

STATISTIQUE ET ANALYSE DES DONNÉES

SOCIÉTÉ FRANÇAISE DE CLASSIFICATION

Résumés - Journées de Statistique, Nice 22-26 mai 1978

Statistique et analyse des données, tome 3, n° 2 (1978), p. 31-44.

http://www.numdam.org/item?id=SAD_1978__3_2_31_0

© Association pour la statistique et ses utilisations, 1978, tous droits réservés.

L'accès aux archives de la revue « Statistique et analyse des données » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/legal.php>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques
<http://www.numdam.org/>

SOCIETE FRANÇAISE
DE CLASSIFICATION

(Résumés - Journées de Statistique, Nice 22-26 MAI 1978)

SUR LES TYPES DE PARTITIONS

J.P. BARTHELEMY

E.N.S.C.M.B.

25030 - BESANCON Cédex

LA BI

U.E.R. DES
UNIVERSITF

L'étude des types de partitions (suite ordonnée des cardinaux des classes), c'est à dire des partages d'un entier peut être envisagée des points de vue suivants :

- Etude mathématique d'un "problème fondamental lié à la démarche du taxinomiste"
- Comparaison de classifications (menées concurremment) selon la taille et le nombre des classes indépendamment des éléments qui constituent ces classes.

C'est dans cette optique que nous nous proposons de présenter quelques résultats récents (essentiellement des propriétés métriques) sur les partages d'un entier.

L'ensemble P_n des partages de n est muni d'une relation d'ordre (non latticielle) pour laquelle il est modulaire. Cette remarque conduit à la construction de métriques qui s'interprètent en termes de graphes, certaines d'entre elles étant liées à la taxinomie et la théorie de l'information.

Les concepts de h
observation lors de dépouillemen
les mettent en évidence les oppo
d'un axe. Alors que les classifi
qui se ressemblent.

Il est souhaitabl
lyse de données les classes d'ob
les classes. Le concept de bi-cl

Une hi-classe est
dire que deux objets (individus,
même classe d'une bi-classe se r
différentes d'une bi-classe s'op

Le principa de la
tableau de données les bi-classe
du type ressemblance et du type

Nous avons dévelo
rechercher une telle structure :

- une adaptation d
 - recherche d'une bi-partition (typ
 - une hi-classific
- pour former des hi-classes

- une méthode qui
fournir des familles de k hi-clas
entre un axe factoriel et une hi-

CONSTRUCTION D'UNE ULTRAMETRIQUE LA PLUS PROCHE
AU SENS DES MOINDRES CARRES APPROXIMATION VS OPTIMISATION

J.L. CHANDON

UNIVERSITE D'AIX MARSEILLE 3
I.A.E.
29, av. Robert Schumann
13670 - AIX en PROVENCE

Soit un ensemble X d'objets à classer et soit un indice de proximité d défini pour l'ensemble IP des paires d'objets de X . d est une application symétrique de IP dans \mathbb{R}_+ , indicatrice du degré de dissimilarité existant entre deux objets.

L'objectif de la classification est de simplifier l'information contenue dans le tableau des proximités en remplaçant les proximités entre paires d'objets par des proximités entre classes d'objets.

Nous nous intéressons à la construction d'une classification ascendante hiérarchique (CAH) sur l'ensemble X . JOHNSON (1967) et BENZECRI (1973) ont démontré que pour obtenir une hiérarchie totale indicée sur X , il suffit de transformer l'indice de dissimilarité d en dissimilarité ultramétrique S .

L'écart entre deux dissimilarités d et S étant mesuré par la norme euclidienne $\|d - S\|$, le problème de la construction d'une ultramétrique optimale S^X minimisant $\|d - S\|$ a été résolu par CHANDON, LEMAIRE, DOUGET (1978). Toutefois, l'algorithme de séparation et évaluation progressive proposé est relativement lent. Il devient impraticable dès que le nombre d'objets est supérieur à 15.

Pour classer un grand nombre d'objets il est nécessaire de disposer d'un algorithme rapide permettant de construire ou d'approximer S^X . Deux algorithmes basés sur le concept de préordonnance sont proposés. Le premier, très rapide, ne garantit pas l'obtention de l'optimum S^X . Néanmoins, il conduit toujours à une bonne approximation et il améliore toujours l'ultramétrique obtenue par l'algorithme de la moyenne de LANCE, WILLIAMS (1967), lorsque celle-ci n'est pas optimale. Le second, moins rapide, améliore encore cette estimation.

Les deux algorithmes sont appliqués aux données classiques de RAO (1952), FITCH, MARGO IASH (1967), KAMEN (1971) et MILLER, NICELY, afin d'illustrer l'amélioration du critère des moindres carrés obtenue par rapport à plusieurs autres méthodes de classification hiérarchique.

RECONNAISSANCE D'OBJETS
D'APPROXIMATION POLYNOMIALE

Lorsque nous désirons approximer une fonction continue $f(x)$ définie sur l'intervalle $[a, b]$, nous cherchons une fonction polynomiale $P(x)$ qui approxime au mieux $f(x)$, au sens où l'erreur $E = \int_a^b (f(x) - P(x))^2 dx$ est minimale. Les plus utilisées sont les polynômes de Tchebychev. Les pentes infinies suivies de pentes faibles, ce qui est nécessaire pour approximer des courbes à pics, nous avons donc choisi une approximation par des splines. Ces fonctions sont définies par des segments dont les limites sont variables. Elles sont plus précises que les polynômes et qu'on ne possède de solution générale puisque la clé du succès d'une approximation est de choisir ces points limites.

Les algorithmes que nous proposons pour la classification automatique (algorithmes de régression) sont basés sur les techniques de régression (norme L_2) mentionnées au sens de Chebyshev (norme L_∞).

Ces algorithmes préconisés sont réalisés aux cas de contours et de dérivabilité (fonctions splines) et appliqués à la fonction approximante globale.

UNE METHODOLOGIE DE CLASSIFICATION :
L'ANALYSE PAR CRITERES ORDONNES SUR UN
ENSEMBLE MUNI D'UNE RELATION DE PROXIMITE.

C. COCHET
IRIA, CEPIA
Domaine de Voluceau
BP 105
78150 LE CHESNAY

Nous proposons, dans le cadre de cette communication, une méthodologie d'approche, "l'analyse par critères ordonnés sur un ensemble muni d'une relation de proximité" (APCO), qui vise des problèmes se posant en des termes complexes à propos du morcellement de certains systèmes. Le déroulement d'une méthode qui permet de discerner des sous parties à l'intérieur d'un système peut prendre deux formes essentielles.

Dans un premier cas il est possible d'identifier un procédé de morcellement dont les résultats seront conformes aux spécifications auxquelles doivent répondre les sous-systèmes discernés. La méthode qu'il faut mettre en oeuvre se réduit alors au procédé qui est ainsi mis en évidence.

Dans tous les autres cas, il n'est pas possible de découvrir une procédure unique qui satisfasse, en fonction des spécificités du domaine étudié, les besoins de l'analyste. L'investigation scientifique auquel est soumis le système sera composé de plus d'un procédé de morcellement.

Nous avons choisi d'établir la synthèse des résultats de morcellement en exploitant les partitions obtenues pour chaque procédé utilisé. Ce moyen autorise l'usage d'une structure mathématique riche : le treillis géométrique des partitions.

Un ensemble de procédures méthodologiques doivent être utilisées pour mener à bonne fin l'analyse des problèmes qui sont soulevés ici, c'est à dire ceux justifiables de différentes analyses de morcellement.

Dans un premier temps, chaque analyse de classification doit être appliquée au champs d'observation du système qui la concerne. Pour développer cette opération, il faut disposer d'une théorie des systèmes dont les concepts soient suffisamment évolués afin de rendre possible cette procédure.

Ensuite, nous avons établi une méthodologie qui permet d'envisager une synthèse des différents morcellements utilisés. Pour atteindre cet objectif, nous avons exprimé les résultats des procédés de morcellement sous forme de "proximités" affectant les points faisant l'objet de cette étude.

Tous les résultats obtenus, sous forme de partition des parties morcellées sont alors synthétisés dans une analyse unique qui est le reflet des différents points de vue présents dans l'élaboration de la stratégie de classification. On débouche alors, par l'exploitation des résultats centraux de cette analyse et des produits intermédiaires, sur des quantifications qui permettent de mesurer l'intervention des éléments de la stratégie et même, à la limite, de mettre en évidence certaines composantes qui sont le résumé d'autres, plus complexe.

ANALYSE ET CLAS
OBTENUES A PARTIR D

CENTRE DE CALCUL I

UER Mathématiques

Il s'agit de la rec
géographique de stèles thessalienn
consiste en deux spirales symétriq
spirales définit un important crit
caractériser les ateliers et de su
sans, en relation avec des situati

Les questions qui n
en Archéologie sont les suivantes
. Peut-on distingue
rales à partir des documents photo
. Si oui, peut-on d
rales et représenter chaque classe

Dans une première p
fecter l'ensemble des spirales à t
"à centres, sommets d'un polygone

Une analyse plus él
rayon de courbure le long de la sp
des différentes modalités de cette
la forme $R = f(\theta)$, où (R, θ) sont
tions ont été retenues :

$R = a\theta + b$, R en es

Nous cherchons main
et à augmenter le nombre de tyne
permettre, d'après nos dernières c
de spirales reconnues.

ANALYSE CANONIQUE DU POINT DE VUE
DE LA CLASSIFICATION AUTOMATIQUE

E. DIDAY

UNIVERSITE PARIS IX-DAUPHINE
Place de Lattre de Tassigny
75775 - PARIS Cédex 16

Quand le tableau des données est de grande taille, il est légitime de chercher des combinaisons linéaires dépendant des tendances locales qui peuvent apparaître dans la population. Il s'agit de détecter ces tendances et simultanément les composantes canoniques qui leur sont le mieux associées.

Suivant que les données sont centrées ou non, on propose plusieurs algorithmes qui tendent à minimiser le critère. Dans le cas où toutes les variables sont qualitatives, le problème se pose en terme d'analyse factorielle des correspondances et revient à minimiser le critère. Dans le cas où toutes les variables sont quantitatives, le problème se pose en terme d'analyse factorielle des correspondances et revient à chercher les classes d'objets qui induisent les plus grands X^2 de contingences entre les variables. Si l'un des deux paquets de variables est formé de variables d'incidence on aboutit à des méthodes intéressantes d'analyse discriminante locale.

LA NOTION DE DISPERSION

UER
UNIV

Alors que généralement
sont fondées sur la notion de dispersion
à traiter, il est proposé ici un

A toute partie A caractérisée par une dispersion D_A caractérisée par des
tativité : $D_{A \cup B} \geq D_A + D_B$ si $A \cap B = \emptyset$.

Différents cas sont
des dispersions avec les types de
dispersion permet de construire
tatifs (notamment la pondération
sible de généraliser des algorithmes
mun.

Enfin est présentée
Méthode Non Hiérarchique Descendante
de dispersion. Cette méthode recense
etc classes sans imposer de structure
elle résulte alors d'une structure
apportée par la méthode.

(+) les programmes de calcul automatisés
1974, par différents laboratoires

R. FAGES 41D rue Phélypeaux

SELECTION ET DISCRETISATION OPTIMALES DE VARIABLES CONTINUES
EN VUE D'UN PROBLEME DE RECONNAISSANCE DE FORMES

R. FAGES

UER de Mathématiques
UNIVERSITE DE LYON I

Lorsque les variables sont en partie ou totalité des variables continues, on se ramène à des variables discrètes (ou de classification) par le choix arbitraire de seuils afin d'utiliser les avantages de processus interrogatifs, comme par exemple les pseudo-questionnaires.

Il est abordé ici le problème du choix optimal pour chacune des variables continues, du nombre de classes et des seuils correspondants, afin de minimiser la probabilité d'erreur de l'identification par la règle de décision de BAYES.

La sélection des variables les plus discriminantes s'en déduit naturellement par le rejet des variables discrétisées par une seule classe. (Cette sélection est étendue aux variables initialement qualitatives).

La technique proposée utilise la majoration la plus fine de la probabilité d'erreur par une mesure d'entropie dérivant du coefficient de BHATTACHARRYA (*).

Un exemple concret est présenté, montrant l'efficacité des sélections et discrétisations obtenues, même dans le cas où les hypothèses nécessaires à la justification du critère d'optimisation ne sont pas vérifiées.

(*) M. TERRENOIRE, D. TOUNISSOUX. "Inequalities using BHATTACHARRYA distance and application to decision process".

3rd International Joint Conference on Pattern Recognition, Coronado
November 1976.

UN NOUVEAU TEST D'UN

ALTER

L. FA

W. FER

Centre Natio
LABORA
POUR LES
31, chemin Joseph

Si $F(x)$ est la fonction unimodale dont le support est inclus entre les points u et v auxquels la densité maximale respectivement se suivent la classification, pour tester l'uniformité de la statistique

$$S = \max_{a < b < 1} [a -$$

où $F^x(\cdot)$ désigne la fonction de répartition

On a déterminé, en utilisant 30 000 échantillons simulés indépendants effectifs égaux à 10, 20 et 40. Des alternatives du type

$$F(z) = \frac{z^k}{a^{k-1}}, \quad z \leq a$$

$$F(z) = 1 - \frac{(1-z)^k}{(1-a)^{k-1}}$$

indiquent que S se compare favorablement

DONNÉES QUANTITATIVES INCOMPLÈTES ET CLASSIFICATION

AMÉLIORA
CLASSIFICA

P.P. FEVRE

IRIA - LABORIA
Domaine de Voluceau
78150 - LE CHESNAY

M
UNIV
Av.
2000

Lorsque des variables sont observées sur une population, il est fréquent que, pour certains individus, des variables ne soient pas relevées.

Pour traiter de telles données, la plupart des auteurs cherchent à "reconstituer" tout d'abord les données non disponibles, afin de pouvoir appliquer, sur ces données complétées, les méthodes usuelles d'analyse.

Le point de vue adopté ici est différent : nous cherchons à mettre en oeuvre directement les méthodes classiques de traitement des données, en ne tenant compte que des observations connues, et sans chercher à reconstituer les observations manquantes.

Pour ce faire, nous calculons, à partir des données disponibles, des approximations des quantités nécessaires au traitement habituel et nous travaillons en nous servant de ces approximations.

En application, afin de montrer ce que permet cette méthodologie, nous montrons comment une méthode de classification, la méthode de Nuées Dynamiques, peut être mise en oeuvre sur des données quantitatives incomplètes.

HARTIGAN (1975)
chique rapide "Quick Tree Lead
tableau des données, qui n'est
méthode de type "leader", const
sage à travers les données un a
sont fixés à priori. Le nombre
pas connu à l'avance.

Considérant un a
certains niveaux sont confondus
l'algorithme en introduisant la
de définir le nombre N de noeud
dimension N par deux tableaux d

Réf. HARTIGAN J.A. Clustering AJ

ANALYSE CLASSIFICATOIRE D'UN TEST SCOLAIRE

R. GRAS

Département de Mathématiques et I.R.E.M. de RENNES
Campus de Beaulieu
35042 RENNES CEDEX

PRESENTATION
DE DISCRETISATION

Cette communication présente quelques résultats didactiques d'une analyse en classification hiérarchique d'un test mathématique proposé à près de 1 100 élèves de 13 à 15 ans (fin de 3ème de C.E.S.). La classification C1, obtenue à l'aide de l'algorithme de la vraisemblance du lien de I.C. Lerman, est comparée à 2 autres classifications :

- classification C2 par rapport au contenu et à la nature de la tâche de l'item
- classification C3 par rapport à une taxinomie d'objectifs cognitifs de R. Gras.

L'hypothèse d'indépendance entre C1 et C2, puis C2 et C3 est rejetée par un test du χ^2 , au seuil de 1‰.

Les 5 classes de C1 conduisent à des interprétations confirmées par une analyse factorielle des correspondances :

- classe α de découverte de règle de production logique
- classe β de fonctions numériques et géométriques
- classe γ de nature numérique, très didactique
- classe δ d'observations de propriétés affines de l'espace
- classe ϵ d'observations de propriétés métriques de l'espace.

Les niveaux cognitifs croissent puis décroissent de α à ϵ , en passant par un maximum dans γ . Les classes $\{\delta, \epsilon\}$ et $\{\alpha, \beta, \gamma\}$ sont relatives aux deux derniers stades de développement cognitif selon Piaget : "opérations concrètes" et "opérations logico-formelles".

Le codage de variab
données pose d'une façon générale
de variation de telles variables e
densité soit unimodale.

Dans le cadre des p
de mélanges de lois de probabilité
définition d'une fonction de densi
larges la séparation en composante

On rend finalement
médical ainsi que sur des mélanges

CLASSIFICATION VISUALISEE DE GRANDS ENSEMBLES
SOUS DOUBLE CONTRAINTE

SELECTION D
UNE METHODE DE

L. LEDART

Y.

C.N.R.S.
CRFDMC - 140, rue du Chevaleret.
75013 - PARIS

Non

C. ROCHE

S.P.A.F.
Direction des Télécommunications

Dépt.

COM

1455 c

La procédure présentée répond aux préoccupations suivantes :

Construire une partition d'objet respectant une contrainte de contiguïté géographique (un zonage) et dont les effectifs des classes soient bornés par une quantité imposée ; faciliter au maximum la compréhension du programme en multipliant les aides à l'interprétation, en introduisant de nouvelles procédures de visualisation, de façon à permettre une utilisation de routine aisée de la procédure ; assurer un encombrement mémoire réduit et une exécution rapide.

L'algorithme de base de construction d'une classification ascendante hiérarchique adaptée au sous-ensemble réactualisé des couples d'objets contigus n'est pas original (cf par exemple la thèse de A. THAURONT, PARIS, 1975).

La matrice de contiguïté ne figure ici que sous la forme de tableau de codage réduit (pour chaque sommet du graphe : adresse des sommets adjacents). C'est sous cette forme qu'elle est actualisée après chaque agrégation, et après chaque intervention du seuil de taille maximale des classes. C'est également à partir de cette forme condensée que cette matrice est soumise à une analyse des correspondances impliquant une technique de diagonalisation particulière, de façon à faire apparaître sur l'imprimante une reconstitution de la carte géographique, sur laquelle seront positionnées les différents noeuds de l'arbre intermédiaire, puis les classes finales. Le principe de cette dernière opération permet d'analyser en quelques secondes des matrices binaires clairsemées d'ordre 1000x1000. Une analyse des correspondances classique effectuée cette fois sur le tableau de données de départ permet également de suivre les évolutions des noeuds et la position des classes non plus dans l'espace géographique, mais dans l'espace des variables. On a ainsi tous les éléments pour suivre et comprendre les mécanismes de formation des classes et le caractère plus ou moins prégnant des contraintes.

L'ensemble des le
phabet et les 10 chiffres. Pour
plusieurs caractères typographiq

Le but de notre a
tre, un caractère de notre ensem
éloigné des caractères des autre

La première étape
Elle consiste à associer à chaque
thème est insoluble si la descri
porte pas les traits pertinents
préalable du codage. L'image de
et 39 lignes, le codage que nous
caractère en intervalles ou en f

La deuxième étape
simple d'affectation décidant, d
ou non d'un nouvel individu à l'
finition de deux fonctions, l'un
tion d'écartement.

ARBRES VALUÉS ET ULTRAMÉTRIQUES

B. TECLERC

C. M. S.
54, bd Raspail
75270 - PARIS Cédex 06

Il est maintenant bien connu que tout arbre (graphe connexe et sans cycle) valué, défini sur un ensemble X de cardinal n , induit naturellement une ultramétrie r sur X , donc une classification hiérarchique sur X . Quelques travaux ont commencé à paraître, cherchant à étudier l'ensemble des ultramétries ainsi définies à partir de tous les arbres valués donnés par un indice de distance d sur X . Inversement, on a aussi posé le problème suivant : soit r une ultramétrie sur X : représenter r par un arbre (Benzécri et Jambu, 1976).

Après avoir rappelé les résultats antérieurs, nous précisons le lien entre arbres valués et ultramétries. On établit d'abord que celles-ci se caractérisent, parmi les indices de distance, par des propriétés ou interviennent uniquement leurs arbres minimaux. Mais, s'il est vrai qu'une ultramétrie r est parfaitement définie par l'un quelconque de ses arbres minimaux (valué par la restriction de r), le nombre $N(r)$ de ceux-ci est compris, dans le cas général, entre $4^{n-1}/n^2$ et $(n-1)!$. Ceci pose le problème du choix d'un arbre particulier pour représenter r , qui n'a pas de réponse évidente que lorsque r a été obtenue à partir d'un arbre valué lisible directement dans les données (c'est le cas dans certaines méthodes classificatoires : lien simple et lien complet).

On s'intéresse ensuite aux propriétés du nombre $N(r)$ qui paraît être un descripteur intéressant de la classification hiérarchique associée à l'ultramétrie r .

CLASSIFICATION DE GRANDS ENSEMBLES

FACULTE UNIVERSITAIRE
FACULTE DE DROIT ET

FACULTE UNIVERSITAIRE
INSTITUT D'ECONOMIE

La méthode "single
est très utilisée en classification
à classer sont grands, une méthode
minimum d'un graphe, exige de $O(n^3)$
où n désigne le nombre d'objets à
à classer sont des distances de n
sans que toutes les dissimilarités
rapide pour ce problème, ainsi qu'
semble de 10 000 étoiles.

H. LERFONDE

Département de Mathématiques
Centre Scientifique et Polytechnique
UNIVERSITE PARIS-NORD (XIII)
Av. J.B. Clément
93430 - VILLETANFISE

Laboratoire de
UNIVERSITE DE R
35031 - RENNES

Le principe général de cette méthode de classification non hiérarchique est de constituer pas à pas des agrégats parmi les éléments à classer.

Pour former un agrégat nous choisissons un élément "pertinent", appelé pôle, auquel nous agrégeons les éléments qui lui sont le plus proche, fonction d'un certain seuil. Cet agrégat constitué, nous recommençons la même opération avec les éléments non encore agrégés, et ceci tant qu'il reste des éléments à classer. A la fin nous obtenons une partition de l'ensemble des éléments, chaque classe de la partition étant l'un des agrégats (pour l'idée de base de cet algorithme, on pourra consulter I.C. LERMAN, reconnaissance et classification de structures finies en analyse des données, Université de RENNES I, 1977, rapport IRISA n° 70).

Nous avons mis en oeuvre cette méthode sous forme d'un algorithme appelé MPAGD (Méthode des Pôles d'Agrégation sur les Distances) où l'on étudie la distribution des distances entre éléments à classer.

La principale difficulté de ce type d'algorithmes est de parvenir à un système permettant d'arrêter la formation d'un agrégat que l'on puisse appliquer indépendamment de la distribution des distances et de la nature des données. Cet algorithme est une tentative dans cette voie.

Différentes tentat

et la Classification ont été prop
M. GONDRAN (1975), M. JAMBU suiv
correspondent en fait, comme nous
"factorielles" du problème de la
tations est en effet grande comot
de l'Analyse Factorielle en Comp
critère de l'inertie expliquée.

Notre but dans cet
pour préciser sa véritable nature
nouvelles. Ce faisant, nous contr
des approches en les situant, dan
aux autres. Ce qui nous permettra
ces différentes tentatives.

Les différents typ
et la Classification se distingue
la recherche d'une classification
che d'un arbre binaire des Classi
la nature du critère optimisé : d
est basé sur l'inertie expliquée
proximité entre parties disjointe
mum). Ces différentes approches s
factorielle retenue : s'agit-il d
l'espace de représentation du qua
trouve défini comme une fonction
Ces approches se distinguent enfi
l'espace engendré par la solution

B. MONJARDET
UNIVERSITE DE PARIS V
Centre de Mathématique Sociale (E.H.E.S.S.)

Il est bien connu que deux méthodes classiques en classification hiérarchique, celle du "lien simple" (ou de l'ultramétrie sous-dominante) et celle du "lien complet" s'interprètent aisément en termes de théories des graphes. Les classes des partitions de l'arbre hiérarchique sont en effet, dans le premier cas, des classes connexes, dans le second cas, des cliques, des graphes "seuils" associés à la dissimilarité considérée. La considération de ces graphes seuils ajoutée à une condition de cohérence, ramène le problème de la classification hiérarchique à celui de la classification des sommets d'un graphe. A cet égard, les classements formés par les classes connexes d'une part, des cliques d'autre part, apparaissent comme deux solutions extrêmes maximisant respectivement un critère de séparation entre classes et un critère d'homogénéité à l'intérieur des classes. Mais en fait, la théorie des graphes permet d'envisager bien d'autres possibilités intermédiaires entre les deux solutions extrêmes. Ces possibilités sont d'ailleurs apparues dans des contextes variés. Par exemple, l'analyse de réseaux sociométriques et la recherche de leur décomposition en groupements homogènes amène à définir des concepts de "cliques généralisées" d'un graphe (Luce, 1950). Inversement, des études théoriques sur le nombre chromatique ou la connectivité d'un graphe conduisent à des notions pouvant être utiles en classification.

Le but de l'exposé est de présenter ces apports de la théorie des graphes à la taxinomie mathématique, qu'il s'agisse ou non de classification hiérarchique. La littérature sur le sujet étant déjà fort vaste (plus de 150 références), on ne tentera pas d'être exhaustif, mais on essaiera d'indiquer les principales contributions, en distinguant les apports conceptuels de ceux plus techniques.

Un problème de re
application d'un ensemble U (l'U
ralement c'est un problème n-dime

En analyse d'image
qui respecterait leur voisinage p
nel en un problème unidimensionne

En classification
siste à partitionner un ensemble
leurs interdistances ou, plus gé
nous trouvons une application T,
nage des points, nous transformon
sionnel.

Une application ré
intéressante en classification. L
(courbe remplissant un carré) peu
le Professeur ALEXANDROV, de LENT

Le principe est le
égaux, qui nous donnerons une div
adressage de ces hypercubes suiv
séquentiellement, en respectant p
de ce type, P et \tilde{P} , l'un translat
les coordonnées sont égales à $\frac{1}{2^p}$
respecter le 2^n -voisinage. Cette p
semble de points en classes, sans

Des exemples sont c
les données IRIS (n=4).

CLASSIFICATION DE MAXIMUM DE VRAISEMBLANCE DE DONNEES BINAIRES

ETUDE DE L'ACR
PAR UNE METHODE

P. ROUSSEMI

Centre de Recherches Mathématiques
UNIVERSITE DE MONTREAL
Caser Postal 6128
MONTREAL (CANADA)

Laborat
H.E.R. d
UNIVERS

Une famille d'ensembles de données binaires est par hypothèse représentée par K modèles logistique-linéaires où chaque modèle représente une sous famille de cette famille.

On définit un algorithme qui trouve simultanément la partition de maximum de vraisemblance en K sous familles et les estimateurs de maximum de vraisemblance des paramètres de chaque modèle.

Cet algorithme est une extension des méthodes de partition itérative selon l'approche de Diday.

On cherche à classer les données de façon à dégager les principaux facteurs de la Vraisemblance du lieu d'O à I des descriptifs ; chaque département

- 1 - surfaces par p
- 2 - importance du
- 3 - structure d'ex

On a défini un indice de (données) d'un tableau de contingence I.C. LERMAN et respectant la méthode cas où les données sont une juxtaposition ayant, tous, le même ensemble de l'application de l'A.V.L. pour des départements français selon un seul caractère deux ou trois caractères réunis,

Nous avons aussi étudié la Classification Ascendante par aggrégation par la maximisation, à la nouvelle classe formée par réunion des algorithmes sont comparés dans les

- a) classification
- et
- b) classification