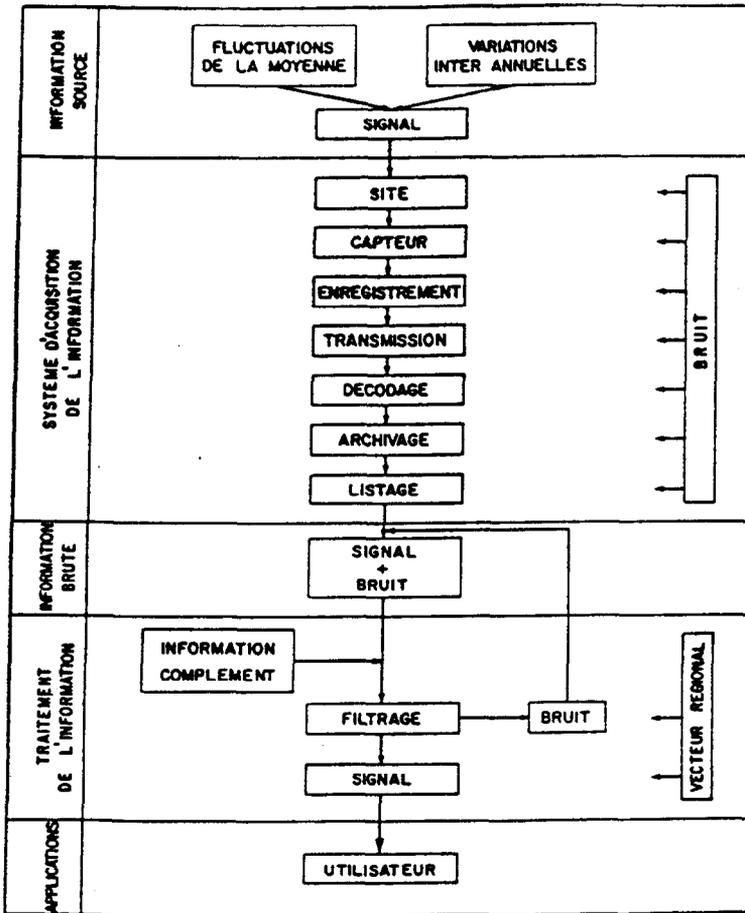


LES SUPPORTS THEORIQUES DU VECTEUR REGIONAL

(Première communication)

DU "SIGNAL" A LA MESURE VRAIE

Les manifestations ponctuelles d'une grandeur physique (par exemple la pluviométrie annuelle à une station) constituent des "signaux" dont les valeurs, pour une année donnée, fluctuent autour d'une moyenne régionale, elle même sujette à des variations interannuelles. En chaque point, chaque année (ou chaque mois), la grandeur physique a une valeur, le signal.



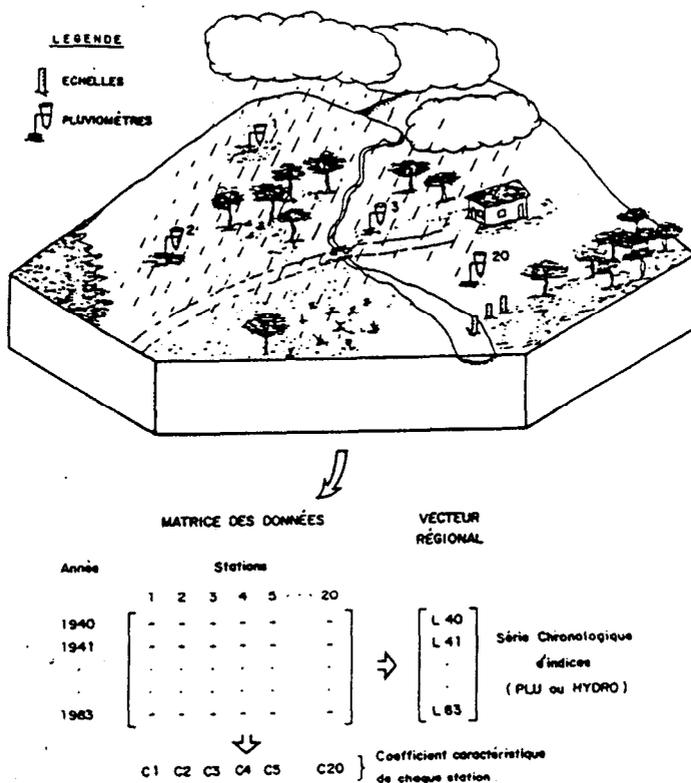
A. S. I. - BUREAU D'HYDROMETRIE - BRUXELLES - 14/17 SEPT 66

DRA. - BELLEGH S. L. 67000

Mais le choix du site et du capteur, le mode d'enregistrement et de transmission, le système de décodage et d'archivage, le type de publication, introduisent chacun des "bruits" qui se superposent au "signal" et l'altèrent.

L'information brute est constituée du couple "signal + bruit" et le traitement de l'information va ajouter à cette information brute toute l'information complémentaire disponible, de toute origine, de sorte d'obtenir après filtrage le "signal" original de la grandeur physique, qui sera fourni à l'utilisateur. Le Vecteur Régional est l'un de ces outils qui permettent de traiter la donnée altérée afin d'en restituer la vérité.

REPRESENTATION NAÏVE DU CONCEPT DE VECTEUR REGIONAL



Le schéma précédent explicite comment une chronique d'observation des données 1940 à 1983, portant sur les 20 stations d'une région homogène, conduit à la matrice des données dont on va chercher à extraire le Vecteur Régional et la liste des coefficients caractéristiques de chaque station.

UNE DEFINITION ET QUELQUES PRINCIPES :

G. HIEZ a ainsi défini son Vecteur Régional :

"Une série chronologique synthétique d'indices pluviométriques ou hydrométriques annuels ou mensuels, provenant de l'extraction de l'information la plus probable contenue dans un ensemble de postes d'observation groupés par région".

Nous rappelons maintenant les deux principes de base, et leur signification mathématique :

- 1) *"Les totaux annuels ou mensuels de stations "voisines" vérifient la règle de la pseudo-proportionnalité.*

Si l'on se place dans le mode de représentation graphique, type "double-cumul", chaque année (i) est représentée par un point de coordonnées X_i et Y_i , correspondant à la valeur de la grandeur physique annuelle (i) considérée (par exemple P_i pour la pluviométrie annuelle). Chaque vecteur élémentaire P_i a alors pour coordonnées $x_i = X_i - X_{i-1}$ et $y_i = Y_i - Y_{i-1}$.

Ecrire qu'il y a pseudo-proportionnalité entre stations voisines revient à dire :

$$\frac{X_i}{Y_i} \approx \frac{X_{i-1}}{Y_{i-1}} \quad \text{ou} \quad \frac{X_i}{X_{i-1}} \approx \frac{Y_i}{Y_{i-1}}$$

- 2) *"L'information la plus probable est celle qui se répète le plus fréquemment"*

Ce principe très oecuménique, est le fondement de toute approche statistique, qu'elle soit pragmatique et événementielle, ou mathématique.

Les principes complémentaire témoignent plus du pragmatisme de l'approche retenue qu'ils ne constituent des principes physiques additionnels. G. HIEZ les énonce simultanément de la façon suivante, malgré les différences de leurs niveaux conceptuels respectifs :

- *"dans un même groupement régional de stations, il n'y a pas de variation sensible des tendances climatiques"*
- *"aucune hypothèse n'est faite sur la distribution statistique des données"*
- *"toute mesure, par nature ponctuelle dans l'espace et dans le temps, est sujette à erreur"*
- *"l'information globale, fournie par un ensemble de stations, contient une valeur estimative des fluctuations plus représentatives que celle délivrée isolément par l'une quelconque des stations de l'ensemble".*

En conséquence, la "philosophie" du Vecteur Régional pourrait être résumée par cette intention:

"Toute l'information contenue dans chacune des stations doit contribuer à l'élaboration du vecteur régional, sans que les données erronées ne puissent avoir une influence sensible sur le résultat".

LE TABLEAU DES DONNEES CONSIDERE COMME UNE MATRICE

Les données observées, les données théoriques ou données réputées "vraies" et les erreurs ou encore les anomalies vont être considérées sous une forme matricielle à n lignes (n années d'observation) et à m colonnes (m stations d'observation) :

$$[A] = [B] + [E]$$

ou $[A]$ = matrice des données observées

$[B]$ = matrice des données théoriques ou réputées vraies

$[E]$ = matrice des erreurs ou anomalies

$[A]$ peut être explicitée, A_{ij} étant par exemple le total annuel pluviométrique observé de l'année i à la station j :

Dans cette approche les m données observées pendant une année (i) représentent les coordonnées d'un vecteur pluie annuel dans l'espace des stations.

$[A] =$

	Station 1	Station 2	...	Station j	...	Station m
Année 1	A_{11}	A_{12}	...	A_{1j}	...	A_{1m}
Année 2	A_{21}	A_{22}	...	A_{2j}	...	A_{2m}
...
Année i	A_{i1}	A_{i2}	...	A_{ij}	...	A_{im}
...
Année n	A_{n1}	A_{n2}	...	A_{nj}	...	A_{nm}

L'application du principe de "pseudo-proportionnalité" permet de définir la matrice $[B]$ par ses éléments $B_{ij} = L_i \cdot C_j$ où L_i est l'indice chronologique de l'année (i) et C_j le coefficient de la station (j) :

$$[B] = [L] \times [C]$$

soit en l'explicitant :

$$\begin{array}{c}
 \left[\begin{array}{c} L_1 \\ L_2 \\ \vdots \\ L_i \\ \vdots \\ L_n \end{array} \right] \\
 \leftarrow \text{VECTEUR COLONNE} \\
 \text{D'INDICES CHRONOLOGIQUES}
 \end{array}
 * \left[\begin{array}{cccc} C_1 & C_2 & \dots & C_j & \dots & C_m \end{array} \right] = \left[\begin{array}{cccc} L_1 * C_1 & \dots & L_1 * C_j & \dots & L_1 * C_m \\ \vdots & & \vdots & & \vdots \\ L_i * C_1 & \dots & L_i * C_j & \dots & L_i * C_m \\ \vdots & & \vdots & & \vdots \\ L_n * C_1 & \dots & L_n * C_j & \dots & L_n * C_m \end{array} \right]$$

VECTEUR-LIGNE
DES COEFFICIENTS DE STATION

La matrice des résidus, [E], sera définie alors comme la différence des matrices [A] et [B] :

$$[E] = [A] - [B]$$

QUE MINIMISER, ET COMMENT, POUR DETERMINER [L] ET [C]

L'approche de la "pseudo-proportionnalité" consiste ainsi à déclarer que l'on reconnaît dans la matrice [A] un ensemble de n vecteurs pluies (i) presque colinéaires. La matrice [B] est une matrice de rang 1, en quelque sorte dégénérée de [A], entièrement définie comme le produit du vecteur ligne [C] des coefficients de station, et du vecteur colonne [L] des indices chronologiques. [B] contient donc toute l'information à caractère linéaire qu'il est possible d'extraire de [A]. La matrice [E] représente le "bruit de fond", dont il était question précédemment, qu'il provienne des fluctuations aléatoires propres à chaque poste, d'origine climatique, ou provoquées par de simples erreurs d'observation. Dans son article de base, G. HIEZ (1977) précise que la "tentation est grande d'imaginer que les éléments de[E] sont linéairement indépendants entre eux", c'est à dire que $[E]^T \cdot [E]$ est une matrice diagonale. Céder à cette tentation revient à chercher à minimiser $\|[A]-[B]\|^2$, c'est à dire engager un processus classique de régression par les moindres carrés, ou encore la recherche des directions principales, donc des vecteurs propres de la matrice des données. La technique des moindres carrés, conséquence de la minimisation de la norme euclidienne, implique la condition d'homocédasticité des résidus qui n'est pas retenue dans les principes complémentaires, afin de conserver une plus grande généralité. Elle suppose aussi que les données ont une distribution gaussienne, ce qui est contraire à une hypothèse initiale. Enfin on sait que cette démarche accorde un poids trop important aux données extrêmes et de ce fait ne permet pas d'extraire des données toute l'information "linéaire". G. HIEZ justifie les réticences qui l'ont conduit aux hypothèses précitées en écrivant :

"Le problème que nous cherchons à résoudre ne consiste pas à "minimiser" les erreurs, et encore moins les erreurs fortes, au profit d'une réduction - somme toute arbitraire - de la variance expliquée, mais bien plutôt à les détecter telles qu'elles existent".

A LA RECHERCHE DE LA NORME PERDUE

En fait on cherche à représenter par un paramètre la tendance centrale de la matrice des données, c'est à dire à minimiser $[E] = [A] - [B]$. Classiquement on procède par l'utilisation d'une norme N_p d'ordre p que l'on minore, de sorte de définir un paramètre m :

$$N_p = \left[\sum (X_i - m)^p \right]^{1/p}$$

où X_i sont les données, prenant toutes les valeurs de $i = 1$ à n ou m .

- si $p = 2$, $\sum(X_i - m)^2$ est minimum lorsque m est la moyenne des X_i ; il s'agit du procédé de minimisation de la distance euclidienne.

De la bibliographie étudiée par G. HIEZ, il résulte que ce procédé n'est applicable que tant que la distribution des X_i est normale. Il n'est notamment pas performant en cas de distribution asymétrique ou encore plurimodale des X_i , voire si l'anomalie d'un groupe de valeurs extrêmes "contamine" la moyenne. Toujours d'après l'étude bibliographique de G. HIEZ, l'écart type est aussi impropre à décrire la dispersion d'une variable autour d'une valeur centrale qui ne serait pas distribuée normalement.

- si $p = 1$, $\sum|X_i - m|$ est minimum lorsque m est la médiane des X_i ; il s'agit d'un procédé de minimisation d'une distance en valeur absolue. Mais G. HIEZ reconnaît que l'estimation de [B] par la recherche des valeurs médianes est délicate, surtout lorsque les séries comportent de nombreuses lacunes.
- G. HIEZ introduit alors le mode qui correspond à la norme d'ordre 0, c'est à dire à la minoration $\sum(X_i - m)^0$, obtenue si m est le mode, ce qui suppose :

$$(X_i - m)^0 = 0, \quad \text{pour } X_i = m$$

$$(X_i - m)^0 = 1, \quad \text{pour } X_i \neq m$$

On peut constater que l'utilisation de la moyenne ou de la médiane pour minimiser $\sum(X_i - m)^2$ ou $\sum|X_i - m|$, équivaut à réduire quantitativement la masse des erreurs. G. HIEZ remarque que l'utilisation du mode vise à réduire numériquement les erreurs et à utiliser une méthode de maximum de vraisemblance dont la solution réside dans la recherche des valeurs modales de la distribution des X_i .

La minimisation de la matrice des erreurs $[E] = [A] - [B]$ peut s'écrire encore $[E] = [A] - [L].[C]$.

L'élément e_{ij} de $[E]$ est alors défini par :

$$e_{ij} = a_{ij} - l_i \cdot c_j$$

puisque dans la définition de la matrice des valeurs vraies [B], [L] et [C] sont respectivement une matrice colonne et une matrice ligne.

L'utilisation de la recherche des valeurs modales comme procédé de minimisation de [E] définissant $[B] = [L] \cdot [C]$ revient à dire que l_i et c_j seront correctement estimés quand la valeur de e_{ij} , la plus fréquente - donc la plus probable ou la plus vraisemblable - dans [E], sera la valeur 0.

Mais il paraît plus intéressant d'opérer en erreur relative, c'est à dire de minimiser les éléments d'une matrice [E] définis comme

$$e_{ij} = \frac{a_{ij}}{l_i \cdot c_j} - 1$$

pour déterminer les l_i et c_j satisfaisants.

La valeur la plus fréquente de a_{ij} sera 0 si la valeur modale de $a_{ij}/l_i c_j$ est 1, ce que G. HIEZ note :

$$\hat{M}(a_{ij}/l_i c_j) = 1$$

Tout le problème est qu'il n'y a pas de méthode usuelle de calcul précis de la valeur modale :

- les données ne sont pas assez nombreuses pour être réparties en classes ; la valeur modale y aurait été la valeur du point médian de la classe de fréquence maximale.
- il faudrait donc définir mathématiquement deux fonctions de densité, "ligne" et "volume" ce qui ne peut se faire a priori.

LE TRAITEMENT L-C

La définition de a_{ij} précédente est claire, il est intéressant de voir par quel procédé on peut obtenir une matrice [E] composée de tels éléments ; il faut pour cela introduire 2 matrices diagonales, dont les éléments diagonaux sont $1/l_i$ et $1/c_j$, matrice que l'on notera $[L^{-1}]$ et $[C^{-1}]$, de dimension n , n et m , m . [E] est alors définie par :

$$[E] = [L^{-1}] \cdot [A] \cdot [C^{-1}] - \quad [1] \text{ est une matrice } n,m \text{ composée de } 1$$

La technique de recherche du mode étant supposée maîtrisée, et son résultat noté

$$\hat{M}\left(\frac{a_{ij}}{l_i c_j}\right)$$

on peut écrire dans [E] :

- pour toute colonne j donnée et quel que soit i

$$\hat{M}\left(\frac{a_{ij}}{l_i c_j}\right) = \frac{1}{c_j} \hat{M}\left(\frac{a_{ij}}{l_i}\right) = 1$$

$$\text{soit } c_j = \hat{M}\left(\frac{a_{ij}}{l_i}\right) \text{ quel que soit } i$$

- pour toute ligne i donnée et quel que soit j

$$\hat{M}\left(\frac{a_{ij}}{l_i c_j}\right) = \frac{1}{l_i} \hat{M}\left(\frac{a_{ij}}{c_j}\right) = 1$$

$$\text{soit } l_i = \hat{M}\left(\frac{a_{ij}}{c_j}\right), \text{ quel que soit } j$$

On peut arriver à ce résultat en restant sous la forme matricielle. Minimiser [E] revient à identifier [A] et [B], or [B] = [L] . [C]. Si l'on fait intervenir $[L^{-1}]$ et $[C^{-1}]$ définies comme auparavant, il vient $[A] \equiv [L] \cdot [C]$ et selon que l'on travaille en [L] ou en [C] :

$[L^{-1}].[A] \equiv [L^{-1}].[L].[C]$
 $[L^{-1}].[L]$ est une matrice colonne de $n \times 1$, le produit $[L^{-1}].[L].[C]$ est alors une matrice n, m répétant n fois la matrice ligne $[C]$.

$[A].[C^{-1}] \equiv [L].[C].[C^{-1}]$ est une matrice ligne de $m \times 1$, le produit $[L].[C].[C^{-1}]$ est alors une matrice n, m , répétant m fois la matrice colonne $[L]$.

Comme $[L^{-1}].[A]$ est aussi une matrice n, m d'éléments $a_{ij}/l_{i.}$, l'identification $[A] \equiv [B]$ conduit à minorer $a_{ij}/l_{i.}$

Comme $[A].[C^{-1}]$ est aussi une matrice n, m d'élément $a_{ij}/c_{.j}$, l'identification $[A] \equiv [B]$ conduit à minorer $a_{ij}/c_{.j}$.

Le choix du mode comme norme conduit bien aux expressions matricielles étendues :

$$[L] = \hat{M}([A].[C^{-1}])$$

$$[C] = \hat{M}([L^{-1}].[A])$$

Cette présentation fournit tous les éléments d'un processus itératif.

Si p est l'ordre des itérations, on écrira :

$$c_{jp} = \hat{M}(a_{ij}/l_{i.p.})$$

$$l_{ip} = \hat{M}(a_{ij}/c_{j.p.})$$

et on initialisera par $l_{i1} = 1$

C'est ce procédé que G. HIEZ nomme "Traitement ligne-colonne ou L.C", procédé qui converge plus ou moins rapidement selon la qualité des données ou le processus d'estimation de la valeur modale. Un "seuil de convergence" S est introduit de sorte d'arrêter l'itération lorsque ce seuil est atteint, qui peut être défini par :

$$v_{ij} = \frac{l_{i.p.} c_{j.p.}}{l_{i.p-1} c_{j.p-1}} - 1 \leq S$$

Une autre approche, après chaque traitement L.C consiste à établir la matrice $[E]$ des erreurs relatives

$$\varepsilon_{ij} = \frac{a_{ij}}{l_{i.} c_{.j}} - 1$$

et à rechercher la valeur de a_{ij} à laquelle correspond l'erreur ε_{ij} maximale. Cette valeur a_{ij} est corrigée par multiplication par le facteur

$$K_{ij} = \frac{1}{1 + \alpha \cdot \varepsilon_{ij}}$$

où α est un facteur d'accélération ou de ralentissement de la convergence. On arrête l'itération lorsque tous les ε_{ij} sont inférieurs à un certain seuil que l'expérience monte de l'ordre de 10^{-2} .

Le produit de ces traitements L.C est donc en premier lieu la creation d'un vecteur L, vecteur de référence régional, dont les composantes L_i traduisent, dans le cas de la pluviométrie, les variations annuelles régionales de cette pluviométrie. En second lieu ces traitement nous fournissent le coefficient, C_j de chacune des stations.

LE SOUS PRODUITS DU TRAITEMENT L.C

Des méthodes d'analyse des totaux pluviométriques annuels classiques s'éclaircissent alors d'un jour nouveau :

Il suffit de comparer au vecteur de référence L les données de chaque poste. On comparera ainsi les données à leur valeur la plus vraisemblable.

Cet écart à la valeur la plus probable sera pour l'année i et la station j :

$$e_{ij} = \frac{a_{ij}}{L_i C_j} - 1$$

que G. HIEZ a nommé EPSI(i, j)

Dans certains cas il est même recommandé de travailler avec le logarithme et de définir

$$EPSI(i, j) = \text{Log} \left(\frac{a_{ij}}{L_i C_j} \right)$$

Deux autres variables peuvent aussi être définies :

SEPSI, valeur cumulée des EPSI pour une année k :

$$\begin{aligned} SEPSI(k) &= \sum_{i=1}^k EPSI(i) \\ &= \frac{1}{C_j} \sum_{i=1}^k \frac{a_{ij}}{L_i} - k \end{aligned}$$

DEPSI, différence première des variables EPSI entre 2 années successives

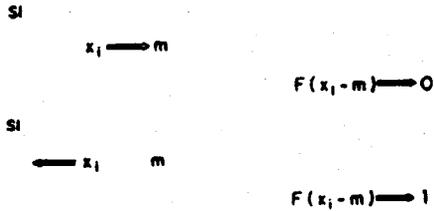
$$DEPSI_i = EPSI_i - EPSI_{i-1}$$

Ces deux variables se révéleront particulièrement performantes dans l'analyse des données dont des exemples circonstanciés sont présentés plus loin.

LA VALEUR MODALE OU LE MODE RETROUVE

G. HIEZ propose une fonction de distribution "naturelle", déjà annoncée dans l'article de 1977 (HIEZ, 1977), qu'il qualifie de plus de "non paramétrique", signifiant sans doute ainsi son caractère discontinu.

VECTEUR REGIONAL
LA VALEUR MODALE



La fonction F est effectivement définie comme une norme d'ordre zéro, analogue donc à un mode. f sera la fonction de distribution de la variable à normer, fonction garantie naturelle et non paramétrique et où l'on retrouve un coefficient de résolution R

$$F = \sum_{i=1}^{i=N} \left[1 - \exp - \left(\frac{x_i - m}{R} \right)^2 \right]$$

$$\hat{M} = \text{MODE} = \min (F)$$

$$f = \sum_{i=1}^{i=N} \exp - \left(\frac{x_i - m}{R} \right)^2$$

$$\hat{M} = \max (f)$$

Un second tableau de G. HIEZ précise la définition de cette fonction de distribution f(x) :

VECTEUR REGIONAL
 FONCTION DE DISTRIBUTION "NATURELLE"

$f \rightarrow$ FONCTION DE DISTRIBUTION.

Il semble bien que l'on attache plus de prix à l'espace temporel des n années, dimension du Vecteur Régional, qu'à l'espace des m stations utilisées.

1^o) $f > 0 \quad \forall x_i$

2^o) $\frac{1}{P} \int_{-\infty}^{+\infty} f \cdot dx = 1$

et $P = NK\sqrt{2\pi}$

Le coefficient de résolution est devenu K, mais on comprend qu'il doit être ajusté à la variabilité de l'échantillon pour ne pas risquer de masquer la structure des données.

$$f(x) = \frac{1}{NK\sqrt{2\pi}} \sum e^{-\frac{1}{2} \left(\frac{x-x_i}{K}\right)^2}$$

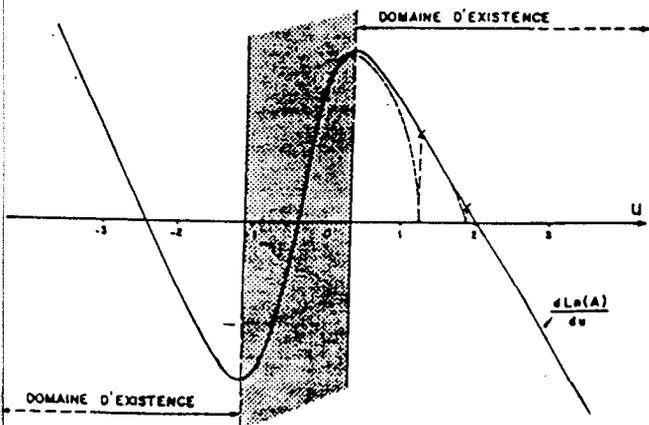
$f(x) =$ FONCTION DE DISTRIBUTION NON PARAMÉTRIQUE

K = COEFFICIENT DE RESOLUTION



STRUCTURE DES DONNÉES

La recherche des modes s'effectue sur une fonction réduite :



$$A = \frac{1}{N\sqrt{2\pi}} \sum_{i=1}^{i=N} \exp -\frac{1}{2} (u-u_i)^2$$

en utilisant la technique de la parabole osculatrice de la fonction dérivée logarithmique de la fonction réduite A.

G. HIEZ rassemble dans un dernier tableau les avantages de la fonction de distribution "naturelle" pour la recherche des modes de la distribution des composantes du vecteur L.

VECTEUR REGIONAL
FONCTION DE DISTRIBUTION "NATURELLE"

PROPRIÉTÉS:

La manipulation des f et des Σ au ni-veau de la fonction de distribution et de son intégrale, n'est pas apparue transparente, au regard de la définition a priori discrète et non continue de l'échantillon de base.

TOUS LES MOMENTS EXISTENT:
LA DISTRIBUTION EST DONC QUELCONQUE
SYMÉTRIQUE OU NON
UNI- OU PLURI-MODALE

K_r INDÉPENDANTS DE R
SAUF $K_2 \longleftrightarrow$ RESOLUTION
 m_3 INDÉPENDANTS DE R
CONSERVATION DE L'ASSYMETRIE

FONCTION CUMULATIVE:

$$FC = \frac{1}{N\sqrt{2\pi}} \int_{-\infty}^u \Sigma e^{-\frac{1}{2}(u-u_i)^2} . du$$

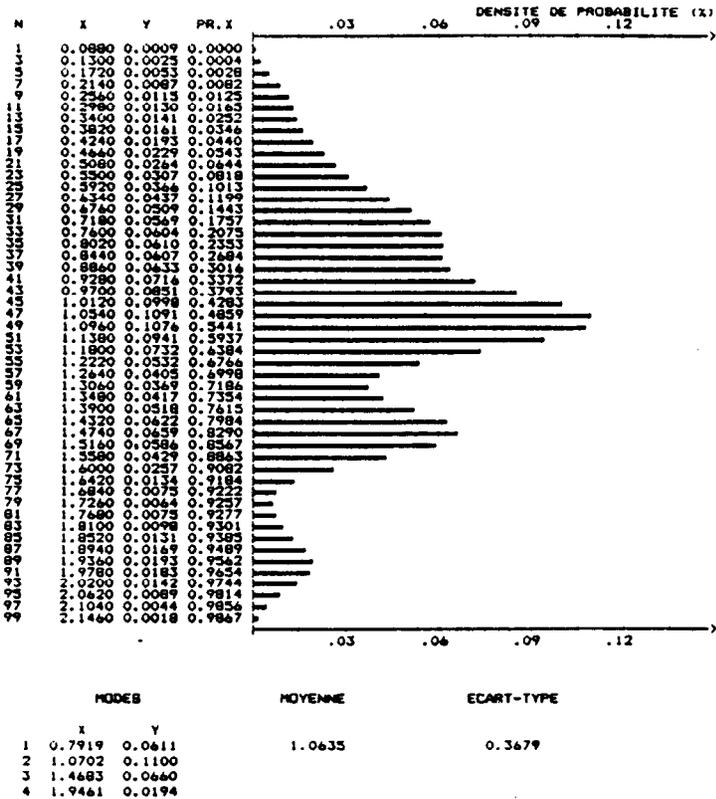
$$= \frac{1}{N\sqrt{2\pi}} \sum_{i=1}^{i=N} \int_{-\infty}^u e^{-\frac{1}{2}(u-u_i)^2} . du$$

CARACTÉRISTIQUE DE TENDANCE — **MODE(S)** —
PEU "CONTAMINÉ" PAR LES VALEURS
DE FAIBLE FRÉQUENCE.

En tout état de cause le pragmatisme conduit à juger une méthode sur ses produits, ce qui peut être esquissé par l'examen du graphe de la visualisation de la distribution du vecteur pluviométrique de la totalité du CEARA (Brésil) pour la période 1911-1985, où apparaissent les modes de cette distribution.

VISUALISATION DE LA DISTRIBUTION :
 VECTEUR PLUVIOMETRIQUE CEARA TOTAL - Période 1911-1985

Densité des probabilités réduites X, Y pour $ALF = .216$
 (paramètres utilisés : $IR = 14$, $IC = 14$ et $IO = 0$)



Dans cet exemple on voit que la distribution temporelle du vecteur est largement plurimodale.