

44

P A C T

10

1984

Cours postgradué européen I
European Postgraduate Course I

Datation-caractérisation des céramiques anciennes

Bordeaux-Talence, 6-18 avril 1981

Cours Intensif européen

organisé dans le cadre du programme intergouvernemental du Conseil de
l'Europe dans le domaine de l'enseignement supérieur et de la recherche

UNIVERSITÉ DE BORDEAUX III

Laboratoire de physique appliquée à l'archéologie du Centre
de recherche interdisciplinaire d'archéologie analytique, associé au CNRS

UNIVERSITÉ DE BORDEAUX I

Laboratoire de cristallographie et physique cristalline, associé au CNRS
Centres de microsondages et de microscopie électroniques

MAISON DES SCIENCES DE L'HOMME D'AQUITAINE

Édité par Tony HACKENS et Max SCHVOERER

Presses du CNRS
22, rue Saint-Amand
PARIS XV^{ème}

Centre Universitaire Européen pour les biens culturels
European University Center for the Cultural Heritage
84010 RAVELLO (Italia) - Villa Rufolo

LE TRAITEMENT DES DONNÉES D'ANALYSE

Résumé

Deux méthodes de traitement des données fournies par l'analyse des céramiques font l'objet d'une présentation détaillée : une méthode de classification, l'analyse de grappes, et une méthode de classement, l'analyse discriminante quadratique. Pour chacune de ces méthodes sont indiqués les principes de base, la présentation des résultats, les limites de validité et les conditions d'emploi. On montre enfin comment ces méthodes interviennent dans le raisonnement archéologique et où se situent les limites de leur compétence.

I. INTRODUCTION

L'utilisation des résultats fournis par l'analyse des céramiques passe généralement par des méthodes informatisées de traitement des données. On distingue traditionnellement les méthodes de classification, qui visent à la constitution de classes, et les méthodes de classement qui cherchent à attribuer des individus à des classes pré-existantes. Il existe une très grande variété de méthodes de classification et de classement. Plutôt que de les passer toutes en revue, ce qui n'aurait pu se faire que d'une manière extrêmement superficielle, il a semblé préférable de s'arrêter plus longuement sur deux d'entre elles, l'analyse de grappes et l'analyse discriminante quadratique. Le choix de ces deux méthodes relève surtout de considérations personnelles, puisque ce sont là les méthodes que nous utilisons systématiquement au Laboratoire de Céramologie de Lyon. Il est vrai que ce sont aussi des méthodes qui sont d'un emploi courant, surtout la première.

Ce cours aurait atteint son but s'il parvenait, au-delà de la complexité superficielle des méthodes, à convaincre que les principes des classifications et des classements sont compréhensibles par tous et qu'il appartient à chacun d'exercer son esprit critique face aux applications archéologiques qui en sont faites. L'ensemble a été conçu comme une initiation assez large aux problèmes de statistique ; on souhaiterait qu'il donne envie d'en savoir plus.

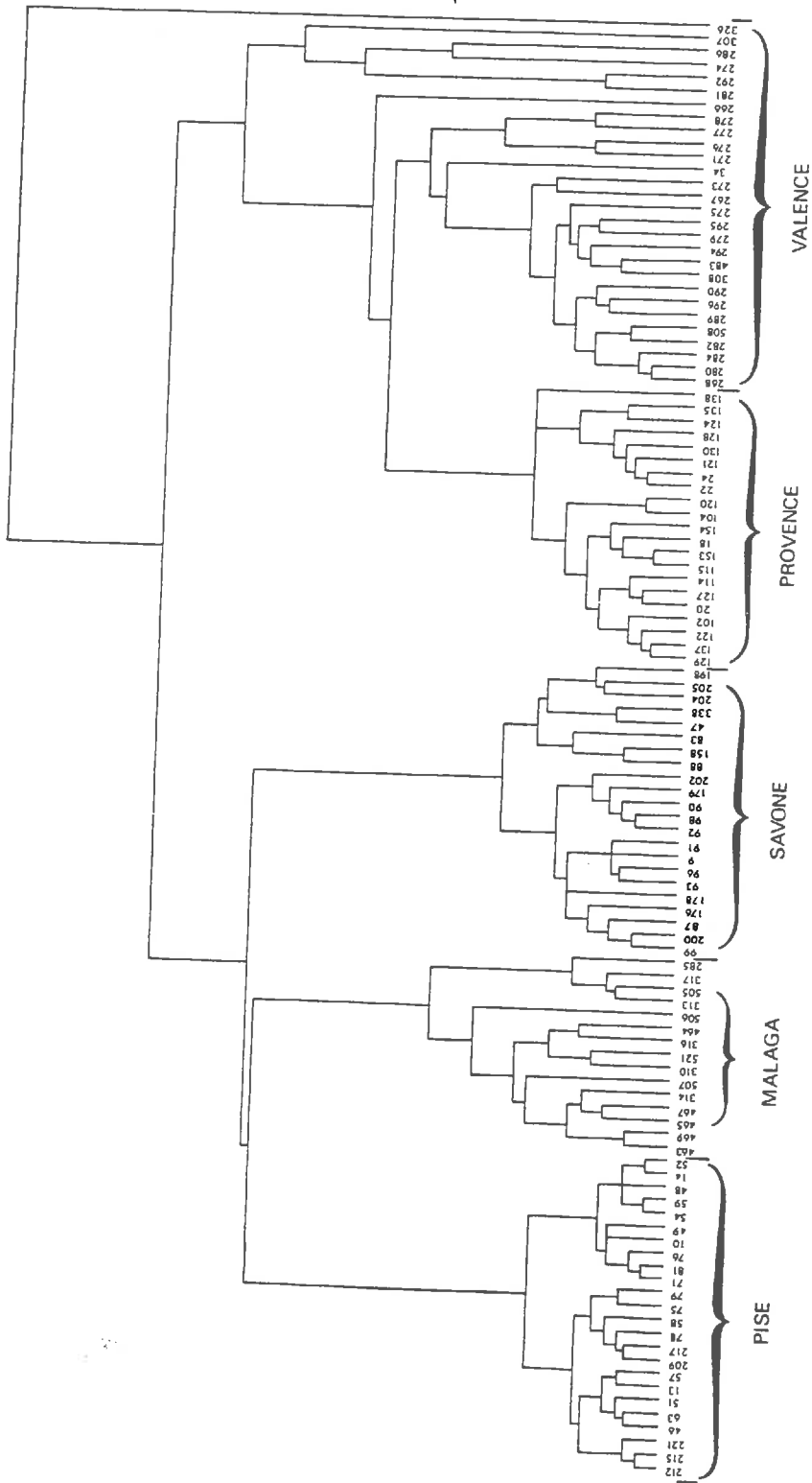


Fig. 1 : Analyse de grappes, par affinité moyenne pondérée sur variables centrées réduites d'un ensemble de céramiques médiévales faisant l'objet d'importations en Provence.

Associant ainsi, de manière progressive, céramiques et pseudo-céramiques dont les compositions chimiques sont les plus proches, on parvient à réunir tous les exemplaires étudiés en une arborescence unique dont les ramifications ou *grappes* matérialisent des *groupes* constitués de céramiques qui se ressemblent d'autant plus que la hauteur de la ramification les réunissant est plus basse. Ainsi, dans l'exemple de la figure 1, le groupe des céramiques de Pise est-il celui dont les compositions sont les plus homogènes, étant constitué de céramiques qui présentent entre elles des ressemblances plus marquées que n'en présentent celles des autres groupes, Savone et Malaga par exemple. Quant au groupe des céramiques de Valence, il est de loin le plus hétérogène! On remarquera sur cette même figure 1 que les groupes se subdivisent eux-mêmes en divers *sous-groupes*. Par exemple le groupe de Pise est constitué d'un premier sous-groupe allant de la céramique 212 à la céramique 79, et d'un second sous-groupe allant de 71 à 52.

Ce sont les différents groupes et sous-groupes résultant du rapprochement des céramiques présentant des compositions voisines que l'on confronte ensuite avec les données archéologiques ou géochimiques, comme on l'a déjà signalé.

C. Exemple de calcul

Essayons à titre d'exemple (et pour mieux être à même de lire une analyse de grappes) de classer, d'après leurs pourcentages de chaux (CaO), 6 exemplaires de céramiques que nous désignerons par a, b, c, d, e et f.

Données initiales

Céramiques à classer :	a	b	c	d	e	f
Pourcentages de chaux :	4	14	12	7	8	17

On choisira, pour déterminer la plus ou moins grande ressemblance des exemplaires pris deux à deux, la valeur absolue de la différence des deux pourcentages de chaux. C'est là une expression de la *distance* des deux exemplaires considérés, mais on pourrait choisir bien d'autres expressions, plus compliquées, pour définir la distance de deux céramiques, ainsi qu'on le verra plus loin. Avec la définition adoptée ici, la distance entre b et e (identique à la distance de e à b) sera de 6, et celle de a à d, de 3. Il est évident que plus la distance séparant deux exemplaires est grande, moins ces exemplaires se ressemblent, ou, si l'on veut, plus est grande leur dissemblance.

On a l'habitude de représenter l'ensemble des distances des céramiques à classer prises 2 à 2 sous forme d'une *matrice des distances*, ou plutôt d'une demi-matrice qui est suffisante pour raisons de symétrie.

Dans le cas pris en exemple cette matrice est la suivante :

Matrice initiale

	a	b	c	d	e	f
a	0	10	8	3	4	13
b	0	2	7	6	3
c	0	5	4	5
d	0	1	10
e	0	9
f	0

On voit clairement sur la matrice des distances que ce sont les céramiques d et e qui fusionnent les premières. Elles sont en effet les plus ressemblantes dans le système de distances que nous avons adopté. Leur distance y est de 1, et nous prendrons cette valeur pour hauteur du pont correspondant à la première fusion (Fig. 2, 1). Quant à la

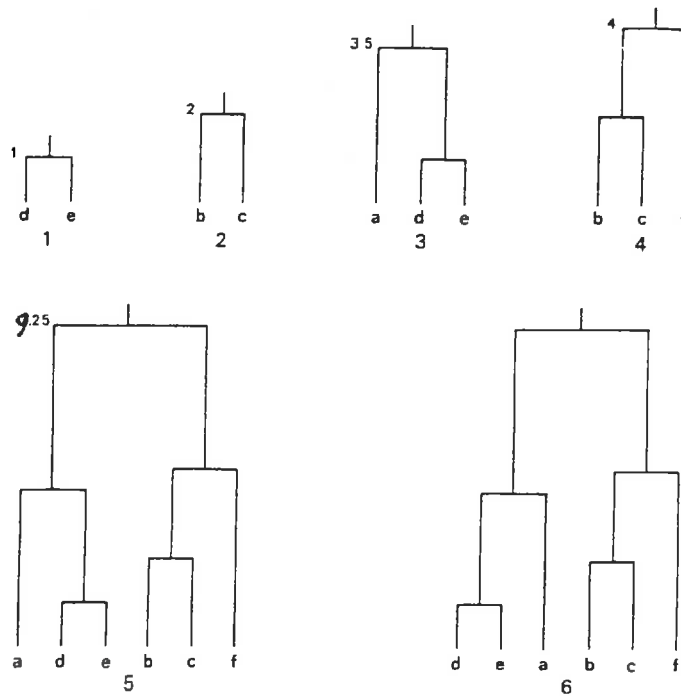


Fig. 2 : Exemple de construction d'un dendrogramme.

pseudo-céramique qui remplacera d et e dans la suite des calculs, nous la désignerons par de et nous déciderons que sa composition sera la moyenne de celles des deux céramiques qu'elle remplace, soit ici 7,5 % de CaO. Nous conserverons la même règle pour toutes les fusions, ce qui se fait le plus souvent, quelle que soit la méthode d'analyse de grappes. Par contre les règles fixant les caractéristiques des pseudo-céramiques varient quant à elles beaucoup d'une méthode à une autre. C'est un point sur lequel nous reviendrons.

Données après la première fusion

Céramiques et pseudo-céramique					
restant à classier :	a	b	c	de	f
Pourcentages de chaux :	4	14	12	7.5	17

Matrice après la première fusion

	a	b	c	de	f
a	0	10	8	3.5	13
b	0	2	6.5	3
c	0	4.5	5
de	0	9.5
f	0

Ce sont à présent les céramiques b et c qui fusionnent, le pont correspondant à cette seconde fusion étant à une hauteur de 2 (Fig. 2,2). On attribue à la pseudo-céramique bc un pourcentage de chaux égal à 13.

Données après la seconde fusion

Céramiques et pseudo-céramiques

restant à classier :

Pourcentages de chaux :

a	bc	de	f
4	13	7.5	17

Matrice après la seconde fusion

	a	bc	de	f
a	0	9	3.5	13
bc	0	5.5	4
de	0	9.5
f	0

Troisième fusion : a et de

Hauteur du pont : 3.5 (Fig. 2,3)

Pourcentage de chaux de ade : 5.75

Données après la troisième fusion

Céramiques et pseudo-céramiques

restant à classier :

Pourcentages de chaux :

ade	bc	f
5.75	13	17

Matrice après la troisième fusion

	ade	bc	f
ade	0	7.25	11.25
bc	0	4
f	0

Quatrième fusion : bc et f

Hauteur du pont : 4 (Fig. 3,4)

Pourcentage de chaux de bcf : 15

Données après la quatrième fusion

Pseudo-céramiques

restant à classier :

Pourcentages de chaux :

ade	bcf
5.75	15

Matrice après la quatrième fusion

	ade	bcf
ade	0	9.25
bcf	0

Dernière fusion : ade et bdf

Hauteur du pont : 9.25 (Fig. 2,5)

Pourcentage de chaux de ade bcf : 10.375

On obtient en fin de compte la séparation de 6 exemplaires étudiés en deux groupes : un groupe pauvre en chaux, ade, et un autre riche en chaux, bcf. Résultat évident à priori pour l'exemple pris ici, mais qui ne l'est guère dès lors qu'on fait intervenir simultanément dans la classification plusieurs constituants chimiques et que le nombre des exemplaires à classier est important.

On notera que le dessin du dendrogramme peut être exécuté pour une même analyse de grappes de différentes manières, en intervertissant à volonté telle ou telle partie. C'est ainsi qu'on utilisera la représentation de droite, au bas de la figure 2, plutôt que celle de gauche, mais cela ne modifie en rien la lecture des classes sur le dendrogramme.

D. Règles de fusion

Nous revenons ici sur les règles qui fixent les caractéristiques des pseudo-céramiques. Nous avons déjà signalé qu'elles varient considérablement d'une méthode d'analyse de grappes à une autre. Dans l'exemple précédent nous avons adopté comme règle celle consistant à doter la pseudo-céramique d'une composition qui soit la *moyenne* de celles des deux céramiques (ou pseudo-céramiques) qui ont fusionné. On dit en ce cas qu'il s'agit d'une analyse de grappes par fusion en *affinité moyenne*. Mais cela mérite d'être regardé de près. Reprenons donc l'exemple précédent, et plus particulièrement la troisième fusion qui nous avait conduit à la pseudo-céramique ade dont le pourcentage de chaux était de 5.75. On notera que ce pourcentage correspond à $1/2 a + 1/4 d + 1/4 e$, ce qui traduit cette propriété générale assez évidente, compte tenu de la règle de fusion adoptée, que dans la composition moyenne qui est attribuée à une pseudo-céramique les compositions individuelles des céramiques ont un rôle d'autant plus restreint que leur fusion est plus ancienne. Il est facile de vérifier par exemple, dans le cas des deux diagrammes de la figure 3, que chacune des compositions individuelles n'intervient dans la composition de la pseudo-céramique finale que pour la fraction indiquée au-dessous. La règle de fusion que nous avons utilisée jusqu'ici est dite en *affinité moyenne pondérée* à cause du *poids* (ou importance) différent donné aux compositions des céramiques entrant dans le calcul de la moyenne.

On aurait pu tout aussi bien définir une moyenne qui soit telle que les compositions de toutes les céramiques intervenant lors d'une fusion quelconque pèsent d'un poids égal dans le calcul de cette moyenne (Fig. 4).

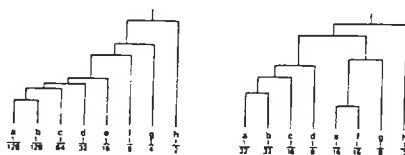


Fig. 3 : Contributions des céramiques intervenant dans la composition de deux pseudo-céramiques, en affinité moyenne pondérée.

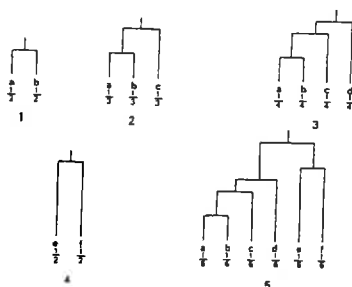


Fig. 4 : Contributions des céramiques lors de la construction d'un dendrogramme en affinité moyenne non pondérée.

Il est instructif alors de comparer les résultats donnés par l'une et l'autre méthode. On pourra refaire à titre d'exercice le calcul en affinité moyenne pondérée et non pondérée de la grappe correspondant aux céramiques $a = 4$, $b = 5$, $c = 7$, $d = 10$ et $e = 15$. Les résultats sont reportés sur la figure 5. On peut observer que des compositions marginales comme d et surtout e paraissent cependant un peu moins marginales en affinité moyenne pondérée (Fig. 5,1) qu'en affinité moyenne non pondérée (Fig. 5,2). Dans ce dernier cas les fusions correspondantes sont en effet un peu plus hautes. C'est là un résultat tout à fait normal puisque les moyennes non pondérées évoluent, au fur et à mesure des fusions, beaucoup moins que les moyennes pondérées. Certes les différences ne sont pas très importantes ici, mais il en est autrement lorsqu'on a affaire à des groupes de composition réels où les moyennes non pondérées se stabilisent très vite après quelques fusions, tandis que les moyennes pondérées peuvent évoluer jusqu'à devenir très éloignées de celles des toutes premières fusions. Par ailleurs il faut savoir qu'existent des lois de fusion qui permettent d'amplifier ces différences de comportement, et qu'on peut donc à volonté obtenir des groupes très restreints ou, tout au contraire, rattacher des éléments fortement marginaux à tel ou tel groupe.

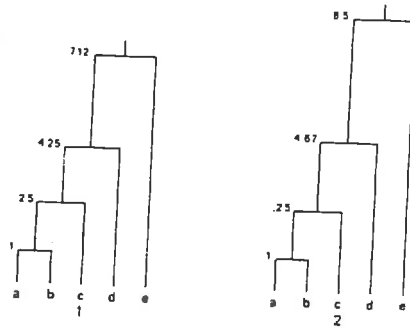


Fig. 5 : Comparaison des dendrogrammes obtenus en affinité moyenné pondérée (1) et non pondérée (2).

Aussi convient-il de ne pas prendre pour un gage d'objectivité le terme de *classification automatique* que l'on donne aux procédés de classification informatisés parmi lesquels se range bien entendu l'analyse de grappes. De même on ne saurait se contenter sans plus d'explications de n'importe quel rapprochement apparaissant sur une analyse de grappes pour conclure à l'identité des origines d'une céramique et d'un groupe, ou de deux groupes. D'ailleurs une similitude de composition, si poussée soit-elle, demeure par principe insuffisante pour décider à elle seule qu'on est dans le cas d'une origine commune. En cette matière la décision, et donc la preuve, ne peuvent résulter que d'une stratégie de recherche plus vaste et plus complexe, que nous évoquerons brièvement au terme de ce cours. On peut toutefois noter dès à présent que l'interpénétration très étroite de deux groupes sur un dendrogramme constitue généralement une sérieuse présomption en faveur de leur identité d'origine. Il n'en est pas de même dans le cas d'un exemplaire unique ; même intégré dans un groupe, il peut toujours être soupçonné de s'y trouver mal classé. Quant aux exemplaires en position marginale dans un groupe, ils doivent être à priori suspectés.

E. Changements de variables

Nous revenons ici sur la notion de distance entre deux céramiques; elle se trouvait définie, dans le premier exemple que nous avons choisi, par la différence de leur pourcentage de chaux CaO. Mais comment définir alors la distance de deux céramiques lorsqu'interviennent, dans le tri des exemplaires en analyse de grappes, plusieurs constituants? Une définition très courante (mais il y en a d'autres) est la *distance euclidienne*. Elle est telle que son carré soit égal à la somme des carrés des différences de concentration de chacun des constituants. Calculons par exemple la distance d_{ab} des deux céramiques a et b dont les compositions sont les suivantes :

	Na ₂ O	K ₂ O	MgO	CaO	MnO	Al ₂ O ₃	Fe ₂ O ₃	SiO ₂	TiO ₂
a :	0.85	5.3	3.1	15.5	0.045	17.1	7.4	50.0	0.55
b :	0.75	6.2	2.7	11.1	0.050	18.6	6.9	52.8	0.65

$$d_{ab}^2 = (0.85 - 0.75)^2 + (5.3 - 6.2)^2 + (3.1 - 2.7)^2 \\ + (15.5 - 11.1)^2 + (0.045 - 0.050)^2 + (17.1 - 18.6)^2 \\ + (7.4 - 6.9)^2 + (50.0 - 52.8)^2 + (0.55 - 0.65)^2$$

$$d_{ab}^2 = (0.1)^2 + (0.9)^2 + (0.4)^2 + (4.4)^2 + (0.005)^2 + (1.5)^2 \\ + (0.5)^2 + (2.8)^2 + (0.1)^2$$

Il est inutile d'aller plus loin dans le calcul pour comprendre que la distance euclidienne ainsi définie, appelée *distance euclidienne sur données brutes*, ne prendra pratiquement jamais en compte les différences de concentration des constituants dont les valeurs moyennes sont les plus faibles dans les céramiques. C'est ainsi que les différences de concentration du manganèse, du titane ou du sodium n'interviendront guère dans la classification, tandis que celles du silicium, du calcium ou de l'aluminium auront un rôle prépondérant. Si l'on veut que les différents constituants aient la possibilité de jouer un rôle comparable lors de la classification, il faut que chacun de ces constituants intervienne d'une manière sensiblement égale dans l'expression de la distance euclidienne.

Une première solution, parmi d'autres, consiste à modifier l'expression des concentrations, de telle sorte que leur valeur moyenne soit la même pour tous les constituants. Supposons par exemple que les moyennes relatives à SiO₂ et à TiO₂ soient respectivement égales à 51.7% et 0.58% pour l'ensemble des exemplaires à classer. On ramènera ces moyennes à 100 en affectant les pourcentages d'un coefficient égal à 100/51.7 pour la silice et à 100/0.58 pour l'oxyde de titane. On obtient alors, pour les céramiques a et b prises en exemple, des taux de 96.7 et 102.1 pour la silice et 94.8 et 112.1 pour l'oxyde de titane. On perçoit à présent que si l'on s'est bien affranchi des difficultés provenant de l'inégalité des valeurs moyennes des divers constituants, d'autres difficultés subsistent cependant.

La plus évidente, au vu des nouvelles valeurs prises par les variables SiO₂ et TiO₂ de l'exemple précédent, tient aux dispersions relatives qui peuvent être très différentes d'un constituant chimique à un autre. Aussi trouve-t-on fréquemment, dans la définition des changements de variables employés en analyse de grappes, l'écart-type σ qui caractérise

la plus ou moins grande dispersion des pourcentages d'un constituant autour de sa valeur moyenne \bar{m} . Les variables centrées réduites, qui sont les plus utilisées, sont dans ce cas; elles sont définies par l'expression $(x - \bar{m}) / \sigma$ où x représente le pourcentage d'un constituant donné dont la moyenne est \bar{m} et l'écart-type σ pour l'ensemble des exemplaires entrant dans l'analyse de grappes. La distance euclidienne sur variables centrées réduites se définit quant à elle de la même manière que pour les données brutes, par une expression dont le carré est égal à la somme des carrés des différences de concentration, ces dernières étant alors exprimées selon les nouvelles variables.

Les inconvénients résiduels résultant de l'emploi des variables centrées réduites sont principalement dus au caractère quelque peu artificiel de l'écart-type σ ; il se rapporte à l'ensemble des céramiques que le hasard a réuni dans la même analyse de grappes, alors qu'il serait plus satisfaisant que l'écart-type soit celui d'un seul groupe. On se rapproche d'une telle situation lorsqu'après un tri préliminaire on teste isolément chacun des groupes par analyse de grappes, ce qui permet beaucoup mieux d'en vérifier l'homogénéité et d'y reconnaître les individus ayant un caractère marginal. L'analyse de grappes, comme bien d'autres méthodes de classification, ne se rapproche de ses conditions optimales de fonctionnement que si l'on n'a plus à faire le tri que d'un groupe unique. Au départ de tout processus de classification l'écart-type est élevé par suite de la présence dans une même analyse de grappes de plusieurs groupes ayant des compositions différentes. Il peut en résulter des perturbations importantes dans certaines classifications, les variables pour lesquelles l'écart-type est particulièrement élevé voyant leur rôle diminuer en proportion inverse. Cela impose que l'on suive une certaine logique dans les classifications en procédant à une réduction progressive des groupes et à l'élimination des individus les plus marginaux.

F. Influence de la dispersion

Il est évident que les problèmes de classification et de classement des céramiques n'existeraient pas si toutes les céramiques produites dans un atelier avaient même composition. En réalité les exemplaires représentatifs de la production d'un atelier présentent, pour un constituant chimique quelconque, des pourcentages variables qui se regroupent autour d'une valeur moyenne (on a déjà signalé d'ailleurs que la dispersion plus ou moins grande des pourcentages d'un constituant autour de sa valeur moyenne \bar{m} s'exprime par l'écart-type σ). C'est à l'existence de cette dispersion et à ses diverses caractéristiques que les problèmes de classification et de classement doivent leur complexité.

Pour illustrer un premier aspect du rôle de la dispersion sur la classification des céramiques on considérera les histogrammes relatifs aux pourcentages de chaux, CaO, d'un atelier 1 ($\bar{m} = 6\%$) auquel appartiennent les céramiques A (% CaO = 10) et B (% CaO = 9), et d'un atelier 2 ($\bar{m} = 15\%$) auquel appartient la céramique C (% CaO = 12) (Fig. 6). Dans l'exemple choisi les histogrammes de la chaux permettent de séparer complètement les productions des ateliers 1 et 2.

Si l'on fait intervenir un second constituant comme la magnésie, MgO, on peut se trouver dans le cas où les histogrammes relatifs à ce constituant aient même moyenne pour l'un et l'autre atelier. Si de plus aucune relation n'existe, pour l'atelier 1, entre

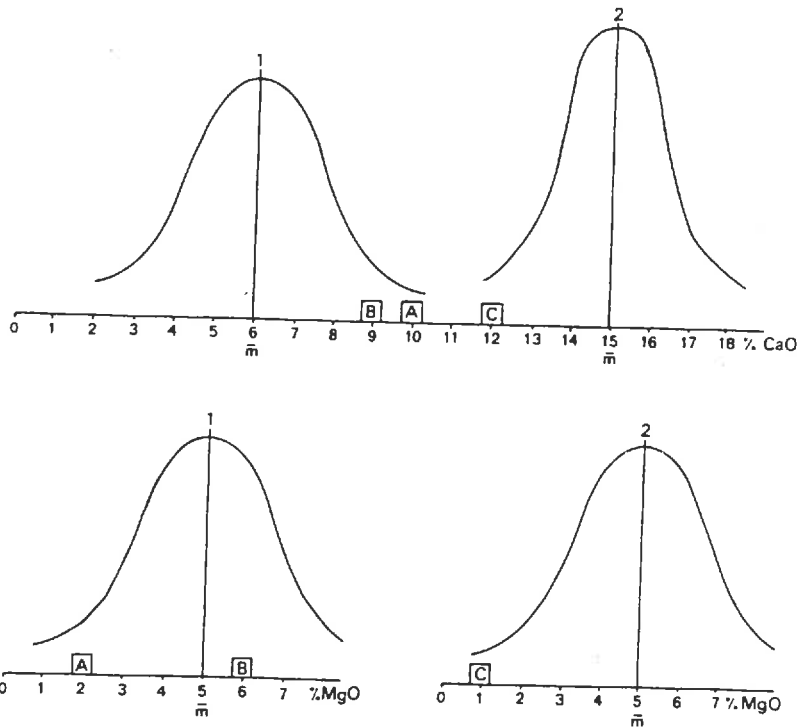


Fig. 6 : Histogrammes partiels de la chaux et de la magnésie des ateliers (1) et (2).

les pourcentages de chaux et de magnésie, ce qui s'exprime en disant qu'on a affaire à des *variables indépendantes*, les céramiques A et B peuvent prendre à priori n'importe quelle valeur dans l'histogramme de la magnésie de l'atelier 1. On peut donc rencontrer une situation comme celle de la figure 6 qui fera de C un voisin plus proche de A que ne l'est B.

En effet si l'on calcule les distances d_{AB} et d_{AC} en utilisant pour simplifier une distance euclidienne sur données brutes on obtient :

$$d_{AB}^2 = (10 - 9)^2 + (2 - 6)^2 = 17 = (4.12)^2$$

$$\times d_{AC}^2 = (10 - 12)^2 + (2 - 1)^2 = 5 = (2.24)^2$$

Cette *inversion des proximités* résulte essentiellement des *fluctuations aléatoires* des variables dans les histogrammes de la chaux et de la magnésie de l'atelier 1. Elle peut dans certains cas, sur lesquels nous reviendrons plus loin, conduire à des rattachements erronés en analyse de grappes. C'est ainsi par exemple que A pourrait se trouver rattaché à l'atelier 2, ou au contraire C à l'atelier 1, dans des analyses de grappes utilisant les distances précédentes d_{AB} et d_{AC} .

Il est certain que si l'écart des moyennes relatives à un ou plusieurs constituants chimiques est considérable, lorsqu'on passe d'un atelier à un autre, les fluctuations aléatoires des variables ne seront jamais en mesure de perturber les classifications. Par

contre il faudra être très attentif aux erreurs éventuelles de classification, lorsque des groupes de composition présentant des caractéristiques peu différentes seront présents dans une même analyse de grappes.

On remarquera toutefois que les inversions de proximités qui résultent des fluctuations des variables ne conduisent pas nécessairement à des rattachements anormaux et donc à des erreurs de classification. Il faut en outre que le *voisinage* des céramiques rapprochées par les fluctuations ne s'oppose pas à leur fusion. C'est ainsi qu'il ne suffit pas, dans l'exemple de la figure 6, que A devienne plus proche de C que de B, pour que A et C fusionnent. Cette fusion peut en effet ne pas se produire si d'autres *individus*, d'autres céramiques, de l'atelier 1, plus proches encore de A que ne l'est C, sont présents dans la grappe.

L'*effet de voisinage* est responsable de nombreuses particularités des analyses de grappes. Nous en montrerons quelques aspects en reprenant sur un diagramme à deux dimensions (Fig. 7) le cas étudié précédemment. Outre les céramiques A, B et C figurées par des points dont les coordonnées reproduisent les compositions, un certain nombre d'autres exemplaires des ateliers 1 et 2 constituent l'*échantillon* (ou échantillonnage) étudié.

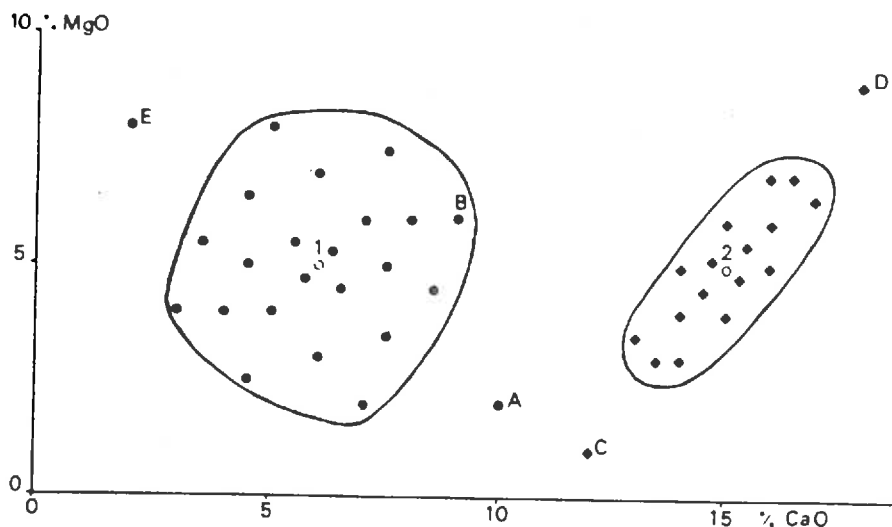


Fig. 7 : Représentation dans un plan (CaO, MgO) d'un échantillon restreint des ateliers (1) et (2).

Imaginons que l'ensemble des exemplaires de la figure 7 fasse l'objet d'une classification par analyse de grappes. On supposera, pour simplifier, qu'on est en affinité moyenne non pondérée, et qu'on utilise des distances euclidiennes sur données brutes ce qui a pour effet de confondre ces distances avec les longueurs qui, sur la figure 7, séparent les points représentatifs des céramiques. L'analyse de grappes réunira d'abord les exemplaires qui sont les plus proches les uns des autres, ceux dont les points représentatifs sont à l'intérieur des deux contours. Ces exemplaires constitueront deux pseudo-céramiques ayant pour composition les compositions moyennes des ateliers 1 et 2, également reportées sur la figure 7. Parvenu à ce stade des regroupements ce sera le

tour des céramiques A et C de fusionner, puis D et E se rattacheront à leurs ateliers respectifs, enfin la pseudo-céramique AC se réunira à la pseudo-céramique la plus proche, c'est-à-dire à celle dont la composition est à peu de chose près la moyenne de l'atelier 2.

L'exemplaire A est donc mal classé, ce qu'on serait quand même parvenu à soupçonner, si l'on avait ignoré son origine, en observant que cet exemplaire a le pourcentage de chaux le plus faible de tout le groupe auquel il se trouve rattaché. De même si AC s'était trouvé rattaché à l'atelier 1 (simplement par l'effet d'un changement de variables) on en viendrait à soupçonner le mauvais classement de C en remarquant cette fois que son pourcentage de chaux est le plus élevé de tout son groupe. D'une manière générale on ne saurait tirer de conclusions d'une analyse de grappes sans une étude attentive des compositions des céramiques de chacun des groupes, et sans réexaminer soigneusement le cas des exemplaires de composition marginale. C'est la raison pour laquelle les méthodes d'analyse de grappes fournissent en complément du dendrogramme la liste des compositions des différentes céramiques, regroupées dans l'ordre où elles se présentent sur ce dendrogramme.

Dans le cas étudié ici les effets de voisinage sont évidents. En effet, c'est parce que les exemplaires A et C se trouvaient isolés des autres céramiques, qu'a été rendue possible la fusion de A et de C, et, par voie de conséquence, le classement erroné de l'un de ces deux exemplaires. Si l'échantillon prélevé sur les deux ateliers avait été plus important, comme c'est le cas sur la figure 8, il ne fait aucun doute que les céramiques A et C eussent été correctement rattachées à leurs ateliers respectifs. De même les céramiques E et D qui devaient paraître très marginales sur la grappe correspondant à la figure 7, perdraient ce caractère sur la grappe issue de la figure 8.

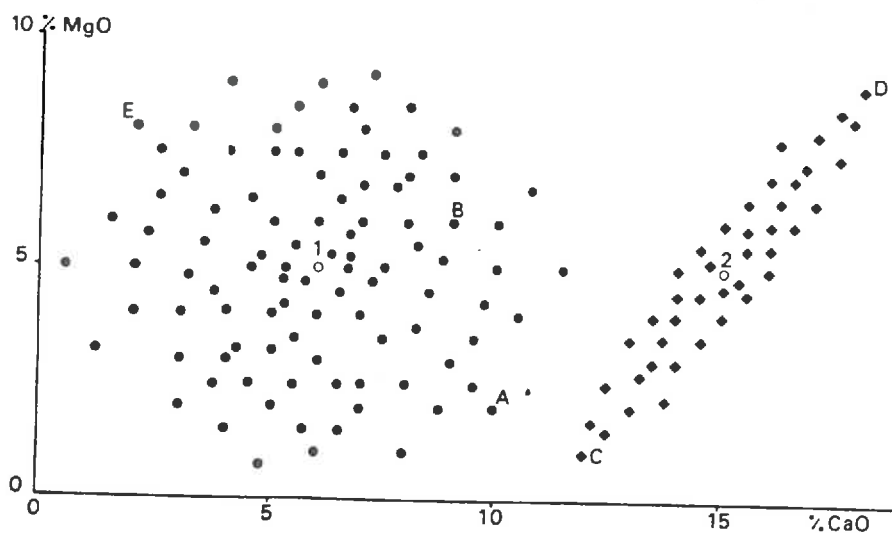


Fig. 8 : Représentation dans un plan (CaO, MgO) d'un échantillon large des ateliers (1) et (2).

La conclusion qui résulte de ces constatations, et qui se généralise à un nombre quelconque de variables, c'est que la séparation par analyse de grappes de deux groupes voisins nécessite un échantillon suffisamment important de l'un et l'autre groupe. Si

l'échantillon prélevé sur chacun des deux ateliers est insuffisant, en sorte que les distances entre individus d'un même atelier soient, dans l'espace à n variables, équivalentes aux distances minimales allant de la frontière d'un atelier à celle de l'autre, il devient impossible d'effectuer une classification satisfaisante des deux ateliers par analyse de grappes. C'est là une des raisons qui incitent à n'utiliser qu'avec beaucoup de prudence l'analyse de grappes comme méthode de classement pour des individus isolés ou pour de petits groupes, ce qu'on peut être tenté de faire en introduisant dans un groupe inconnu des individus d'origine connue ou en procédant à l'inverse. Mais ces réserves ne concernent pas les groupes importants, et l'on a déjà signalé que l'interpénétration très étroite de deux groupes sur un dendrogramme constitue généralement une sérieuse présomption en faveur de leur identité d'origine.

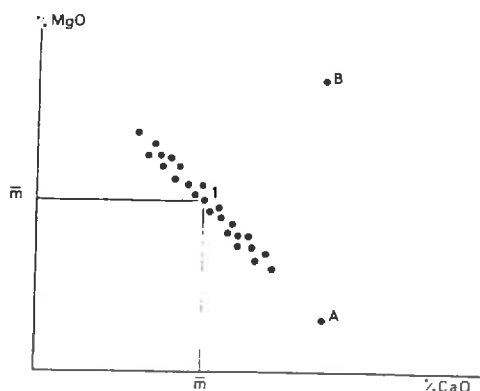


Fig. 9 : Positions de deux exemplaires A et B par rapport à la moyenne d'un groupe présentant une forte corrélation CaO/MgO.

Une autre raison qui incite à se méfier des attributions et donc des classements faits par analyse de grappes, c'est le caractère non géométrique des distances qui y sont employées, alors que les points représentatifs des céramiques sont situés dans un espace à n dimensions (ou n variables). Nous illustrerons cette dernière difficulté en revenant une fois de plus aux deux variables précédentes : CaO et MgO, qui définissent un plan où se trouvent les points représentatifs d'un groupe 1 qui présente la particularité d'être très allongé (Fig. 9). Cette particularité correspond au fait que les pourcentages de chaux et de magnésie de ce groupe ne varient pas d'une manière indépendante ; dans le cas présent les pourcentages de l'un des constituants diminuent lorsqu'augmentent ceux de l'autre. On dit qu'on a affaire à des *variables non indépendantes* et qu'il existe une *corrélation* entre ces deux variables (ce qui était aussi le cas pour l'atelier 2, sur les figures 7 et 8).

Lorsqu'en analyse de grappes par affinité moyenne non pondérée on sera parvenu à rassembler en une pseudo-céramique unique, ayant la composition moyenne du groupe 1, toutes les céramiques à l'exception des exemplaires A et B, il est clair que ces deux exemplaires sont alors à égale distance de la pseudo-céramique 1. Or ceci ne traduit absolument pas le fait que l'exemplaire A présente une certaine probabilité d'appartenir au groupe 1 tandis que B, par suite de la forme même du groupe, n'en a aucune. Il y a là une limitation sérieuse pour l'analyse de grappes, imposée par l'emploi des distances euclidiennes ; elle ne concerne pas seulement les analyses de grappes servant aux classements, mais aussi les classifications. De tels risques pourraient être graves si l'on ne disposait avec l'analyse discriminante quadratique des moyens de correction nécessaires.

III. ANALYSE DISCRIMINANTE QUADRATIQUE

A. Objectifs

L'analyse discriminante quadratique sert à attribuer une céramique donnée au groupe de composition auquel elle appartient, ce groupe étant en principe constitué de céramiques ou d'argiles ayant une même origine, connue, présumée ou inconnue. L'opération consiste à sélectionner le groupe auquel appartient la céramique, parmi un ensemble de groupes préexistants, la sélection se fondant sur les caractéristiques de composition de la céramique et des groupes utilisés en référence. Il s'agit là d'une opération de classement puisqu'on attribue la céramique étudiée à une certaine classe ou groupe de composition. Cette attribution se fait *individuellement* et s'exprime en termes probabilistes.

Comme les groupes de référence ont le plus souvent une origine connue, l'analyse discriminante quadratique peut sembler constituer à tort la partie essentielle des déterminations d'origine. En fait elle n'en est que l'instrument, et nous y reviendrons dans la dernière partie de ce cours.

B. Principes

L'analyse discriminante quadratique est une méthode de calcul qui permet de déterminer des probabilités d'appartenance; elle ne se traduit donc par aucune représentation graphique. Elle consiste dans son principe à extraire d'une base de données un certain nombre de groupes parmi lesquels on suppose que se trouve, ou pourrait se trouver, celui auquel appartient la céramique étudiée. On cherche ensuite par le calcul quel est le groupe qui présente le plus d'affinité avec la céramique étudiée. Pratiquement on détermine la probabilité qu'a la céramique d'appartenir à chacun des groupes considérés, la somme de ces différentes probabilités d'appartenance étant égale à 1.

L'attribution se fait naturellement au groupe qui présente la probabilité d'appartenance la plus élevée. Mais il existe de nombreux cas où toute attribution demeure impossible, en particulier si le groupe auquel appartient la céramique étudiée ne figure pas sur la liste des groupes retenus pour le calcul! La difficulté en ce cas tient au fait qu'il existe toujours un groupe de la liste, qui est plus proche de la céramique étudiée que les autres. Or il est essentiel de pouvoir distinguer les proximités qui ne sont qu'accidentelles, de celles qui résultent du fait que la céramique appartient bien au groupe qui a été sélectionné. A cet égard les cas particulièrement critiques ne sont pas ceux où les ressemblances sont toutes fort lointaines; on peut les éliminer facilement. Mais avec des ressemblances plus marquées le risque existe d'attribuer la céramique à un groupe dont elle ne provient pas.

Pour minimiser ce risque il semblerait souhaitable de procéder à des calculs de probabilité à deux niveaux. Le premier niveau est celui de la *probabilité inter-groupes* qui vient d'être exposé; il sélectionne le groupe le plus probable. Le second niveau serait celui de la *probabilité intra-groupe* qui préciserait la plus ou moins grande probabilité qu'a la céramique d'appartenir au groupe sélectionné. On déterminerait donc dans un premier temps le groupe auquel la céramique a le *plus de chance* d'appartenir, et, dans

un second temps, si la céramique a réellement *beaucoup de chance* d'appartenir au groupe le plus probable, *une chance moyenne* ou *peu de chance*.

Le principe de la probabilité intra-groupe tient à cette constatation que plus on s'éloigne des caractéristiques d'un groupe, plus on a de risque de se trouver dans le cas d'une ressemblance accidentelle. Malheureusement, en l'absence d'une solution mathématique rigoureuse, on ne peut que procéder à une évaluation sommaire des risques. Le procédé utilisé pour cette évaluation fait appel à la *distance généralisée* ou *distance de Mahalanobis* qui intervient déjà dans le calcul de la probabilité inter-groupes. Elle mesure la distance d'une céramique à un groupe d'une manière beaucoup plus satisfaisante que ne le ferait par exemple la distance euclidienne existant entre cette céramique et la moyenne du groupe. Appliquée aux céramiques constituant le groupe lui-même, elle mesure leur plus ou moins grand éloignement des caractéristiques de ce groupe. On peut ainsi ordonner les distances des différentes céramiques constitutives du groupe entre une limite inférieure et une limite supérieure, et comparer la distance de la céramique étudiée aux distances précédentes. Si cette distance est voisine de la limite inférieure, on peut considérer que la probabilité d'appartenance de la céramique au groupe est élevée. Si cette distance est au contraire proche de la limite supérieure, voire plus grande, la probabilité d'appartenance de la céramique au groupe est faible ou nulle. La position de la céramique étudiée, parmi l'ensemble des distances de Mahalanobis du groupe, s'exprime par un coefficient r qui vaut 100 lorsque la distance de la céramique au groupe est plus courte ou égale à la limite inférieure des distances du groupe, et qui vaut 0 lorsque cette distance est plus grande ou égale à la limite supérieure. C'est ce coefficient r qui permet, à défaut d'un véritable calcul de la probabilité intra-groupe, une évaluation sommaire des risques d'attribution erronée.

C. Calculs

Nous donnons ici, à toutes fins utiles, les formules fondamentales de l'analyse discriminante quadratique.

La probabilité d'appartenance d'une céramique x au groupe i , parmi k groupes considérés, est P_i telle que :

$$P_i = d_i(x) / \sum_{j=1}^{j=k} d_j(x)$$

avec $d_i(x) = (2\pi)^{-n/2} \cdot |C_i|^{-1/2} \cdot \exp \left[-1/2 (x - \bar{m}_i)' C_i^{-1} (x - \bar{m}_i) \right]$

$d_i(x)$ est la densité de probabilité relative à la céramique x et au groupe i .

$|C_i|$ est le déterminant de la matrice de dispersion ou matrice des variantes et covariances C_i . Cette dernière est telle que :

$$C_i = \begin{pmatrix} \sigma_1 \sigma_1 & \sigma_1 \sigma_2 r_{12} & \sigma_1 \sigma_3 r_{13} & \dots & \sigma_1 \sigma_n r_{1n} \\ \sigma_2 \sigma_1 r_{21} & \sigma_2 \sigma_2 & \sigma_2 \sigma_3 r_{23} & \dots & \sigma_2 \sigma_n r_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ \sigma_n \sigma_1 r_{n1} & \sigma_n \sigma_2 r_{n2} & \sigma_n \sigma_3 r_{n3} & \dots & \sigma_n \sigma_n \end{pmatrix}$$

$(x - \bar{m}_i) = (x_1 - \bar{m}_1, x_2 - \bar{m}_2, x_3 - \bar{m}_3, \dots, x_n - \bar{m}_n)$,
 $(x - \bar{m}_i)'$ est la matrice transposée de $(x - \bar{m}_i)$,

$x_1, x_2, x_3, \dots, x_n$ sont les pourcentages des constituants chimiques 1, 2, 3, ..., n, dans la céramique x,

$\bar{m}_1, \bar{m}_2, \bar{m}_3, \dots, \bar{m}_n$ sont les valeurs moyennes des constituants chimiques 1, 2, 3, ..., n, pour le groupe i,

$\sigma_1, \sigma_2, \sigma_3, \dots, \sigma_n$ sont les écarts-types de ces mêmes constituants pour le groupe i,

r_{12}, r_{13}, \dots sont les facteurs de corrélation relatifs au groupe i.

La distance généralisée ou distance de Mahalanobis d'une céramique x au groupe i est telle que son carré soit :

$$d_{\text{mah}}^2 = (x - \bar{m}_i)' C_i^{-1} (x - \bar{m}_i)$$

Enfin le coefficient r a pour expression $r = 100l/L$ où l est le nombre de céramiques du groupe qui ont une distance de Mahalanobis supérieure à celle de la céramique étudiée, L étant le nombre total de céramiques du groupe.

D. Conditions d'emploi

On notera d'abord que les relations précédentes supposent que les variables utilisées aient une *répartition normale* (ou de Laplace-Gauss) pour les différents groupes considérés. Cela signifie, simplement, que les histogrammes relatifs à chacun de ces groupes doivent dessiner pour chaque variable une courbe en cloche (ou courbe de Gauss). On peut toutefois montrer qu'en présence d'une répartition un peu différente, s'écartant donc de la *normalité*, ces relations demeurent utilisables, moyennant quelques précautions que nous indiquerons au passage.

L'intérêt majeur de l'analyse discriminante quadratique résulte de l'introduction de la distance généralisée de Mahalanobis. Celle-ci prend en compte, dans l'évaluation des ressemblances, la répartition spatiale, dans l'espace à n variables, des céramiques constituant les groupes de référence, ou, si l'on veut, tient compte des lois internes régissant les compositions des groupes (voir par exemple le cas de la figure 9). Il en résulte que la plupart des incertitudes ou des erreurs qui sont inhérentes à l'analyse de grappes se corrigent en analyse discriminante, celle-ci constituant alors un puissant moyen de contrôle des classifications. Mais c'est quand même à l'occasion d'attributions de céramiques à des groupes préalablement définis que l'analyse discriminante quadratique se révèle irremplaçable. Ce qui ne signifie pas, pour autant, que son emploi soit sans risque.

Outre le fait évident qu'on ne saurait faire d'attributions correctes qu'à des groupes correctement constitués, diverses causes peuvent entraîner des erreurs d'attribution si l'on n'y porte pas attention. Ces erreurs se rencontrent en analyse discriminante, comme en analyse de grappes, seulement lorsqu'on est en présence de groupes ayant des caractéristiques voisines. Lorsque la céramique étudiée présente par exemple un pourcentage qui s'écarte anormalement de la moyenne de son groupe, cela peut entraîner une attri-

bution à un autre groupe que le sien, surtout dans le cas où la variation de l'écart-type entre les deux groupes accentue encore les effets de cette fluctuation. Prenons le cas par exemple d'un groupe A dont le pourcentage moyen d'oxyde de titane et l'écart-type correspondant sont 0.90 ± 0.05 , et un groupe B pour lequel ces valeurs sont 0.95 ± 0.05 . Une céramique appartenant à A, et ayant par suite d'une fluctuation accidentelle un pourcentage d'oxyde de titane de 1.02, se situera à 2.4σ de A et 1.4σ de B, soit une différence d'un écart-type, sans doute assez facile à rattraper par l'effet des autres constituants. Mais si les caractéristiques du titane des groupes A et B sont respectivement 0.90 ± 0.02 et 0.95 ± 0.07 , un pourcentage de 1.02 se trouvera à 6σ de A et à 1σ de B, soit une différence de 5σ qui est généralement impossible à rattraper par les autres constituants et aura pour effet d'attribuer la céramique au groupe B. Ces risques, et ceux qu'on peut facilement imaginer sur un modèle semblable, ne sont toutefois pas très difficiles à prévoir, et donc à éviter, en comparant les compositions moyennes et les écarts-type des différents groupes.

D'une autre nature sont les erreurs de classement qui résultent d'écarts trop importants à la normalité. Bien que les relations permettant le calcul des probabilités d'appartenance aient un assez large domaine de validité, ce que nous avons déjà signalé, on ne saurait s'écarter sans risque de la normalité, au point que la moyenne des groupes ne correspondent plus du tout aux maxima de fréquences, ou que les coefficients de corrélation et les écarts-type n'aient plus de signification concrète. Par ailleurs il ne paraît pas souhaitable de faire intervenir des groupes dont les dispersions soient trop différentes, ce qu'on peut contrôler en comparant les déterminants de leurs matrices des variances et covariances. Si la dispersion d'un groupe est trop élevée, il est préférable, dans la mesure du possible, de scinder ce groupe en plusieurs populations mieux regroupées.

Signalons enfin pour mémoire que lorsqu'une céramique se trouve à égale distance de deux groupes, les relations définissant les probabilités d'appartenance indiquent que c'est toujours le groupe dont le discriminant est le plus faible auquel se fait l'attribution. En pratique on ne tient pas compte de telles attributions, ni, d'une manière plus générale, de toutes celles où le partage des probabilités entre deux ou plusieurs groupes a pour effet de ramener la probabilité la plus élevée à des valeurs inférieures à 0.9 (la probabilité maximale étant égale à l'unité). Les écarts, même faibles, de normalité des groupes, les fluctuations de l'échantillonnage et les variations aléatoires des variables font qu'il suffit alors de très peu de chose pour que le partage des probabilités se fasse en faveur du groupe qui arrive en seconde position. Dans ces conditions mieux vaut laisser en suspens l'attribution.

E. *Présentation des résultats*

Pour se familiariser un peu plus avec l'analyse discriminante quadratique on a reproduit sur la figure 10 une présentation de résultat correspondant à une céramique de type italique, GRA 73, trouvée à La Graufesenque. Les deux premières colonnes de gauche donnent la liste des groupes de référence utilisés; ils correspondent à des ateliers ayant produit ce type de céramique. Les groupes de référence sont rangés d'après les valeurs décroissantes des probabilités d'appartenance relatives à la céramique GRA 73 (dernière colonne à droite); ils sont repérés par un numéro de code et par leur nom en clair. On

peut à l'impression garder la liste complète des groupes qui sont intervenus dans le calcul, ou ne conserver que les groupes pour lesquels la céramique GRA 73 a les probabilités d'appartenance les plus élevées. Dans la colonne suivante n désigne le nombre d'exemplaires ayant servi à constituer chacun des groupes de référence, ce qui permet de se rendre compte des risques de sous-représentativité de certains groupes (ne pas confondre ce n avec celui qui apparaît dans les relations de définition précédentes; n est ici l'équivalent de L dans la relation de définition du coefficient r). La quatrième colonne donne la densité de probabilité de la céramique GRA 73 pour chacun des groupes (avec un luxe de chiffres correspondant à une précision parfaitement illusoire). Enfin les cinquième et sixième colonnes donnent les valeurs du coefficient r , et la distance de Mahalanobis de la céramique GRA 73 à chacun des groupes.

La lecture rapide des résultats se fait en commençant par la probabilité d'appartenance la plus élevée, en haut de la colonne de droite, ce qui sélectionne donc le groupe le plus probable de la liste. Si cette probabilité est inférieure à 0.90 on considère généralement que le classement n'est pas possible. Dans le cas contraire la lecture se continue par le coefficient r du groupe de plus grande probabilité; il détermine si l'on peut ou non accorder quel crédit à la sélection qui vient d'être faite. Pour la céramique GRA 73 le coefficient r est élevé; l'attribution de cette céramique au groupe arezzo-2 est donc hautement probable (les autres valeurs fournies par le calcul sont utilisées pour une lecture plus approfondie des résultats).

GRA 73

code	nom du groupe	n	dens. prob.	r	dmah	prob
414	arezzo-2	91	0.2958188E+06	75.00	2.22	0.94
403	cincelli	18	0.1879290E+05	0.00	3.77	0.06
415	pise-ateius-2	134	0.9463095E+00	2.00	5.22	0.00
413	campanien-3	56	0.3504153E+00	1.00	5.29	0.00
411	pouzzoles	26	0.1032520E-03	0.00	7.08	0.00
409	padan-1-2	74	0.5230292E-18	0.00	10.63	0.00
412	campanien-2	20	0.0000000E+00	0.00	17.99	0.00
410	padan-3	17	0.0000000E+00	0.00	21.03	0.00
402	muette-1	97	0.0000000E+00	0.00	13.82	0.00
400	muette-2	28	0.0000000E+00	0.00	33.07	0.00

Fig. 10 : Exemple de sortie d'une analyse discriminante quadratique.

IV. LOGIQUE ET TRAITEMENT DES DONNÉES

On a vu jusqu'ici différents aspects du traitement des données fournies par l'analyse : certaines méthodes, leur principe, leurs limites aussi, et l'on pourrait envisager de clore là cette initiation. Ce faisant il semblerait toutefois qu'on accrédi-terait une version, presque inexacte à force d'être partielle, du traitement des données d'analyse. Il faudrait en effet bien se garder de croire que les problèmes de classification et de classement des céramiques se ramènent à l'emploi raisonné de quelques méthodes, judicieusement choisies et intelligemment exploitées. Nous ne prendrons ici qu'un seul exemple, facilement généralisable à d'autres cas.

Que l'on utilise des méthodes de classification ou des méthodes de classement, l'objectif que l'on se fixe le plus souvent est d'arriver à déterminer l'origine des céramiques. Or on ne peut traiter de ces problèmes (par des techniques d'analyse) qu'en comparant des céramiques ou des argiles d'origine connue et des céramiques d'origine inconnue. Si nous avons affaire par exemple à une céramique inconnue, x , que l'on compare à une céramique ou à un groupe de céramiques d'origine connue, A , les méthodes de classification et de classement nous permettront d'apprécier, voire de claculer la plus ou moins grande ressemblance existant entre x et A ; elles ne nous permettront jamais à elles seules de dire que x et A ont même origine. Ce serait demander aux méthodes de traitement des données autre chose que ce qu'elles peuvent fournir, qui n'est rien de plus que la mise en évidence d'une certaine ressemblance entre x et A . Le passage de la constatation d'une ressemblance à l'affirmation d'une communauté d'origine ne relève pas du traitement des données. Mais la tentation est souvent grande de laisser croire que ce sont les méthodes de classification ou de classement qui autorisent des conclusions relevant en fait d'une autre logique, fréquemment plus exigeante et difficile à mettre en œuvre.

Si la ressemblance entre x et A ne permet pas de conclure à une même origine, c'est qu'on doit admettre le *postulat fondamental* suivant lequel il est à priori possible de trouver une autre source d'argile, différente de celle de A , dont la ressemblance avec x soit plus forte que celle de A , ou pour le moins égale. Le seul moyen qui permette de s'affranchir des graves conséquences de ce postulat consiste à ne plus se contenter de la comparaison de x et de A , mais à étendre les comparaisons à un très grand nombre de matériaux dont les origines, connues, tissent au sol un réseau suffisamment dense et suffisamment étendu. Si, parmi tous les points du *réseau de renseignements localisés* ainsi constitué le lieu d'origine de A est celui dont les caractéristiques sont les plus proches de x , on est en droit de conclure à la communauté d'origine de x et de A .

Le plus souvent on en viendra à définir, plutôt qu'une origine ponctuelle, une *zone d'incertitude*, voire même une *zone de conjonction*. Par ailleurs il est bien rare, même lorsqu'on a affaire à des *problèmes très localisés*, qu'on puisse établir un réseau suffisamment continu pour que l'attribution de x à A ne laisse place au moindre doute. Cela devient plus évident encore dans le cas des *problèmes peu localisés* où les déterminations d'origine feront nécessairement appel à un certain nombre de *probabilités à priori*. On retrouve ici les différentes notions qui ont fait l'objet du cours sur les « problèmes de détermination de l'origine des céramiques ».

Dans tous les cas une détermination d'origine devient une opération qui nécessite un grand nombre de renseignements localisés réunis en un réseau cohérent, donc une opération longue et difficile à mettre en œuvre. Il serait évidemment bien plus commode de croire que tout peut se résoudre par une comparaison ne faisant intervenir que x et A , quitte à faire endosser aux méthodes de traitement des données la responsabilité de conclusions qui ne les concernent plus. Remarquons, et ce n'est pas une consolation, que les méthodes de traitement des données ne sont pas les seules à être abusivement sollicitées! Que ne dit-on des méthodes d'analyse des céramiques où la course aux listes les plus longues possibles de constituants chimiques a eu pour objet principal de se persuader qu'on parviendrait ainsi à décider de l'attribution de x à A en faisant l'économie des réseaux de renseignements localisés. C'est la théorie du « fingerprint », de l'empreinte digitale, qui exprime bien qu'on pensait ainsi mettre en échec le postulat fondamental des déterminations d'origine. On sait à présent que cet espoir est vain. Il faut donc veiller, dans le cas des analyses comme dans celui du traitement des données, à ne pas faire porter aux méthodes le poids qui, normalement, devrait revenir aux renseignements localisés des réseaux.

BIBLIOGRAPHIE

- P. DAGNELLE, *Théorie et méthodes statistiques*, 2 vol., Les Presses Agronomiques de Gembloux, Gembloux, 1975.
ID., *Analyse statistique à plusieurs variables*, Les Presses Agronomiques de Gembloux, Gembloux, 1977.
P. LAFFITTE, *Traité d'informatique géologique*, Masson, Paris, 1972.
J.-M. ROMODER, *Méthodes et programmes d'analyse discriminante*, Dunod, Paris, 1973.

M. Maurice PICON
Centre de recherches archéologiques CNRS
URA 3 Université de Lyon
Maison de l'Orient méditerranéen ancien
1, rue Raulin
69365 Lyon Cedex 2 (France)
Tél. (78) 69 24 45 - (78) 72 02 53