

# Classification d'un tableau de contingence et modèle probabiliste

G rard Govaert\*, Mohamed Nadif\*\*

\*Heudiasyc, UMR CNRS 6599, Universit  de Technologie de Compi gne,  
BP 20529, 60205 Compi gne Cedex, France  
gerard.govaert@utc.fr

\*\*LITA, Universit  de Metz, Ile du Saulcy, 57045 Metz Cedex, France  
mohamed.nadif@univ-metz.fr

**R sum .** Ces derni res ann es, la classification crois e ou classification par blocs, c'est- -dire la recherche simultan e d'une partition des lignes et d'une partition des colonnes d'un tableau de donn es, est devenue un outil tr s utilis  en fouille de donn es. Dans ce domaine, l'information se pr sente souvent sous forme de tableaux de contingence ou tableaux de co-occurrence croisant les modalit s de deux variables qualitatives. Dans cet article, nous  tudions le probl me de la classification crois e de ce type de donn es en nous appuyant sur un mod le de m lange probabiliste. En utilisant l'approche vraisemblance classifiante, nous proposons un algorithme de classification crois e bas  sur la maximisation altern e de la vraisemblance associ e   deux m langes multinomiaux classiques et nous montrons alors que sous certaines contraintes restrictives, on retrouve les crit res du Chi2 et de l'information mutuelle. Des r sultats sur des donn es simul es et des donn es r elles illustrent et confirment l'efficacit  et l'int r t de cette approche.

## 1 Introduction

La classification automatique, comme la plupart des m thodes d'analyse de donn es peut  tre consid r e comme une m thode de r duction et de simplification des donn es. Dans le cas o  les donn es mettent en jeu deux ensembles  $I$  et  $J$ , ce qui est le cas le plus fr quent, la classification automatique en ne faisant porter la structure recherch e que sur un seul des deux ensembles, agit de fa on dissym trique et privil gie un des deux ensembles, contrairement par exemple   l'analyse factorielle des correspondances qui obtient simultan ment des r sultats sur les deux ensembles ; il est alors int ressant de rechercher *simultan ment* une partition des deux ensembles. Ce type d'approche a suscit  r cemment beaucoup d'int r t dans divers domaines tels que celui des biopuces o  l'objectif est de caract riser des groupes de g nes par des groupes de conditions exp rimentales ou encore celui de l'analyse textuelle o  l'objectif est de caract riser des classes de documents par des classes de mots. Notons que dans ce domaine, les donn es se pr sentent g n ralement sous forme d'un tableau de contingence o  chaque cellule correspond au nombre d'occurrences d'un mot dans un document.