'un jeu de données issu du *marketing*

- Suivi des achats d'un panel de consommateurs mis à disposition dans le cadre du SLDS 2009³
 - 10 068 clients
 - 2 marchés de biens de consommation
 - 6 marques de produits
- Période analysée :
 - 14 mois (du 09 juillet 2007 jusqu'au 08 septembre 2008)
- Champs descriptifs :

Champ	Signification
ident	identification client
date	premier jour de la semaine de l'achat
marché	marché concerné (marche_1, marche_2)
marque	marque du produit acheté (A, B, C, D, E, F)
valeur	valeur des achats

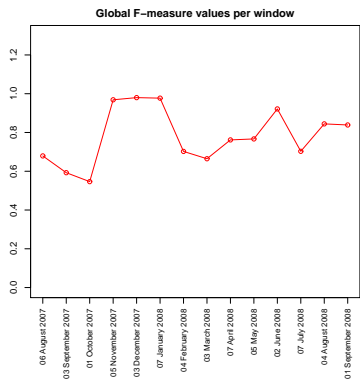
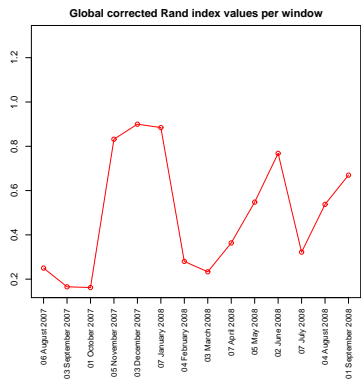
3. <http://www.ceremade.dauphine.fr/SLDS2009>

Analyse d'un jeu de données issu du *marketing*

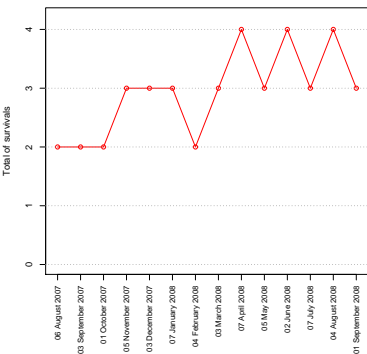
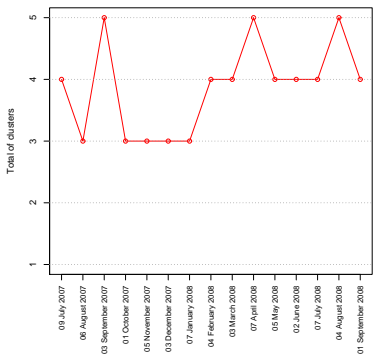
Paramètres

- méthode de classification : SOM (Kohonen, 1995) initialisée par une ACP suivie d'une CAH avec le critère de Ward, puis du critère basé sur les dérivées
- distance : Euclidienne
- nombre de dimensions : 6 (marques de produit)
- type de fenêtre : temporelle de taille égale à 1 mois
- données d'entrée : tableau croisé *client x marques de produit*
- total d'individus : 262 215

- Valeurs obtenues par l'indice de Rand corrigé (à gauche) et par la F-mesure (à droite)
- Quelques périodes d'**instabilité** : mois de *Octobre 2007*, *Mars 2008* et *Juillet 2008*.



- D'*Octobre 2007* jusqu'à *Janvier 2008* : le nombre total de clusters est égal à 3.
- Ces mois correspondent à une période de **grande stabilité** selon montrent les valeurs élevées de l'**indice de Rand** et de la **F-measure**



Trois premières variables les plus significatives des clusters de comportement issus du mois d'*octobre 2007*.

id cluster	variable	$\frac{B_l^j}{T_j}$	$\frac{B_l^j}{B_l}$	$\frac{B^j}{T_j}$	$> \frac{B}{T}$
1	F	2.23e-006	0.44	3.56e-005	oui
1	B	2.19e-006	0.43	3.11e-005	oui
1	A	2.18e-007	0.04	1.78e-006	non
2	F	3.32e-005	0.93	3.56e-005	oui
2	A	1.31e-006	0.03	1.78e-006	non
2	E	4.24e-007	0.01	1.04e-006	non
3	B	2.88e-005	0.93	3.11e-005	oui
3	D	8.01e-007	0.02	1.38e-006	non
3	E	4.69e-007	0.01	1.04e-006	non

Conclusion de l'analyse de données du *marketing*

- Possibilité d'application de l'approche de détection et suivi des changements sur d'autres domaines

- Pour des stratégies de CRM, le profil d'un client pourrait être interprété comme suit :
 - Si le client reste dans un même cluster au cours du temps : client *fidèle* aux marques de produit
 - Si le client change considérablement de cluster au cours du temps : client *zappeur* entre les marques de produit

- ① Introduction
- ② Analyses exploratoires des stratégies de classification
- ③ Approche de classification pour la détection et le suivi des changements sur des données évolutives
- ④ Expérimentation
- ⑤ Conclusion**

Apports

- ✓ Deux modélisations des données d'usage du Web
 - 1 Caractérisation du **mode de navigation** de l'internaute
 - 2 Caractérisation du **centre d'intérêt** de l'internaute
- ✓ Une méthodologie de **générations** de données d'usage et de **simulation** de changements
- ✓ Analyse théorique et expérimentale des méthodes de classification MND, MNDS et **MNDSO**.

Apports

- ✓ Une approche de classification permettant la détection et le suivi des changements sur des données évolutives
 - Opère sur de fenêtrés non recouvrantes de taille préfixée
 - Indépendante de la méthode de classification
 - Applique des indices de comparaison de partition basés sur l'extension
 - Intègre l'interprétation des changements repérés

- ✓ Applicabilité de l'approche à d'autres domaines
Ex. : Surveillance de matériels (tâche WP3.2, projet MIDAS)

- ✓ Outil d'analyse mis en ligne :
 - <http://atwueda.gforge.inria.fr/>

Publications majeures

- DA SILVA, Alzenny, LECHEVALLIER, Yves, ROSSI, Fabrice, DE CARVALHO, Francisco. [Clustering Dynamic Web Usage Data](#). In Nadia Nedjah, Luiza Mourelle and Janusz Kacprzyk (Editors), Springer Berlin/Heidelberg, Series : Studies in Computational Intelligence, *Innovative Applications in Data Mining*, Vol. 169, pages 71-82, ISBN : 978-3-540-88044-8, 2009.
- DA SILVA, Alzenny, LECHEVALLIER, Yves, DE CARVALHO, Francisco. [Monitoring Data Changes through a Clustering Approach](#). In *International Federation of Classification Societies (IFCS 2009)*. Dresden, March 13-18, 2009.
- DA SILVA, Alzenny. [Diverses approches permettant l'introduction du temps dans la fouille de données d'usage du Web](#). Editeurs : Chantal Reynaud et Gilles Venturini. *Numéro spécial sur la fouille du Web de la Revue des Nouvelles Technologies de l'Information (RNTI W-1)*, pages 35-55, ISBN 978.2.85428.793.6, cépaduès éditions, 2007.
- DA SILVA, Alzenny. [Analyzing the Evolution of Web Usage Data](#). In *Special issue on Data Stream Analysis of MODULAD* (Monde des Utilisateurs de L'Analyse de Données), numéro 36, pages 75-84, May, 2007.
- DA SILVA, Alzenny, LECHEVALLIER, Yves, ROSSI, Fabrice, DE CARVALHO, Francisco. [Construction and Analysis of Evolving Data Summaries : an Application on Web Usage Data](#). In Luiza Mourelle, Nadia Nedjah and Janusz Kacprzyk editors, *VII IEEE International Conference on Intelligent Systems Design and Applications (ISDA 2007)*, Pages 377-380, ISBN : 978-0-7695-2976-9, IEEE Computer Society, Rio de Janeiro, Brazil, 22-24 October, 2007.
- DA SILVA, Alzenny, DE CARVALHO, Francisco, LECHEVALLIER, Yves, TROUSSE, Brigitte. [Mining Web Usage Data for Discovering Navigation Clusters](#). In : *XI IEEE Symposium on Computers and Communications (ISCC 2006)*, pages 910-915, ISBN ISSN :1530-1346, 0-7695-2588-1, Mining Web Usage Data for Discovering Navigation Clusters, IEEE Computer Society, Pula-Cagliari, Italy, 2006.
- DA SILVA, Alzenny, DE CARVALHO, Francisco, LECHEVALLIER, TROUSSE, Brigitte. [Characterizing Visitor Groups from Web Data Streams](#). In : *IEEE International Conference on Granular Computing (GrC 2006)*, Atlanta, USA, 2006, pages 389-392, ISBN : 1-4244-0134-8, IEEE Computer Society, 2006.

Comparaison de deux partitions

Tableau de contingence entre les partitions P_1^t et P_2^t :

clusters de P_1^t	clusters de P_2^t					
	C_1	\dots	C_j	\dots	C_k	
C_1	n_{11}	\dots	n_{1j}	\dots	n_{1k}	$n_{1.}$
\vdots						
C_i	n_{i1}	\dots	n_{ij}	\dots	n_{ik}	$n_{i.}$
\vdots						
C_m	n_{m1}	\dots	n_{mj}	\dots	n_{mk}	$n_{m.}$
	$n_{.1}$		$n_{.j}$		$n_{.k}$	$n_{..} = n$

Indices basés sur l'extension : F-mesure

Rappel et Précision

$$\text{rappel}(C_i, C_j) = \frac{n_{ij}}{\sum_{j=1}^k n_{ij}} ; \text{precision}(C_i, C_j) = \frac{n_{ij}}{\sum_{i=1}^m n_{ij}}$$

F-mesure entre deux clusters

$$F_{\text{measure}}(C_i, C_j) = \frac{2 \text{precision}(C_i, C_j) \text{rappel}(C_i, C_j)}{\text{precision}(C_i, C_j) + \text{rappel}(C_i, C_j)}$$

F-mesure entre deux partitions

$$F(P_1^t, P_2^t) = \frac{1}{n} \sum_{i=1}^m n_i \cdot \max_{j=1, \dots, k} F_{\text{measure}}(C_i, C_j)$$

Indices basés sur l'extension : indice corrigé de Rand

Indice corrigé de Rand (RC) entre deux partitions

$$RC(P_1^t, P_2^t) = \frac{\sum_{i=1}^m \sum_{j=1}^k \binom{n_{ij}}{2} - \binom{n}{2}^{-1} \sum_{i=1}^m \binom{n_{i.}}{2} \sum_{j=1}^k \binom{n_{.j}}{2}}{\frac{1}{2} \left[\sum_{i=1}^m \binom{n_{i.}}{2} + \sum_{j=1}^k \binom{n_{.j}}{2} \right] - \binom{n}{2}^{-1} \sum_{i=1}^m \binom{n_{i.}}{2} \sum_{j=1}^k \binom{n_{.j}}{2}}$$

Variables descriptives des navigations

Le nombre d'accès aux pages Web (page hit count) a long temps été utilisé comme indicateur majeur des préférences des internautes (Yan *et al.*, 1996).

N°	Variable	Signification
1	NbRequests	Nombre de clics effectués durant la navigation
2	NbRequests_OK	Nombre de requêtes réussies (statut = 200) dans la navigation
3	NbRequests_Bad	Nombre de requêtes échouées (statut ≠ 200) dans la navigation
4	PRequests_OK	Pourcentage de requêtes réussies (= $NbRequests_OK / NbRequests$)
5	NbRepetitions	Nombre de requêtes répétées dans la navigation
6	PRepetitions	Pourcentage de répétitions (= $NbRepetitions / NbRequests$)
7	DureeTotale	Durée totale de la navigation (en secondes)
8	MDuree	Moyenne de la durée des requêtes (= $DureeTotale / NbRequests$)
9	MDuree_OK	Moyenne de la durée des requêtes réussies (= $DureeTotale_OK / NbRequests_OK$)
10	NbRequests_Sem	Nombre de requêtes liées aux pages de la structure sémantique du site
11	PRequests_Sem	Pourcentage de requêtes liées aux pages de la structure sémantique du site (= $NbRequests_Sem / NbRequests$)
12	TotalSize	Total d'octets transférés durant la navigation
13	MSize	Moyenne d'octets transférés durant la navigation (= $TotalSize / NbRequests_OK$)
14	DureeMax_OK	Durée maximale parmi les durées des requêtes réussies

Expérimentation sur des données artificielles

Méthodologie de génération de données

Paramètre	Valeurs utilisées lors des expérimentations
p	10
nbWindow	2
nbClicMIN	10
nbClicMAX	50
K	5
totalIndividuals	10 000
c'	1 (changement d'effectifs) et 2 (déplacement)
β_1 (facteur de rétrécissement)	{0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}
β_2 (facteur de grossissement)	{0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}
β_3 (facteur de déplacement)	{0.1, 0.2, 0.3, 0.4, 0.5, 0.6}

Deux cas de figures :

- ① Classes bien séparées
- ② Classes recouvrantes

References I



AGGARWAL, CHARU C. 2005.

On change diagnosis in evolving data streams.

IEEE transactions on knowledge and data engineering, **17**(5),
587–600.



AGGARWAL, CHARU C., HAN, JIAWEI, WANG, JIANYONG,
& YU, PHILIP S. 2003.

A framework for clustering evolving data streams.

*Pages 81–92 of : Vldb '2003 : Proceedings of the 29th
international conference on very large data bases.*

VLDB Endowment.

References III

New York, NY, USA : Springer-Verlag New York, Inc.



CELEUX, G., DIDAY, E., GOVAERT, G., LECHEVALLIER, Y., & RALAMBONDRAINY, H. 1989.

Classification automatique des données.

Dunod, Paris.



CHEN, XIAODONG, & PETROUNIAS, ILIAS. 1999.

Mining temporal features in association rules.

Pkdd'99 : Proceedings of the third european conference on principles of data mining and knowledge discovery, 295–300.



CIAMPI, A., & LECHEVALLIER, Y. 2000.

Clustering large, multi-level data sets : An approach based on kohonen self organizing maps.

Principles of data mining and knowledge discovery, springer berlin / heidelberg, 1910/2000, 161–177.

References IV



COOLEY, ROBERT, MOBASHER, BAMSHAD, & SRIVASTAVA, JAIDEP. 1999.

Data preparation for mining world wide web browsing patterns.

Journal of knowledge and information systems, 1(1), 5–32.



CORVAISIER, FRANÇOISE, MILLE, ALAIN, & PINON, JEAN MARIE. 1997.

Information retrieval on the world wide web using a decision making system.

Pages 284–295 of : Proceedings of the computer-assisted searching on the internet (ria0 97).

References V



CSERNEL, B. 2008.

Résumé généraliste de flux de données.

Ph.D. thesis, ENST Paris.



DIDAY, E. 1971.

Une nouvelle méthode en classification automatique et reconnaissance des formes la méthode des nuées dynamiques.

Revue de statistique appliquée, **19**(2), 19–33.



DIDAY, E. 1975.

Classification automatique séquentielle pour grands tableaux.

Revue française d'automatique, informatique et recherche opérationnelle (rairo), 29–61.

References VI



DUBES, RICHARD C. 1987.

How many clusters are best?—an experiment.

Pattern recogn., **20**(6), 645–663.



ELEMENTO, O. 1999.

Apport de l'analyse en composantes principales pour l'initialisation et la validation de cartes topologiques de kohonen.

In : Actes des 7èmes journées de la société francophone de classification (sfc'99).

References VII



GANTI, VENKATESH, GEHRKE, JOHANNES, RAMAKRISHNAN, RAGHU, & LOH, WEI-YIN. 1999.
A framework for measuring changes in data characteristics.
Pages 126–137 of : In pods.
ACM Press.



GANTI, VENKATESH, GEHRKE, JOHANNES, & RAMAKRISHNAN, RAGHU. 2000.
Demon : Mining and monitoring evolving data.
Pages 439–448 of : lee transactions on knowledge and data engineering.

References X



LEBART, L., MORINEAU, A., & PIRON, M. 1995.
Statistique exploratoire multidimensionnelle.
 Dunod.



LIU, BING, MA, YIMING, & LEE, RONNIE. 2001.
 Analyzing the interestingness of association rules from the
 temporal dimension.
*Icdm '01 : Proceedings of the 2001 ieee international
 conference on data mining, 377–384.*



MALEK, MARIA, & KANAWATI, RUSHED. 2001.
 Cobra : A cbr-based approach for predicting users actions in a
 web site.
*Case-based reasoning research and development : 4th
 international conference on case-based reasoning (iccbbr 2001),
 336–346.*

