



Mesurer l'intérêt des règles d'association

Benoît Vaillant*, Patrick Meyer***, Elie Prudhomme**,
Stéphane Lallich**, Philippe Lenca*, Sébastien Bigaret*

*GET ENST Bretagne / Département LUSSI – UMR 2872 CNRS TAMCIC
{*prenom.nom*}@enst-bretagne.fr

**Laboratoire ERIC - Université Lumière - Lyon 2
lallich@univ-lyon2.fr

***Faculté de Droit, d'Economie et de Finance, Université du Luxembourg,
patrick.meyer@uni.lu

Atelier DKQ – EGC 05, 18 Janvier 2005



Résumé/Contexte

- ▶ Action Spécifique STIC Fouille de Bases de Données (GAFODONNÉES), animée par R. Cicchetti et M. Sebag (2002)
- ▶ Groupe de travail sur les Mesures de Qualité (GAFOQUALITÉ), animé par F. Guillet (2002)
- ▶ laboratoires LUSI (ENST Bretagne) et ERIC (Université Lyon 2) engagent une collaboration sur le thème de l'intérêt des règles d'association

Synthèse des travaux réalisés sur ce thème par les deux parties.



Plan

Synthèse des travaux réalisés sur ce thème par les deux parties.

- ▶ Vingt mesures, huit propriétés.
- ▶ Contrôle du risque.
- ▶ Etude formelle des mesures selon les propriétés.
- ▶ Etude de comportement, en situation (plate-forme HERBS)
- ▶ Confrontation des typologies formelle et expérimentale.
- ▶ Assistance au choix d'une mesure de qualité.



Pour des mesures d'intérêt des règles d'association

Qu'est-ce qu'une règle d'association ?

- ▶ focalise sur les coprésences [Agrawal et al., 1993]
- ▶ règle d'association, implication, équivalence et corrélation

Règles admissibles :

- ▶ Support
- ▶ Confiance

Extraction des règles :

- ▶ But : extraire les règles admissibles en totalité
- ▶ Algorithme : Apriori et dérivés

Problème : trop de règles, certaines sans intérêt

Une nécessité : mesurer l'intérêt des règles d'association

Mesures

Soit $n = |E|$, le nombre total d'enregistrements

Pour $A \rightarrow B$, on note :

$n_a = |A|$, le nombre d'enregistrements vérifiant A.

$n_b = |B|$, le nombre d'enregistrements vérifiant B.

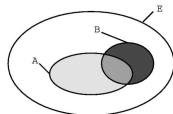
$n_{ab} = |A \cap B|$, le nombre d'exemples de la règle.

$n_{a\bar{b}} = |A \cap \bar{B}|$, le nombre de contre-exemples à la règle.

$A \rightarrow B$ est évaluée à l'aide de mesures généralement **monotones décroissantes** en fonction de $n_{a\bar{b}} = n_a - n_{ab}$. $A \rightarrow B$ est jugée intéressante selon la mesure μ lorsque $\mu(A \rightarrow B) \geq \alpha$, α devant être fixé par l'utilisateur.

Pour $X \subset E$, on remplace n_X/n par p_X lorsque l'on considère les fréquences relatives plutôt que les fréquences absolues.

Légende :





Mesures

Mesure	Abréviation	Définition
support	SUP	$\frac{n_a - n_{a\bar{b}}}{n}$
confidence	CONF	$1 - \frac{n_{a\bar{b}}}{n_a}$
linear correlation coefficient	R	$\frac{nn_{ab} - n_a n_b}{\sqrt{nn_a n_b n_{\bar{a}} \cdot n_{\bar{b}}}}$
centred confidence	CONF _{CEN}	$\frac{nn_{ab} - n_a n_b}{n_a n_{\bar{b}}}$
conviction	CONV	$\frac{nn_{ab}}{n_a n_{\bar{b}}}$
Piatetsky-Shapiro	PS	$\frac{1}{n} \left(\frac{nn_{ab}}{n_a n_{\bar{b}}} - n_{a\bar{b}} \right)$
Loevinger	LOE	$1 - \frac{n_{a\bar{b}}}{n_a n_{\bar{b}}}$
information gain	GI	$\log\left(\frac{nn_{ab}}{n_a n_b}\right)$
Sebag-Schoenauer	SEB	$\frac{n_a - n_{a\bar{b}}}{n_a n_{\bar{b}}}$
lift	LIFT	$\frac{nn_{ab}}{n_a n_b}$



Mesures

Mesure	Abréviation	Définition
Laplace	LAP	$\frac{n_{ab}+1}{n_a+2}$
least contradiction	MoCo	$\frac{n_{ab}-n_{a\bar{b}}}{n_b}$
odd multiplier	MC	$\frac{(n_a-n_{a\bar{b}})n_{\bar{b}}}{n_b n_{a\bar{b}}}$
example and counter example rate	TEC	$\frac{n_a-2n_{a\bar{b}}}{n_a-n_{a\bar{b}}}$
Kappa	IQC	$2 \frac{nn_a-hn_{a\bar{b}}-n_a n_b}{nn_a+nn_b-2n_a n_b}$
Zhang	ZHANG	$\frac{nn_{a\bar{b}}-n_a n_b}{\max\{n_{ab}n_{\bar{b}}, n_b n_{a\bar{b}}\}}$
implication index	-INDIMP	$\frac{nn_{a\bar{b}}-n_a n_b}{\sqrt{nn_a n_b}}$
intensity of implication	INTIMP	$P\left[\text{poisson}\left(\frac{n_a n_{\bar{b}}}{n}\right) \geq n_{a\bar{b}}\right]$
entropic intensity of implication	IIE	$\left\{ \left[(1-h_1\left(\frac{n_{a\bar{b}}}{n}\right))^2 \right] \times \left[(1-h_2\left(\frac{n_{a\bar{b}}}{n}\right))^2 \right] \right\}^{1/4}_{\text{INTIMP}}$
probabilistic discriminant index	IPD	$P\left[\mathcal{N}(0,1) > \text{INDIMP}^{\text{CR}/\mathcal{B}}\right]$



Propriétés

Distinguer $A \rightarrow B$ et $A \rightarrow \bar{B}$.

Expression des mesures en fonction de la confiance.

Relation d'équivalence *classer comme ...*

	Sémantique	Valeurs	Compétence
g_1	non symétrie	2	E_a
g_2	décroissance avec n_b	2	E_a
g_3	situation à l'indépendance	2	E_a
g_4	situation à la règle logique	2	E_a
g_5	non-linéarité autour de 0^+	3	E_r
g_6	prise en compte de n	2	E_r
g_7	facilité à fixer un seuil	2	E_a et E_r
g_8	intelligibilité	3	E_r

[Lallich, 2002], [Gras et al., 2004], [Lenca et al., 2004]



Contrôle du risque multiple,

[Teytaud et Lallich, 2001], [Lallich et al., 2004]

- ▶ règle significative : $p_{b/a}$ significativement supérieur à p_b
- ▶ on teste $H_0 : p_{b/a} = p_b$ contre $H_1 : p_{b/a} > p_b$
- ▶ statistique de test : $r(A, B)$, où $nr^2 = \chi^2(A, B)$
- ▶ sous $H_0 : r(A, B) \approx N(0, 1/\sqrt{n})$
- ▶ règle de décision : rejet H_0 au risque $\alpha_0 \iff r(A, B) \times \sqrt{n} > u_{1-\alpha_0}$
- ▶ problème : inflation de fausses découvertes à 5%

Nécessité : contrôle du risque multiple !

Pour 10 000 règles, même si aucune n'est pertinente, on doit s'attendre à sélectionner 500 règles à tort



Contrôle du risque multiple

Réa. \ Déc.	Acc.	Rej.	Tot.		Réa. \ Déc.	Acc.	Rej.	Tot.
H_0 vraie	$1 - \alpha$	α	1	\implies	H_0 vraie	U	V	m_0
H_1 vraie	β	$1 - \beta$	1		H_1 vraie	T	S	m_1
					total	W	R	m

- ▶ nouveau critère : V_0 fausses découvertes au risque δ
 $FWER = P(V > 0) \longrightarrow UAFWER = P(V > V_0)$,
- ▶ algorithme : BS_FD fondé sur le bootstrap pour évaluer la marge de fluctuation de la mesure et assurer $UAFWER = \delta$
- ▶ complexité : en $\mathcal{O}(lmn)$, où l constante due au bootstrap
- ▶ expériences : jusqu'à 50% de règles éliminées suivant la base

Perspectives

Règles différentiellement exprimées.



HERBS : une plate-forme d'expérimentation

Développements à poursuivre . . .

- ▶ HERBS V1.0 [Vaillant, 2002], GAFoQUALITÉ
- ▶ HERBS V2.0 [Vaillant et al., 2005], démonstration à EGC 05

Deux types d'analyse :

- ▶ intra-mesure
- ▶ inter-mesures

Post-analyse :

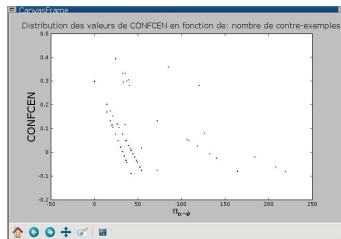
- ▶ données et règles sont des entrées
- ▶ format C4.5, APRIORI et CSV



HERBS : une plate-forme d'expérimentation

Etude d'une mesure :

- ▶ évaluation des objets au moyen de quelques grandeurs : nombre de cas et de règles, taux de couverture, indice de recouvrement, et nombre de règles *particulières*
- ▶ sélection de l'ensemble des N meilleures règles
- ▶ tracé de la distribution des valeurs prises par la mesure

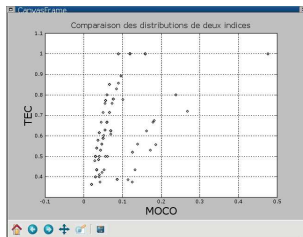




HERBS : une plate-forme d'expérimentation

Comparaison de mesures :

- ▶ extraction de l'ensemble des règles classées k fois parmi les N meilleures par p mesures.
- ▶ comparaison des préordres induits par deux mesures.
- ▶ tracé des distributions croisées des valeurs de deux mesures.





Typologie formelle

Classification formelle des mesures :

- ▶ construction d'une matrice de distance entre les mesures (matrice de décision sous forme disjonctive complète – carré des distances euclidiennes)
- ▶ classification ascendante hiérarchique avec le critère de WARD

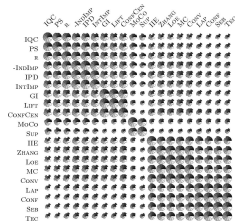
Quatre classes principales :

- ▶ {PS, IQC, GI, CONF CEN, LIFT, R, -INDIMP, IPD}
- ▶ {INTIMP, IIE, LOE, ZHANG, MC, CONV}
- ▶ {CONF, SEB, TEC}
- ▶ {LAP, MoCo, SUP}



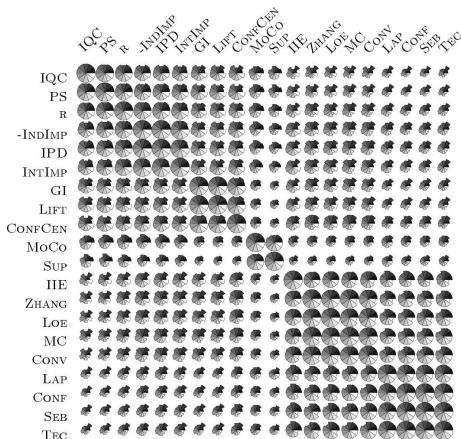
Typologie expérimentale

- ▶ coefficients d'accords entre préordres
- ▶ 10 bases (UCI Repository of machine learning databases)
- ▶ réarrangement de l'ordre des lignes et des colonnes afin de mieux mettre en évidence les structures de blocs (AMADO, [Chauchat et Risson, 1998])
- ▶ classification expérimentale des mesures de qualité





Typologie expérimentale





Confrontation

Du point de vue expérimental, seule la classe 3 présente de forts désaccords avec les résultats formels :

For. \ Exp.	Classe 1	Classe 2	Classe 3	Classe 4
Classe 1	PS, IQC, GI CONFCEN, LIFT, R -INDIMP, IPD			
Classe 2	INTIMP	IIE	LOE, ZHANG, MC, CONV	
Classe 3			CONF, SEB, TEC	
Classe 4			LAP	MoCo, SUP



Assistance au choix des mesures

Sélectionner les *bonnes* règles c'est aussi sélectionner les *bonnes* mesures [Lenca et al., 2002], [Lenca et al., 2003] :

- ▶ Aide Multicritère à la Décision
- ▶ six éléments définissant le contexte : l'ensemble de données, l'ensemble de règles, l'ensemble de mesures, l'ensemble de propriétés des mesures, l'ensemble de préférences de l'utilisateur, l'ensemble de critères de décision [Lenca et al., 2004]
- ▶ diverses méthodes : sous-ensemble de bonnes mesures, classement des mesures

Assistance au choix des mesures

	g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8		
SUP	0	0	0	0	1	0	1	2		
CONF	1	0	0	1	1	0	1	2		
R	0	1	1	0	1	0	1	1		
CONF CEN	1	1	1	0	1	0	1	2		
PS	0	1	1	0	1	1	1	1		
LOE	1	1	1	1	1	0	1	1		
ZHANG	1	1	1	1	2	0	0	0		
- INDIMP	1	1	1	0	1	1	1	0	g_1	non symétrie
LIFT	0	1	1	0	1	0	1	1	g_2	décroissance avec n_b
MOCO	1	1	0	0	1	0	1	1	g_3	situation à l'indépendance
SEB	1	0	0	1	0	0	1	1	g_4	situation à la règle logique
MC	1	1	1	1	0	0	1	2	g_5	non-linéarité autour de 0^+
CONV	1	1	1	1	0	0	1	1	g_6	prise en compte de n
TEC	1	0	0	1	2	0	1	1	g_7	facilité à fixer un seuil
IQC	0	1	1	0	1	0	1	0	g_8	intelligibilité
GI	0	1	1	0	2	0	1	0		
INTIMP	1	1	1	1	2	1	1	0		
IIE	1	1	1	1	2	1	0	0		
IPD	1	1	1	0	1	1	1	0		
LAP	1	0	0	0	1	0	1	0		

Matrice de décision



Assistance au choix des mesures

- ▶ deux scénarios utilisateur (tolérance **–Sc1–** ou non **–Sc2–** de l'apparition de contre-exemples dans les règles); [Gras et al., 2004]
- ▶ rangements totaux pour **Sc1** et **Sc2** – poids égaux PROMETHEE [Brans et Mareschal, 1994]

Rang:	1	2	3	4	5	6	7
Sc1:	INTIMP	LOE	MC	CONF	CONV	-INDIMP, IPD	
Sc2:	MC	CONV	LOE	CONF	INTIMP	-INDIMP, IPD	
Rang:	8	9	10	11	12	13	14
Sc1:	IIE, ZHANG		PS	TEC	CONF	GI	R, LIFT
Sc2:	PS	SEB	CONF	R, LIFT		MoCo	IIE
Rang:	15	16	17	18	19	20	
Sc1:		MoCo	SEB	IQC	SUP	LAP	
Sc2:	ZHANG	IQC	TEC	SUP	GI	LAP	



Conclusion

- ▶ Double approche à la fois formelle et expérimentale.
- ▶ Diverses extensions à développer :
 - ▶ capacité discriminante d'une mesure
 - ▶ résistance au bruit
 - ▶ autres mesures . . .
- ▶ Une collaboration fructueuse !



Références



Agrawal, R., Imielinski, T., et Swami, A. (1993).

Mining association rules between sets of items in large databases.

In Buneman, P. et Jajodia, S., editors, *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, pages 207–216, Washington, D.C.



Brans, J. et Mareschal, B. (1994).

The PROMETHEE-GAIA decision support system for multicriteria investigations.

Investigation Operativa, 4(2):102–117.



Chauchat, J.-H. et Risson, A. (1998).

Visualization of Categorical Data, chapter 3, pages 37–45.

Blasius J. & Greenacre M. ed.

New York : Academic Press.



Gras, R., Couturier, R., Bernadet, M., Blanchard, J., Briand, H., Guillet, F., Kuntz, P.,

Lehn, R., et Peter, P. (2004).

Quelques critères pour une mesure de qualité de règles d'association - un exemple : l'intensité d'implication.

Revue des Nouvelles Technologies de l'Information.

A paraître.



Lallich, S. (2002).

Mesure et validation en extraction des connaissances à partir des données.

Habilitation à Diriger des Recherches – Université Lyon 2.

