

**NOTICE D'INSTALLATION ET D'UTILISATION DE PROGRAMMES BASES SUR
L'ALGORITHME DE KOHONEN ET DEDIES A L'ANALYSE DE DONNEES**

3/10/2001 (V8.2)

**Patrick Letrémy
MATISSE-SAMOS UMR CNRS 8595
Université Paris 1
pley@univ-paris1.fr**

*Algorithme d'apprentissage de Kohonen (SOM)
Algorithme KACP (Analyse d'un tableau de données)
Algorithme KFAST (Kohonen à zéro voisin)
Algorithme KBATCH (Kohonen dans la version déterministe)
Algorithme KORRESP (Analyse d'un tableau de contingence)
Algorithme KACM(j=1,2) (Analyse d'une table de Burt ou d'un tableau Disjonctif complet)*

I - CONTENU DE LA DISQUETTE (OU DU REPERTOIRE), INSTALLATION

La distribution contient les trois fichiers suivants :

*AVERTI8.TXT : le fichier d'annonce
DABOR8.TXT : le fichier d'introduction
KSET_V8.ZIP : le fichier à décompresser*

L'installation par DEFAULT comporte 2 étapes:

*1) - Décompresser, avec WinZip, le fichier KSET_V8.ZIP de la disquette (ou du répertoire) sur c:\
2) - Charger et soumettre (Submit ou F3) dans SAS l'exécutable KINSTAL.SAS situé dans le répertoire C:\K8.*

Pour plus de détails concernant l'installation par défaut, lire le paragraphe II "DESCRIPTION DE LA PROCEDURE D'INSTALLATION" de la présente notice.

L'installation terminée, il est conseillé de lire le paragraphe III "CONTENU DES PROCEDURES ET REFERENCES" de la présente notice.

Pour une information complète sur l'utilisation des programmes, il faut consulter les paragraphes IV, V, VI, VII et VIII du présent fichier C:\K8\LADOK8.DOC.

II - DESCRIPTION DE LA PROCEDURE D'INSTALLATION

Vous avez suivi les indications de la section I, qui consistent en deux étapes:

*- Décompression du fichier KSET_V8.ZIP de la disquette (ou du répertoire) sur c:\
- Chargement et lancement dans SAS de l'exécutable KINSTAL.SAS situé dans le répertoire c:\k8.*

Vous trouverez ci-dessous la liste des fichiers créés au cours de l'installation.

L'étape 1) génère sur le disque dur C: l'arborescence suivante

*C:\K8
C:\K8\ABRI
C:\K8\ABOUL
C:\K8\AKOR
C:\K8\AMAC*

Contenu du répertoire K8 :

*1. LADOK8.DOC le présent fichier.
2. les batch de lancement des algorithmes :
BKOHN.SAS, DKACP.SAS, DKFAST.SAS, DKBATCH.SAS, DKORR.SAS et DKACM.SAS
3. KOPTIONS.SAS est une macro qui définit les chemins d'accès des catalogues et des jeux de données.*

4. *KOLCLUS.SAS* donne un aperçu des 22 couleurs disponibles pour les Clusters (Super Classes).
5. *KOLCAM.SAS* donne un aperçu des 13 couleurs disponibles pour les secteurs de Camemberts
6. *KOMAC.SAS* permet la création du catalogue permanent de macros *SASMACR.sas7bcat* dans *C:\K8\AKOR* (cf. le paragraphe VIII " TRAITEMENTS COMPLEMENTAIRES ").
7. *KINSTAL.SAS* déjà cité.
8. *DSTAT.SAS* est un utilitaire qui permet de calculer, pour une classification obtenue à partir de variables quantitatives, les moyennes conditionnelles, l'analyse de la variance et certaines statistiques de tests multidimensionnels.
9. *DTAIL.SAS* est un utilitaire qui permet, à partir d'un tableau de réponses, de créer le tableau disjonctif complet et la table de Burt sous forme de deux tables sas; dans le cas où les réponses sont codées sous forme numérique, l'utilitaire fournit une troisième table qui correspond à un recodage des réponses sous forme de libellés.
10. *DMISS.SAS* est un utilitaire qui permet de contrôler la présence de données manquantes.
11. *DVIATION.SAS* est un programme de calcul de déviations et d'écartés pondérés qu'il faut actualiser à chaque utilisation (voir la fin du paragraphe VIII).

Contenu du répertoire *K8\ABRI* : 188 fichiers de programme qui seront compilés dans des catalogues IML permanents

Contenu du répertoire *K8\ABOUL* : jeux de données (8 tables sas) = *BLAYO.SAS7BDAT*, *BLAYO2.SAS7BDAT*, *BLAYO3.SAS7BDAT*, *BOPAY3.SAS7BDAT*, *SUPPL96.SAS7BDAT*, *MONUMENT.SAS7BDAT*, *CHIENS.SAS7BDAT* et *PIRON .SAS7BDAT*.

Contenu du répertoire *K8\AKOR* : au niveau de l'étape 1) il est vide

A l'étape 2) le batch *KINSTAL.SAS* permet la création et le stockage dans le répertoire *C:\K8\AKOR* de 8 catalogues IML permanents :

1. *KOUPEL.sas7bcat* (9 entrées) utilisé par *KOHONEN*
2. *KOMPA.sas7bcat* (37 entrées) qui est commun à tous les algorithmes
3. *KOKAR.sas7bcat* (31 entrées) qui est commun à l'utilisation de *KACP*, *KFAST* et *KBATCH*
4. *KOUL.sas7bcat* (14 entrées) qui est propre à l'utilisation de *KORRESP*
5. *KALM.sas7bcat* (19 entrées) qui est commun à l'utilisation de *KACM*, *KACM1* et *KACM2*
6. *KOUTO.sas7bcat* de 28 entrées (pour les traitements complémentaires)
7. *KABA.sas7bcat* de 26 entrées (pour les traitements complémentaires)
8. *KANIF.sas7bcat* de 24 entrées (pour les traitements complémentaires)

A la fin de l'exécution de *KINSTAL*, pour vérification, apparaissent dans la fenêtre *OUTPUT* de *SAS* : le nom et le contenu de ces 8 catalogues IML.

Contenu du répertoire *K8\AMAC* : 19 fichiers de macro-commandes = *MCELL.SAS*, *MCONT.SAS*, *MDIMA.SAS*, *MDIMA12.SAS*, *MDIST.SAS*, *MFIND.SAS*, *MFMOD.SAS*, *MGVAR.SAS*, *MKAMF1.SAS*, *MKAMF2.SAS*, *MKAMG.SAS*, *MKLUF1.SAS*, *MKLUF2.SAS*, *MKLUG.SAS*, *MSCAL.SAS*, *MSCOL.SAS*, *MSTAT.SAS*, *MSUPP.SAS* et *MVC3D.SAS* qui seront compilés et stockés dans le catalogue *C:\K8\AKOR\SASMACR.sas7bcat* après la soumission de *KOMAC.SAS*.

III - CONTENU DES PROCEDURES ET REFERENCES

Thème : Programmes conversationnels pour :

- l'algorithme d'apprentissage de *KOHONEN* (*SOM*)
- l'algorithme *KACP* (version Kohonen de L'ACP)
- l'algorithme *KFAST* (version Kohonen à zéro voisin)
- l'algorithme *KBATCH* (version déterministe de Kohonen)
- l'algorithme *KORRESP* (version Kohonen de L'AFC)
- l'algorithme *KACM(j=1,2)* (versions Kohonen de L'ACM)
- traitements complémentaires à ces algorithmes

- Le paragraphe IV décrit l'exécution de l'algorithme de KOHONEN pour des entrées qui sont des points de R^2 choisis au hasard dans un carré.
 - Le paragraphe V décrit le déroulement de KACP (KFAST et KBATCH) sur un jeu de données.
 - Le paragraphe VI décrit le déroulement de KORRESP sur un jeu de données.
 - Le paragraphe VII décrit le déroulement de KACM($j=1,2$) sur un jeu de données.
- Le paragraphe VIII décrit des traitements complémentaires.

Références :

- F. Blayo et M. Verleysen, *Les réseaux de neurones artificiels*, Collection Que sais-je ?, vol 3042, PUF, 1996.
- M. Cottrell et P. Letrémy, Classification et analyse des correspondances au moyen de l'algorithme de Kohonen : application à l'étude de données socio-économiques, Prépublication du SAMOS n° 42 janvier 95, *Actes de Neuro Nîmes, 1994*.
- M. Cottrell et S. Ibbou, Multiple Correspondence Analysis of a crosstabulation matrix using the Kohonen algorithm, Prépublication du SAMOS n° 49, octobre 95, *Proc. of ESANN'95*, Editions D Facto, Bruxelles, 1995.
- M. Cottrell et E. de Bodt, A Kohonen map representation to avoid misleading interpretation, Prépublication du SAMOS n°62, juin 96, *Proc. of ESANN'96*, Editions D Facto, Bruxelles, 1996.
- Demartines, Analyse de données par réseaux de neurones auto-organisés, Thèse de Doctorat, Laboratoire TIRF, Institut National Polytechnique de Grenoble.
- Classification, Analyse des Correspondances et Méthodes Neuronales, Thèse de S. Ibbou soutenue le 20/1/98.
- Applications des Algorithmes d'Auto-Organisation à la Classification et à la Prévion, Thèse de P. Rousset soutenue le 3/12/99

IV - DEMONSTRATION

Algorithme d'apprentissage de KOHONEN pour des points dans un carré de R^2

Charger et soumettre dans SAS le fichier BKOHN.SAS qui se trouve dans $c:\backslash k8$. Ce programme de démonstration ne nécessite aucune donnée. **TOUTE LA SUITE EST CONVERSATIONNELLE**

Une fenêtre "CHOIX" apparaît, elle correspond à une boucle sans fin dont on ne peut sortir que par choix : A (Arrêt définitif). Sinon l'utilisateur peut choisir entre un réseau de dimension 1, (ficelle choix : F) ou un réseau de dimension 2 (grille choix : G). Pour toutes les fenêtres, le choix retenu est validé par F3.

- Si le choix est : F, une fenêtre "FICELLE" s'ouvre pour demander la taille de la ficelle (n), le nombre maximum d'itérations ($tmax$) et la périodicité de l'affichage ($periode$).

Par exemple : $n = 20$, $tmax = 200^1$ et $periode = 25$ suivi de F3 pour valider.

Alors, $200/25+1=9$ écrans graphiques seront produits et stockés provisoirement dans le catalogue WORK.GSEG.

A la fin des 200 itérations, l'utilisateur peut revenir sur l'un des 9 graphiques (à l'aide de "CATALOG WORK.GSEG") pour une éventuelle impression.

Dans le catalogue WORK.GSEG, les 9 graphiques sont nommés IMLG à IMLG8.

On peut toujours revenir à la fenêtre "CHOIX" avec la commande Window.

- Si le choix est : G, une fenêtre "GRILLE" s'ouvre pour demander la taille de la grille ($m =$ nombre de lignes et $n =$ nombre de colonnes), le nombre maximum d'itérations ($tmax$) et la périodicité de l'affichage ($periode$).

Par exemple : $n = 7$, $m = 7$, $tmax = 500^2$ et $periode = 50$ suivi de F3 pour valider.

¹ Dans cette démonstration, on peut prendre pour $tmax$ une valeur de l'ordre d'au moins 10 fois la taille du réseau

² id.

Alors, $500/50+1=11$ graphiques nommés IMLG à IMLG10 seront produits et stockés dans le catalogue WORK.GSEG.

Pour éviter un mélange de numérotation entre les graphiques "ficelles" et "grilles", on doit avant d'opter pour un nouveau choix (F ou G) détruire ou renommer le catalogue WORK.GSEG.

Il est donc possible de boucler indéfiniment sur choix : F ou G. Pour y mettre fin, il suffit de choisir: A suivi de plusieurs F3 pour revenir dans la fenêtre " PROGRAM EDITOR ".

V - L'ALGORITHME KACP

Algorithme KACP : Analyse en Composantes Principales version KOHONEN

Il s'agit d'une classification en classes liées par une structure de voisinage, qui utilise l'algorithme de Kohonen. On peut dire que cette analyse s'apparente à une analyse en composantes principales.

Pour illustrer le déroulement d'une session, nous utiliserons les données fournies par F.BLAYO qui concernent 53 pays en 1984 et qui sont présentées dans le Que sais-je ? vol 3042 "LES RESEAUX DE NEURONES ARTIFICIELS" de F. BLAYO et M. VERLEYSSEN.

Pour chaque pays on dispose des valeurs de 7 variables :

1. PAYS (nom du pays),
2. ANCRX (croissance annuelle de la population),
3. TXMORT (taux de mortalité infantile),
4. TXANAL (taux d'illettrisme),
5. SCOL2 (fréquentation scolaire du 2ème degré),
6. PIBH (PIB par habitant),
7. CRXPIB (croissance annuelle du PIB).

Ces données sont dans une table sas (BLAYO.sas7bdat) qui est placée dans le répertoire c:\k8\aboul. La table C:\K8\ABOUL\BOPAY3.sas7bdat (96 pays en 1996) est une mise à jour de la table BLAYO.sas7bdat, avec introduction de nouvelles variables comme les taux de chômage et d'inflation, ainsi qu'une variable qualitative indiquant le niveau de l'IDH (Indice du Développement Humain).

DEBUT DE LA SESSION KACP

Charger et soumettre dans SAS : DKACP.SAS qui se situe dans c:\k8.

Toute la suite est CONVERSATIONNELLE, des fenêtres s'ouvrent et guident l'utilisateur dans ses choix. Pour toutes les fenêtres, les choix retenus seront validés par F3.

Fenêtre n°1 : "CHEMIN".

L'utilisateur précise le chemin d'accès aux données (par défaut c:\k8\aboul) .

Exemple : path: C:\K8\ABOUL suivi de F3.

Dans le cas où le répertoire proposé est vide ou inexistant, un message d'erreur apparaît dans la fenêtre LOG, la commande WINDOW permet le passage de la fenêtre CHEMIN à celle du LOG.

Fenêtre n°2 : "CHOIXDAT".

A partir des tables sas situées dans le répertoire choisi dans "CHEMIN", l'utilisateur sélectionne (par X) la table retenue pour l'analyse.

Exemple : X BLAYO suivi de F3.

Fenêtre n°3 : "CHOIXVAR".

A partir de la liste des variables de la table sélectionnée dans "CHOIXDAT", l'utilisateur indique (par C) la variable qui identifie les observations (cette variable sera traitée comme étant de type Caractère) et (par N) les variables Numériques à retenir pour l'analyse.

Exemple :

CHOIX	VARIABLE
C	PAYS
N	ANCRX
N	TXMORT
N	TXANAL
N	SCOL2
N	PIBH
N	CRXPIB

suivi de F3.

Remarque: En cas de données manquantes un message apparaît dans l'OUTPUT indiquant leur nombre et le programme s'arrête.

Fenêtre n°4 : "STRUCTUR".

L'utilisateur choisit son type de réseau : F pour Ficelle, G pour Grille.

Exemple : choix : G suivi de F3.

Fenêtre n°5 : "PARAM_FI" ou "PARAM_GR".

Pour une Ficelle, "PARAM_FI" demande sa taille (n) et le nombre maximum d'itérations (tmax).

Pour une Grille, "PARAM_GR" demande le nombre de lignes (m), le nombre de colonnes (n) et le nombre maximum d'itérations (tmax).

Les nu unités du réseau (nu = n pour une ficelle ou nu = m×n pour une grille) seront numérotées de 1 à nu. Par exemple si nu=12 pour la ficelle de n=12, on aura : 1 2 11 12.

Pour la grille de m = 3 et n = 4 on aura :

1 4 7 10
2 5 8 11
3 6 9 12

On peut prendre pour tmax une valeur de l'ordre de 5 à 6 fois le nombre d'observations; ceci revient à dire qu'en moyenne chaque observation sera présentée 5 à 6 fois durant l'apprentissage. Notons que ce rapport peut être diminué si le nombre d'observations est très élevé.

Dans "PARAM_FI" comme dans "PARAM_GR" deux choix supplémentaires sont offerts.

Le premier choix (O/N) permet d'initialiser le générateur de nombre au hasard.

O : le point de départ est fixe, ce qui rend les résultats reproductibles (valeur par défaut).

N: le point de départ est calé sur l'horloge de la machine, ce qui rend les résultats NON reproductibles.

Le deuxième choix (O/N) propose le calcul et la représentation éventuels de la fonction "énergie". Cette "énergie" ou "potentiel" généralise la notion de "variance intra" en l'étendant aux plus proches voisins. Passé les trois quarts de tmax (à zéro voisin) les deux notions coïncident. En fin d'itérations, l'énergie doit se situer sur un minimum (local).

Le choix négatif (N) est la valeur par défaut; si le nombre d'itérations est élevé (dès 500), un choix affirmatif (O) peut s'avérer coûteux en temps de calcul et implique la création d'un graphique de nom : NRJ stocké dans un catalogue graphique permanent (cf. la fenêtre n°6 : "INFORM").

Exemple :

m = 8
n = 8
tmax = 300
CHOIX : O
CHOIX : O suivi de F3.

Fenêtre n°6 : "INFORM".

Dans cette fenêtre, l'utilisateur doit renseigner le champ Nom (fixé par défaut à `_TEMPOR_`) et décider (choix O ou N) de l'éventuelle visualisation et stockage d'un graphique de type "dx-dy" (cf. : la thèse de P. DEMARTINES dont la référence est donnée dans la section III).

Le champ Nom sera utilisé pour nommer des catalogues (fichier d'extension `.sas7bcat`) et des tables sas (fichier d'extension `.sas7bdat`). Il doit comporter au moins 5 caractères; dans le cas contraire, le programme complète la réponse par des X pour obtenir un champ Nom de 5 caractères.

Le choix : N est proposé par défaut; en cas de réponse positive (O), un graphique de nom `DX_DY` sera créé dans le catalogue graphique. C'est un nuage de points dans un carré de côté 1, où l'on compare les distances (normalisées à 1) théoriques entre les unités gagnantes avec les distances (normalisées à 1) euclidiennes entre les vecteurs poids associés. La situation idéale correspondant à une organisation parfaite serait celle où tous les points (représentés sur le graphique par des cercles) sont situés sur la diagonale du carré.

Exemple :

Nom : B8G3CPAY
choix : O suivi de F3.

Ces réponses impliquent la création de 4 tables sas et de 2 catalogues qui seront placés dans `C:\K8\ABOUL` (cf.: le champ path de la fenêtre n°1 : "CHEMIN").

Description des 4 tables sas :

1. la table `B8G3C_CL` (valeur par défaut : `_TEMP_CL`) donne, pour chaque unité gagnante, le contenu de la classe correspondante.
2. la table `B8G3C_UG` (valeur par défaut : `_TEMP_UG`) donne les valeurs moyennes des variables brutes pour chaque unité gagnante.
3. La table `B8G3C_WS` (valeur par défaut : `_TEMP_WS`) donne pour chaque unité (gagnante ou pas) sa position dans la grille : ligne, colonne, son effectif (zéro si l'unité n'est pas gagnante) et son "vecteur poids final " ou "vecteur représentant ".
4. La table `B8G3CPAY` (valeur par défaut : `_TEMPOR_`) donne pour chaque modalité (ici PAYS) de la variable (C) qui identifie les observations (cf. la fenêtre n°3 : "CHOIXVAR") son unité gagnante (sa classe) : `_codage_`, sa position dans la grille : ligne, colonne et ses valeurs pour les variables brutes ainsi que pour les variables éventuellement transformées (cf.: la fenêtre n°8 "PREPROC").

Description des 2 catalogues :

1. Le catalogue `iml B8G3CPAY` (valeur par défaut : `_TEMPOR_`) contient tous les intermédiaires de calcul (matrices de poids initiaux et finaux, liste des unités gagnantes , etc.) qui seront ultérieurement utilisés dans les traitements complémentaires (cf. le paragraphe VIII).
2. Le catalogue graphique `GACPB8G3` (valeur par défaut : `GACP_TEM`) contient 8 graphiques (au plus). Chaque graphique possède son nom et sa description qui correspond aux 40 premiers caractères de son titre.

Les 8 graphiques du catalogue graphique `GACPB8G3` :

Nom	Description
<code>DX_DY</code>	<code>dx_dy</code>
<code>G_DIMA</code>	Distances (M) avec les plus proches voisins
<code>G_PAVAGE</code>	KACP : grille 8x8 et 300 itérations
<code>HIS_GRI</code>	Représentants des classes (Poids Finaux)
<code>HIS_GRII</code>	Moyennes des classes (valeurs normalisées)
<code>LIN_GRI</code>	Représentants des classes (Poids Finaux)
<code>LIN_GRII</code>	Moyennes des classes (valeurs normalisées)
<code>NRJ</code>	Potentiel

Les graphiques, dans le catalogue, sont classés en ordre alphabétique alors qu'à l'affichage l'ordre sera

pour une grille :

G_PAVAGE, G_DIMA, LIN_GRI, HIS_GRI, LIN_GRII, HIS_GRII,DX_DY et NRJ

pour une ficelle :

F_PAVAGE, F_DIMA, LIN_FIC, HIS_FIC, LIN_FIC1, HIS_FIC1,DX_DY et NRJ

*G_DIMA pour une Grille (ou F_DIMA pour une Ficelle) permet d'apprécier pour chaque unité (classe) son effectif et les distances de **Mahalanobis** normalisées avec ses plus proches voisins (au plus 8 pour une grille et au plus 2 pour une ficelle). A l'exception des unités des bords, si l'unité est très proche de ses (8 ou 2) voisins, son polygone sera très proche des bords du carré. (d'après: A Kohonen map representation to avoid misleading interpretation de M. COTTRELL et E. DE BODT, ESANN'96).*

G_PAVAGE pour une Grille (ou F_PAVAGE pour une Ficelle) donne l'illustration graphique du contenu des classes (unités) du réseau (cf.: la table B8G3C_CL).

HIS_GRI pour une Grille (HIS_FIC pour une Ficelle) est un pavage d'histogrammes (diagramme en bâtons) des représentants des classes (poids finaux).

HIS_GRII pour une Grille (HIS_FIC1 pour une Ficelle) est un pavage d'histogrammes des moyennes normalisées pour les classes non vides.

LIN_GRI pour une Grille (LIN_FIC pour une Ficelle) est un pavage de courbes des représentants des classes (poids finaux).

LIN_GRII pour une Grille (LIN_FIC1 pour une Ficelle) est un pavage de courbes des moyennes normalisées pour les classes non vides.

Il n'y a pas de graphique de pavage pour une ficelle de plus de 50 unités

Fenêtre n°7 : "INITO".

L'utilisateur indique la façon d'initialiser les vecteurs poids initiaux :

E : s'ils sont pris au hasard entre le min et le max (choix par défaut).

H : s'ils sont pris au hasard parmi les entrées.

Fenêtre n°8 : "PREPROC".

L'utilisateur choisit un éventuel type de prétraitement de ses données :

A : pour Aucun

C : pour Centrer les variables

N : pour centrer et réduire les variables (Normer)

K : pour transformer les lignes de la matrice des données en pourcentages de somme 1 (profils lignes) et utiliser la distance du Khi deux sur ces profils.

Le choix N (normer) est proposé par défaut

Exemple : choix : N suivi de F3.

PATIENCE KACP travaille pour vous !!!!

Avant d'afficher les graphiques du catalogue graphique IML, les 5 autres fichiers permanents sont créés.

LA SESSION KACP EST TERMINEE

L'ALGORITHME KFAST

Algorithme KFAST : Version stochastique de l'algorithme de FORGY

Il s'agit d'une classification qui utilise l'algorithme de Kohonen à zéro voisin sur toutes les étapes du processus.

L'ALGORITHME KBATCH

Cet algorithme est la version déterministe de l'algorithme de KOHONEN.

Le déroulement d'une session KFAST (ou KBATCH) est similaire à celle de KACP : il suffit de charger et soumettre dans SAS : DKFAST.SAS (ou DKBATCH.SAS) qui se situe dans c:\k8.

VI - L'ALGORITHME KORRESP

Algorithme KORRESP : Version Kohonen de l'analyse d'un tableau de contingence

On analyse un tableau de contingence croisant deux variables qualitatives, au moyen d'un algorithme dérivé de l'algorithme de Kohonen. cf.: la référence M.Cottrell et P.Létrémy, dans les Actes de Neuro Nîmes, 1994 .

Pour illustrer le déroulement d'une session, nous utiliserons un tableau de contingence qui croise les monuments historiques en France, suivant leur catégorie (au nombre de p=11) et leur type de propriétaire (au nombre de q=6).

Ce tableau est stocké dans la table sas de nom MONUMENT qui possède p=11 observations et q+1=7 variables soit :

Une variable caractère nommée MONU dont les valeurs sont les p=11 catégories et q=6 variables numériques qui ont pour nom les q=6 types de propriétaires.

Détaillons les p=11 catégories de monuments :

preh (antiquités préhistoriques),

hist (antiquités historiques),

chat (châteaux),

mili (architecture militaire),

cath (cathédrales),

egli (églises),

chap (chapelles),

mona (monastères),

ecpu (édifices civils publics),

ecpr (édifices civils privés)

div (divers).

et les q=6 types de propriétaires :

COMM (commune),

PRIV (privé),

ETAT (état),

DEPA (département),

ETPU (établissement public)

NDET (non déterminé).

L'allure de la table MONUMENT (placée dans le répertoire c:\k8\aboul) est la suivante:

MONU	COMM	PRIV	NDET
preh	244	790	144
hist	246	166	31
.....
ecpr	224	909	4
div	967	242	9

DEBUT DE LA SESSION KORRESP

Charger et soumettre dans SAS : DKORR.SAS qui se trouve dans c:\k8.

Toute la suite est CONVERSATIONNELLE, des fenêtres s'ouvrent et guident l'utilisateur dans ses choix. Pour toutes les fenêtres, les choix retenus seront validés par F3.

Fenêtre n°1 : "CHEMIN".

L'utilisateur précise le chemin d'accès aux données (par défaut c:\k8\aboul).

Exemple : path: C:\K8\ABOUL suivi de F3.

Dans le cas où le répertoire proposé est vide ou inexistant, un message d'erreur apparaît dans la fenêtre LOG, la commande WINDOW permet le passage de la fenêtre CHEMIN à celle du LOG.

Fenêtre n°2 : "CHOIXDAT".

A partir des tables sas situées dans le répertoire choisi dans "CHEMIN", l'utilisateur sélectionne (par X) la table retenue pour l'analyse.

Exemple : X MONUMENT suivi de F3.

Fenêtre n°3 : "SELECT".

L'utilisateur doit sélectionner les éléments du tableau de contingence .

C : Pour l'identificateur des lignes du tableau de contingence.

N : Pour les colonnes (variables numériques) du tableau de contingence.

Exemple :

CHOIX	VARIABLE	
C	MONU	
N	COMM	
N	PRIV	
N	ETAT	
N	DEPA	
N	ETPU	
X	NDET	suivi de F3.

Remarque: En cas de données manquantes un message apparaît dans l'OUTPUT indiquant leur nombre et le programme s'arrête.

En absence de donnée manquante , le programme calcule et affiche dans la fenêtre "OUTPUT" la statistique du khi_deux et sa "p-value". Si la p-value est > 5% alors un message d'avertissement apparaît dans la fenêtre "OUTPUT" qui explique en quoi l'analyse du tableau de contingence n'est pas très pertinente mais le programme KORRESP continue.

Fenêtre n°4 : "STRUCTUR".

L'utilisateur choisit son type de réseau : F pour Ficelle, G pour Grille.

Exemple : choix : G suivi de F3.

Fenêtre n°5 : "PARAM_FI" ou "PARAM_GR".

Pour une FIcelle, "PARAM_FI" demande sa taille (n) et le nombre maximum d'itérations (tmax).

Pour une GRille, "PARAM_GR" demande le nombre de lignes (m), le nombre de colonnes (n) et le nombre maximum d'itérations (tmax).

Les nu unités du réseau (nu = n pour une ficelle ou nu = m×n pour une grille) seront numérotées de 1 à nu. Par exemple si nu=12 pour la ficelle de n=12, on aura : 1 2 11 12.

Pour la grille de $m = 3$ et $n = 4$ on aura :

1 4 7 10
2 5 8 11
3 6 9 12

Dans "PARAM_FI" comme dans "PARAM_GR" deux choix supplémentaires sont offerts.
Le premier choix (O/N) permet d'initialiser le générateur de nombre au hasard.

O : le point de départ est fixe, ce qui rend les résultats reproductibles (valeur par défaut).

N: le point de départ est calé sur l'horloge de la machine, ce qui rend les résultats NON reproductibles.

Le deuxième choix (O/N) propose le calcul et la représentation éventuels de la fonction "énergie". Cette "énergie" ou "potentiel" généralise la notion de "variance intra" en l'étendant aux plus proches voisins (à zéro voisin les deux notions coïncident). En fin d'itérations, elle doit se situer sur un minimum (local).

Le choix négatif (N) est la valeur par défaut; si le nombre d'itérations est élevé (dès 500), un choix affirmatif (O) peut s'avérer coûteux en temps de calcul et implique la création de 2 graphiques (nommés NRJ et NRJ1) stockés dans un catalogue graphique permanent (cf. la fenêtre n°6 : "INFORM").

Exemple :

$m = 5$
 $n = 5$
 $tmax = 300$
CHOIX : O
CHOIX : O suivi de F3.

Si le tableau de contingence a $p \times q$ modalités, on peut prendre pour :

- une ficelle, n au moins égal à $\max(p,q)$
- une grille, $m = n$, où n est tel que $n \times n$ est immédiatement supérieur à $2 \max(p,q)$.

Ici $p=11$ et $q=6$, on prend une grille 5×5 puisque 25 est le carré immédiatement supérieur à 22.

Fenêtre n°6 : "INFORM".

Dans cette fenêtre, l'utilisateur doit renseigner le champ Nom (fixé par défaut à `_ TEMPOR _`) et décider (choix O/N) de l'éventuel visualisation et stockage de deux graphiques de type "dx-dy" (cf. : la thèse de P. DEMARTINES dont la référence est donnée dans la section III).

Le champ Nom sera utilisé pour nommer des catalogues (fichier d'extension `.sas7bcat`) et des tables sas (fichier d'extension `.sas7bdat`). Il doit comporter au moins 5 caractères; dans le cas contraire, le programme complète la réponse par des X pour obtenir un champ Nom de 5 caractères.

Le choix : N est proposé par défaut; en cas de réponse positive (O), deux graphiques de noms `DX_DY` et `DX_DY1` seront créés et stockés dans le catalogue graphique. Un graphique de type "dx-dy" est un nuage de points dans un carré de côté 1, où l'on compare les distances (normalisées à 1) théoriques entre les unités gagnantes avec les distances (normalisées à 1) euclidiennes entre les vecteurs poids associés. La situation idéale correspondant à une organisation parfaite serait celle où tous les points (représentés sur le graphique par des cercles) sont situés sur la diagonale du carré.

Exemple :

Nom : MO3C5G
choix : O suivi de F3.

Ces réponses impliquent la création de 3 tables sas et de 2 catalogues qui seront placés dans C:\K8\ABOUL (cf.: le champ path de la fenêtre n°1 : "CHEMIN").

Description des 3 tables sas :

1. La table MO3C5_CL (valeur par défaut : _TEMP_CL) donne, pour chaque unité gagnante, les modalités du croisement dont elle est la référence (contenu de la classe).
2. La table MO3C5_WS (valeur par défaut : _TEMP_WS) donne pour chaque unité (gagnante ou pas) sa position dans la grille : ligne, colonne, son effectif (zéro si l'unité n'est pas gagnante) et son "vecteur poids final concaténé " ou "vecteur représentant ".
3. La table MO3C5G (valeur par défaut : _TEMPOR_) donne pour chaque modalité du croisement son unité gagnante (sa classe) :_codage_ , sa position dans la grille : ligne, colonne. Dans le prolongement des modalités du croisement (ici 11+6=17 valeurs) se place la table de BURT.

Description des 2 catalogues :

1. Le catalogue iml MO3C5G (valeur par défaut : _TEMPOR_) contient tous les intermédiaires de calcul (matrices de poids initiaux et finaux, liste des unités gagnantes , etc.) qui seront ultérieurement utilisés dans les traitements complémentaires (cf. le paragraphe VIII).
2. Le catalogue graphique GKORMO3C (valeur par défaut : GKOR_TEM) contient 10 graphiques (au plus). Chaque graphique possède son nom et sa description qui correspond aux 40 premiers caractères de son titre.

Les 10 graphiques du catalogue graphique GKORMO3C :

Nom	Description
DX_DY	Dx_dy pour les colonnes
DX_DYI	Dx_dy pour les lignes
G_DIST	Distances (E) avec les plus proches voisins
G_PAVAGE	KORRESP : grille 5x5 et 300 itérations
HIS_GRI	Représentants (Poids Finaux) pour les colonnes
HIS_GRII	Représentants (Poids Finaux) pour les lignes
LIN_GRI	Représentants (Poids Finaux) pour les colonnes
LIN_GRII	Représentants (Poids Finaux) pour les lignes
NRJ	Potentiel pour les colonnes
NRJI	Potentiel pour les lignes

Les graphiques, dans le catalogue, sont classés en ordre alphabétique alors qu'à l'affichage l'ordre sera

pour une grille: G_PAVAGE, G_DIST, HIS_GRI, LIN_GRI, HIS_GRII, LIN_GRII, DX_DY, DX_DYI, NRJ et NRJI

pour une ficelle: F_PAVAGE, F_DIST, HIS_FIC, LIN_FIC, HIS_FICI, LIN_FICI, DX_DY, DX_DYI, NRJ et NRJI

G_DIST pour une Grille (ou F_DIST pour une Ficelle) permet d'apprécier pour chaque unité (classe) son effectif et les distances euclidiennes normalisées avec ses plus proches voisins (8 pour une grille et 2 pour une ficelle). A l'exception des unités des bords, si l'unité est très proche de ses (8 ou 2) voisins, son polygone sera très proche des bords du carré. (cf. : A Kohonen map representation to avoid misleading interpretation de M. COTTRELL et E. DE BODT, SAMOS n°62).

G_PAVAGE pour une Grille (ou F_PAVAGE pour une Ficelle) donne l'illustration graphique du contenu des classes (unités) du réseau (cf. : la table MO3C5_CL).

HIS_GRI pour Grille (ou HIS_FIC pour une Ficelle) est un pavage d'histogrammes des vecteurs "poids finaux" ou "représentants" associés aux colonnes du tableau de contingence.

HIS_GRII pour Grille (ou HIS_FIC1 pour une Ficelle) est un pavage d'histogrammes des vecteurs "poids finaux" ou "représentants" associés aux lignes du tableau de contingence.

LIN_GRI pour Grille (ou LIN_FIC pour une Ficelle) est un pavage de courbes des vecteurs "poids finaux" ou "représentants" associés aux colonnes du tableau de contingence.

LIN_GRII pour Grille (ou LIN_FIC1 pour une Ficelle) est un pavage de courbes des vecteurs "poids finaux" ou "représentants" associés aux lignes du tableau de contingence.

Il n'y a pas de graphique de pavage pour une ficelle de plus de 50 unités.

PATIENCE KORRESP travaille pour vous !!!!

Avant d'afficher les graphiques du catalogue graphique IML, les 4 autres fichiers permanents sont créés.

LA SESSION KORRESP EST TERMINEE

VII - LES ALGORITHMES KACM, KACM1 ET KACM2

Algorithmes KACM(j) (j=1,2) : Version Kohonen de l'Analyse des Correspondances Multiples

Il s'agit d'une analyse des relations entre plusieurs variables qualitatives. cf.: la référence M.Cottrell et S.Ibbou, dans Proc.ESANN'95.

Pour illustrer le déroulement d'une session, nous utiliserons les données de BREFORT,1982 (cité dans l'ouvrage de G. SAPORTA "Probabilités, analyse des données et statistique").

On dispose pour 27 races de chiens de 7 variables qualitatives :

- 1. TAILLE dont les 3 modalités sont : taille1, taille2 et taille3*
- 2. POIDS dont les 3 modalités sont : poids1, poids2 et poids3*
- 3. VELOCITE dont les 3 modalités sont : veloce1, veloce2 et veloce3*
- 4. INTELLIG(ence) dont les 3 modalités sont : malin1, malin2 et malin3*
- 5. AFFECTIO(n) dont les 2 modalités sont : affec1 et affec2*
- 6. AGRESSION(ité) dont les 2 modalités sont : agres1 et agres2*
- 7. FONCTION dont les 3 modalités sont : chasse, compagnie et garde (soit, au total, 19 modalités différentes).*

Ces données sont dans la table sas CHIENS qui est placée dans le répertoire 'c:\k8\aboul, et dont l'allure est la suivante:

RACE	TAILLE	POIDS	VELOCITE	INTELLIG	AFFECTIO	AGRESSION	FONCTION
BEAUCERON	taille3	poids2	veloce3	malin3	affec2	agres2	garde
.....							

DEBUT DE LA SESSION KACM(j)

Charger et soumettre dans SAS : DKACM.SAS qui se situe dans c:\k8.

Toute la suite est CONVERSATIONNELLE, des fenêtres s'ouvrent et guident l'utilisateur dans ses choix. Pour toutes les fenêtres, les choix retenus seront validés par F3.

Fenêtre n°1 : "CHEMIN".

L'utilisateur précise le chemin d'accès aux données (par défaut c:\k8\aboul) .

Exemple : path: C:\K8\ABOUL suivi de F3.

Dans le cas où le répertoire proposé est vide ou inexistant, un message d'erreur apparaît dans la fenêtre LOG, la commande WINDOW permet le passage de la fenêtre CHEMIN à celle du LOG.

Fenêtre n°2 : "CHOIXDAT".

A partir des tables sas situées dans le répertoire choisi dans "CHEMIN", l'utilisateur sélectionne (par X) la table retenue pour l'analyse.

Exemple : X CHIENS suivi de F3.

Fenêtre n°3 : "CHOIX".

En premier lieu, l'utilisateur précise la nature des données :

- Dans le cas où la table sas correspond à un tableau de Contingence. par : C
(cette option permet de comparer les résultats obtenus avec KORRESP)
- Dans le cas où la table sas correspond aux Réponses d'individus. par : R
- Dans le cas où la table sas correspond à un tableau Disjonctif complet. par : D
- Dans le cas où la table sas correspond à un tableau de Burt. par : B

Ensuite, l'utilisateur précise le type des observations :

- par : O si elles sont "anonymes" (valeur par défaut),
seules les modalités seront prises en compte : **KACM**
- par : N si les observations sont "NON anonymes",
leurs identifiants seront traitées avec les modalités : **KACM1** ou **KACM2**

Dans les cas C (tableau de contingence) ou B (tableau de Burt), la question ne se posant pas, l'utilisateur est contraint de garder la réponse par défaut (O).

Enfin, l'utilisateur identifie la nature des variables :

- par : K
- par : X

Le choix de K est unique et il identifie une variable considérée comme étant de type caractère.

Dans les cas R (tableau de réponses) ou D (disjonctif complet), K correspond à l'identificateur des observations (individus, répondants).

Dans les cas C (contingence) ou B (Burt), K correspond à la variable qui identifie les lignes ou les modalités du tableau.

Plusieurs choix sont possibles pour X :

Dans les cas C (contingence), D (disjonctif) ou B (Burt), l'utilisateur est contraint de prendre toutes les variables restantes dans la liste.

Dans le cas R (tableau de réponses), le choix de 2 variables est l'option minimum.

Exemple :

CHOIX	R	
CHOIX	N	
CHOIX	VARIABLE	
K	RACE	
X	TAILLE	
X	POIDS	
X	VELOCITE	
X	INTELLIG	
X	AFFECTIO	
X	AGRESSIV	
X	FONCTION	suivi de F3.

Si les observations sont "NON anonymes", la Fenêtre n°4 : "KELKACMJ". propose à l'utilisateur de choisir entre KACM1 (Choix : A) et KACM2 (Choix : B)

Pour KACM1 l'apprentissage et le classement initial s'effectuent sur les individus à l'aide du tableau disjonctif complet "corrigé" puis les modalités sont présentées comme données supplémentaires à l'aide de la table de Burt "corrigée".

Pour KACM2 l'apprentissage et le classement initial s'effectuent sur les modalités à l'aide de la table de Burt "corrigée" puis les individus sont présentés comme données supplémentaires à l'aide du tableau disjonctif "corrigé".

A l'instar de l'ACM les modalités sont classées de façon identiques pour KACM et KACM2.

Exemple : choix : A suivi de F3.

Remarque: En cas de données manquantes un message apparaît dans l'OUTPUT indiquant leur nombre et le programme s'arrête.

Comme pour KORRESP, si KACM traite un tableau de contingence (C), le programme calcule et affiche dans la fenêtre "OUTPUT" la statistique du khi2 et sa "p-value". Si la p-value est > 5% alors un message d'avertissement apparaît dans la fenêtre "OUTPUT" qui explique en quoi l'analyse du tableau de contingence n'est pas très pertinente mais le programme KACM continue.

Fenêtre n°5 : "STRUCTUR".

L'utilisateur choisit son type de réseau : F pour Ficelle, G pour Grille.

Exemple : choix : G suivi de F3.

Fenêtre n°6 : "PARAM_FI" ou "PARAM_GR".

Pour une Ficelle, "PARAM_FI" demande sa taille (n) et le nombre maximum d'itérations (tmax).

Pour une Grille "PARAM_GR" demande le nombre de lignes (m), le nombre de colonnes (n) et le nombre maximum d'itérations (tmax).

Les nu unités du réseau (nu = n pour une ficelle ou nu = m×n pour une grille) seront numérotées de 1 à nu. Par exemple si nu=12 pour la ficelle de n=12, on aura : 1 2 11 12.

Pour la grille de m = 3 et n = 4 on aura :

```
1 4 7 10
2 5 8 11
3 6 9 12
```

Dans "PARAM_FI" comme dans "PARAM_GR" deux choix supplémentaires sont offerts.

Le premier choix (O/N) permet d'initialiser le générateur de nombre au hasard.

O : le point de départ est fixe, ce qui rend les résultats reproductibles (valeur par défaut).

N: le point de départ est calé sur l'horloge de la machine, ce qui rend les résultats NON reproductibles.

Le deuxième choix (O/N) propose le calcul et la représentation éventuels de la fonction "énergie". Cette "énergie" ou "potentiel" généralise la notion de "variance intra" en l'étendant aux plus proches voisins (à zéro voisin les deux notions coïncident). En fin d'itérations, elle doit se situer sur un minimum (local).

Le choix négatif (N) est la valeur par défaut, si le nombre d'itérations est élevé (dès 500), un choix affirmatif (O) peut s'avérer coûteux en temps de calcul et implique la création d'un graphique de nom : NRJ stocké dans un catalogue graphique permanent (cf. la fenêtre n°6 : "INFORM").

Exemple : m = 5
 n = 5
 tmax = 300
 CHOIX : O
 CHOIX : O suivi de F3.

Si q est le nombre total de modalités et r le nombre d'individus, on peut prendre pour une ficelle, n de l'ordre de $q/2$ et pour une grille où $m = n$, n tel que $n \times n$ est immédiatement supérieur à q .
Pour t_{max} une valeur d'au moins 10 fois le nombre d'individus

Ici $q=19$, $r=27$ donc on prend une grille 5×5 , puisque 25 est le carré immédiatement supérieur à 19 et $t_{max} = 300 > 10 \times 27$.

Fenêtre n°7 : "INFORM".

Dans cette fenêtre, l'utilisateur doit renseigner le champ Nom (fixé par défaut à `_TEMPOR_`) et décider (choix O ou N) de l'éventuelle visualisation et stockage d'un graphique de type "dx-dy" (cf. : la thèse de P. DEMARTINES dont la référence est donnée dans la section III).

Le champ Nom sera utilisé pour nommer des catalogues (fichier d'extension `.sas7bcat`) et des tables sas (fichier d'extension `.sas7bdat`). Il doit comporter au moins 5 caractères; dans le cas contraire le programme complète la réponse par des X pour obtenir un champ Nom de 5 caractères.

Le choix : N est proposé par défaut; en cas de réponse positive (O), un graphique de nom `DX_DY` sera créé dans le catalogue graphique. C'est un nuage de points dans un carré de côté 1, où l'on compare les distances (normalisées à 1) théoriques entre les unités gagnantes avec les distances (normalisées à 1) euclidiennes entre les vecteurs poids associés. La situation idéale correspondant à une organisation parfaite serait celle où tous les points (représentés sur le graphique par des cercles) sont situés sur la diagonale du carré.

Exemple : Nom : D3C5G
 choix : O suivi de F3.

Ces réponses impliquent la création de 3 tables sas et de 2 catalogues qui seront placés dans `C:\K8\ABOUL` (cf.: le champ path de la fenêtre n°1 : "CHEMIN").

Description des 3 tables sas :

1. La table `D3C5G_CL` (valeur par défaut : `_TEMP_CL`) donne, pour chaque unité gagnante, les modalités (caractéristiques des chiens) ainsi que les identifiants des individus (race des chiens) dont elle est la référence, dans le cas d'un tableau de Réponses ou d'un tableau Disjonctif portant sur des individus Non anonymes (CHOIX R ou D et N dans la fenêtre n°3 : "CHOIX").
2. La table `D3C5G_WS` (valeur par défaut : `_TEMP_WS`) donne pour chaque unité (gagnante ou pas) sa position dans la grille : ligne, colonne, son effectif (zéro si l'unité n'est pas gagnante) et son "vecteur poids final" ou "vecteur représentant".
3. La table `D3C5G` (valeur par défaut : `_TEMPOR_`) donne pour chaque modalité son unité gagnante (sa classe) : `_codage_`, sa position dans la grille : ligne, colonne. Dans le prolongement des modalités (ici 19 valeurs) se place la table de BURT. Dans le cas d'un tableau de Réponses ou d'un tableau Disjonctif portant sur des individus Non anonymes, la table donne en plus, pour chaque modalité identifiant un individu son unité gagnante (sa classe) : `_codage_`, sa position dans la grille : ligne, colonne. Dans le prolongement des identifiants des individus se place le tableau disjonctif complet.

Description des 2 catalogues :

1. Le catalogue iml `D3C5G` (valeur par défaut : `_TEMPOR_`) contient tous les intermédiaires de calcul (matrices de poids initiaux et finaux, liste des unités gagnantes, etc.) qui seront ultérieurement utilisés dans les traitements complémentaires (cf. le paragraphe VIII).

2. Le catalogue graphique GACMD3C5 (valeur par défaut : GACM_TEM) contient 7 graphiques (au plus). Chaque graphique possède son nom et sa description qui correspond aux 40 premiers caractères de son titre.

Les 7 graphiques du catalogue graphique GACMD3C5 :

Nom	Description
CELLW	Valeur moyenne et représentant des 25 cell
DX_DY	dx_dy
G_DIMA	Distances (M) avec les plus proches voisins
G_PAVAG1	KACM1 : grille 5x5 et 300 itérations
G_PAVAGE	KACM1 : grille 5x5 et 300 itérations
LIN_GRI	Représentants des classes (Poids Finaux)
NRJ	Variance intra etendue aux voisins

Les graphiques, dans le catalogue, sont classés en ordre alphabétique alors qu'à l'affichage l'ordre sera

pour une grille : G_PAVAGE, G_PAVAG1, CELLW, G_DIMA, LIN_GRI, DX_DY et NRJ
 pour une ficelle : F_PAVAGE, F_PAVAG1, CELLW, F_DIMA, LIN_FIC, DX_DY et NRJ

Si le nombre de modalités > 50, il n'y a pas de graphique LIN_FIC ou LIN_GRI.

Le graphique CELLW ne figure que dans les sorties de KACM1, il permet de comparer l'allure du vecteur représentant de l'unité gagnante (classe non vide) de la structure, avec le vecteur moyen des entrées. (l'individu moyen de cette classe non vide).

G_DIMA pour une Grille (ou F_DIMA pour une Ficelle) permet d'apprécier pour chaque unité (classe) son effectif et les distances de **Mahalanobis** normalisées avec ses plus proches voisins (8 pour une grille et 2 pour une ficelle). A l'exception des unités des bords, si l'unité est très proche de ses (8 ou 2) voisins, son polygone sera très proche des bords du carré. (d'après : A Kohonen map representation to avoid misleading interpretation de M. COTTRELL et E. DE BODT, Proc.ESANN' 96).

G_PAVAGE pour une Grille (ou F_PAVAGE pour une Ficelle) donne l'illustration graphique du contenu des classes (les modalités pour KACM et les individus en plus pour KACM1 et KACM2) du réseau (cf.: D3C5G_CL).

G_PAVAG1 pour une Grille (ou F_PAVAG1 pour une Ficelle) donne l'illustration graphique des seules modalités pour KACM1 et KACM2.

LIN_GRI pour une Grille (LIN_FIC pour une Ficelle) est un pavage de courbes des représentants des classes (poids finaux).

Il n'y a pas de graphique de pavage pour une ficelle de plus de 50 unités.

PATIENCE KACM travaille pour vous !!!!

Avant d'afficher les graphiques du catalogue graphique IML, les 4 autres fichiers permanents sont créés.

LA SESSION KACM EST TERMINEE

VIII - TRAITEMENTS COMPLEMENTAIRES

A) GENERALITES CONCERNANT LES MACROS

Les 19 macros (MKLUG.SAS, MKLUF1.SAS, MKLUF2.SAS, MSTAT.SAS, MCONT.SAS, MCELL.SAS, MSCAL.SAS, MSCOL.SAS, MSUPP.SAS, MGVAR.SAS, MVC3D.SAS, MKAMG.SAS, MKAMF1.SAS, MKAMF2.SAS, MDIMA.SAS, MDIMA12.SAS, MDIST.SAS, MFMOD.SAS, MFIND.SAS) sont des traitements qui peuvent être envisagés à la suite des programmes KACP, KFAST, KBATCH, KORRESP ou KACM(j=1,2).

La soumission de KOMAC.SAS, situé dans C:\K8, permet de les compiler et de les stocker dans le catalogue permanent SASMACR.sas7bcat placé dans C:\K8\AKOR.

L'appel d'une macro s'obtient en soumettant l'instruction **%nom_de_la_macro;** dans la fenêtre PROGRAM EDITOR (exemple : %MKLUG;).

Chaque macro est conversationnelle et s'ouvre sur une fenêtre où est indiqué son domaine d'application : COMPLEMENT(S) TOUT PROGRAMME signifie qu'elle est applicable aux sorties des 7 programmes KACP, KFAST, KBATCH, KORRESP ou KACM(j=1,2).

C'est par le biais du champ " type du catalogue " que l'utilisateur précise le type de programme (gacp pour KACP, KFAST, et KBATCH gkor pour KORRESP ou gacm pour KACM(j=1,2)).

Le vocable : TOUTE STRUCTURE signifie que le traitement peut s'appliquer aussi bien à une ficelle qu'à une grille.

Par exemple pour la macro MKLUG, il est indiqué dans la fenêtre " SURCLASS " :
COMPLEMENTS TOUT PROGRAMME POUR UNE GRILLE

Suite au domaine d'application, il est donné une brève description de la nature du traitement.

Par exemple pour la macro MKLUG il est indiqué :

=>>> PAVAGE EN SUPER CLASSES (CLUSTERS)

Deux champs sont communs à toutes ces macros : " librairie " et " nom commun ".

La valeur par défaut de " librairie " est C:\K8\ABOUL, ce champ est équivalent au champ "path " de la fenêtre " CHEMIN "; il correspond au chemin d'accès des données et de certains catalogues (comme le catalogue graphique).

La valeur par défaut de " nom commun " est _TEMPOR_; ce champ est équivalent au champ " Nom " de la fenêtre " INFORM "; il correspond au nom utilisé pour nommer des tables sas et le catalogue graphique.

Pour abandonner l'exécution de la macro, il suffit de taper O dans le champ " STOP ".

En bas de la fenêtre de chaque macro il est stipulé que la touche F3 valide les choix et qu'en cas de problèmes de cohérence dans la saisie des réponses, l'utilisateur reste dans la fenêtre de la macro avec les réponses par défaut, il doit consulter les messages d'erreurs de la fenêtre LOG avant de reprendre la saisie dans la fenêtre de la macro.

B) DESCRIPTION ET EXEMPLE D'APPLICATION POUR CHAQUE MACRO

MKLUG : (fenêtre " SURCLASS ", domaine : *tout programme pour une grille*).

Création de Super-Classes via une Classification Ascendante Hiérarchique (cf.: PROC CLUSTER avec la méthode de Ward) appliquée aux vecteurs poids finaux des unités de la grille.

Cette macro produit au plus 5 graphiques rajoutés au catalogue graphique déjà existant (cf.: le champ "nom commun ") et une table sas placée dans le répertoire des données (cf.: le champ "librairie").

Description des 5 graphiques :

Le graphique nommé G_DENDRO qui représente le dendrogramme associé à la classification hiérarchique ascendante, est précédé de l'historique de la classification qui apparaît dans la fenêtre OUTPUT et permet de confirmer a posteriori le choix du nombre de super-classes.

Le graphique nommé CLUSPAV représente la grille où dans chaque unité figure le vecteur poids final (le représentant) avec le numéro de son cluster; si le nombre de modalités (KORRESP, KACM(j=1,2)) ou si le nombre de variables (KACP, KFAST, KBATCH) est strictement supérieur à 50, alors seuls les numéros de cluster figurent dans le graphique nommé CLUSPAZ.

Les trois autres graphiques nommés : CLUSCOUL, CLUSCOLT, CLUSDIMA ou CLUSDIST ne seront produits que dans la mesure où l'utilisateur a choisi (cf.: le champ " nombre de super-classes " de la fenêtre " SURCLASS ") un nombre de clusters qui ne dépasse pas 22 qui est le nombre de couleurs disponibles (cf. : le programme KOLCLUS.SAS dans C:\K8).

Dans le graphique nommé CLUSCOUL, on colorie CLUSPAV (on remplace des numéros de cluster par des couleurs); si le nombre de modalités (KORRESP, KACM(j=1,2)) ou si le nombre de variables (KACP, KFAST, KBATCH) est strictement supérieur à 50, alors seules les couleurs de cluster figurent dans le graphique nommé CLUSCOZ.

Dans le graphique nommé CLUSCOLT, figure les "valeurs" associées aux unités de la grille ainsi que les couleurs des clusters.

Par "valeurs" on entend :

- *les individus pour KACP, KFAST ou KBATCH*
- *les modalités pour KORRESP ou KACM*
- *les modalités et les individus pour KACMj (j=1,2)*

Dans le graphique nommé CLUSDIMA, on colorie le graphique G_DIMA avec les couleurs associées aux clusters (super-classes).

Pour KORRESP correspond le graphique nommé CLUSDIST; dans ce graphique on colorie le graphique G_DIST avec les couleurs associées aux clusters (super-classes).

La table sas nommée CLUSTn (où n correspond au nombre de clusters retenu par l'utilisateur) donne pour chaque unité de la grille son numéro de cluster ainsi que son vecteur poids final. Cette table sas sera utilisée par les macros MSTAT, MCONT, MSCOL, MVC3D, MKAMG, MFINN et MFMOD qui effectuent des représentations graphiques ou des calculs impliquant les clusters.

Exemple :

<i>librairie :</i>	<i>C:\K8\ABOUL</i>	
<i>nom commun :</i>	<i>B8G3CPAY</i>	
<i>type du catalogue :</i>	<i>GACP</i>	
<i>nombre de super-classes :</i>	<i>5</i>	<i>suivi de F3</i>

Ces réponses impliquent la création du fichier C:\K8\ABOUL\CLUST5.sas7bdat, et portent à 13 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MKLUF1 : (*fenêtre “ SURFIC10 ”, domaine : tout programme pour une ficelle (≤ 10)*).

MKLUF2 : (*fenêtre “ SURFICEL ”, domaine : tout programme pour une ficelle (> 10 et ≤ 50)*).

*Ces deux macros effectuent sur une **ficelle** un traitement qui est équivalent à celui de MKLUG pour une grille.*

Les graphiques :

G_DENDRO, CLUSPAV, CLUSCOUL, CLUSTCOLT, CLUSDIMA, CLUSDIST, CLUSPAZ et CLUSCOZ deviennent respectivement :

Pour MKLUF1:

F_DENDRO, CLUFPA10, CLUFCO10, CLUFCT10, CLUFDM10, CLUFDE10, CLUFAZ10 et CLUFOZ10

Pour MKLUF2:

F_DENDRO, CLUFPAV, CLUFCOUL, CLUFCTXT, CLUFDMAH, CLUFDEUC, CLUFPAZ et CLUFCOZ

La table sas nommé CLUSFn (l'équivalent de la table CLUSTn) donne pour chaque unité de la ficelle son numéro de cluster ainsi que son vecteur poids final. Cette table sera utilisée par les macros MSTAT, MCONT, MSCOL, MVC3D, MKAMFj ($j=1,2$), MFINd et MFMOD qui effectuent des représentations graphiques ou des calculs impliquant les clusters.

MSTAT : (*fenêtre “ STATCLUS ”, domaine : **KACP, KFAST, KBATCH** et toute structure*).

Cette macro doit faire suite à MKLUG ou MKLUFj ($j=1,2$), elle considère les super-classes (la variable cluster) comme une variable de classification et affiche dans la fenêtre OUTPUT les moyennes conditionnelles des variables brutes (cf. la sélection de la fenêtre “ CHOIXVAR ” dans KACP). Pour chaque variable figure dans L'OUTPUT la décomposition de la variance, la statistique de Fisher et la P-value associée. Enfin, y apparaissent les tests multidimensionnels de Wilks et Hotelling avec leur valeur approchée en terme de khi-deux ainsi que leur p-value.

La macro MSTAT produit :

Dans le cas d'une grille la table sas nommée CLASGn (où n correspond au nombre de clusters retenu dans MKLUG).

Dans le cas d'une ficelle la table sas nommée CLASFn (où n correspond au nombre de clusters retenu dans MKLUFj ($j=1,2$)).

La table CLASGn ou CLASFn est placée dans le répertoire des données (cf.: le champ “librairie”). Elle indique pour chaque observation (ici chaque PAYS) ses numéros de cluster et de classe (unité gagnante) ainsi que ses valeurs pour les variables brutes.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	
nombre de clusters :	5	suivi de F3

Outre les sorties dans la fenêtre OUTPUT, ces réponses impliquent la création du fichier C:\K8\ABOUL\CLASG5.sas7bdat.

MCONT : (fenêtre “ CONTENT ”, domaine : **tout programme et toute structure**).

Cette macro doit faire suite à MKLUG ou MKLUFj (j=1,2); elle ajoute un graphique nommé CONCLU au catalogue graphique déjà existant.

Ce graphique représente sous forme de pavage les vecteurs poids finaux des unités du réseau (grille ou ficelle) associés à un même cluster (le CONTenu des CLUsters). Il permet d'évaluer l'homogénéité des clusters (super-classes).

Si le nombre de clusters, choisi par l'utilisateur, ne dépasse pas le nombre de couleurs disponibles, le graphique CONCLU est colorié (une même couleur par contenu de cluster).

Il faut noter que, si le nombre de modalités (KORRESP, KACM(j=1,2)) ou si le nombre de variables (KACP, KFAST, KBATCH) est strictement supérieur à 50, alors la macro MCONT ne produit aucun graphique.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	
nombre de clusters :	5	suivi de F3

Ces réponses portent à 14 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MCELL : (fenêtre “ CELLS ”, domaine : **KACP, KFAST, KBATCH et toute structure**).

Cette macro ajoute deux graphiques nommés CELL et CELLW au catalogue graphique déjà existant.

Le graphique CELL présente sous forme de pavage le contenu des unités (les valeurs numériques des variables, éventuellement transformées selon le choix de la fenêtre “ PREPROC ” et qui correspondent aux individus associés à une même unité gagnante) du réseau (grille ou ficelle).

Le graphique CELLW présente sous forme de pavage le vecteur poids final (représentant) de l'unité du réseau ainsi que la moyenne des variables dans le cas où l'unité est gagnante (classe non vide). Il permet de comparer l'allure du vecteur représentant d'une unité gagnante avec son vecteur moyen. (l'individu moyen de cette classe non vide).

Exemple :

<i>librairie :</i>	<i>C:\K8\ABOUL</i>	
<i>nom commun :</i>	<i>B8G3CPAY</i>	<i>suivi de F3</i>

Ces réponses portent à 16 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MSCAL : (*fenêtre “ ESCALING ”, domaine : tout programme et toute structure*).

Utilisation de la technique Multi Dimensional Scaling (cf.: PROC MDS qui permet une représentation graphique des unités gagnantes du réseau (grille ou ficelle) à partir des distances euclidiennes entre les différents représentants (poids finaux) de ces unités gagnantes.

Cette macro produit un graphique nommé ESCALING qui est rajouté au catalogue graphique déjà existant et une table sas nommée par défaut _TEMP_EG.sas7bdat qui est placée dans le répertoire des données. Cette table sas fournit les distances euclidiennes entre les vecteurs poids finaux des unités gagnantes du réseau.

Exemple :

<i>librairie :</i>	<i>C:\K8\ABOUL</i>	
<i>nom commun :</i>	<i>B8G3CPAY</i>	
<i>type du catalogue :</i>	<i>GACP</i>	<i>suivi de F3</i>

Ces réponses impliquent la création du fichier C:\K8\ABOUL\B8G3C_EG.sas7bdat, et portent à 17 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MSCOL : (*fenêtre “ ESCALCLU ”, domaine : tout programme et toute structure*).

Cette macro doit faire suite à MKLUG ou MKLUFj (j=1,2); elle effectue le même traitement que la macro MSCAL en ajoutant à chaque unité gagnante la couleur associée à son cluster.

Cette macro produit un graphique nommé ESCALCLU (version couleur de ESCALING) qui est rajouté au catalogue graphique déjà existant et une table sas nommée par défaut _TEMP_EG.sas7bdat qui est placée dans le répertoire des données. Cette table sas fournit les distances euclidiennes entre les vecteurs poids finaux des unités gagnantes du réseau.

Si le nombre de clusters dépasse 22 il faut utiliser la macro MSCAL.

Exemple :

<i>librairie :</i>	<i>C:\K8\ABOUL</i>	
<i>nom commun :</i>	<i>B8G3CPAY</i>	
<i>type du catalogue :</i>	<i>GACP</i>	
<i>nombre de clusters :</i>	<i>5</i>	<i>suivi de F3</i>

Ces réponses portent à 18 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MSUPP : (fenêtre “ EN_PLUS ”, domaine : **KACP, KFAST, KBATCH et toute structure**).

Cette macro effectue le traitement des données (observations) supplémentaires (avec éventuellement des valeurs manquantes).

Elle produit un graphique nommé **G_PAVAG1** ou **F_PAVAG1** (selon que le réseau est une grille ou une ficelle) qui est rajouté au catalogue graphique déjà existant.

Le graphique **G_PAVAG1** pour une grille (ou **F_PAVAG1** pour une ficelle) représente la position des données (individus) supplémentaires sur le réseau.

Elle produit aussi une table sas dont le nom (de 8 lettres) est construit à partir des 5 premières lettres du nom de la table contenant les données supplémentaires suivies de “_SG” ou “_SF”(selon que le réseau est une grille ou une ficelle); cette table est placée dans le répertoire des données.

Cette table donne pour chaque modalité de la variable qui identifie les données (observations) supplémentaires son numéro d'unité gagnante (sa classe), sa position dans le réseau et ses valeurs (éventuellement manquantes) pour les variables éventuellement transformées par le traitement initial.

Elle produit enfin un catalogue iml des résultats intermédiaires de même nom que celui de la table sas précédemment créée.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	suivi de F3

Suite à cette soumission, s'ouvre la fenêtre “ WHAT_NEW ” qui permet de sélectionner la table qui contient les données supplémentaires.

Exemple :	CHOIX	Nom de la Table	
	X	BLAYO3	suivi de F3.

Puis s'ouvre la fenêtre “ QUEL_VAR ” qui permet la sélection par C de la variable qui identifie les observations supplémentaires et par N des variables numériques déjà retenues dans l'analyse.

Exemple :	CHOIX	VARIABLE	
	C	PAYS	
	N	ANCRX	
	N	TXMORT	
	N	TXANAL	
	N	SCOL2	
	N	PIBH	
	N	CRXPIB	suivi de F3.

Ces réponses impliquent la création de la table sas **BLAYO_SG.sas7bdat** et du catalogue iml **BLAYO_SG.sas7bcat** et portent à 19 le nombre de graphiques du catalogue graphique **GACPB8G3.sas7bcat** placé dans **C:\K8\ABOUL**.

MGVAR : (fenêtre “ GRAF_VAR ”, domaine : **tout programme et toute structure**).

Représentations graphiques (profils) des variables ou des modalités (mais pas des individus) à travers la structure du réseau (grille ou ficelle).

Pour éviter d'alourdir le contenu du catalogue graphique déjà existant, cette macro crée son propre catalogue graphique (par défaut GVAR_TEM.sas7bcat) qui est placé dans le répertoire des données et contient un nombre total de graphiques égal à deux fois le nombre de variables (ou de modalités) retenues dans l'étude.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	suivi de F3

Ces réponses impliquent la création du catalogue graphique GVARB8G3.sas7bcat placé dans C:\K8\ABOUL; ce nouveau catalogue graphique contient $2 \times 6 = 12$ graphiques.

MVC3D : (fenêtre “ G3DVACLU ”, domaine : **tout programme et toute structure**).

Cette macro doit faire suite à MKLUG ou MKLUFj (j=1,2). Comme la macro MGVAR, elle présente l'influence des variables ou des modalités (mais pas des individus) à travers le réseau (grille ou ficelle), mais en plus, elle tient compte de la répartition en super-classes (clusters) en faisant apparaître les couleurs des différents clusters.

Si le nombre de clusters dépasse 22 il faut utiliser la macro MGVAR.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	
nombre de clusters :	5	suivi de F3

Ces réponses rajoutent 6 graphiques au catalogue graphique GVARB8G3.sas7bcat placé dans C:\K8\ABOUL et portent à 18 son nombre total de graphiques.

MKAMG : (fenêtre “ CAMGRID ”, domaine : **KACP, KFAST, KBATCH pour une grille**).

Il s'agit d'une méthode qui permet de croiser une variable qualitative avec les variables quantitatives utilisées dans le cadre de KACP, KFAST ou KBATCH (cf. la thèse P. ROUSSET dont la référence est donnée dans la section III).

Cette macro doit faire suite à MKLUG, elle produit au plus 4 graphiques rajoutés au catalogue graphique déjà existant.

Description et chronologie des 4 graphiques :

Le graphique nommé CLUSPAC représente sur la grille, pour chaque unité gagnante, la répartition sous forme sectorielle (camembert) des modalités d'une variable qualitative .

Suit le graphique nommé CLUSCAM; ici on colorie CLUSPAC en remplaçant les valeurs numériques (modalités) de la variable qualitative par des couleurs.

Puis vient le graphique nommé CANCLU qui produit un camembert de la variable qualitative pour chaque cluster.

Arrive enfin, le graphique nommé CAMCLU; là encore on colorie le graphique précédent (CANCLU).

Les graphiques, en couleurs, nommés CLUSCAM et CAMCLU ne seront produits que dans la mesure où le nombre de modalités de la variable qualitative et le nombre de clusters ne dépassent pas respectivement 13 et 22 (cf. : les programmes KOLCLUS.SAS et KOLCAM.SAS dans C:\K8).

Dans le cas où la variable qualitative n'est pas numérique, le programme la recode et affiche dans la fenêtre OUTPUT le tableau de conversion.

La table temporaire sas WORK.GG fournit, pour chaque observation, son unité gagnante (variable _codage_), sa super-classe (variable cluster) et la valeur de la variable qualitative (variable var_qual).

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	
nombre de clusters :	5	suivi de F3

Suite à cette soumission, s'ouvre la fenêtre " QUELDSN " qui permet de sélectionner la table qui contient la variable qualitative.

Exemple :	CHOIX	Nom de la Table	
	X	BLAYO2	suivi de F3.

Puis, s'ouvre la fenêtre " QUEL_VAR " qui permet la sélection par V de la variable qualitative et par O de la variable qui identifie les observations.

Exemple :	CHOIX	VARIABLE	
	O	PAYS	
	V	IDH	suivi de F3.

IDH correspond à 6 niveaux (codés de 1 à 6) de l'Indice du Développement Humain.

Ces réponses portent à 23 le nombre de graphiques du catalogue graphique GACPB8G3.sas7bcat placé dans C:\K8\ABOUL.

MKAMF1 : (fenêtre " CAMFIC10 ", domaine : **KACP** ou **KFAST** pour une ficelle (≤ 10)).

MKAMF2 : (fenêtre " CAMFICEL ", domaine : **KACP** ou **KFAST** pour une ficelle ($>10 \leq 50$)).

Ces deux macros doivent faire suite à MKLUF_j (j=1,2). Elles réalisent sur une ficelle un traitement qui est équivalent à celui de MKAMG pour une grille.

Les graphiques CLUSPAC, CLUSCAM, CANCLU et CAMCLU deviennent respectivement :

Pour MKAMF1: CANBF110, CACOF110, CANCLU et CAMCLU

Pour MKAMF2: CANBFIC, CACOLFIC, CANCLU et CAMCLU

La table temporaire sas WORK.GG (pour MKAMG) devient la table temporaire WORK.F1 (pour MKAMF1) et la table temporaire WORK.F2 (pour MKAMF2).

MDIMA : (fenêtre “ DIMAWS ”, domaine : **KACP, KFAST, KBATCH, KACM et toute structure**).

Cette macro produit une table sas nommée par défaut *_TEMP_DM.sas7bdat* qui est placée dans le répertoire des données. Elle fournit les distances de Mahalanobis entre tous les vecteurs poids finaux des unités du réseau (grille ou ficelle).

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	B8G3CPAY	
type du catalogue :	GACP	suivi de F3

Ces réponses impliquent la création du fichier C:\K8\ABOUL\B8G3C_DM.sas7bdat

MDIMA12 : (fenêtre “ DIMA12WS ”, domaine : **KACM1, KACM2 et toute structure**).

Cette macro est identique à MDIMA.

MDIST : (fenêtre “ DISTWS ”, domaine : **tout programme et toute structure**).

Cette macro produit une table sas nommée par défaut *_TEMP_DE.sas7bdat* qui est placée dans le répertoire des données. Elle fournit les distances euclidiennes entre tous les vecteurs poids finaux des unités du réseau (grille ou ficelle).

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	D3C5G	suivi de F3

Ces réponses impliquent la création du fichier C:\K8\ABOUL\D3C5G_DE.sas7Bdat.

MFMOD : (fenêtre "FILEMOD", domaine : **KORRESP ou KACM(j=1,2) pour toute structure**).

Cette macro doit faire suite à MKLUG ou MKLUFj (j=1,2); elle produit une table sas nommée *MODALn.sas7bdat* où n correspond au nombre de clusters retenu dans MKLUG ou dans MKLUFj (j=1,2). Ce fichier sera placée dans le répertoire des données.

La table *MODALn* fournit pour chaque modalité (variable *idobs*) son numéro de super-classes (variable *cluster*) et son unité gagnante (variable *_codage_*) et sa position (variables: ligne, colonne) si le réseau est une grille.

Exemple :

librairie :	C:\K8\ABOUL	
nom commun :	D3C5G	
type du catalogue :	GACM	
nombre de clusters :	5	suivi de F3

Ces réponses impliquent la création du fichier *MODAL5.sas7bdat* dans le répertoire C:\K8\ABOUL.

MFIND : (fenêtre “ FICKACMj ”, domaine : **KACMj** (j=1,2) pour toute structure).

Cette macro doit faire suite à **MKLUG** ou **MKLUFj** (j=1,2), elle produit une table sas nommée **INDIVn.SAS7BDAT** où n correspond au nombre de clusters retenu dans **MKLUG** ou dans **MKLUFj** (j=1,2).

La table **INDIVn**, placée dans le répertoire des données, ne sera créée que si les individus ont été déclarés non anonymes (**KACMj** pour j=1,2); elle fournit pour chaque individu (variable **idobs**) son numéro de super-classes (variable **cluster**), son unité gagnante (variable **_codage_**) et ses modalités de réponses (les autres variables de la table).

Supposons que l'on ait soumis au préalable la macro **MKLUG** avec les renseignements suivants :

librairie : C:\K8\ABOUL
nom commun : D3C5G
type du catalogue : GACM
nombre de super-classes : 5.

Exemple :

librairie : C:\K8\ABOUL
nom commun : D3C5G
nombre de clusters : 5 suivi de F3

suite à cette soumission, s'ouvre la fenêtre "CHOIXREP" qui permet de sélectionner le nom de la table qui contient les réponses des individus aux questions posées.

Exemple : CHOIX Nom de la Table
X CHIENS suivi de F3.

Puis, s'ouvre la fenêtre “ CHOIXQST ” qui permet de sélectionner par I la variable identifiant les individus et Q les variables qui correspondent aux questions posées.

Exemple :

CHOIX QUESTION
I RACE
Q TAILLE
Q POIDS
Q VELOCITE
Q INTELLIG
Q AFFECTIO
Q AGRESSIV
Q FONCTION suivi de F3.

Ces réponses impliquent la création du fichier **INDIV5.sas7bdat** dans le répertoire **C:\K8\ABOUL**.

Le programme DVIATION.SAS situé dans c:\k8 et mis à jour comme suit, permet de contrôler le bon placement des modalités au regard des super-classes.

```
options nodate nocenter ps=40 ls=125;
libname a 'c:\k8\aboul';
%let nc=5;
%let qdeb=taille;
%let qfin=fonction;
proc freq data=a.indiv&nc;
tables (&qdeb--&qfin)*cluster / norow deviation cellchi2;
run;
```

Ce programme effectue les tris croisés entre la variable cluster (des super-classes) et les variables correspondants aux questions posées .

Pour chaque modalité de réponse :

Une déviation positive indique une "attraction" entre la modalité et la super-classe (la modalité est bien placée).

Une déviation négative indique une "répulsion" entre la modalité et la super-classe (la modalité est mal placée).

L'importance de la déviation est mesurée par l'écart pondéré (la contribution au khi-deux) : cellchi2.

Les pourcentages en colonnes comparés aux pourcentages marginaux du tri croisé donnent une idée des valeurs tests.