

REVUE DE STATISTIQUE APPLIQUÉE

MICHEL ARNAUD

XAVIER EMERY

CHANTAL DE FOUQUET

MARINUS BROUWERS

MICHEL FORTIER

L'analyse krigéante pour le classement d'observations spatiales et multivariées

Revue de statistique appliquée, tome 49, n° 2 (2001), p. 45-67

http://www.numdam.org/item?id=RSA_2001__49_2_45_0

© Société française de statistique, 2001, tous droits réservés.

L'accès aux archives de la revue « *Revue de statistique appliquée* » (<http://www.sfds.asso.fr/publicat/rsa.htm>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques

<http://www.numdam.org/>

L'ANALYSE KRIGEANTE POUR LE CLASSEMENT D'OBSERVATIONS SPATIALES ET MULTIVARIÉES

Michel Arnaud¹, Xavier Emery², Chantal de Fouquet²,
Marinus Brouwers¹, Michel Fortier¹

1. CIRAD Avenue du Val de Montferrand, BP5035, 34032 Montpellier Cedex

2. Ecole des Mines de Paris, Centre de Géostatistique, 35, rue Saint-Honoré,
77305 Fontainebleau Cedex

RÉSUMÉ

Lorsque l'ingénieur, le chercheur... désire classer des observations multivariées, un certain nombre de méthodes statistiques sont mises à sa disposition. Lorsque, en outre, les observations sont disposées dans un espace géographique et que l'information spatiale est riche de sens et doit donc être prise en compte, les méthodes adéquates sont beaucoup moins nombreuses. L'article propose une méthodologie d'analyse pour obtenir des groupes de sites ayant les mêmes caractéristiques. Basée sur les techniques d'analyses de données et de géostatistique elle permet également à l'utilisateur de choisir le niveau de perception (échelle) pour le classement.

Mots-clés : Classification, Géostatistique, Analyse en Composantes Principales, Analyse krigéante.

ABSTRACT

There are various statistical methods that consultants, researchers, etc can use to classify multivariate observations. However, if their observations concern a specific geographic area and the corresponding spatial information is highly relevant and thus needs to be taken into account, the choice of appropriate methods is much more limited. This article proposes an analysis method aimed at identifying groups of sites with the same characteristics. It is based on data analysis and geostatistical techniques, and also enables users to select a perception level (scale) for classification.

Keywords : Classification, Geostatistics, Principal Component Analysis, Multivariate Factorial Kriging.

1. Introduction

Cartographier une zone géographique \mathcal{D} à partir d'observations multivariées et géoréférencées est un objectif courant dans de nombreuses disciplines traitant des phénomènes naturels qui s'étendent dans l'espace : météorologie, agronomie, environnement, géographie... On retrouve cet objectif en sciences du sol pour des études détaillées destinées par exemple à :

- l'analyse des causes attribuables au sol et pouvant expliquer la variabilité spatiale constatée dans les paramètres du rendement au sein d'une parcelle;
- la mise en place d'essais agronomiques afin que la variabilité de l'essai, attribuable aux variations spatiales existant au sein de la parcelle dans les caractéristiques du sol, soit la plus faible possible;
- mieux connaître l'extension spatiale du ou des types de sols au sein de la parcelle sur lesquels on fera en un endroit déterminé des multiples mesures dans le temps destinées à mieux comprendre la croissance de la culture, son bilan hydrique...

Lorsqu'il désire cartographier une zone déterminée, le pédologue prélève des échantillons de sol et relève avec précision leur position sur le terrain. Ces échantillons seront, par la suite, analysés au laboratoire et caractérisés par un ensemble de valeurs quantitatives (les variables physico-chimiques : granulométrie, CEC, bases échangeables, pH, N, C...) ainsi que par leurs coordonnées dans la zone. Pour chaque variable, une première «cartographie» sommaire de la zone, peut s'obtenir, sans avoir recours à des outils compliqués. Un moyen, parmi d'autres, consiste à découper l'histogramme de la variable en plusieurs classes de valeurs (faible, moyenne et forte) et à identifier, par des symboles (ou des couleurs) différents les positions des échantillons selon leur classe d'appartenance. Il est possible alors soit de repérer des secteurs «riches», des secteurs «pauvres», des gradients ou, au contraire, d'obtenir une répartition homogène des symboles sur toute l'étendue de la zone sans lien avec la position des échantillons. Le problème se complique lorsqu'on veut utiliser simultanément l'ensemble des variables pour fournir une cartographie globale et synthétique.

La cartographie multivariée consiste, à partir d'observations ponctuelles pour lesquelles on a mesuré plusieurs variables et relevé leurs coordonnées, à synthétiser par un nombre réduit de nouvelles variables, appelées ici « facteurs », à la fois les valeurs observées sur les emplacements échantillonnés et les dépendances spatiales dues à la proximité, au sens géographique, entre les sites. En tenant compte de la position des sites de prélèvement dans l'espace, on admet implicitement que les sites géographiquement proches ont des valeurs voisines et que les différences, lorsqu'elles existent, ne sont dues qu'aux fluctuations d'échantillonnage. Cette hypothèse aura pour effet de simplifier la visualisation des cartes et de supprimer une grande partie de l'effet de « mitage ».

État de l'art

Quels outils et quelles méthodes statistiques ou géostatistiques sont disponibles pour résoudre cette question ?

En première approche, on peut utiliser les outils classiques bien connus tels ceux de l'analyse factorielle : Analyse en Composantes Principales (ACP) ou Analyse Factorielle des Correspondances (AFC) (Cailliez et Pagès, 1976), en liaison avec les techniques de classification comme les Nuées dynamiques (Diday, 1972) ou la Classification Ascendante Hiérarchique (Jambu, 1978). Malheureusement, on obtient souvent une cartographie de type mosaïque où les classes obtenues apparaissent très fragmentées et très dispersées dans l'espace. Ce « mitage » provient de l'absence de prise en compte de l'information spatiale par ces techniques.

Pour cette raison, la dimension spatiale a été introduite de diverses façons. D'abord en modélisant chaque variable en fonction des coordonnées géographiques. Dans un cadre monovarié, il s'agit de découper l'espace géographique par un arbre dichotomique. L'arbre est construit par détection de seuils de coupure, selon les deux axes de coordonnées, de telle manière que l'homogénéité des observations dans chaque unité spatiale soit maximum (Breiman *et al*, 1984). D'autres auteurs (Ambroise *et al*, 1997) proposent d'utiliser l'algorithme « EM », où le problème de classification est formulé en un problème d'optimisation. Ils introduisent l'information spatiale en ajoutant au critère à maximiser une quantité ayant une signification géographique.

Dans une autre optique (analyse des données à la française), certains auteurs ont introduit la contrainte de contiguïté entre les observations. Mesurée par le coefficient de Geary (Lebart, 1969), elle est utilisée dans la Classification Hiérarchique avec des contraintes de voisinages (Lebart, 1978). Ces travaux ont été poursuivis en introduisant des opérateurs de voisinage dans l'analyse des données spatio-temporelles (Méot *et al*, 1993) et avec Chessel et Sabatier (1993) qui proposent l'étude de deux ensembles de variables tenant compte de la proximité des observations mais dans une optique explicative. Enfin Thioulouse *et al* (1995) réalisent des ACP et des AFC qui maximisent les structures globales, locales et totales.

La géostatistique (Matheron, 1962) offre une autre perspective. En général, ces techniques permettent de répondre à des problèmes d'estimation locale ou globale de paramètres mesurés à l'intérieur du domaine géographique \mathcal{D} . Le caractère multivarié est pris en compte dans les techniques de cokrigage. La géostatistique permet ainsi de construire des blocs diagrammes à trois dimensions ou de dessiner des cartes d'isovaleurs. Cependant, ces produits finaux sont monovariés : ils ne cartographient qu'une variable à la fois.

A partir de l'analyse géostatistique et des méthodes de classification, Oliver et Webster (1989) proposent d'appliquer un algorithme de classification non hiérarchique sur une matrice de dissimilarité particulière. Cette procédure considère, non seulement les différences entre les valeurs mesurées sur les variables mais également la structure spatiale par l'intermédiaire des variogrammes modélisés.

Le caractère multivarié et l'information spatiale peuvent aussi être analysés dans une méthodologie utilisant de façon combinée : géostatistique, analyse multivariée et classification (Arnaud et Pichot, 1996). Après découpage de l'espace géographique en cellules et cokrigage par bloc de ces cellules, les valeurs estimées sont analysées par ACP et classification pour obtenir une cartographie de l'espace géographique. Cette méthodologie n'est pas sans critique en raison, notamment, de la non indépendance spatiale des facteurs de l'ACP (Wackernagel, 1998) utilisés dans la classification.

Plus récemment Allard *et al* (1999) proposent, dans le cas monovarié, une méthodologie de classification où la composante spatiale est prise en compte par l'estimation de la variance intra-groupe lorsque le nombre de groupes est fixé *a priori*.

Dans cet article, nous proposons d'utiliser l'analyse krigeante (Matheron, 1982; Wackernagel, 1998) qui considère simultanément les caractères multivariés et géographiques des observations. Nous verrons comment cette méthode permet de répondre à un problème de classification et comment elle l'enrichit en donnant des

indications sur le caractère plus ou moins régulier du phénomène régionalisé, les différents niveaux d'échelle, la détection d'éventuelles anisotropies...

2. Présentation des méthodes géostatistiques et de l'analyse krigéante

La géostatistique a été développée pour résoudre le problème de l'estimation des variables régionalisées (Matheron, 1962). Une variable régionalisée $z(s)$, définie pour s variant dans le domaine \mathcal{D} , est interprétée comme une réalisation, un tirage particulier, d'une fonction aléatoire $Z(s)$. L'aspect erratique de $z(s)$ est pris en compte par un modèle probabiliste, l'aspect régionalisé par la corrélation spatiale ou plus généralement par la loi spatiale de la fonction aléatoire $Z(s)$.

On peut généraliser cette approche au cas multivarié : chaque variable régionalisée $z_i(s)$, $i = 1, \dots, N$, est considérée comme une réalisation d'un processus $\{Z_i(s), s \in \mathcal{D}\}$. Deux outils, les fonctions de covariance et les variogrammes, reliés entre eux dans certaines conditions, vont rendre possible l'analyse de la structure spatiale *conjointe* des variables $Z_i(s)$. L'introduction d'hypothèses de stationnarité va être nécessaire : la stationnarité des deux premiers moments des variables considérées (stationnarité conjointe du second ordre) ou la stationnarité de deux premiers moments des *accroissements* des variables (hypothèse intrinsèque conjointe). Ces hypothèses vont rendre possible l'inférence statistique à partir des valeurs expérimentales pour estimer les fonctions de covariance ou les variogrammes. Leur modélisation se situe, en général, dans le cadre du modèle linéaire de corégionalisation. Celui-ci va permettre de décomposer les fonctions aléatoires en facteurs dont les propriétés de non corrélation spatiale vont donner les moyens de cartographier la zone d'étude.

2.1. Fonctions de covariance et variogrammes : propriétés et inférence

On dit que les fonctions aléatoires Z_i sont **conjointement stationnaires du second ordre** si leurs deux premiers moments (espérance et covariance) existent et sont invariants par translation :

$$\begin{cases} E[Z_i(s)] = m_i & \forall i = 1, \dots, N \text{ et } \forall s \in \mathcal{D} \\ \text{cov}[Z_i(s + \mathbf{h}), Z_j(s)] = C_{ij}(\mathbf{h}) & \forall i, j = 1, \dots, N \text{ et } \forall s, s + \mathbf{h} \in \mathcal{D} \end{cases}$$

Les $C_{ij}(\mathbf{h})$ sont les fonctions de covariances simples ($i = j$) et croisées ($i \neq j$). Ces fonctions ne sont pas quelconques et doivent vérifier certaines contraintes mathématiques.

Les fonctions aléatoires Z_i sont **conjointement intrinsèques** si leurs accroissements sont conjointement stationnaires du second ordre. On peut résumer cette condition par les relations suivantes :

$$\begin{cases} E[Z_i(s + \mathbf{h}) - Z_i(s)] = 0 & \forall i = 1, \dots, N \text{ et } \forall s, s + \mathbf{h} \in \mathcal{D} \\ \text{cov}[Z_i(s + \mathbf{h}) - Z_i(s), Z_j(s + \mathbf{h}) - Z_j(s)] & \\ = 2\gamma_{ij}(\mathbf{h}) & \forall i, j = 1, \dots, N \text{ et } s, s + \mathbf{h} \in \mathcal{D} \end{cases}$$

Les $\gamma_{ij}(\mathbf{h})$ sont les variogrammes simples ($i = j$) et croisés ($i \neq j$). On montre que, pour tout vecteur \mathbf{h} , la matrice $\Gamma(\mathbf{h}) = [\gamma_{ij}(\mathbf{h})]_{i,j=1\dots N}$ est symétrique et de type positif.

Les fonctions de covariance et les variogrammes simples et croisés fournissent une description élémentaire de la dépendance entre les valeurs prises entre deux points séparés de \mathbf{h} par les différentes variables. Ils modélisent la « structure » spatiale du phénomène multivarié étudié.

L'hypothèse stationnaire conjointe implique l'hypothèse intrinsèque conjointe, mais la réciproque n'est pas vraie. On qualifie de strictement intrinsèques des fonctions aléatoires intrinsèques mais non stationnaires. Dans le cas où l'hypothèse stationnaire conjointe est vérifiée, on a la relation :

$$\forall i, j = 1, \dots, N, \forall \mathbf{h}, \gamma_{ij}(\mathbf{h}) = C_{ij}(\mathbf{0}) - \frac{1}{2} [(C_{ij}(-\mathbf{h}) + C_{ij}(\mathbf{h}))]$$

qui montre que le variogramme croisé ne contient que la partie paire de la fonction de covariance croisée, et renferme donc moins d'information structurale. En dépit de cette restriction, on préfère, le plus souvent, utiliser les variogrammes simples et croisés, car ils sont définis dans le cadre intrinsèque, plus général que le cadre stationnaire du second ordre.

A partir des données expérimentales, on peut estimer les fonctions de covariance et les variogrammes. Notons s_α^i ($\alpha = 1 \dots n_i$) les sites de mesure de la variable z_i . Ces sites peuvent être différents d'une variable à l'autre (cas d'*hétérotopie*, partielle ou totale). La fonction de covariance simple ou croisée $C_{ij}(\mathbf{h})$ est habituellement estimée par :

$$\hat{C}_{ij}(\mathbf{h}) = \frac{1}{|N_{ij}(\mathbf{h})|} \sum_{N_{ij}(\mathbf{h})} [Z_i(s_\alpha^i) - \bar{Z}_i] [Z_j(s_\beta^j) - \bar{Z}_j]$$

$$\text{où } \bar{Z}_i = \frac{1}{n_i} \sum_{\alpha=1}^{n_i} Z_i(s_\alpha^i) \text{ et } \bar{Z}_j = \frac{1}{n_j} \sum_{\beta=1}^{n_j} Z_j(s_\beta^j)$$

$$N_{ij}(\mathbf{h}) = \{(\alpha, \beta) \text{ tel que } s_\alpha^i - s_\beta^j = \mathbf{h}\}$$

$|N_{ij}(\mathbf{h})|$ est le nombre de paires distinctes de l'ensemble $N_{ij}(\mathbf{h})$.

Quant au variogramme $\gamma_{ij}(\mathbf{h})$, il est estimé par :

$$\hat{\gamma}_{ij}(\mathbf{h}) = \frac{1}{2|N_{ij}(\mathbf{h})|} \sum_{N_{ij}(\mathbf{h})} [Z_i(s_\alpha^i) - Z_i(s_\beta^i)] [Z_j(s_\alpha^j) - Z_j(s_\beta^j)]$$

$$\text{où } N_{ij}(\mathbf{h}) = \{(\alpha, \beta) \text{ tel que } s_\alpha^i = s_\alpha^j, s_\beta^i = s_\beta^j \text{ et } s_\alpha^i - s_\beta^i = s_\alpha^j - s_\beta^j = \mathbf{h}\}$$

Le calcul de $\hat{\gamma}_{ij}(\mathbf{h})$ (mais pas celui de $\hat{C}_{ij}(\mathbf{h})$) est impossible en cas d'hétérotopie totale, auquel cas l'ensemble $N_{ij}(\mathbf{h})$ est vide.

2.2. Le modèle linéaire de corégionalisation

Les variogrammes (ou covariances) simples et croisés d'un ensemble de variables $Z_1 \dots Z_N$ ne peuvent être modélisés indépendamment, car il existe entre eux des contraintes mathématiques. En pratique, pour satisfaire ces contraintes, on a recours au *modèle linéaire de corégionalisation*, qui permet de faire face à la plupart des situations. On peut en donner une définition dans les deux cas : stationnaire et intrinsèque. L'interprétation de ce modèle sera détaillée dans le paragraphe 2.3 concernant l'analyse krigéante.

N fonctions aléatoires $Z_i(\mathbf{s})$ conjointement stationnaires d'ordre deux obéissent au modèle linéaire de corégionalisation lorsque leurs covariances sont des combinaisons linéaires des mêmes covariances de base $\rho_u(\mathbf{h})$:

$$\forall i, j = 1, \dots, N, \forall \mathbf{h}, C_{ij}(\mathbf{h}) = \sum_{u=1}^S b_{ij}^u \rho_u(\mathbf{h})$$

en notation matricielle : $\mathbf{C}(\mathbf{h}) = \sum_{u=1}^S \mathbf{B}_u \rho_u(\mathbf{h})$

où $\mathbf{C}(\mathbf{h})$ est la matrice des covariances : $\mathbf{C}(\mathbf{h}) = [C_{ij}(\mathbf{h})]_{i,j=1,\dots,N}$

$\rho_u(\mathbf{h})$ est une covariance telle que $\rho_u(\mathbf{0}) = 1$

$\mathbf{B}_u = (b_{ij}^u)$ est appelée matrice de corégionalisation.

On obtient un modèle cohérent en imposant que chaque matrice \mathbf{B}_u soit symétrique et de type positif. En particulier, on aura la relation :

$$\forall i, j = 1 \dots N, \forall u = 1 \dots S, |b_{ij}^u| \leq \sqrt{b_{ii}^u b_{jj}^u} \quad (1)$$

Cette inégalité implique que si, pour un modèle ρ_u quelconque, la valeur b_{ij}^u associée à la covariance croisée est non nulle, les valeurs b_{ii}^u et b_{jj}^u relatives aux deux covariances simples correspondantes sont nécessairement strictement positives.

N fonctions aléatoires $Z_i(\mathbf{s})$ conjointement intrinsèques obéissent au modèle linéaire de corégionalisation lorsque les variogrammes sont des combinaisons linéaires des mêmes variogrammes de base $g_u(\mathbf{h})$:

$$\forall i, j = 1 \dots N, \forall \mathbf{h}, \gamma_{ij}(\mathbf{h}) = \sum_{u=1}^S b_{ij}^u g_u(\mathbf{h})$$

en notation matricielle : $\Gamma(\mathbf{h}) = \sum_{u=1}^S \mathbf{B}_u g_u(\mathbf{h})$

où $\Gamma(\mathbf{h})$ est la matrice des variogrammes,

$g_u(\mathbf{h})$ est un variogramme de base,

$\mathbf{B}_u = (b_{ij}^u)$ est appelée matrice de corrégalisation.

Comme $\Gamma(\mathbf{h})$ doit être une matrice symétrique de type positif pour tout \mathbf{h} , on impose (condition suffisante) que *chacune* des matrices de corrégalisation \mathbf{B}_u soit symétrique de type positif.

La recherche des coefficients des matrices \mathbf{B}_u se fait par ajustement à partir des covariances ou variogrammes expérimentaux. Dans le cas bivarié, il s'effectue d'abord sur les variogrammes simples puis se poursuit par celui du variogramme croisé en utilisant la relation (1). Dans le cas de plus de deux fonctions aléatoires, Goulard *et al* (1992) proposent une procédure itérative qui fournit des matrices \mathbf{B}_u ayant les bonnes propriétés, mais la convergence de l'algorithme n'a pu être démontrée.

2.3. L'analyse krigeante

L'analyse krigeante s'effectue en deux étapes successives : la première étape réalise la décomposition des fonctions aléatoires en une somme de *facteurs* indépendants; la seconde procède à l'estimation, par cokrigeage, de ces facteurs en chaque point $s \in \mathcal{D}$.

2.3.1. Cas d'une seule variable $Z(s)$

2.3.1.1. Décomposition de la variable $Z(s)$

a) cas d'une fonction aléatoire stationnaire

Nous allons considérer dans un premier temps une fonction aléatoire $Z(s)$ stationnaire d'ordre deux, de moyenne (espérance) m et de fonction de covariance « gigogne », c'est-à-dire composée de plusieurs modèles élémentaires :

$$\forall \mathbf{h} : C(\mathbf{h}) = \sum_{u=1}^S b_u \rho_u(\mathbf{h})$$

où $\forall u = 1 \dots S, b_u > 0$

$\forall u = 1 \dots S, \rho_u(\mathbf{h})$ est une covariance telle que $\rho_u(\mathbf{0}) = 1$

On peut décomposer $Z(s)$ en une somme de S *composantes spatiales* stationnaires X^u :

$$\forall s \in \mathcal{D}, Z(s) = \sum_{u=1}^S X^u(s) + m$$

- d'espérances nulles : $\forall u = 1 \dots S, \forall s \in \mathcal{D}, E[X^u(s)] = 0$

- spatialement non corrélées et de fonctions de covariance $b_u \rho_u(\mathbf{h})$:

$$\forall u, v = 1 \dots S, \forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \text{cov}[X^u(\mathbf{s} + \mathbf{h}), X^v(\mathbf{s})] = \begin{cases} 0 & \text{si } u \neq v \\ b_u \rho_u(\mathbf{h}) & \text{si } u = v \end{cases}$$

La variable $Z(\mathbf{s})$ apparaît comme la superposition de composantes « indépendantes » qui agissent à des échelles spatiales différentes (chaque covariance est associée à une dimension caractéristique, donnée par sa *portée*, *i.e.* la distance à partir de laquelle elle est identiquement nulle).

A titre d'exemple, supposons que la fonction aléatoire stationnaire d'ordre deux $Z(\mathbf{s})$ admette une covariance gigogne composée de deux modèles sphériques de portées 15 m et 70 m :

$$C(\mathbf{h}) = 3C_{\text{sph}}(\mathbf{h}, 15 \text{ m}) + 5C_{\text{sph}}(\mathbf{h}, 70 \text{ m}).$$

On rappelle que le modèle de covariance sphérique de portée a s'écrit :

$$C_{\text{sph}}(r, a) = \begin{cases} 1 - \frac{3}{2} \frac{r}{a} + \frac{1}{2} \frac{r^3}{a^3} & \text{pour } 0 \leq r \leq a \\ 0 & \text{pour } r \geq a \end{cases}$$

Il s'agit d'un des modèles de covariance les plus couramment utilisés.

$Z(\mathbf{s})$ peut alors se décomposer de la façon suivante :

$$Z(\mathbf{s}) = X^1(\mathbf{s}) + X^2(\mathbf{s}) + m$$

où $X^1(\mathbf{s})$ et $X^2(\mathbf{s})$ sont deux fonctions aléatoires stationnaires d'ordre deux, sans corrélation spatiale mutuelle, de moyenne nulle et de covariances respectives $C_1(\mathbf{h}) = 3 \text{sph}(\mathbf{h}, 15 \text{ m})$ et $C_2(\mathbf{h}) = 5 \text{sph}(\mathbf{h}, 70 \text{ m})$.

b) cas d'une fonction aléatoire intrinsèque stricte

Considérons à présent une fonction aléatoire $Z(\mathbf{s})$ obéissant au modèle linéaire de corégionalisation avec un variogramme gigogne composé de $S - 1$ variogrammes stationnaires (*i.e.* qui se stabilisent à l'infini autour d'une valeur limite appelée *palier*) et d'un variogramme intrinsèque strict (non borné) :

$$\forall \mathbf{h}, \gamma(\mathbf{h}) = \sum_{u=1}^{S-1} b_u g_u(\mathbf{h}) + b_S g_S(\mathbf{h}).$$

$Z(\mathbf{s})$ peut être décomposée en une somme de S composantes spatiales X^u , où $X^1 \dots X^{S-1}$ sont stationnaires et X^S strictement intrinsèque, de variogrammes respectifs $b_u g_u(\mathbf{h})$ et $b_S g_S(\mathbf{h})$:

$$\forall \mathbf{s} \in \mathcal{D}, Z(\mathbf{s}) = \sum_{u=1}^{S-1} X^u(\mathbf{s}) + X^S(\mathbf{s}).$$

Les composantes stationnaires X^u et les accroissements de X^S sont d'espérances nulles et sans corrélation spatiale.

Ainsi donc, toute fonction aléatoire $Z(\mathbf{s})$ qui admet une covariance ou un variogramme *gigogne* peut être décomposée en plusieurs composantes spatiales spatialement non corrélées ou à accroissements spatialement non corrélés, agissant à différentes échelles.

2.3.1.2. Estimation des composantes spatiales par krigeage

Lorsqu'on considère le modèle intrinsèque décrit ci-dessus¹, on estime les valeurs des composantes spatiales $X^u(\mathbf{s})$ par krigeage à partir des mesures expérimentales $Z(\mathbf{s}_\alpha)$. On pose, pour tout $u = 1 \dots S$:

$$\widehat{X}_u(\mathbf{s}) = \sum_{\alpha} \omega_{\alpha}^u Z(\mathbf{s}_{\alpha})$$

où les ω_{α}^u sont des pondérateurs réels. Ils sont déterminés par le système d'équations suivant :

$$\begin{cases} \sum_{\beta} \omega_{\beta}^u \gamma(\mathbf{s}_{\alpha} - \mathbf{s}_{\beta}) - \mu_u = b_u g_u(\mathbf{s}_{\alpha} - \mathbf{s}) \text{ pour } \alpha = 1 \dots n \\ \sum_{\beta} \omega_{\beta}^u = \begin{cases} 0 & \text{pour } u = 1 \dots S - 1 \text{ (} X^u \text{ stationnaire)} \\ 1 & \text{pour } u = S \text{ (} X^S \text{ intrinsèque stricte)} \end{cases} \end{cases}$$

où μ_u est un multiplicateur de Lagrange.

Ce système diffère du krigeage usuel par l'apparition dans le membre de droite des termes $b_u g_u(\mathbf{s}_{\alpha} - \mathbf{s})$ spécifiques de la composante $X^u(\mathbf{s})$ et, dans le cas où $X^u(\mathbf{s})$ est stationnaire, par la condition sur la somme des poids (Wackernagel, 1998).

2.3.2. Décomposition de plusieurs fonctions aléatoires en corrélation intrinsèque

2.3.2.1. Décomposition des variables $Z_i(\mathbf{s})$

a) Cas stationnaire du second ordre

N fonctions aléatoires conjointement stationnaires d'ordre deux $Z_1(\mathbf{s}) \dots Z_N(\mathbf{s})$, de moyennes respectives $m_1 \dots m_N$ sont en *corrélation intrinsèque* si leurs fonctions de covariance sont proportionnelles entre elles :

$$\forall i, j = 1 \dots N, \forall \mathbf{h}, C_{ij}(\mathbf{h}) = \sigma_{ij} \rho_0(\mathbf{h})$$

où $\rho_0(\mathbf{h})$ est une covariance telle que $\rho_0(\mathbf{0}) = 1$

$V = (\sigma_{ij})$ est la matrice de variance-covariance des $Z_i(\mathbf{s})$

¹ On se limite à une seule structure intrinsèque stricte car l'estimation par krigeage de X^u est impossible pour plus de deux composantes strictement intrinsèques.

On montre que les Z_i peuvent être décomposées en un ensemble de N facteurs spatiaux $Y_p(\mathbf{s})$:

$$\forall i = 1 \dots N, \forall \mathbf{s} \in \mathcal{D}, Z_i(\mathbf{s}) = \sum_{p=1}^N a_p^i Y_p(\mathbf{s}) + m_i \quad (2)$$

- d'espérances nulles : $\forall p = 1 \dots N, \forall \mathbf{s} \in \mathcal{D} : E[Y_p(\mathbf{s})] = 0$
- spatialement non corrélés et de fonction de covariance $\rho_0(\mathbf{h})$:

$$\forall p, q = 1 \dots N, \forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \text{cov}[Y_p(\mathbf{s} + \mathbf{h}), Y_q(\mathbf{s})] = \begin{cases} 0 & \text{si } p \neq q \\ \rho_0(\mathbf{h}) & \text{si } p = q \end{cases}$$

Les facteurs spatiaux $Y_p(\mathbf{s})$ sont définis par la relation :

$$(\sqrt{\lambda_1} Y_1(\mathbf{s}), \dots, \sqrt{\lambda_N} Y_N(\mathbf{s})) = (Z_1(\mathbf{s}) - m_1, \dots, Z_n(\mathbf{s}) - m_N) \times \mathbf{Q}$$

obtenue lorsque l'on effectue l'analyse en composantes principales de la matrice de variance-covariance \mathbf{V} . On a alors en notation matricielle :

$$\mathbf{V} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^t \text{ avec } \mathbf{Q}^t \mathbf{Q} = \mathbf{I} \text{ (matrice identité) et } \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$$

où \mathbf{Q} est une matrice orthogonale de vecteurs propres de \mathbf{V} et $\mathbf{\Lambda}$ la matrice diagonale des valeurs propres associées.

Chaque facteur spatial Y_p est associé à une valeur propre positive λ_p de \mathbf{V} qui mesure la « quantité » d'information contenue dans ce facteur. La comparaison des valeurs propres permet de classer les facteurs par ordre d'importance.

Lorsque les $Z_i(\mathbf{s})$ sont exprimées dans des unités différentes, il est préférable de travailler sur les variables standardisées (centrées et réduites) :

$$\tilde{Z}_i(\mathbf{s}) = \frac{Z_i(\mathbf{s}) - m_i}{\sqrt{\sigma_{ii}}}$$

qui deviennent alors comparables (car sans dimensions) et l'analyse en composantes principales est réalisée sur la matrice des corrélations au lieu de celle de variance-covariance.

A titre d'exemple, supposons que les deux fonctions aléatoires stationnaires d'ordre deux, $Z_1(\mathbf{s})$ et $Z_2(\mathbf{s})$ soient en corrélation intrinsèque avec les covariances simples et croisés suivantes :

$$C_{11}(\mathbf{h}) = 4C_{\text{sph}}(\mathbf{h}, 70 \text{ m})$$

$$C_{12}(\mathbf{h}) = C_{\text{sph}}(\mathbf{h}, 70 \text{ m})$$

$$C_{22}(\mathbf{h}) = 9C_{\text{sph}}(\mathbf{h}, 70 \text{ m})$$

soit en notation matricielle :

$$\begin{pmatrix} C_{11}(\mathbf{h}) & C_{12}(\mathbf{h}) \\ C_{21}(\mathbf{h}) & C_{22}(\mathbf{h}) \end{pmatrix} = \begin{pmatrix} 4 & 1 \\ 1 & 9 \end{pmatrix} C_{\text{sph}}(\mathbf{h}, 70 \text{ m}).$$

On obtient alors la décomposition des $Z_i(\mathbf{s})$ en une somme de facteurs spatiaux $Y_i(\mathbf{s})$:

$$\begin{pmatrix} Z_1(\mathbf{s}) \\ Z_2(\mathbf{s}) \end{pmatrix} = \begin{pmatrix} 1.92 & 0.57 \\ -0.37 & 2.98 \end{pmatrix} \begin{pmatrix} Y_1(\mathbf{s}) \\ Y_2(\mathbf{s}) \end{pmatrix} + \begin{pmatrix} m_1 \\ m_2 \end{pmatrix}$$

Les valeurs propres associées sont $\lambda_1 = 9.19$ et $\lambda_2 = 3.81$. Le facteur Y_1 contient donc 70 % de l'information totale.

b) Cas intrinsèque

N fonctions aléatoires intrinsèques $Z_1(\mathbf{s}) \dots Z_N(\mathbf{s})$ sont en **corrélacion intrinsèque** si leurs variogrammes sont proportionnels entre eux :

$$\forall i, j = 1 \dots N, \forall \mathbf{h}, \gamma_{ij}(\mathbf{h}) = p_{ij} g_0(\mathbf{h})$$

où $g_0(\mathbf{h})$ est un variogramme de base,

$P = (p_{ij})$ est appelée matrice des paliers.

Les $Z_i(\mathbf{s})$ peuvent être décomposées en un ensemble de N facteurs spatiaux $Y_p(\mathbf{s})$:

$$\forall i = 1, \dots, N, \forall \mathbf{s} \in \mathcal{D}, Z_i(\mathbf{s}) = \sum_{p=1}^N a_p^i Y_p(\mathbf{s}).$$

Les Y_p ont pour variogramme $g_0(\mathbf{h})$; leurs accroissements sont d'espérances nulles et sans corrélation spatiale mutuelle. Ils sont définis par la relation :

$$\forall \mathbf{s} \in \mathcal{D}, (\sqrt{\lambda_1} Y_1(\mathbf{s}), \dots, \sqrt{\lambda_N} Y_N(\mathbf{s})) = (Z_1(\mathbf{s}), \dots, Z_n(\mathbf{s})) \times \mathbf{Q} \quad (3)$$

où \mathbf{Q} est une matrice orthogonale de vecteurs propres de \mathbf{P} et $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ la matrice diagonale de valeurs propres associées :

$$\mathbf{P} = \mathbf{Q} \Lambda \mathbf{Q}^t \text{ avec } \mathbf{Q}^t \mathbf{Q} = \mathbf{I}.$$

Comme dans le cas stationnaire, la comparaison des valeurs propres $\lambda_1, \dots, \lambda_N$ permet de hiérarchiser les facteurs par ordre d'importance.

2.3.2.2. Estimation des facteurs spatiaux $Y_p(\mathbf{s})$ en un point $\mathbf{s} \in \mathcal{D}$

a) si le point \mathbf{s} coïncide avec l'un des points de données \mathbf{s}_α où toutes les variables Z_j sont connues, alors $Y_p(\mathbf{s})$ sera évalué par l'une des relations (2) ou (3) (celle correspondant à l'hypothèse adoptée : stationnaire d'ordre deux ou intrinsèque).

b) si \mathbf{s} n'est pas un point de données, on estimera $Y_p(\mathbf{s})$ par cokrigage à partir des $Z_j(\mathbf{s}_\alpha)$, ou, ce qui revient au même, à partir des $Y_j(\mathbf{s}_\alpha)$ précédemment calculés. Or les facteurs $Y_j(\mathbf{s})$ étant en corrélation intrinsèque, leur cokrigage, en cas d'homotopie, se réduit à leur krigeage séparé, ce qui simplifie considérablement les calculs.

2.3.3. Cas de plusieurs fonctions aléatoires obéissant au modèle linéaire de corégionalisation

2.3.3.1. Décomposition des N variables $Z_i(\mathbf{s})$

a) Cas stationnaire du second ordre : dans un premier temps, en appliquant les résultats du paragraphe 2.3.1 et en supposant que les $Z_i(\mathbf{s})$ obéissent au modèle linéaire de corégionalisation, on peut décomposer chaque fonction aléatoire $Z_i(\mathbf{s})$ en un ensemble de composantes spatiales $\{X_i^u(\mathbf{s})\}_{u=1\dots S}$:

$$\forall i = 1 \dots N, \forall \mathbf{s} \in \mathcal{D}, Z_i(\mathbf{s}) = \sum_{u=1}^S X_i^u(\mathbf{s}) + m_i$$

Les $X_i^u(\mathbf{s})$ sont spatialement non corrélées si $u \neq v$; si $u = v$, leur covariance est donnée par la matrice $\mathbf{B}_u \rho_u(\mathbf{h})$.

Dans un deuxième temps, comme les composantes spatiales $\{X_i^u(\mathbf{s})\}_{i=1\dots N}$ sont en corrélation intrinsèque, elles peuvent être à leur tour décomposées en **facteurs spatiaux** $Y_p^u(\mathbf{s})$ non corrélés spatialement (paragraphe 2.3.2). On obtient alors le résultat général suivant :

N fonctions aléatoires conjointement stationnaires d'ordre deux $Z_1(\mathbf{s}) \dots Z_N(\mathbf{s})$, de moyennes respectives $m_1 \dots m_N$ et de covariances de la forme :

$$\forall i, j = 1 \dots N, \forall \mathbf{h}, C_{ij}(\mathbf{h}) = \sum_{u=1}^S b_{ij}^u \rho_u(\mathbf{h})$$

peuvent être décomposées en un ensemble de facteurs spatiaux $Y_p^u(\mathbf{s})$:

$$\forall i = 1 \dots N, \forall \mathbf{s} \in \mathcal{D}, Z_i(\mathbf{s}) = \sum_{u=1}^S \sum_{p=1}^N a_{up}^i Y_p^u(\mathbf{s}) + m_i$$

- d'espérances nulles :

$$\forall u = 1 \dots S, \forall p = 1 \dots N, \forall \mathbf{s} \in \mathcal{D}, E[Y_p^u(\mathbf{s})] = 0$$

- spatialement non corrélés :

$$\forall p, q = 1 \dots N, \forall u, v = 1 \dots S, p \neq q \text{ ou } u \neq v, \forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \\ \text{cov}[Y_p^u(\mathbf{s} + \mathbf{h}), Y_q^v(\mathbf{s})] = 0$$

- de fonction de covariance $\rho_u(\mathbf{h})$:

$$\forall p = 1 \dots N, \forall u = 1 \dots S, \forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \text{cov}[Y_p^u(\mathbf{s} + \mathbf{h}), Y_p^u(\mathbf{s})] = \rho_u(\mathbf{h})$$

Les facteurs spatiaux $Y_p^u(\mathbf{s})$ sont définis par la relation :

$$\forall \mathbf{s} \in \mathcal{D} : \left(\sqrt{\lambda_1^u} Y_1^u(\mathbf{s}), \dots, \sqrt{\lambda_N^u} Y_N^u(\mathbf{s}) \right) \\ = (X_1^u(\mathbf{s}) - m_1, \dots, X_N^u(\mathbf{s}) - m_N) \times \mathbf{Q}_u \quad (4)$$

où \mathbf{Q}_u est une matrice orthogonale de vecteurs propres de \mathbf{B}_u et $\Lambda_u = \text{diag}(\lambda_1^u, \dots, \lambda_N^u)$ la matrice diagonale de valeurs propres associées :

$$\mathbf{B}_u = \mathbf{Q}_u \Lambda_u \mathbf{Q}_u^t \text{ avec } \mathbf{Q}_u^t \mathbf{Q}_u = \mathbf{I}.$$

A titre d'exemple, supposons que les deux fonctions aléatoires stationnaires d'ordre deux, $Z_1(\mathbf{s})$ et $Z_2(\mathbf{s})$ admettent les covariances gigognes simples et croisées :

$$C_{11}(\mathbf{h}) = 4 \text{sph}(\mathbf{h}, 5 m) + 2 \text{sph}(\mathbf{h}, 50 m) + 15 \text{sph}(\mathbf{h}, 500 m) \\ C_{12}(\mathbf{h}) = \text{sph}(\mathbf{h}, 5 m) - 3 \text{sph}(\mathbf{h}, 50 m) + 16 \text{sph}(\mathbf{h}, 500 m) \\ C_{22}(\mathbf{h}) = 9 \text{sph}(\mathbf{h}, 5 m) + 10 \text{sph}(\mathbf{h}, 50 m) + 25 \text{sph}(\mathbf{h}, 500 m)$$

soit, en notation matricielle :

$$\begin{pmatrix} C_{11}(\mathbf{h}) & C_{12}(\mathbf{h}) \\ C_{21}(\mathbf{h}) & C_{22}(\mathbf{h}) \end{pmatrix} = \begin{pmatrix} 4 & 1 \\ 1 & 9 \end{pmatrix} \text{sph}(\mathbf{h}, 5 m) + \begin{pmatrix} 2 & -3 \\ -3 & 10 \end{pmatrix} \text{sph}(\mathbf{h}, 50 m) \\ + \begin{pmatrix} 15 & 16 \\ 16 & 25 \end{pmatrix} \text{sph}(\mathbf{h}, 500 m).$$

On obtient alors la décomposition des $Z_i(\mathbf{s})$ en une somme de facteurs spatiaux $Y_i(\mathbf{s})$ que l'on peut écrire matriciellement :

$$\begin{pmatrix} Z_1(\mathbf{s}) \\ Z_2(\mathbf{s}) \end{pmatrix} = \begin{pmatrix} 1.92 & 0.57 \\ -0.37 & 2.98 \end{pmatrix} \begin{pmatrix} Y_1^1(\mathbf{s}) \\ Y_1^2(\mathbf{s}) \end{pmatrix} + \begin{pmatrix} 0.95 & -1.05 \\ 0.32 & 3.15 \end{pmatrix} \begin{pmatrix} Y_2^1(\mathbf{s}) \\ Y_2^2(\mathbf{s}) \end{pmatrix} \\ + \begin{pmatrix} 1.45 & 3.59 \\ -1.07 & 4.89 \end{pmatrix} \begin{pmatrix} Y_1^3(\mathbf{s}) \\ Y_2^3(\mathbf{s}) \end{pmatrix} + \begin{pmatrix} m_1 \\ m_2 \end{pmatrix}$$

b) Cas intrinsèque : le même raisonnement peut être effectué pour N fonctions aléatoires intrinsèques $Z_i(\mathbf{s})$, obéissant au modèle linéaire de corégionalisation :

$$\forall i, j = 1 \dots N, \forall \mathbf{h}, \gamma_{ij}(\mathbf{h}) = \sum_{u=1}^{S-1} b_{ij}^u g_u(\mathbf{h}) + b_{ij}^S g_S(\mathbf{h})$$

où

- $g_u(\mathbf{h})$, $u = 1 \dots S - 1$, sont des variogrammes stationnaires (bornés) de palier unité, *i.e.* tels que $g_u(\mathbf{h}) \xrightarrow[|\mathbf{h}| \rightarrow \infty]{} 1$;

- $g_S(\mathbf{h})$ est un variogramme strictement intrinsèque (non borné).

Les $Z_i(\mathbf{s})$ peuvent être décomposées en un ensemble de facteurs spatiaux $Y_p^u(\mathbf{s})$:

$$\forall i = 1 \dots N, \forall \mathbf{s} \in \mathcal{D}, Z_i(\mathbf{s}) = \sum_{u=1}^S \sum_{p=1}^N a_{up}^i Y_p^u(\mathbf{s})$$

où les facteurs stationnaires ($u = 1 \dots S - 1$) et les accroissements des facteurs intrinsèques ($u = S$) sont d'espérances nulles et spatialement non corrélés.

Les facteurs ainsi définis ne contiennent pas d'information redondante. A chaque facteur $Y_p^u(\mathbf{s})$ est associée une valeur propre λ_p^u de la matrice de corégionalisation \mathbf{B}_u correspondante, qui mesure son importance. Il faut remarquer que la matrice \mathbf{B}_S , correspond à une structure intrinsèque stricte et n'a pas le même sens que les autres; les facteurs non stationnaires doivent être privilégiés par rapport aux facteurs stationnaires.

2.3.3.2. Estimation des facteurs spatiaux $Y_p^u(\mathbf{s})$

Quel que soit le point $\mathbf{s} \in \mathcal{D}$, on estime la valeur des facteurs spatiaux $Y_p^u(\mathbf{s})$ par cokrigage à partir des mesures expérimentales. Pour chaque fonction Z_i , on utilise les sites voisins échantillonnés les plus proches de \mathbf{s} et l'estimateur du facteur $Y_p^u(\mathbf{s})$ en un point $\mathbf{s} \in \mathcal{D}$ est de la forme :

$$\hat{Y}_p^u(\mathbf{s}) = \sum_{i=1}^N \sum_{\alpha=1}^{n_i} Z_i(\mathbf{s}_i^\alpha).$$

Dans le cas où toutes les variables sont stationnaires, de moyennes m_i inconnues, on utilise les fonctions de covariance et le système de cokrigage s'écrit :

$$\begin{cases} \sum_{j=1}^N \sum_{\beta=1}^{n_j} \omega_\beta^j C_{ij}(\mathbf{s}_\alpha^i - \mathbf{s}_\beta^j) + \mu_i = a_{up}^i \rho_u(\mathbf{s}_\alpha^i - \mathbf{s}) & \text{pour } i = 1 \dots N, \alpha = 1 \dots n_i \\ \sum_{\beta=1}^{n_j} \omega_\beta^j = 0 & \text{pour } i = 1 \dots N \end{cases}$$

Dans le cas plus général où les Z_i vérifient l'hypothèse intrinsèque conjointe avec une seule structure strictement intrinsèque (on ne peut réaliser le cokrigage des facteurs dans le cas de plus d'une structure non stationnaire), il est également possible d'estimer par cokrigage les différents facteurs spatiaux. On obtient le système de cokrigage :

$$\begin{cases} \sum_{j=1}^N \sum_{\beta=1}^{n_j} \omega_{\beta}^j \gamma_{i_j} (\mathbf{s}_{\alpha}^i - \mathbf{s}_{\beta}^j) - \mu_i = a_{u_p}^i g_u (\mathbf{s}_{\alpha}^i - \mathbf{s}) \text{ pour } i = 1 \dots N, \alpha = 1 \dots n_i \\ \sum_{i=1}^N a_{u_q}^i \sum_{\beta=1}^{n_j} \omega_{\beta}^j = \begin{cases} 0 & \text{pour } u = 1 \dots S - 1 (Y_p^u \text{ stationnaire}) \\ \delta_p^q & \text{pour } u = S (Y_p^S \text{ intrinsèque strict)} \end{cases} \text{ pour } q = 1 \dots N \end{cases}$$

où δ_p^q est le symbole de Kronecker.

2.3.4. Conclusion sur l'analyse krigeante : atouts et limitations

Le lecteur aura pu noter des similitudes entre l'analyse krigeante et l'analyse en composantes principales (ACP). Nous allons les préciser. L'ACP a pour but la recherche des vecteurs propres de la matrice de variance-covariance des variables Z_i indépendamment de la position des sites $(\mathbf{s}_i)_{i=1 \dots N}$ d'observation dans l'espace géographique. Elle fournit des facteurs $Y_p(\mathbf{s})$ qui, s'ils ont la propriété d'être orthogonaux point à point :

$$\forall \mathbf{s} \in \mathcal{D}, \text{cov}[Y_p(\mathbf{s}), Y_q(\mathbf{s})] = 0 \text{ si } p \neq q$$

ne sont en général pas spatialement décorrélés (Wackernagel, 1998) :

$$\forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \text{cov}[Y_p(\mathbf{s} + \mathbf{h}), Y_q(\mathbf{s})] \neq 0 \text{ si } p \neq q.$$

En revanche, l'analyse krigeante, dans le cadre de l'hypothèse de stationnarité d'ordre deux, décompose la matrice de variance-covariance \mathbf{V} en une somme de S matrices de corégionalisation ($\mathbf{V} = \sum_{u=1}^S \mathbf{B}_u$) et les analyse séparément pour extraire de chacune d'elles des facteurs qui, eux, sont spatialement non corrélés :

$$\forall u, v = 1 \dots S, \forall \mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}, \text{cov}[Y_p^u(\mathbf{s} + \mathbf{h}), Y_q^v(\mathbf{s})] = 0 \text{ si } p \neq q \text{ ou } u \neq v.$$

En outre, l'analyse krigeante peut se réaliser sous l'hypothèse intrinsèque stricte, alors que la matrice de variance-covariance n'existe pas (les variances sont infinies). L'analyse krigeante apparaît donc plus riche que l'ACP car elle appréhende simultanément les caractéristiques structurales de la corégionalisation et le caractère multivarié et spatialisé des observations.

L'estimation par cokrigage des facteurs spatiaux en chaque point de l'espace \mathcal{D} va fournir les éléments pour établir leur cartographie.

Les cartes obtenues pour une structure particulière correspondent à une échelle déterminée, donnée par exemple par la portée de la fonction de covariance. On peut, par exemple, ne s'intéresser qu'aux variations du phénomène à grande distance et vouloir mettre en évidence les éventuelles dérives (« tendances ») dans la zone étudiée. On a la possibilité également de choisir le rang du facteur dans la structure. Comme en ACP, il est caractérisé par une valeur propre et par conséquent son taux d'inertie donne une mesure de la quantité d'information qu'il contient.

La cartographie d'un facteur spatial synthétise l'information multivariée et s'appuie sur la structure spatiale du phénomène en intégrant le contenu de l'information des points de données voisins. A partir des facteurs spatiaux les plus significatifs (ceux associés aux plus grandes valeurs propres), on peut appliquer des algorithmes de classification qui proposeront des zones homogènes.

Néanmoins, une des limitations de l'analyse krigéante repose sur le relatif arbitraire concernant le choix des structures (variogrammes ou covariances) que l'on identifie au cours de l'analyse variographique. En effet, les facteurs spatiaux dérivent directement de ces structures. A partir d'un ensemble de données l'ajustement d'un modèle est une opération subjective. Pour réduire cet arbitraire, il convient d'intégrer les informations qualitatives que l'expérimentateur est à même de connaître (genèse du phénomène, différents processus en synergie ou en opposition, fiabilité et erreurs des mesures).

3. Exemple sur des données de sol

Nous allons illustrer les concepts de l'analyse krigéante sur des données de sol. Il s'agit de l'analyse de 110 prélèvements de terre à trois niveaux (0-20 cm, 40-60 cm et 80-100 cm) faits à la tarière, répartis sur une grille régulière carrée de six mètres de côté. Pour étudier la variabilité à courte distance, trois croix de sondage ont été implantées avec des sites distants de un mètre. La figure 1 localise l'implantation des prélèvements dans la zone étudiée. La zone étudiée fait partie d'une plaine peu large polygénique parcourue d'un cours d'eau cadré de reliefs calcaires d'âge secondaire. Elle présente une topographie plane, faiblement inclinée vers la droite où se trouve le cours d'eau. Diverses variables pédologiques ont été mesurées mais nous ne retiendrons pour l'exemple que le pourcentage de terre fine (fraction de taille < 2 mm), le pourcentage d'argile + limon et l'humidité au point de flétrissement permanent (1.6 MPa), par convention, humidité retenue à une tension telle qu'elle est indisponible pour les plantes. La structure conjointe des trois variables étudiées a été ajustée selon le modèle linéaire de corégionalisation, à l'aide de trois structures de base : un effet de pépite, un schéma exponentiel et un modèle puissance. Les deux dernières structures présentent une anisotropie géométrique, dont les directions principales sont orientées de 40° et -50° par rapport à l'axe est-ouest. On trouve les différents paramètres dans le tableau 1, les variogrammes expérimentaux et leurs modèles dans la figure 2.

L'analyse krigéante est menée en trois étapes :

- La première étape consiste à estimer sur une grille régulière que l'on détermine *a priori* les composantes spatiales $X_i^u(s)$ associées à chaque variable ($i = 1, 2, 3$) et à chaque structure ($u = 1, 2, 3$).

- On calcule ensuite les matrices Q_u et Λ_u permettant de passer des composantes spatiales $X_i^u(s)$ aux facteurs spatiaux $Y_u^p(s)$. En règle générale, les variables

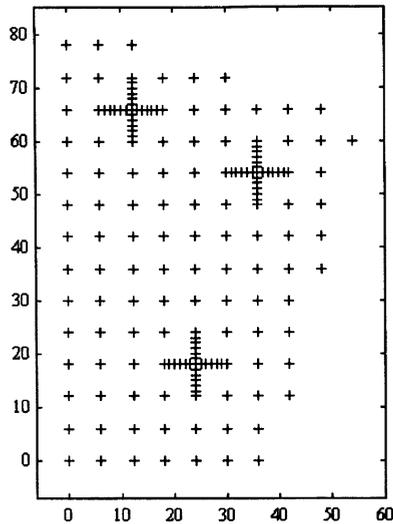


FIGURE 1

Carte d'implantation des sites échantillonnés (axes gradués en mètres)

TABLEAU 1

Matrices des corégionalisations pour les trois schémas respectifs

Structure	Matrice de corégionalisation
Structure 1 : Effet de pépité $\gamma(\mathbf{h}) = \begin{cases} 0 & \text{si } \mathbf{h} = 0 \\ 1 & \text{sinon} \end{cases}$	$\mathbf{B}_1 = \begin{pmatrix} 20.00 & 15.00 & 1.56 \\ 15.00 & 13.00 & 1.00 \\ 1.56 & 1.00 & 0.33 \end{pmatrix}$
Structure 2 : Schéma exponentiel $\gamma(\mathbf{h}) = 1 - \exp(-3 \mathbf{h} /a)$ portée pratique : $a = 40$ m (direction d'angle $+40^\circ$) $a = 25$ m (direction d'angle -50°)	$\mathbf{B}_2 = \begin{pmatrix} 5.00 & 10.00 & 1.50 \\ 10.00 & 20.00 & 3.00 \\ 1.50 & 3.00 & 0.62 \end{pmatrix}$
Structure 3 : Schéma puissance $\gamma(\mathbf{h}) = (\mathbf{h} /a)^{1.8}$ facteurs d'échelle : $a = 5$ (direction d'angle $+40^\circ$) $a = 1$ (direction d'angle -50°)	$\mathbf{B}_3 = \begin{pmatrix} 0.1100 & 0.0500 & -0.0160 \\ 0.0500 & 0.0240 & -0.0080 \\ -0.0160 & -0.0080 & 0.0032 \end{pmatrix}$

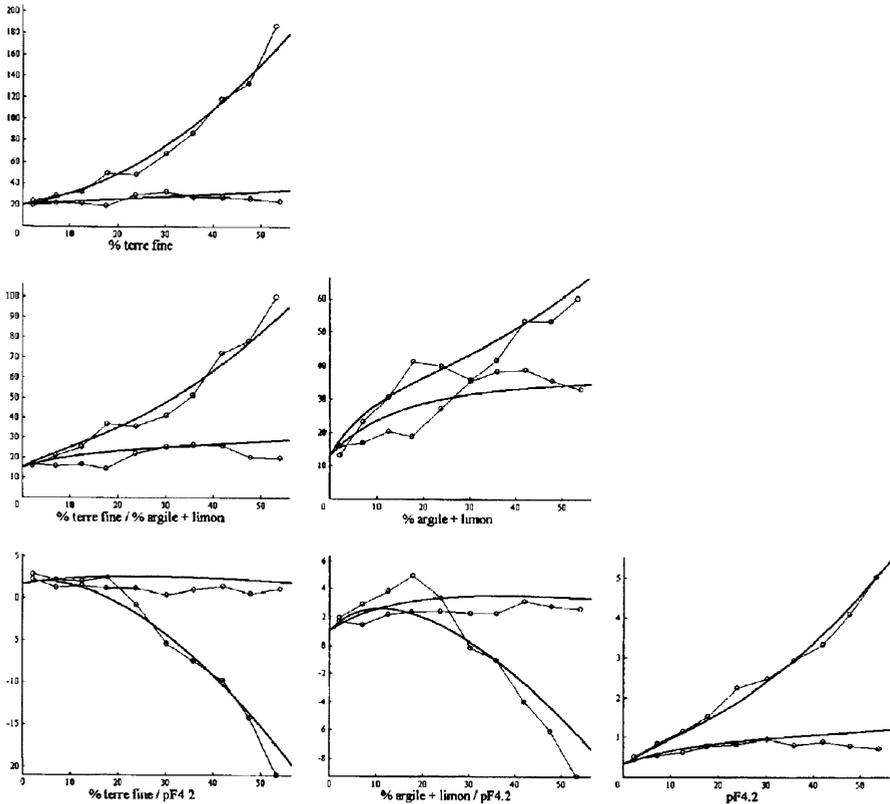


FIGURE 2

Variogrammes expérimentaux simples (diagonale) et croisés et modèles associés selon les deux directions d'anisotropie pour les trois variables analysées (en abscisse : la séparation en mètres; en ordonnée : les variogrammes simples et croisés)

sont standardisées lorsqu'elles sont d'unités différentes, c'est-à-dire qu'on les a divisées par leur écart-type. Si dans le cas stationnaire, les variances *a priori* sont finies, il n'en est pas de même en présence d'une structure strictement intrinsèque comme c'est le cas ici du schéma puissance. Les variances *a priori* sont alors infinies. Il convient donc de trouver d'autres coefficients de standardisation. Ils peuvent être fournis par les coefficients diagonaux de la matrice de corégionalisation B_3 relative au schéma puissance, qui indiquent le degré d'« infinité » des variances *a priori*. Cela revient en fait à effectuer l'analyse krigeante avec les variables non standardisés, mais en recherchant des facteurs spatiaux orthogonaux par rapport à la métrique

$$M = \begin{pmatrix} b_{11}^3 & 0 & 0 \\ 0 & b_{22}^3 & 0 \\ 0 & 0 & b_{33}^3 \end{pmatrix} = \begin{pmatrix} 0.1100 & 0 & 0 \\ 0 & 0.0240 & 0 \\ 0 & 0 & 0.0032 \end{pmatrix}$$

Les matrices Q_u et Λ_u cherchées sont alors solutions du problème de valeurs propres généralisé :

$$Q_u^t B_u Q_u = \Lambda_u \text{ avec } Q_u^t M Q_u = I.$$

Dans le tableau 2, on trouve les matrices de vecteurs et valeurs propres relatives à la diagonalisation des matrices de corégionalisation.

TABLEAU 2

Matrices des vecteurs et valeurs propres pour les trois schémas respectifs

Structure	Diagonalisation	
	Matrice des vecteurs propres	Matrice des valeurs propres
Structure 1 : effet de pépîte	$Q_1 = \begin{pmatrix} -1.45 & -2.61 & 0.42 \\ -5.49 & 2.72 & -2.04 \\ -3.83 & 4.76 & 16.59 \end{pmatrix}$	$\Lambda_1 = \begin{pmatrix} 735.7 & 0 & 0 \\ 0 & 14.1 & 0 \\ 0 & 0 & 76.8 \end{pmatrix}$
Structure 2 : Schéma exponentiel	$Q_2 = \begin{pmatrix} 2.94 & -0.63 & -0.27 \\ -1.47 & -5.79 & -2.44 \\ 0.00 & -6.87 & 16.29 \end{pmatrix}$	$\Lambda_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1027 & 0 \\ 0 & 0 & 45.5 \end{pmatrix}$
Structure 3 : Schéma puissance	$Q_3 = \begin{pmatrix} 1.74 & -1.76 & -1.72 \\ 3.81 & 5.05 & -1.30 \\ -9.98 & 3.90 & -14.1 \end{pmatrix}$	$\Lambda_3 = \begin{pmatrix} 2.827 & 0 & 0 \\ 0 & 0.017 & 0 \\ 0 & 0 & 0.156 \end{pmatrix}$

• La troisième étape est l'estimation des facteurs spatiaux $\hat{Y}_1^u(s)$ ($p = 1, 2, 3$) aux nœuds de la grille à partir de la relation

$$(\sqrt{\lambda_1^u} \hat{Y}_1^u(s), \sqrt{\lambda_2^u} \hat{Y}_2^u(s), \sqrt{\lambda_3^u} \hat{Y}_3^u(s)) = (\hat{X}_1^u(s), \hat{X}_2^u(s), \hat{X}_3^u(s)) \times Q_u$$

où les λ_i^u sont les éléments diagonaux de Λ_u .

Pour chaque structure, on obtient trois facteurs spatiaux, hiérarchisés selon les valeurs propres associées. Nous n'allons retenir que le premier facteur des structures relatives aux modèles exponentiel et puissance (les facteurs relatifs à l'effet de pépîte n'ont aucun intérêt car ils sont nuls partout sauf aux points d'observations). La figure 3 visualise ces facteurs.

Le premier facteur relatif au schéma puissance, qui correspond à une structure intrinsèque stricte, affiche une nette dérive dans la direction nord-ouest/sud-est. En quelque sorte, il prend en charge la non-stationnarité du phénomène. Quant au premier facteur du schéma exponentiel, il peut être vu comme un résumé de la structure sous-jacente, lorsque le phénomène est débarrassé de sa dérive. Il correspond, par construction, à une structure stationnaire dont la portée s'étend entre 25 et 40 mètres selon les deux directions principales d'anisotropie.

Il est intéressant d'interpréter les cartes obtenues en termes pédologiques pour s'assurer si elles confirment bien les connaissances qualitatives que l'on a sur la zone

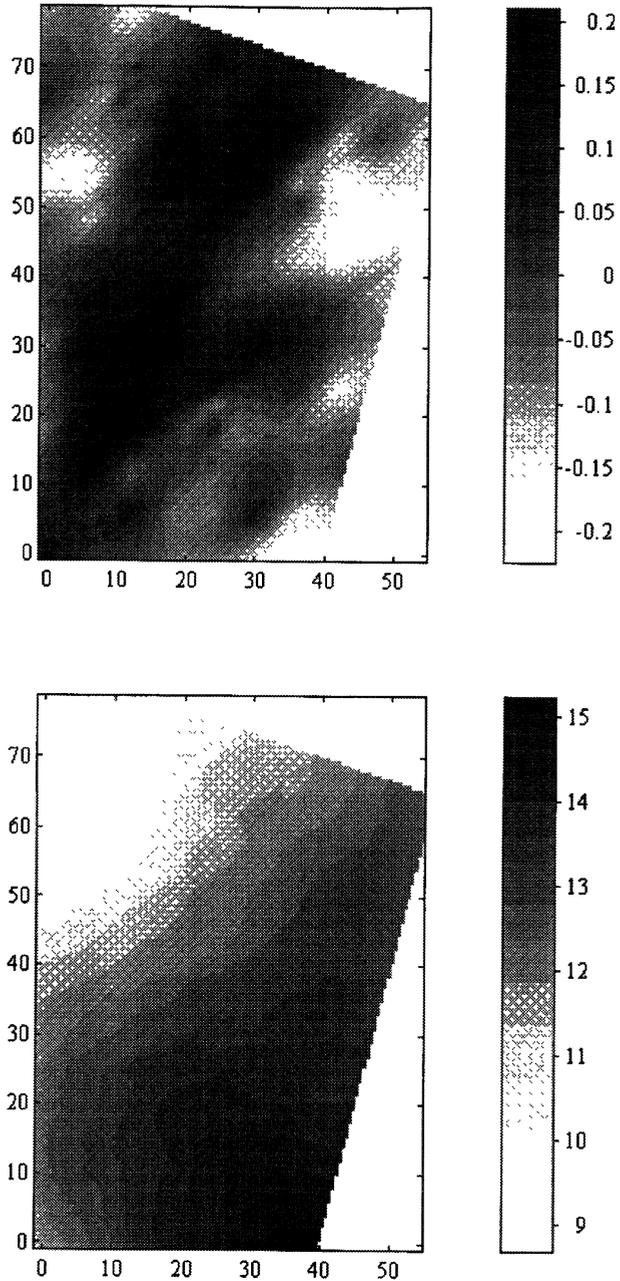


FIGURE 3

*Visualisation du premier facteur relatif au modèle exponentiel (en haut),
et au modèle puissance (en bas)
(axes gradués en mètres)*

étudiée ou si elles sont cohérentes avec d'autres informations que l'on a pu obtenir par d'autres méthodes. A cet effet, simultanément au travail d'échantillonnage et d'analyses au laboratoire, une carte pédologique sommaire a été établie sur cette même zone lors de la réalisation des sondages à partir d'observations qualitatives visuelles et tactiles comme la couleur, la texture de la terre fine de la terre extraite et de sa charge en éléments grossiers. Ce travail de terrain a permis d'identifier quatre unités pédologiques (figure 4), à savoir :

- colluvions (*C*) matériaux **assez argileux à charge grossière** pouvant être importante venant du versant adjacent (à gauche sur la carte) sur toute la profondeur du sondage (110 cm environ);
- colluvions sur tuf calcaire (*C/T*) dépôt lacustre à texture **limono-argileuse**, riche en **calcaire**, du pleistocène ancien;
- colluvions sur limon (*C/L*);
- limon (*L*) alluvions récentes essentiellement formées de **limons**.

On peut mettre en parallèle la carte pédologique avec celle du premier facteur de la structure exponentielle. Visuellement, on retrouve, sur cette dernière, les trois unités « colluvions sur tuf » (*C/T*), « colluvions » (*C*) et « limon » (*L*). Les cartes des facteurs sont donc en cohérence avec des études plus qualitatives et confirment le bien fondé de la technique.

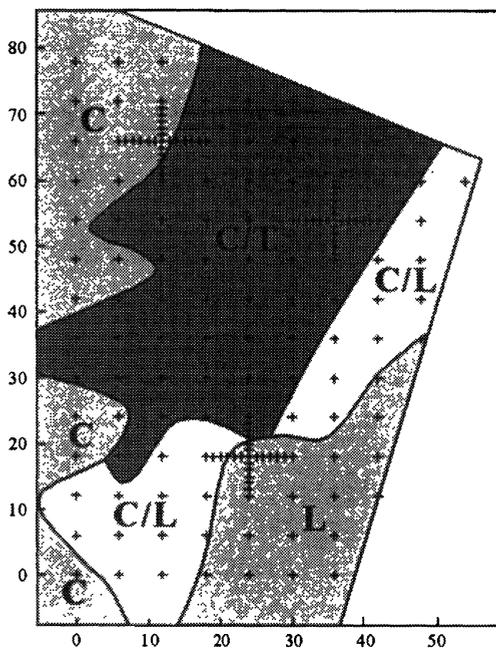


FIGURE 4
Carte pédologique
(axes gradués en mètres)

Conclusion

L'analyse krigeante a permis de prendre en compte le caractère spatial et multivarié des observations pour cartographier une zone géographique. Elle fournit des facteurs non corrélés spatialement et hiérarchisés en fonction de la quantité d'information qu'ils contiennent. Chaque facteur est associé à une échelle spatiale particulière. On peut ainsi soit choisir une échelle d'observation et cartographier les facteurs correspondants. On peut aussi combiner l'analyse des facteurs les plus significatifs, associés à des échelles différentes, et appliquer un algorithme de classification pour aboutir à une cartographie de la zone.

Bibliographie

- ALLARD, D. & MONESTIEZ, P. *Classification par critère de variance en contexte spatial*. XXXI^{èmes} Journées de Statistiques, 17-21 mai 1999, Grenoble, France, pp. 945-948.
- AMBOISE, C., DANG, M. & GOVAERT, G. *Clustering of spatial data by the EM algorithm*. In A. Soares *et al* (eds.), *geoENV I - Geostatistics for Environmental Applications*, Kluwer Academic Publishers, Dordrecht, 1997, pp. 493-504.
- ARNAUD, M. & PICHOT, J. P. *Poster : A methodology to build a classification of spatialised and multivariate data*. Third International Conference/Workshop on Integrating GIS and Environmental Modeling. January 21-25, 1996. Santa FE, New Mexico, USA.
- BREIMAN, L., FRIEDMAN, J. H., OLSHEN R. A., STONE, C. H., *Classification and Regression Trees*. 1984 The Wadsworth Statistics / Probabilities Series. Monterey California. 358 p.
- CAILLIEZ, F. & PAGES, J.-P. *Introduction à l'analyse des données*, 1976 (616p). SMASH.
- CHEssel, D. & SABATIER, R. *Couplage de triplets statistiques et graphes de voisinage*. *Biométrie et analyse de données spatio-temporelles* (B. Asselain et Coll., eds). Société Française de Biométrie, ENSA, Rennes, 1993, pp. 28-37.
- DIDAY, E. *Une nouvelle méthode en classification automatique et reconnaissance des formes. La méthode des nuées dynamiques*. *Analyse des données en architecture et urbanisme*. centre de méthodologie mathématique et informatique. Institut de l'environnement. Nov. 1972, pp. 97-111.
- GOULARD, M. & VOLTZ, M. *Linear Coregionalization Model : Tools for Estimation and Choice of Cross-Variogram Matrix*. *Mathematical Geology*, Vol 24, N° 3, 1992, pp. 269- 286.
- JAMBU, M. *Classification automatique pour l'analyse des données 1- méthodes et algorithmes*, Dunod, 1978, 310 p.
- LEBART, L. *Analyse statistique de la contiguïté*. Publication de l'ISUP, XVIII, 1969, pp. 81-112.

- LEBART, L. *Programme d'agrégation avec contraintes*. Les cahiers de l'analyse des données. Dunod. Vol III, 1978, pp. 275-287.
- MATHERON, G. *Traité de géostatistique appliquée*. 1962, Technip, Paris
- MATHERON G. *Pour une Analyse Krigeante des données Régionalisées*. Publication N-732, Centre de Géostatistique, Fontainebleau, France, 1982, 22 p.
- MÉOT, A., CHESSEL, D. & SABATIER, R. *Opérateurs de voisinage et analyse des données spatio-temporelles*. In : Biométrie et environnement (Lebreton, J.-D. & Asselin B. Eds), 1993, Masson, Paris. pp. 45-71.
- OLIVER, M. A. WEBSTER, R. *A geostatistical Basis for Spatial Weighting in Multivariate Classification*. Mathematical Geology, 1989, Vol. 21, N° 1.
- THIOULOUSE, J. CHESSEL, D. & CHAMPELY, S. *Multivariate analysis of spatial patterns : a unified approach to local and global structures*. Environmental and Ecological Statistics 2, pp. 1-14, 1995.
- WACKERNAGEL, H. Cours de géostatistique multivariable. (4^{ème} édit.), 1993, Ecole des Mines de Paris, Fontainebleau (France) [C-146], 80 p.
- WACKERNAGEL, H., *Multivariate Geostatistics : an introduction with applications*, Springer-Verlag, Berlin, 1998, 291 p.
- WACKERNAGEL, H., SANGUINETTI, H. Gold prospecting with factorial cokriging in Limousin, France. In J.C. Davis & U. Herzfeld, *Computers in Geology : 25 years of progress*, O.U.P., Oxford, 1993, 33-43 (Studies in Math. Geol. 5).