

L'extraction de règles de dépendance bien définies entre ensembles de variables multivaluées

Frédéric Pennerath*,**

*UMI 2958 GeorgiaTech - CNRS

**Supélec

2 rue Édouard Belin, 57070 Metz

Résumé. Cet article étudie la faisabilité et l'intérêt de l'extraction de règles de dépendance entre ensembles de variables multivaluées en comparaison du problème bien connu de l'extraction des règles d'association fréquentes. Une règle de dépendance correspond à une dépendance fonctionnelle approximative caractérisée principalement par l'entropie conditionnelle associée. L'article montre comment établir une analogie formelle entre les deux familles de règles et comment adapter à l'aide de cette analogie l'algorithme « Eclat » afin d'extraire d'un jeu de données les règles de dépendance dites bien définies. Une étude expérimentale conclut sur les forces et inconvénients des règles de dépendance bien définies vis-à-vis des règles d'association fréquentes.

1 Introduction

Dans de nombreux problèmes de fouille de motifs, les données à analyser forment une collection d'échantillons décrits par différents attributs multivalués appelés *variables* dans ce qui suit. Un tel jeu de données est illustré par l'exemple simplifié de la figure 1(a). Les motifs

	V ₁	V ₂	V ₃	V ₄
D ₁	a	1	oui	1
D ₂	c	2	oui	1
D ₃	b	1	non	1
D ₄	c	2	non	1

(a) Avec variables multivaluées

	V ₁ = a	V ₁ = b	V ₁ = c	V ₂ = 1	V ₂ = 2	V ₃ = yes	V ₃ = no	V ₄ = 1
D ₁	×			×		×		×
D ₂			×		×	×		×
D ₃		×		×			×	×
D ₄			×		×		×	×

(b) Avec items ou attributs monovalués

FIG. 1 – *Interprétation mono ou multivaluée d'un même jeu de données.*

les plus souvent recherchés dans ces données sont des *itemsets* : un itemset est un ensemble d'attributs monovalués appelés *items* représentant chacun une des valeurs possibles d'une des variables multivaluées. Un jeu de données défini sur un ensemble d'attributs multivalués peut être représenté de façon équivalente par des items comme illustré sur la figure 1(b). Un itemset *couvre* une donnée si tous les items de l'itemset sont aussi des items de la donnée. La *fréquence* d'un itemset est la proportion des données couvertes par celui-ci. Une question consiste alors