

# STATISTIQUE ET ANALYSE DES DONNÉES

J. L. MALLET

## **Proposition pour un coefficient de corrélation entre individus**

*Statistique et analyse des données*, tome 5, n° 1 (1980), p. 19-31

[http://www.numdam.org/item?id=SAD\\_1980\\_\\_5\\_1\\_19\\_0](http://www.numdam.org/item?id=SAD_1980__5_1_19_0)

© Association pour la statistique et ses utilisations, 1980, tous droits réservés.

L'accès aux archives de la revue « Statistique et analyse des données » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/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/>

PROPOSITION POUR UN COEFFICIENT DE CORRELATION ENTRE INDIVIDUS

par J.L. MALLET

Ecole Nationale Supérieure de Géologie de Nancy  
Centre de Recherches Pétrographiques et Géochimiques

1 - INTRODUCTION

1.1 - Définitions préliminaires

Métrique sur  $R^V$

Soit  $[C]$  une matrice carrée symétrique définie positive à  $v$  lignes et  $v$  colonnes donnée. Dans ce qui suit, nous supposerons que  $R^V$  est muni de la métrique  $\| \cdot \|$  telle que :

$$\| z \|^2 = z^t \cdot [C]^{-1} \cdot z \quad \forall z \in R^V$$

Fonction  $f(x, y)$

Soit  $\sigma(x)$  une fonction positive donnée définie sur  $R^V$  :

$$\sigma(x) > 0 \quad \forall x \in R^V$$

Soit d'autre part  $K$  la constante positive telle que :

$$K = \frac{1}{(2\pi)^{V/2} \cdot \{\det [C]\}^{1/2}}$$

Nous désignerons par  $f(x, \omega)$  la fonction positive définie sur  $(R^V \times R^V)$  telle que pour tout  $x \in R^V$  et tout  $\omega \in R^V$  on ait :

$$f(x, \omega) = \frac{K}{\sigma^V(x)} \cdot \exp \left\{ - \frac{1}{2 \cdot \sigma^2(x)} \cdot \| x - \omega \|^2 \right\}$$

On vérifie facilement que  $f(x, \omega)$  n'est autre que la densité d'une loi normale définie sur  $R^V$ , de moyenne  $x$  et de matrice de covariance  $[\Sigma]$  telle que :

$$[\Sigma] = \sigma^2(x) \cdot [C]$$

Il s'ensuit en particulier que l'on a :

$$\left[ \int_{R^V} \tilde{f}(x, \omega) \cdot d\omega = 1 \quad \forall x \in R^V \right.$$

*Mesure  $\mu$*

Soit S un ensemble de N points  $s^i$  appartenant à  $R^V$  :

$$\left[ \begin{array}{l} S = \{ s^1, \dots, s^i, \dots, s^N \} \\ \text{avec : } s^i \in R^V \quad \forall i \end{array} \right.$$

Par définition, nous désignerons par  $\mu$  la mesure définie sur la tribu borélienne  $B^V$  et admettant une densité  $f_\mu(\omega)$  par rapport à la mesure de Lebesgue telle que :

$$f_\mu(\omega) = \frac{1}{N} \cdot \sum_{s^i \in S} f(s^i, \omega)$$

Compte tenu que l'intégrale de  $f(s^i, \omega)$  par rapport à  $\omega$  sur le domaine  $R^V$  est égale à 1 pour tout  $s^i \in R^V$ , on est assuré que :

$$\int_{R^V} f_\mu(\omega) \cdot d\omega = 1$$

De plus, puisque  $\tilde{f}(s^i, \omega)$  est positif pour tout  $s^i \in R^V$  et tout  $\omega \in R^V$ , on est assuré que :

$$f_\mu(\omega) > 0 \quad \forall \omega \in R^V$$

On en conclut que  $\mu$  est une mesure positive normée sur  $B^V$  et par conséquent,  $\mu$  définit sur  $R^V$  une répartition continue de masse telle que :

$$\left\{ \begin{array}{l} \mu(B) = \text{masse de } B \quad \forall B \in B^V \\ \mu(R^V) = 1 \end{array} \right.$$

*Exemple*

Pour illustrer les définitions énoncées ci-dessus, supposons que  $V = 2$  et que S est constitué par les  $N = 41$  points de  $R^V$  représentés sur la figure 1. Si l'on se donne une constante  $\sigma > 0$  et si l'on pose ...

$$\left\{ \begin{array}{l} 1^\circ) [C] = \text{matrice des covariances des masses } 1/N \\ \quad \text{placées aux } N \text{ points } s^i \in S \\ 2^\circ) \sigma(x) = \sigma \quad \forall x \in R^V \end{array} \right.$$

... alors, suivant la valeur affectée à la constante  $\sigma$ , on obtient pour  $\mu$  des densités du type de celles représentées sur les figures 2, 3 et 4. On notera en particulier sur ces figures que  $\sigma$  apparaît comme un coefficient d'agrégation en ce sens que lorsque  $\sigma$  augmente, la densité  $f_\mu$  tend à devenir unimodale.

## 1.2 - Présentation du problème posé

*Remarque préliminaire*

Soit  $(\Omega, A, P)$  un espace probabilisé et soit  $\dot{X}_x$  une classe d'équivalence de fonction aléatoire (en abrégé CEFA) définie sur  $[(\Omega, A, P) \times \mathbb{R}^V]$  telle que :

$$\left. \begin{array}{l} \dot{X}_x \in L^2(\Omega, A, P) \\ \|\dot{X}_x\|_{L^2} = 1 \end{array} \right\} \forall x \in \mathbb{R}^V$$

On réalise ainsi une cartographie de  $\mathbb{R}^V$  sur une partie de la sphère unité  $S_{L^2}$  de  $L^2(\Omega, A, P)$  car, à tout point  $x \in \mathbb{R}^V$ , on peut associer le point image  $\dot{X}_x$  situé sur  $S_{L^2}$  puisque  $\|\dot{X}_x\| = 1$ . Si nous désignons par  $R(x^1, x^2)$  la fonction d'autocovariance de  $\dot{X}_x$  et par  $\langle \cdot | \cdot \rangle_{L^2}$  le produit scalaire dans  $L^2(\Omega, A, P)$ , alors on a :

$$R(x^1, x^2) = \langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2} = E(\dot{X}_{x^1} \cdot \dot{X}_{x^2}) \quad \forall \left\{ \begin{array}{l} x^1 \\ x^2 \end{array} \right\} \in \mathbb{R}^V$$

Comme le suggère la figure 5, il s'ensuit qu'un observateur situé à l'origine  $\dot{\theta}_{L^2}$  de  $L^2(\Omega, A, P)$  voit deux points  $(\dot{X}_{x^1}, \dot{X}_{x^2})$  images de  $(x^1, x^2)$  sous un angle  $\alpha(x^1, x^2)$  tel que :

$$\cos \{ \alpha(x^1, x^2) \} = R(x^1, x^2)$$

Quand  $x$  parcourt  $\mathbb{R}^V$ , son image  $\dot{X}_x$  parcourt une partie  $D_\mu$  de la sphère unité  $S_{L^2}$ . Si, de plus, on suppose que ...

$$R(x^1, x^2) \geq 0 \quad \forall \left\{ \begin{array}{l} x^1 \\ x^2 \end{array} \right\} \in \mathbb{R}^V$$

... alors, comme le suggère la figure 6, il existe dans  $L^2(\Omega, A, P)$  un cône  $C_\mu$  ayant pour sommet  $\dot{\theta}_{L^2}$  et possédant un angle au sommet égal à  $\pi/2$  tel que :

$$\dot{X}_x \in \{C_\mu \cap S_{L^2}\} \quad \forall x \in \mathbb{R}^V$$

*Objectif souhaité*

Considérons une population  $S$  de  $N$  individus de référence (étalons) plongés dans  $\mathbb{R}^V$  et engendrant une distribution de masse  $\mu$  suivant le procédé décrit au paragraphe 1.1. On souhaite associer à chaque point  $x \in \mathbb{R}^V$  une CEFA  $\dot{X}_x$  telle que  $R(x^1, x^2)$  appartienne à  $[0, 1]$  et que l'on ait :

- a)  $R(x^1, x^2) =$  fonction décroissante de  $\|x^1 - x^2\|$
- b)  $R(x^1, x^2)$  d'autant plus élevé que  $x^1$  et  $x^2$  appartiennent à un même mode (bosse) de  $f_\mu$

On notera que la condition (b) traduit la préoccupation communément admise qui veut que l'on interprète les modes de  $f_\mu$  comme des "indicatrices floues" des sous-populations à distinguer dans  $S$ . Pour quantifier la condition (b), nous conviendrons de la remplacer par la condition (b\*) suivante :

- b\*)  $R(x^1, x^2) =$  fonction croissante de  $f_\mu(x^{12})$  lorsque  $x^{12}$  est le barycentre de  $x^1$  et  $x^2$  affectés respectivement des masses  $1/\sigma^2(x^1)$  et  $1/\sigma^2(x^2)$

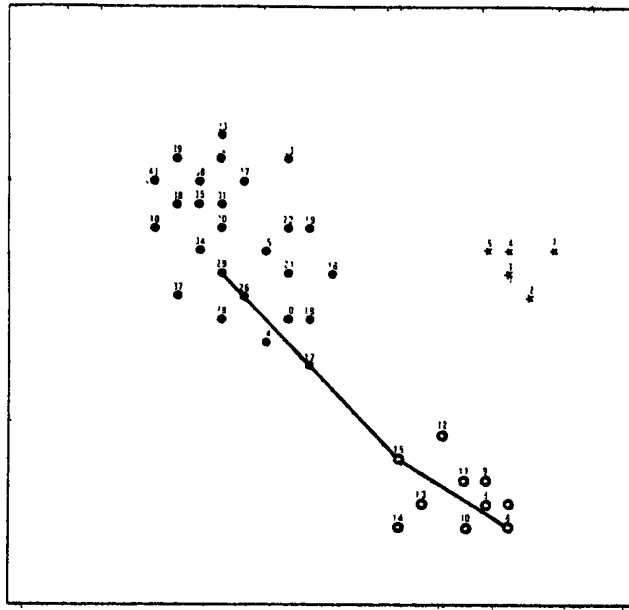


FIGURE 1

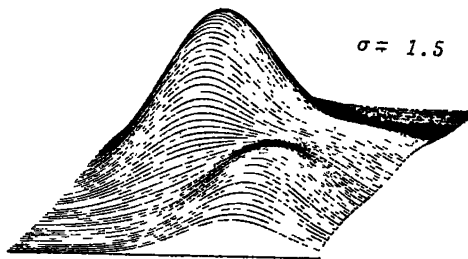


FIGURE 2

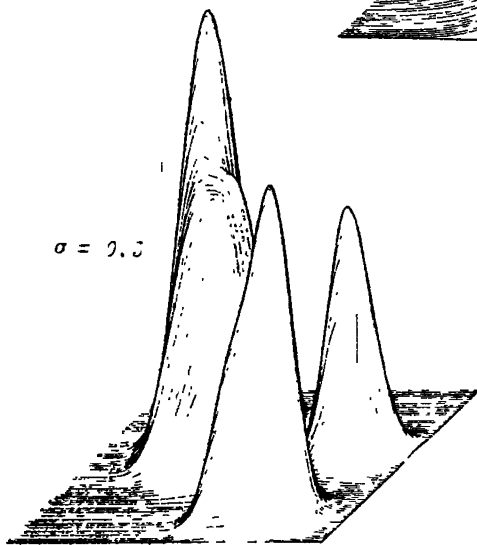


FIGURE 4

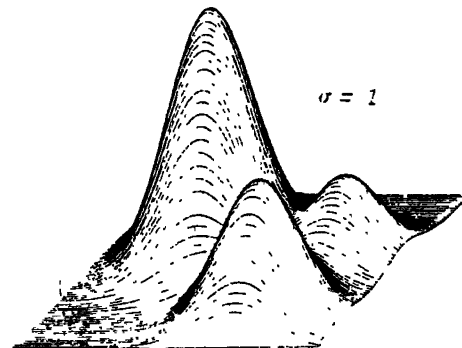


FIGURE 3

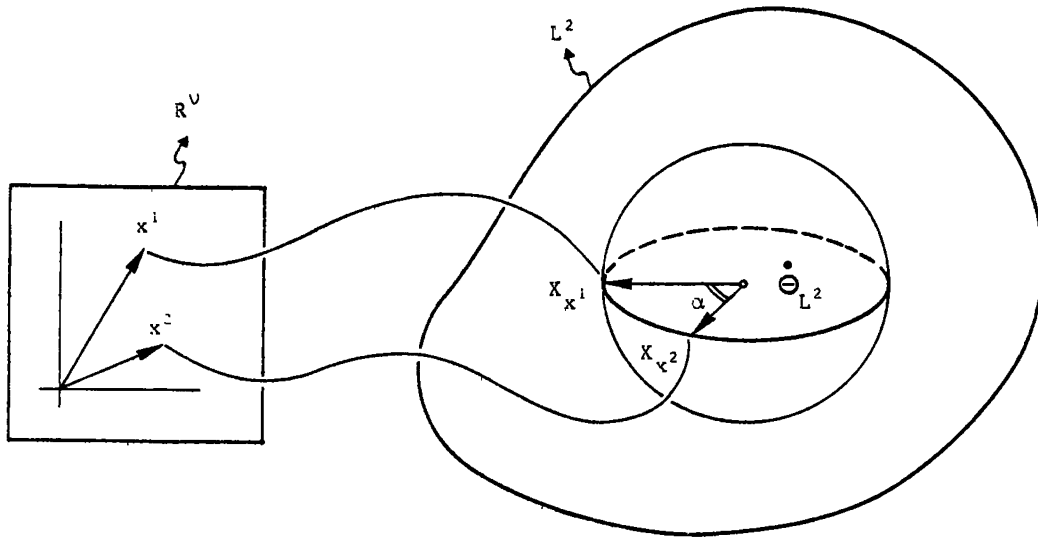


FIGURE 5

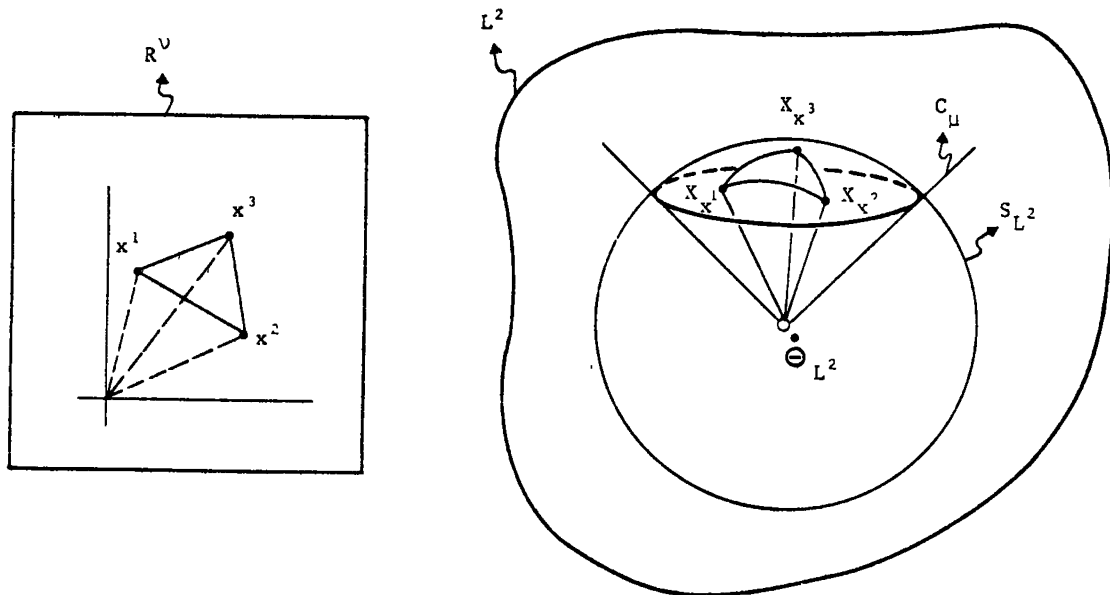


FIGURE 6

### Motivations

Le coefficient  $R(x^1, x^2)$  apparaît comme étant un indice de similarité possédant les propriétés d'un coefficient de corrélation. L'intérêt essentiel de ce coefficient est de permettre l'utilisation des méthodes d'analyse des variables aléatoires pour étudier les similarités entre points de  $R^V$  compte tenu de la population de référence  $S$ .

## 2 - ESPACE $(\Omega, A, P)$ ET CEFA $\dot{X}_x$ ASSOCIÉS A $\mu$

### 2.1 - Théorème fondamental

Compte tenu des définitions présentées au paragraphe 1.1, considérons l'espace probabilisé  $(\Omega, A, P)$  suivant :

$$\begin{cases} \Omega = R^V \\ A = \text{tribu borélienne de } \Omega \\ P = \mu \end{cases}$$

Soit d'autre part  $\gamma(x^1, x^2)$  la fonction définie sur  $(R^V \times R^V)$  telle que pour tout couple de points  $(x^1, x^2)$  appartenant à  $R^V$  on ait :

$$\left| \begin{aligned} \gamma(x^1, x^2) &= \frac{K^2}{N} \cdot \sum_{s^i \in S} \left( \frac{1}{a^{12} \cdot \sigma^2(s^i) + b^{12}} \right)^{V/2} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{a^{12}}{a^{12} \cdot \sigma^2(s^i) + b^{12}} \cdot \|x^{12} - s^i\|^2 \right\} \\ \text{avec : } \begin{cases} a^{12} &= \sigma^2(x^1) + \sigma^2(x^2) \\ b^{12} &= \sigma^2(x^1) \cdot \sigma^2(x^2) \\ x^{12} &= \text{barycentre de } x^1 \text{ et } x^2 \text{ affectés des masses } 1/\sigma^2(x^1) \text{ et } 1/\sigma^2(x^2) \end{cases} \end{aligned} \right.$$

Dans ces conditions, la fonction  $X_x$  définie de la façon suivante pour tout  $x \in R^V$  et tout  $\omega \in \Omega \dots$

$$X_x(\omega) = \frac{f(x, \omega)}{\sqrt{\gamma(x, x)}}$$

... est une fonction aléatoire. De plus, si l'on désigne par  $\dot{X}_x$  la classe d'équivalence de  $X_x$ , alors pour tous points  $x^1, x^2$  et  $x$  appartenant à  $R^V$ , on a :

$$\left| \begin{aligned} 1^\circ) \quad \dot{X}_x &\in L^2(\Omega, A, P) \\ 2^\circ) \quad \|\dot{X}_x\|_{L^2} &= 1 \\ 3^\circ) \quad \langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2} &= \frac{\gamma(x^1, x^2)}{\sqrt{\gamma(x^1, x^1) \cdot \gamma(x^2, x^2)}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{\|x^1 - x^2\|^2}{\sigma^2(x^1) + \sigma^2(x^2)} \right\} \end{aligned} \right.$$

### 2.2 - Corollaire

*Inoncé*

Le produit scalaire  $\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2}$  est tel que :

$$0 \leq \langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2} \leq 1 \quad \forall \left\{ \begin{matrix} x^1 \\ x^2 \end{matrix} \right\} \in \mathbb{R}^V$$

De plus, si l'on désigne par  $x^{12}$  le barycentre des points  $x^1$  et  $x^2$  affectés respectivement des masses  $1/\sigma^2(x^1)$  et  $1/\sigma^2(x^2)$  et si l'on désigne par  $B(s^i, \sqrt{2} \cdot \sigma(s^i))$  la boule de centre  $s^i \in S$  et de rayon  $\sqrt{2} \cdot \sigma(s^i)$  alors, toutes choses égales par ailleurs, le produit scalaire  $\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2}$  est d'autant plus élevé qu'il y a plus de points  $s^i \in S$  vérifiant la relation suivante :

$$x^{12} \in B(s^i, \sqrt{2} \cdot \sigma(s^i))$$

*Interprétation*

Ce corollaire exprime le fait que le produit scalaire  $\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2}$  est d'autant plus élevé qu'il y a plus de points  $s^i \in S$  entre  $x^1$  et  $x^2$  ou, plus précisément, au voisinage du barycentre  $x^{12}$  de ces deux points. Par ailleurs, compte tenu que  $\|\dot{X}_x\|_{L^2} = 1$  pour tout  $x \in \mathbb{R}^V$ , on est assuré que :

$$\|\dot{X}_{x^1} - \dot{X}_{x^2}\|^2 = 2 \cdot \left( 1 - \langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2} \right) \quad \forall \left\{ \begin{matrix} x^1 \\ x^2 \end{matrix} \right\} \in \mathbb{R}^V$$

Il s'ensuit que, toutes choses égales par ailleurs, la distance entre les points images  $\dot{X}_{x^1}$  et  $\dot{X}_{x^2}$  est d'autant plus petite qu'il y a plus de points  $s^i \in S$  entre  $x^1$  et  $x^2$  (et réciproquement). Compte tenu de la façon dont a été définie la mesure  $\mu$ , on peut également considérer que la distance  $\|\dot{X}_{x^1} - \dot{X}_{x^2}\|$  est d'autant plus faible que la densité  $f_\mu(x)$  est plus élevée entre  $x^1$  et  $x^2$ .

### 2.3 - Interprétation de $X_x$ et de $\mu$ et choix de $\sigma(x)$

*Interprétation de  $X_x$*

Soit  $x$  fixé dans  $\mathbb{R}^V$  ; si nous reprenons la définition

$$X_x(\omega) = \frac{f(x, \omega)}{\sqrt{\gamma(x, x)}} \quad \forall \omega \in \Omega$$

... alors on constate que  $\dot{X}_x$  est une gaussienne de moyenne  $x$ , d'écart type  $\sigma(x)$  et dont le maximum (atteint pour  $\omega = x$ ) vaut  $K / \{\sigma^V(x) \cdot \sqrt{\gamma(x, x)}\}^{\frac{1}{2}}$ . Si nous considérons maintenant deux points  $x^1$  et  $x^2$  fixés dans  $\mathbb{R}^V$ , alors comme le suggère la figure 7 (dans le cas où  $v = 1$ ), le produit scalaire  $\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2}$  défini par la relation

$$\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2} = \int_{\Omega} X_{x^1}(\omega) \cdot X_{x^2}(\omega) \cdot t_\mu(\omega) \cdot d\omega$$

... apparaît comme le produit scalaire des gaussiennes  $\dot{X}_{x^1}$  et  $\dot{X}_{x^2}$  pondérées par  $t_\mu$  ; on peut noter en passant que le seul rôle joué par  $\gamma(x, x)$  dans la définition de  $X_x$  est d'assurer que la condition  $\|\dot{X}_x\|_{L^2} = 1$  soit remplie pour tout  $x \in \mathbb{R}^V$ .



Interprétation de  $\mu$ 

Soit  $\nu$  une mesure positive normée définie sur  $B^V$ . Si l'on tire  $N$  points  $s^i$  au hasard dans  $R^V$  suivant la loi  $\nu$ , alors dès que  $N$  est "suffisamment" élevé, on constate expérimentalement que la mesure  $\mu$  définie au paragraphe 1.1 à partir de  $S = \{s^1, \dots, s^N\}$  constitue une approximation de la mesure  $\nu$  lorsque  $\sigma(x)$  est "bien choisie". En particulier, si  $\nu$  admet une densité  $f_\nu$  par rapport à la mesure de Lebesgue, alors dans ces conditions,  $f_\mu$  est une approximation de  $f_\nu$ .

Conséquence : proposition pour un "bon" choix de  $\sigma(x)$ 

Pour atteindre l'objectif souhaité présenté au paragraphe 1.2, le corollaire 2.2 ainsi que la figure 7 nous suggèrent de choisir une fonction  $\varphi$  monotone non décroissante sur  $[0, 1]$  et de poser :

$$\left[ \sigma(x) = \varphi \left( f_\mu(x) \right) \right]$$

Premier exemple de fonction  $\varphi$ 

L'exemple le plus simple de fonction  $\varphi$  est celui pour lequel, une constante positive  $\sigma^*$  étant donnée, on pose :

$$\varphi(s) = \sigma^* \quad \forall s \in [0, 1]$$

On obtient alors pour  $\sigma(x)$  la fonction telle que :

$$\sigma(x) = \sigma^* \quad \forall x \in R^V$$

L'expérience montre que, lorsque  $[C]$  est la matrice de covariance de  $S$ , le choix de  $\sigma^*$  entre 0.5 et 1.0 conduit en général à des résultats satisfaisants comme on pourra s'en convaincre en considérant les figures 1, 2, 3 et 4.

Le "pouvoir séparateur" des images  $\hat{x}_i$  est alors d'autant plus fort que  $\sigma^*$  est faible (et inversement).

Deuxième exemple de fonction  $\varphi$ 

Supposons que l'on se soit donné  $\sigma^* \in [0.5, 1]$  et que, par le procédé décrit au paragraphe 1.1, on ait ainsi obtenu une mesure  $\mu^*$  de densité  $f_{\mu^*}$  associée à :

$$\sigma(x) = \sigma^* \quad \forall x \in R^V$$

$$[C] = \text{matrice de covariance de } S$$

Comme le montre l'exemple présenté sur les figures 1, 2, 3 et 4, on constate que la forme de  $f_{\mu^*}(x)$  est relativement peu sensible au choix de  $\sigma^*$ , c'est pourquoi on peut songer en première approximation à chercher à déterminer  $\varphi$  sous forme d'une fonction positive croissante sur  $[0, 1]$  telle que :

$$\left[ \sigma(x) = \varphi \left( f_{\mu^*}(x) \right) \right]$$

Si l'on décide de choisir pour  $\varphi(s)$  une fonction linéaire de  $s$ , alors on peut par exemple poser :

$$\left[ \begin{array}{l} \sigma(x) = (\bar{\sigma} / \bar{f}_{\mu^*}) \cdot f_{\mu^*}(x) \\ \text{avec : } \left\{ \begin{array}{l} \bar{f}_{\mu^*} = \frac{1}{N} \sum_{s^1 \in S} f_{\mu^*}(s^1) \\ \bar{\sigma} = \text{constante positive donnée} \end{array} \right. \end{array} \right.$$

L'expérience montre que le choix de  $\bar{\sigma}$  entre 0.5 et 1.0 conduit en général à des résultats satisfaisants.

### 2.4 remarque

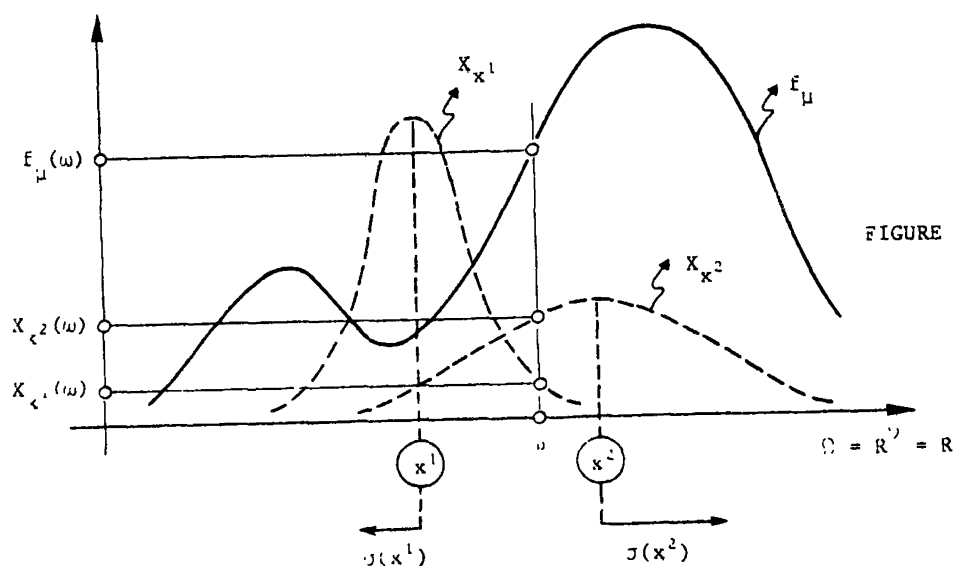
Compte tenu de ce qui vient d'être dit au paragraphe 2, on peut considérer que le coefficient ...

$$R(x^1, x^2) = \langle \dot{X}_{x^1} \mid \dot{X}_{x^2} \rangle_{L^2}$$

... est un indice de similarité entre  $x^1 \in R^V$  et  $x^2 \in R^V$  possédant des propriétés intéressantes du point de vue de la classification. En effet, compte tenu du corollaire 2.1 et compte tenu de la technique de choix de  $\sigma(x)$  présentée ci-dessus, on vérifie facilement que les propriétés souhaitées pour  $R(x^1, x^2)$  au paragraphe 1.2 sont effectivement vérifiées.

Ceci étant dit, bien que  $\dot{X}_{x^1}$  et  $\dot{X}_{x^2}$  ne soient (par construction) jamais centrées, par abus de langage, dans l'étude des fonctions aléatoires on a souvent coutume de dire que  $R(x^1, x^2)$  est le "coefficient d'autocorrélation" entre  $\dot{X}_{x^1}$  et  $\dot{X}_{x^2}$ . Par un deuxième abus de langage, nous proposons de dire que  $R(x^1, x^2)$  est le "coefficient de corrélation" entre  $x^1$  et  $x^2$ , cette dernière appellation étant justifiée par le fait que :

$$\left\{ \begin{array}{l} 0 \leq R(x^1, x^2) \leq 1 \\ R(x^1, x^2) = 1 \iff x^1 = x^2 \\ R(x^1, x^2) = 0 \iff \|x^1 - x^2\| = \infty \end{array} \right.$$



## 3 - APPLICATIONS

3.1 - Analyse d'une partie finie T de  $R^V$ 

Définition de T et de la CEVA  $\dot{X}_T$  associée

Soit T une partie finie de  $R^V$  (éventuellement on peut avoir  $T \equiv S$ ) constituée de n points  $t^i \in R^V$  :

$$T = \{ t^1, \dots, t^i, \dots, t^n \}$$

Soit d'autre part  $\mu$  une mesure positive normée du type de celle définie au paragraphe 1.1 et soit  $\dot{X}_x$  la CEVA associée définie au paragraphe 2.1.

Par définition, nous désignerons par  $\dot{X}_T$  la classe d'équivalence de variable aléatoire (en abrégé CEVA) vectorielle telle que :

$$\dot{X}_T = \begin{bmatrix} \dot{X}_{t^1} \\ \vdots \\ \dot{X}_{t^n} \end{bmatrix}$$

Matrice  $E_{X_{T_1}, X_{T_2}}$

Pour tout couple  $(T_1, T_2)$  de parties finies de  $R^V$ , nous poserons :

$$E_{X_{T_1}, X_{T_2}} = E(\dot{X}_{T_1}, \dot{X}_{T_2}^t)$$

Si l'on désigne par  $R(x^1, x^2) = \langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle_{L^2}$ , la fonction d'autocorrélation de  $\dot{X}_x$ , alors le terme  $(E_{X_{T_1}, X_{T_2}})_{ij}$  situé sur la ième ligne et la jème colonne de  $E_{X_{T_1}, X_{T_2}}$  est tel que :

$$(E_{X_{T_1}, X_{T_2}})_{ij} = R(t_1^i, t_2^j) \text{ avec : } \begin{cases} t_1^i \in T_1 \\ t_2^j \in T_2 \end{cases}$$

Compte tenu de la formule définissant  $\langle \dot{X}_{x^1} | \dot{X}_{x^2} \rangle$  au paragraphe 2.1, on peut montrer que  $E_{X_{T_1}, X_{T_2}}$  est toujours inversible si les points  $t^i \in T$  sont tous distincts.

Conséquence

En conséquence de tout ce qui vient d'être dit, il est tout à fait naturel de songer à appliquer aux CEVA  $\dot{X}_T$  associés à des parties finies de  $R^V$  les techniques classiques d'analyse comme par exemple :

- l'analyse de régression de  $\dot{X}_{t^0}$  en fonction de  $\dot{X}_T$
- l'analyse en composantes principales de  $\dot{X}_T$
- l'analyse canonique du couple  $(\dot{X}_{T_1}, \dot{X}_{T_2})$
- etc.

De même, on peut définir des coefficients analogues aux :

- coefficient de corrélation multiple  $r(X_{t^0} | X_T)$  entre  $\dot{X}_{t^0}$  et  $\dot{X}_T$

- coefficient RV ( $\hat{X}_{T_1}, \hat{X}_{T_2}$ ) d'Escouffier entre  $\hat{X}_{T_1}$  et  $\hat{X}_{T_2}$
- etc.

Compte tenu que  $E(\hat{X}_x)$  ne présente pas un intérêt particulier dans toutes ces techniques, il est recommandé d'agir formellement comme si :

$$E(\hat{X}_x) = 0 \quad \forall x \in R^V$$

$$R(x_1, x_2) = \text{coefficient de corrélation entre } \hat{X}_{x_1} \text{ et } \hat{X}_{x_2}$$

Par exemple, on obtient ainsi les formules utiles suivantes :

$$r(X_{t^0} | X_T) = \sqrt{\frac{E_{X_{t^0} X_T} \cdot E_{X_T X_T}^{-1} \cdot E_{X_T X_{t^0}}}{n_1 \cdot n_2}} \quad \forall t^0 \in R^V$$

$$RV^2(X_{T_1}, X_{T_2}) = \frac{\text{trace} \left\{ \begin{matrix} E_{X_{T_1} X_{T_2}} & E_{X_{T_2} X_{T_1}} \\ E_{X_{T_2} X_{T_1}} & E_{X_{T_1} X_{T_2}} \end{matrix} \right\}}{n_1 \cdot n_2} \quad \text{avec : } \left\{ \begin{matrix} n_1 = \text{card}(T_1) \\ n_2 = \text{card}(T_2) \end{matrix} \right.$$

Notons en passant que ces deux formules nous fournissent respectivement un indice de similarité d'une part entre un point  $t^0 \in R^V$  et une partie finie  $T$  de  $R^V$  et d'autre part, entre deux parties finies  $T_1$  et  $T_2$  de  $R^V$  ; on peut d'ailleurs remarquer que ces deux indices sont cohérents avec l'indice  $R(t^1, t^2)$  en ce sens que :

$$\left. \begin{matrix} t^1 = t^0 = T_1 \\ t^2 = T = T_2 \end{matrix} \right\} \Rightarrow R(t^1, t^2) = r(X_{t^0} | X_T) = RV(X_{T_1}, X_{T_2})$$

### 3.2 - Etude d'un exemple particulier

#### *Introduction*

Nous nous proposons de reprendre dans ce qui suit, l'exemple présenté sur les figures 1, 2, 3 et 4, en supposant que  $V = 2$  et que :

$$\left\{ \begin{array}{l} S = \text{ensemble des } N = 41 \text{ points représentés sur la figure 1} \\ \sigma(x) = \sigma^* \quad \forall x \in R^V \\ [C] = \text{matrice des covariances des masses } 1/N \text{ placées aux points } s^i \in S \end{array} \right.$$

Par ailleurs, en nous reportant à la figure 1, nous poserons :

$$\left\{ \begin{array}{l} T_1 = \text{groupe des } n_1 = 5 \text{ étoiles noires} \\ T_2 = \text{groupe des } n_2 = 10 \text{ étoiles blanches cerclées de noir} \\ T_3 = \text{groupe des } n_3 = 26 \text{ points noirs} \end{array} \right.$$

Matrice  $E_{X_T X_T}$

Suivant la valeur du paramètre d'agrégation  $\sigma^*$ , la matrice des "coefficients de corrélation"  $E_{X_T X_T}$  varie de façon sensible.

En particulier, pour  $\sigma^* = 0,5$ , cette matrice est quasiment partitionnée de la façon suggérée par la figure 8.

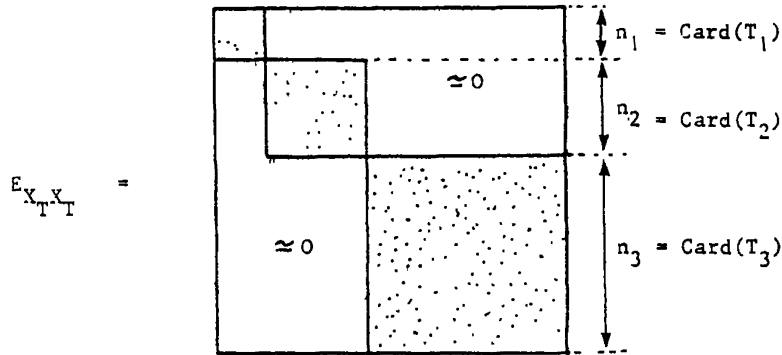


FIGURE 8

Pour se faire une idée de la variation de ces coefficients de corrélation en fonction de  $\sigma^*$ , on peut considérer le tableau représenté sur la figure 9.

$\sigma^*$	$R(t^6, t^{15})$	$R(t^{15}, t^{17})$	$R(t^{17}, t^{29})$
0,50	0,407	0,246	0,383
0,75	0,654	0,582	0,652
1,00	0,781	0,767	0,790

Fig. 9

*Effet de "chaînage-répulsion"*

D'après le corollaire 2.2, on peut s'attendre à observer un effet de "chaînage" des images  $(\dot{X}_{t^\alpha}, \dot{X}_{t^\beta})$  des couples de points  $(t^\alpha, t^\beta)$  appartenant à T par des points  $s^i \in S$  situés au voisinage du barycentre  $t^{\alpha\beta}$  de  $t^\alpha$  et  $t^\beta$ , ou au contraire une "répulsion" de  $\dot{X}_{t^\alpha}$  et  $\dot{X}_{t^\beta}$  s'il n'y a pas (ou peu) de points  $s^i \in S$  au voisinage de  $t^{\alpha\beta}$ . En d'autres termes, nous devons avoir :

$$\left\{ \begin{array}{l} R(t^\alpha, t^\beta) \text{ élevé s'il existe des points } s^i \in S \text{ au voisinage de } t^{\alpha\beta} \\ R(t^\alpha, t^\beta) \text{ faible s'il n'existe pas de points } s^i \in S \text{ au voisinage de } t^{\alpha\beta} \end{array} \right.$$

On pourra se faire une idée de l'amplitude de ce phénomène en se reportant au tableau représenté sur la figure 9 et en remarquant que les points  $\{t^6, t^{15}, t^{17}, t^{29}\}$  sont quasiment alignés et équidistants

$$|| t^6 - t^{15} || = || t^{15} - t^{17} || = || t^{17} - t^{29} ||$$

S'il n'y avait pas d'effet de chaînage-répulsion, on devrait donc avoir  $R(t^6, t^{15}) = R(t^{15}, t^{17}) = R(t^{17}, t^{29})$ . On notera en particulier sur la figure 9 que cet effet de chaînage-répulsion est d'autant plus sensible que le paramètre d'agrégation  $\sigma^*$  est plus faible.

## 4 - BIBLIOGRAPHIE

- CACOULOS, T. (1966). - Estimation of a multivariate density. Annals of the Inst. of Statistical Mathematics, Japan, vol. 18, p. 179-189.
- PARZEN, E. (1962). - On estimation of a probability density function and mode. Ann. Math. Statistics, vol. 33, p. 1065-1076.
- ROSENBLATT, M. (1971). - Curve estimates. The annals of Mathematical Statistics, vol. 42, p. 1815-1842.