

NOTION D'AGROSYSTEME

PREMIERE PARTIE

NOTION DE BASE ET NOTIONS DERIVEES

METHODE D'ETUDE

Chapitre 2

CARACTERISATION D'UN AGROSYSTEME

par

B. BONZON et J. DEJARDIN

H= 54870

~~N AG 1983/1/1
N° enote qu'a
Noumba~~



Fonds Documentaire IRD
Cote: B * 25730 Ex: *uniqua*

NOTE SUR UN PROJET D'OUVRAGE COLLECTIF
CONCERNANT LA NOTION D'AGROSYSTEME.

La notion d'agrosystème a été élaborée en 1972 à Adiopodoumé par J. DEJARDIN, C. FILLONNEAU et B. BONZON en réponse aux difficultés rencontrées à l'époque par les agronomes de Côte d'Ivoire dans l'organisation de l'interprétation de leurs données sur les "Interaction sol-plantes fourragères en milieu tropical humide".

Le concept et les notions qu'ils en ont dérivées/^{ont}été utilisés en premier par J.C. TALINEAU pour différencier l'évolution de la matière organique du sol sous les graminées et les légumineuses étudiées dans le cadre de ce programme.

Ils ont été mis en oeuvre ensuite par P. de BOISSEZON et B. BONZON pour étudier de façon systématique un certain nombre des relations sol-techniques culturales-manioc dans les conditions d'une étude expérimentale conduite elle aussi à Adiopodoumé (travail qui doit être présenté prochainement).

Ils sont utilisés actuellement de façon très systématique à Nouméa par B. DENIS et B. BONZON dans le cadre du programme multidisciplinaire sur la fertilité et l'évolution sous culture des sols de Nouvelle-Calédonie.

L'expérience acquise par les uns et les autres montrent finalement que ces concepts permettent une plus grande finesse au niveau de l'interprétation scientifique après en avoir facilité l'organisation préalable.

A la suite d'un certain nombre de discussions entre B. DENIS, P. de BOISSEZON et B. BONZON à Nouméa, il est apparu à ces chercheurs qu'il serait certainement intéressant, tant au plan conceptuel qu'aux plans méthodologique et scientifique, de présenter enfin officiellement la notion d'agrosystème dans un ouvrage collectif.

Contactés fin Juin 1983 par B. BONZON pour participer à ce travail, J. DEJARDIN et J.C. TALINEAU ont déjà donné leur accord de principe. C. FILLONNEAU contacté lui aussi à la même époque en Côte d'Ivoire n'a peut être pas encore reçu cette correspondance.

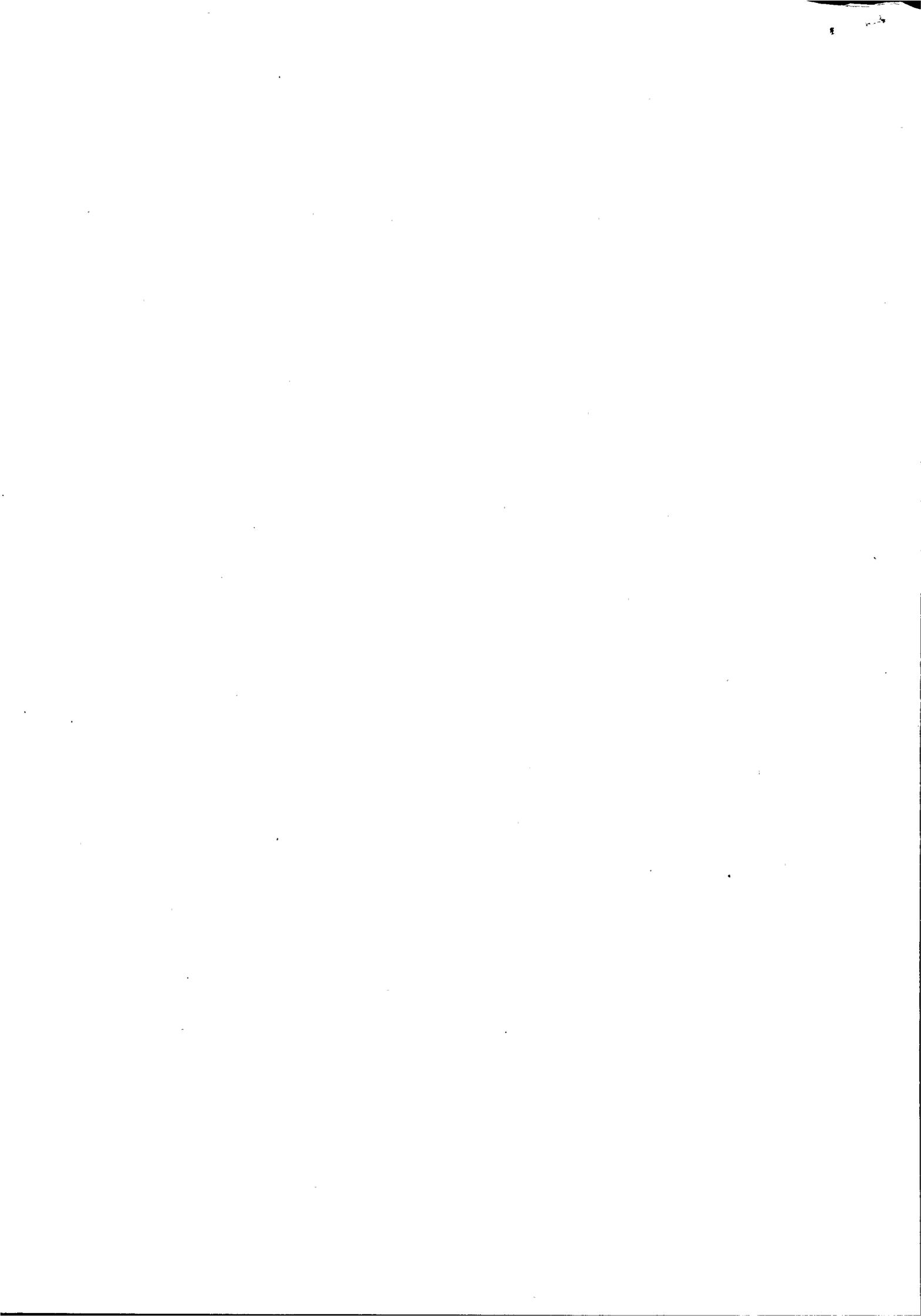
Plus précisément, le projet pourrait comporter deux parties. La première, théorique, exposerait la notion d'agrosystème, les notions que l'on peut y rattacher et les principales méthodes et techniques d'analyses statistiques nécessaires.

La seconde illustrerait la première en s'appuyant sur des résultats tirés des études conduites à Adiopodoumé et Nouméa.

Nombre des chapitres de l'ouvrage pourraient faire l'objet, par ailleurs, de publications préalables.

Le plan général de la première partie est le plus avancé. Elle pourrait comporter quatre chapitres et une annexe.

Le premier chapitre présenterait la notion d'agrosystème et les notions que l'on peut y rattacher.



Les suivants traiteraient successivement :

- . le chapitre II , de la caractérisation d'un agrosystème ;
- . le chapitre III , de la comparaison de deux ou plusieurs agrosystèmes ;
- . le chapitre IV , de l'effet de facteurs de variations contrôlées sur un agrosystème.

L'annexe regrouperait les informations nécessaires à l'application des principaux tests statistiques utilisés pour caractériser un agrosystème, comparer deux ou plusieurs agrosystèmes, étudier l'influence de facteurs de variations sur un agrosystème, l'accent étant mis sur les possibilités d'exécutions automatiques de ces tests.

Rédigée à Nouméa le 06/09/83 et adressée pour critique et compléments à
J. DEJARDIN, J.C. TALINEAU, C. FILLONNEAU.

Adressée en même temps pour info à MM. P. FRANQUIN, B. DABIN, J. FORESTIER.

Documents de référence actuels.

1. BONZON B., DEJARDIN J., FILLONNEAU C., GERI M., 1974 .

Programme multilocal d'étude des interactions sol plantes fourragères en milieu tropical humide. Organisation générale de l'analyse statistique des résultats. ORSTOM - Adiopodoumé, multig., 14 p, 4 tableaux.

2. TALINEAU J.C., BONZON B., FILLONNEAU C., HAINNAUX G., 1980.

Contribution à l'étude d'un agrosystème praival dans le milieu tropical humide de la Côte d'Ivoire. 2. Analyse des données relatives à l'état de la matière organique. Cah. ORSTOM, sér. Pédol., vol. XVIII, n° 1, 1980-81 : 29-47.

Pièces jointes : premières ébauches des chapitres 2 et 3 et de l'annexe de la première partie (par B. BONZON, J. DEJARDIN et C. FILLONNEAU si ce dernier est intéressé par le projet).



CARACTERISATION PARTIELLE APPROCHEE D'UN AGROSYSTEME

ASPECTS STATISTIQUES

-§-

	<u>Pages</u>
1. NOTION DE CORRELATION ET PROBLEMES A RESOUDRE	3
1.1. Notions de liaison fonctionnelle, de corrélation et d'indépendance	3
1.2. Problèmes à résoudre	4
2. CORRELATION DANS LE CAS OU LES DEUX VARIABLES SONT DISTRI- BUEES NORMALEMENT ET LIEES LINEAIREMENT	5
2.1. Loi de référence	5
2.2. Estimations des coefficients de corrélation et de régression	7
2.3. Conditions d'utilisation des estimations des coefficients de corrélation et de régression	9
2.3.1. Vérification de la normalité des deux variables	10
2.3.1.1. Principe du test	10
2.3.1.2. Aspects particuliers de la conduite du test	11
2.3.1.2.1. Conduite du test de KOLMOGOROV et SMIRNOV sur des distributions classées	12
2.3.1.2.2. Problème du nombre de classes k à retenir	13
2.3.1.2.3. Problème de la détermination de la fréquence relative cumulée théorique dans le cas d'une réalisation entièrement automatique du test	14
2.3.1.2.4. Réalisation en série du test dans le cas où l'on ne peut calculer la fréquence relative cumulée théorique	14
2.3.2. Signification des coefficients de corrélation et de régression	15
2.3.2.1. Première méthode	15
Réalisation du test : table des valeurs limites de $ r $ à différents seuils α et réalisation automatique du test ...	16
2.3.2.2. Deuxième méthode	17

2.3.2.2.1.	Intervalle de confiance du coefficient de corrélation	17
2.3.2.2.2.	Comparaison de r à une valeur théorique r_0	
2.3.2.2.3.	Remarques sur une meilleur approximation de r et sur le nombre minimum de couples	18
2.3.3.	Vérification de la linéarité des régressions	18
2.3.3.1.	Principe du test	18
2.3.3.2.	Conduite du test	19
2.3.3.3.	Remarque concernant la réalisation automatique du test	21
3.	CORRELATION DANS LE CAS OU LES DEUX VARIABLES SONT DISTRIBUEES NORMALEMENT MAIS NE SONT PAS LIEES LINEAIREMENT	21
3.1.	Rapport de corrélation de X sur U	21
3.2.	Rapport de corrélation de U sur X	23
3.3.	Linéarisation des relations $X = f(U)$ et $U = g(X)$	23
3.3.1.	Principe d'une démarche générale	23
3.3.2.	Fonctions de linéarisation utilisées couramment	24
4.	CORRELATION DANS LE CAS OU LES DEUX VARIABLES NE SONT PAS DISTRIBUEES NORMALEMENT	24
4.1.	Fonctions de transgénération classiques	25
4.2.	Règle empirique de choix d'une fonction de transgénération stabilisant les variances	27
4.3..	Remarques sur certains aspects de la non-linéarité ...	27
5.	ORGANISATION GENERALE DES CALCULS	28
6.	CONCLUSION	28

CARACTERISATION PARTIELLE APPROCHEE D'UN AGROSYSTEME,

ASPECTS STATISTIQUES,

-§-

Soit donc un agrosystème A dont on désire caractériser la partie à p éléments observés indépendamment les uns des autres sur n sites répartis de façon aléatoire à la surface du champ de référence.

Chaque élément est défini par un "paramètre caractéristique", ou "élémentaire", ou "de base", ou encore plus simplement, par une "caractéristique", vocable plus général car pouvant inclure des aspects qualitatifs. Pour la suite, on ne considérera cependant, ici, que des "caractéristiques mesurables".

Il y a donc p caractéristiques

U, V, W, ..., X, ..., Z

et chaque caractéristique fait l'objet de n observations, par exemple pour X :

$x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n$

Lorsqu'une observation de X est faite sur chacun des n sites, une autre observation est faite, simultanément ou non, de chacune des (p-1) autres caractéristiques sur chacun des mêmes n sites.

L'ensemble de ces informations peut être rassemblé en une matrice de données \hat{A} de dimensions (p.n) :

$$\hat{A} = \begin{pmatrix} u_1 & u_2 & \dots & u_i & \dots & u_n \\ v_1 & v_2 & \dots & v_i & \dots & v_n \\ : & & & & & \\ x_1 & x_2 & \dots & x_i & \dots & x_n \\ : & & & & & \\ z_1 & z_2 & \dots & z_i & \dots & z_n \end{pmatrix}$$

Chaque colonne de \hat{A} caractérise un site, chaque ligne une caractéristique.

Si l'on connaissait la fonction de liaison de chacun des $p.(p-1) / 2$ couples de paramètres (X,U) , (X,V) , ... etc, et si les conditions de milieu des sites d'observations possibles étaient parfaitement homogènes, le système constitué par la partie de A à p éléments serait défini à partir d'une seule série d'observations des caractéristiques U,V,\dots etc, c'est-à-dire à partir d'un seul site. Les n vecteurs-colonnes de \hat{A} seraient, en conséquence, identiques.

Dans la réalité, l'intervention de facteurs de variations aléatoires sur l'aire des sites d'observations possibles fait osciller chaque caractéristique X autour de sa valeur moyenne réelle que l'on ne peut, évidemment, qu'estimer en répétant les observations un certain nombre de fois, n ici en l'occurrence.

Dans ces conditions d'ailleurs, si l'on peut estimer μ_X , moyenne vraie de X , par

$$\bar{X} = (1/n) \cdot \sum_{i=1}^n x_i \quad (1)$$

on sait qu'il faut aussi préciser, en même temps, l'étendue du champ de variation de \bar{X} en estimant l'écart-type de sa distribution, $s_{\bar{X}}$, à partir de

$$s_{\bar{X}}^2 = \{ 1 / n.(n-1) \} \cdot \sum_{i=1}^n (x_i - \bar{X})^2 \quad (2)$$

Quant à la détermination des fonctions de liaison des $p.(p-1) / 2$ couples de paramètres (X,U) , $(X,V),\dots$ etc, il s'agit-là du problème fondamental de la caractérisation d'un agrosystème.

Pour chaque couple de paramètres (X,U) deux situations peuvent se présenter :

- ou bien l'on connaît a priori la fonction théorique de la liaison qui les unit, et le problème est alors celui de l'ajustement de cette fonction aux n valeurs observées des couples (X,U) . D'une façon générale, il s'agit d'un problème de modélisation ;

- ou bien l'on ne possède aucune information sur cette fonction et il s'agit dans ce cas, plus modestement, d'estimer l'intensité de la liaison puis de préciser une fonction ou un graphique d'estimation de X à partir de U , ou de l'inverse.

D'une façon générale les problèmes à résoudre relèvent alors de la corrélation.

Comme dans l'immense majorité des cas on se trouve dans la seconde situation on rappellera seulement ici l'essentiel de la théorie et surtout des techniques d'analyse de la corrélation nécessaires à la caractérisation approchée de la partie observée de A.

Remarque

Le fait pour chacune des p caractéristiques de pouvoir être liée simultanément aux $(p-1)$ autres, ou du moins à un certain nombre d'entre elles, pose naturellement le problème des liaisons multiples et partielles susceptibles de les unir. D'un point de vue théorique, deux autres séries de problèmes devraient donc être examinées ensuite :

- celle des problèmes relatifs à la caractérisation des liens multiples et partiels entre X d'une part, U, V, W, \dots etc, d'autre part ;

- celle des problèmes concernant la caractérisation des liens existant entre les ensembles de deux groupes de caractéristiques : par exemple ceux existant entre X et U d'une part, V, W, \dots etc, d'autre part (corrélations canoniques).

Comme les conditions d'application de l'analyse multidimensionnelle de p caractéristiques et la conduite des calculs peuvent reposer sur les conditions d'application et les résultats de l'analyse bi-dimensionnelle de ces p caractéristiques prises deux à deux, on se limitera ici, pour l'instant, aux problèmes de la corrélation simple.

S'agissant, enfin, de rappels très orientés, le lecteur aura intérêt à consulter les ouvrages de base de l'analyse statistique et de la statistique multidimensionnelle, notamment ceux de P. DAGNELIE (1969-70, 1975) et J. LEFEVRE (1976).

1 - NOTION DE CORRELATION ET PROBLEMES A RESOUDRE.

1.1. Notions de liaison fonctionnelle, de corrélation et d'indépendance.

Soient donc deux caractéristiques X et U prises parmi les p caractéristiques du système :

1°/ si X et U sont rigidement liées l'une à l'autre, c'est-à-dire si à toute valeur de U l'on sait aussitôt faire correspondre une valeur et une seule de X ,

et réciproquement, on dira qu'elles sont liées de façon fonctionnelle : par exemple la surface S d'un cercle et son rayon R .

Dans la réalité et graphiquement, cela se traduira par un nuage de points très étroit à l'allure générale de courbe : dans l'exemple ci-dessus, une parabole que l'on retrouvera aux erreurs près commises sur S et R ;

2°/ Si X et U sont liées de façon moins rigide, en d'autres termes si l'on observe pour un niveau donné de U une certaine variation de X , et réciproquement, les deux caractéristiques seront dites en liaison stochastique ou en corrélation. On dira aussi que leurs répartitions sont liées.

Dans la réalité et graphiquement, cela se traduira par un nuage de points d'autant plus étiré que la liaison sera proche d'une liaison fonctionnelle.

Un diagramme des moyennes liées pourra être tracé à travers le nuage de points. La relation que l'on mettra ainsi en évidence en classant X sur U $\{ X = f(U) \}$ n'est cependant pas la même que celle que l'on obtiendra en classant U sur X $\{ U = g(X) \}$.

Les deux courbes de régression ainsi dégagées seront d'autant plus éloignées l'une de l'autre et différentes que la liaison s'éloignera d'une liaison du type fonctionnelle et linéaire, un cas particulier - très important - étant celui où les lignes de régression sont des droites ;

3°/ si X et U ne sont pas liées du tout la répartition des valeurs de X sera la même quelle que soit la valeur de U , et réciproquement.

Graphiquement, cette situation devrait se traduire par des droites de régression $X = f(U)$ et $U = g(X)$ parallèles aux axes des U et des X .

Un point à noter au passage : aucune des deux variables n'a été privilégiée.

1.2. Problèmes à résoudre.

Devant caractériser la liaison (X,U) deux premières questions se posent immédiatement. Comment peut-on :

- . caractériser l'intensité de cette liaison, son étroitesse ?
- . estimer l'une des variables à partir de l'autre ?

Comme les réponses à ces questions dépendent de la normalité des distributions de X et de U d'une part, de la linéarité des fonctions de régression $X = f(U)$ et $U = g(X)$ d'autre part, deux nouveaux problèmes apparaissent. Peut-on et comment :

- . vérifier la normalité des distributions des variables ?
- . vérifier la linéarité de leur régressions ?

Si l'une des deux régressions n'est pas linéaire - et a fortiori si les deux ne le sont pas - une troisième série de questions se pose :

- . est-il possible dans ces conditions de caractériser l'intensité de la liaison ?
- . ou : peut-on transformer l'une ou les deux variables pour linéariser la ou les régressions ?

Enfin, si la distribution de X (ou celle de U) n'est pas normale, une septième question devra être examinée :

- . est-il possible de transformer X (ou U) pour que sa distribution apparaisse normale ?

Ces questions sont, à l'évidence, très étroitement liées les unes aux autres.

2. CORRELATION DANS LE CAS OU X ET U SONT DISTRIBUEES NORMALEMENT ET LIEES LINEAIREMENT.

2.1. Loi de référence.

. Si l'on fait l'hypothèse que X et U sont distribuées normalement et liées linéairement, on se place dans une situation où l'on dispose d'une loi de référence classique, la loi normale à deux variables.

. Cette loi dépend de cinq paramètres :

- la moyenne μ_X et l'écart-type σ_X de X ;
- la moyenne μ_U et l'écart-type σ_U de U ;
- une constante ρ , appelée coefficient de corrélation de X et U, qui est un nombre compris entre -1 et +1.

. Elle jouit des propriétés suivantes :

- les distributions marginales de X et de U (c'est-à-dire celles de l'ensemble des observations de X d'une part, de U de l'autre) sont normales, de même que les distributions de X à U constant ou de U à X constant ;
- les courbes de régression sont des droites ;
- les écart -types dans les sous-populations liées sont réduits dans la proportion de 1 à $\sqrt{1-\rho^2}$ par rapport aux écart -types σ_x et σ_u .

. Le coefficient de corrélation ρ résume à lui seul le degré de liaison des deux variables :

- si $\rho = 0$, X et U sont indépendantes : les droites de régression sont des parallèles aux axes de coordonnées et la variabilité dans les tranches à U ou à X constant est la même qu'au niveau des distributions marginales ;

- si $\rho = \pm 1$, X et U sont en liaison fonctionnelle linéaire ; les deux droites de régression sont confondues et la variabilité dans les tranches est nulle ;

- si $-1 < \rho < +1$, X et U sont en dépendance stochastique, ou en corrélation, plus ou moins étroite.

. Les relations suivantes existent entre ces cinq paramètres :

- l'équation de la droite de régression de X sur U (ou de X en U) qui permet d'estimer X à partir de U, X étant considérée comme variable dépendante de U, est

$$x(u) - \mu_x = \rho \cdot \frac{\sigma_x}{\sigma_u} \cdot (u - \mu_u) \quad (3)$$

la quantité

$$\rho \cdot \frac{\sigma_x}{\sigma_u} = \beta_{xu} \quad (4)$$

désignée par le vocable de "coefficient de régression de X sur U", étant la pente de la droite en question ;

- l'équation de la droite de régression de U sur X (ou de U en X) qui permet d'estimer U à partir de X, U étant considérée comme variable dépendante de X est

$$u(x) - \mu_U = \rho \cdot \frac{\sigma_U}{\sigma_X} \cdot (x - \mu_X) \quad (5)$$

la quantité

$$\rho \cdot \frac{\sigma_U}{\sigma_X} = \beta_{UX} \quad (6)$$

appelée de la même façon "coefficient de régression de U sur X" étant la pente de cette deuxième droite ;

- les distributions de X dans les sous-populations à U constant ont pour écart-type

$$\sigma_{X(U)} = \sqrt{1 - \rho^2} \cdot \sigma_X \quad (7)$$

- les distributions de U dans les sous-populations à X constant ont de la même façon pour écart-type

$$\sigma_{U(X)} = \sqrt{1 - \rho^2} \cdot \sigma_U \quad (8)$$

- le coefficient de corrélation ρ est égal à

$$\rho = \frac{M_{XU}}{\sigma_X \cdot \sigma_U} \quad (9)$$

où M_{XU} représente le moment centré d'ordre 1 par rapport à X et U ;

- les coefficients de régression et de corrélation sont liés entre eux par :

$$\rho^2 = \beta_{XU} \cdot \beta_{UX} \quad (10)$$

. Les deux droites de régression passent par le point M (\bar{x}, \bar{u})

2.2. Estimation des coefficients de corrélation et de régression.

Disposant de n couples (x_i, u_i) d'observations de X et U on peut donc estimer σ_X par s_X calculé à partir de

$$s_x^2 = \{1/ (n-1)\} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 \quad (11)$$

σ_u par s_u calculé à partir de

$$s_u^2 = \{1/ (n-1)\} \cdot \sum_{i=1}^n (u_i - \bar{u})^2 \quad (12)$$

M_{xu} par M_{xu} calculé à partir de

$$M_{xu} = \{1/ (n-1)\} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (u_i - \bar{u}) \quad (13)$$

ρ par r_{xu} tel que

$$r_{xu} = \frac{M_{xu}}{\sqrt{s_x^2 \cdot s_u^2}} \quad (14)$$

β_{xu} par b_{xu} tel que

$$b_{xu} = M_{xu} / s_u^2 \quad (15)$$

β_{ux} par b_{ux} tel que

$$b_{ux} = M_{xu} / s_x^2 \quad (16)$$

Remarque.

Si l'on applique la méthode des moindres carrés à l'estimation des coefficients des droites de régression de X sur U et de U sur X on retrouve les mêmes estimations de b_{xu} et b_{ux} . En effet, soit par exemple

$$x(u) = a + b_{xu} \cdot u \quad (17)$$

l'équation de la droite de régression de X sur U.

Si liaison il y a entre X et U, l'une des positions les plus probables de cette droite peut être celle qui minimise, globalement, les écarts entre les

valeurs observées x_i et les valeurs estimées $x(u)_i$ données par l'équation de la droite ci-dessus, c'est-à-dire, en fait, qui minimise la quantité

$$Q = \sum_{i=1}^n \{x_i - x(u)_i\}^2 = \sum_{i=1}^n (x_i - a - b_{xu} \cdot u_i)^2 \quad (18)$$

Comme x_i et u_i sont connus, le minimum peut être obtenu si l'on a simultanément :

$$\frac{\delta Q}{\delta a} = \frac{\delta Q}{\delta b_{xu}} = 0 \quad (19)$$

c'est-à-dire :

$$\left\{ \begin{array}{l} \sum_{i=1}^n (x_i - a - b_{xu} \cdot u_i) = 0 \\ \sum_{i=1}^n (x_i - a - b_{xu} \cdot u_i) \cdot u_i = 0 \end{array} \right. \quad (20)$$

$$\left\{ \begin{array}{l} \sum_{i=1}^n (x_i - a - b_{xu} \cdot u_i) = 0 \\ \sum_{i=1}^n (x_i - a - b_{xu} \cdot u_i) \cdot u_i = 0 \end{array} \right. \quad (21)$$

système de deux équations du premier degré à deux inconnues qui permet d'estimer les deux inconnues a et b_{xu} par

$$a = \bar{x} - b_{xu} \cdot \bar{u} \quad (22)$$

$$b_{xu} = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (u_i - \bar{u})}{\sum_{i=1}^n (u_i - \bar{u})^2} \quad (23)$$

La droite de régression ainsi définie est bien confondue avec celle donnée par l'équation (3) : elle a la même pente et passe par le même point $M(\bar{x}, \bar{u})$.

2.3. Conditions d'utilisation des estimations de r_{xu} , b_{xu} et b_{ux}

Les quantités r_{xu} , b_{xu} et b_{ux} peuvent toujours être calculées. Leurs estimations reposent néanmoins sur deux hypothèses :

- . la normalité des distributions marginales de x et de u ;
- . la linéarité des régressions.

La vérification de la première est toujours possible pourvu que le nombre n de couples d'observations soit au moins égal à 5.

Celle de la seconde, qui implique une mise en classes des observations, n'est réalisable, de ce fait, que si n est suffisamment grand (20 à l'extrême minimum).

Un autre problème se pose d'autre part :

. celui du seuil de signification du r_{xu} , b_{xu} ou b_{ux} .

En toute logique cette question doit d'ailleurs être réglée avant celle de la linéarité : si r_{xu} n'est pas significativement différent de zéro, un lien non-linéaire peut néanmoins exister entre X et U dont la signification peut être testée si X et U sont normales.

2.3.1. Vérification de la normalité de X et de U .

Tester la normalité de X et de U peut consister en la comparaison de leurs distributions réduites à la distribution normale réduite.

Cette comparaison, qui se conduit de la même façon sur les deux variables, gagne beaucoup en simplicité à s'appuyer sur le test de KOLMOGOROV et SMIRNOV (cf. in DAGNELIE, Théorie et Méthodes Statistiques. Applications Agronomiques, Vol. II, paragraphe 12.2.4, pages 70 et suivantes).

2311 - Principe du test.

Soit U la variable à tester et s_u l'écart-type de sa distribution.

Rangeons les n valeurs observées de U .

Soient u_m la plus petite, u_M la plus grande, u_i la c^{eme} .

La fréquence relative cumulée "observée" de u_i est

$$fco(u_i) = c/n \quad (24)$$

Sa valeur réduite correspondante est :

$$t(u_i) = (u_i - \bar{u})/s_u \quad (25)$$

Si la distribution de U suit une loi normale, comme on sait faire correspondre une fréquence cumulée théorique $f_{cth}(u_i)$ à toute valeur $t(u_i)$ observée de la variable réduite (cette fréquence est donnée par la fonction de répartition Φt , cf. le paragraphe 1 et la table 1 de l'annexe 1) on peut donc calculer pour chaque valeur rangée de U l'écart absolu des fréquences relatives cumulées théoriques et observées.

Pour la c^{eme} cet écart est

$$e_c = |f_{cth}(u_i) - f_{co}(u_i)| \quad (26)$$

KOLMOGOROV et SMIRNOV ont montré que, pour que la courbe de répartition observée puisse être considérée comme normale, ces écarts absolus ne doivent pas dépasser des valeurs limites maximales. Ces valeurs sont fonctions des seuils de probabilité retenus et des effectifs de la population.

D'après DAGNELIE (cf référence ci-dessus), pour des distributions normales définies par leurs moyennes et leur écart-types estimés, \bar{u} et s_u , les valeurs critiques relatives aux seuils 0,05 et 0,01 sont approximativement pour $n > 5$

$$e_{0,05} = 0,886/\sqrt{n+1,5} \quad (27)$$

$$e_{0,01} = 1,031/\sqrt{n+1,5} \quad (28)$$

La comparaison des écarts absolus observés e_c à ces limites $e_{0,05}$ et $e_{0,01}$ permet donc d'accepter ou non l'hypothèse de la normalité de U.

2312 - Aspects particuliers de la conduite du test.

Le test de KOLMOGOROV et SMIRNOV peut être, comme on vient de le voir, conduit sur les valeurs individuelles rangées de la distribution observée. L'importance des opérations de rangement et des calculs à effectuer (calculs de $t(u_i)$, $f_{cth}(u_i)$, etc) d'une part, la nécessité, de l'autre, de classer ensuite les données pour le test de linéarité, font que l'on a généralement avantage, chaque fois que le test de linéarité est applicable, à réaliser le test précédent au niveau des limites de classes de la distribution observée.

et les limites supérieures réduites de classes

$$tu_1, tu_2, \dots, tu_c, \dots, tu_k$$

telles que

$$tu_c = (u_c - \bar{u}) / s_u \quad (38)$$

On peut alors appliquer le test de KOLMOGOROV et SMIRNOV en déterminant la fréquence cumulée théorique $fcthu_c$ correspondant à chaque limite supérieure de classe réduite tu_c , et en calculant

$$e_c = | fcthu_c - fcou_c | \quad (39)$$

que l'on compare aux limites théoriques données par les formules (27) et (28)

2312.2 - Problème du nombre de classes k à retenir.

Le nombre de classes à retenir peut être approximativement défini par l'une des deux formules suivantes :

$$k_1 = 1/2 \cdot \{n/10 - \log_{10}(n)\} \quad (41)$$

$$k_2 = 10 \cdot \log_{10}(n/10) \quad (42)$$

Ces formules n'aboutissent pas aux mêmes résultats.

Formules possibles	Nombres de classes donnés par les formules (chiffres supérieurs) et à prendre (chiffres inférieurs) en fonction de quelques effectifs.											
	20	26	32	40	51	64	80	100	200	330	400	500
k_1	0,3 -	0,5 -	0,8 -	1,1 -	1,7 -	2,2 2	3,0 3	4,0 4	8,8 8	15,2 15	18,7 18	23,6 23
k_2	3,0 3	4,1 4	5,0 5	6,0 6	7,0 7	8,0 8	9,0 9	10,0 10	13,0 13	15,0 15	16,0 16	16,9 16

La seconde, comme on le remarque sur le tableau ci-dessus, semble convenir aux populations de faibles effectifs. Si on l'adopte, on peut, en principe, essayer de mettre en classe des populations de 20 observations. Les nombre de classes à prendre sont alors de :

3	pour	$20 \leq n < 26$
4	"	$26 < n < 32$
5	"	$32 < n < 40$
6	"	$40 < n < 51$
7	"	$51 < n < 64$
8	"	$64 < n < 80$
9	"	$80 < n < 100$
etc...		

Remarque.

Une troisième formule empirique donnée par STURGES $k_3 = 1 + (10/3) \cdot \log_{10} n$ aboutit à des nombres encore plus élevés de classes que la seconde pour les petits échantillons.

Elle ne semble pas adaptée aux problèmes de l'étude d'un agrosystème.

2312.3 - Problème de la détermination de la fréquence relative cumulée théorique dans le cas d'une réalisation entièrement automatique du test.

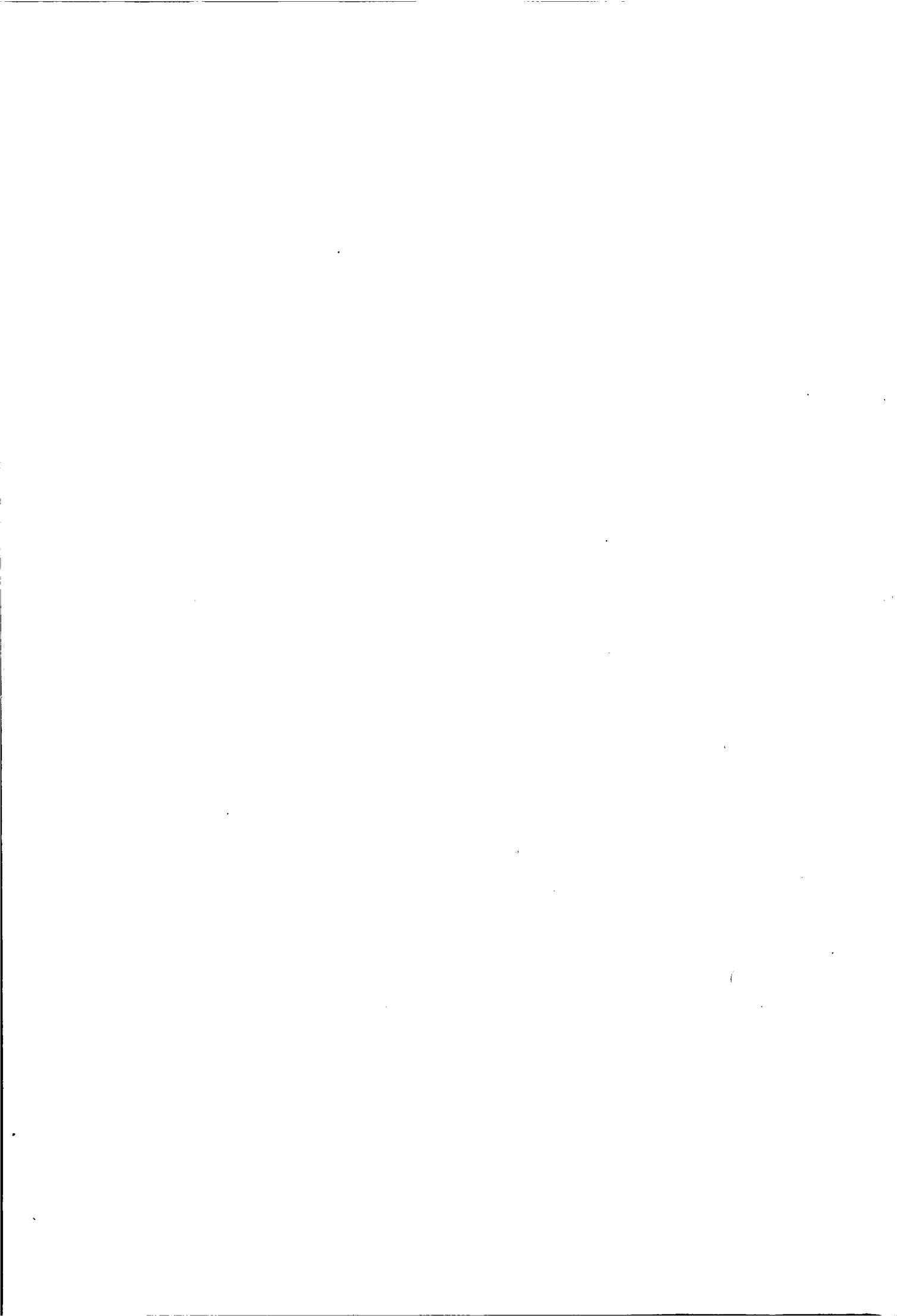
La fréquence relative cumulée théorique f_{cth} est donnée, classiquement, par les tables de la fonction de répartition :

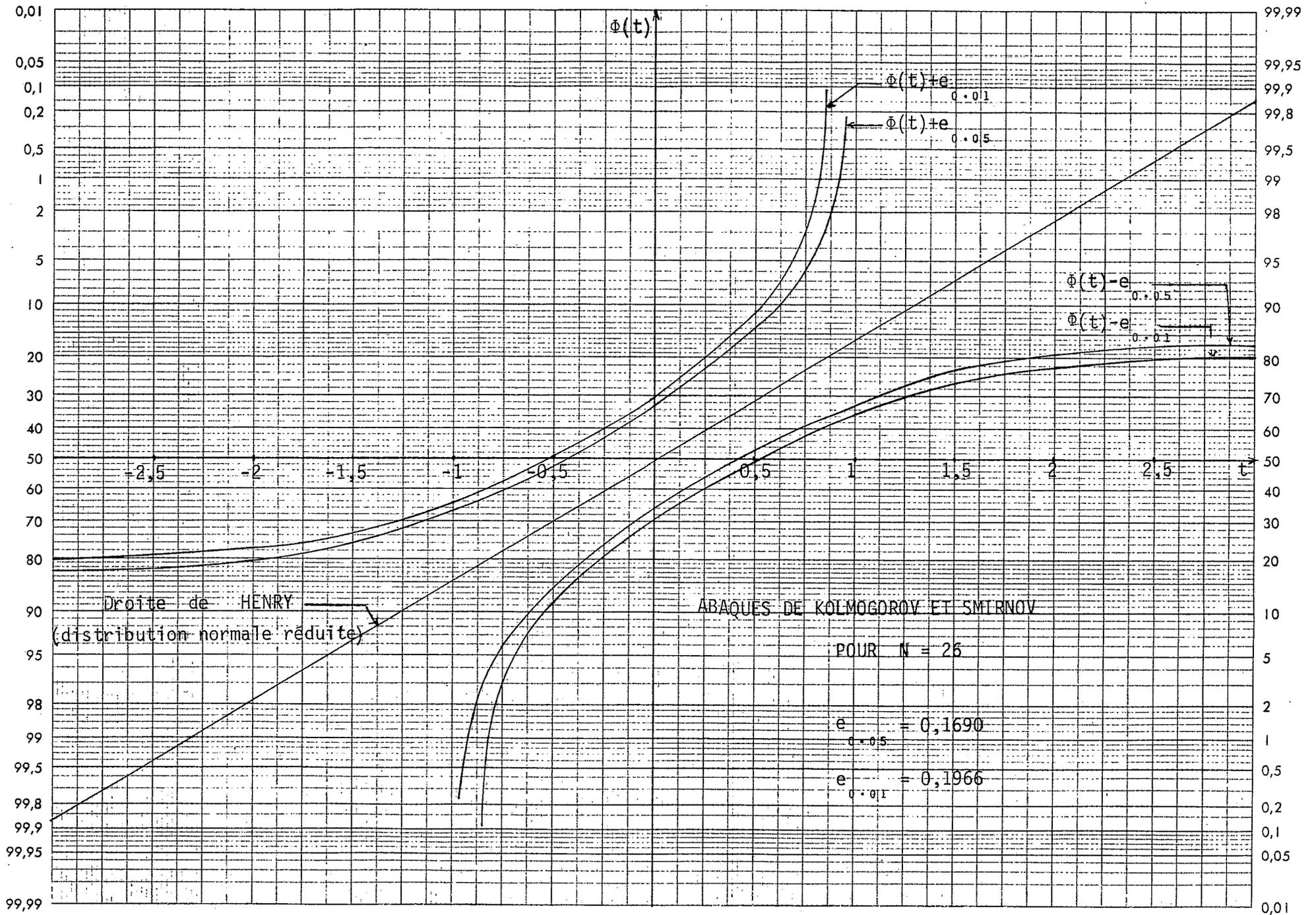
$$f_{cth} = \Phi(t) = \int_{-\infty}^t (1/\sqrt{2\pi}) \cdot e^{-t^2/2} \cdot dt \quad (43)$$

Dans le cas où les moyens de traitement des données le permettent, la "meilleure approximation de HASTINGS" donne directement une valeur suffisamment précise de $\Phi(t)$ (cf annexe 1, paragraphe 1).

2312.4. - Réalisation en série du test dans le cas où l'on ne peut calculer la fréquence relative cumulée théorique.

Dans le cas où les moyens de traitement des données ne permettent pas de calculer les quantités f_{cth} par la meilleure approximation de HASTINGS, mais où





L'on a un nombre assez important de distributions de même effectif n à tester, il peut être avantageux d'établir préalablement un "diagramme de KOLMOGOROV et SMIRNOV" sur du papier probit.

Pour mémoire, l'échelle des ordonnées de ce papier est telle que les fonctions de répartition des distributions normales sont représentées par des droites (droites de HENRY).

A partir de la droite représentative de la fonction de répartition de la loi normale réduite et des valeurs critiques théoriques $e_{0,05}$ et $e_{0,01}$, il est possible, en effet, de tracer sur ce papier les deux courbes, la première située à une distance $+e$ de la droite, la seconde à une distance $-e$, qui délimitent l'aire du diagramme à l'intérieur de laquelle doivent se situer tous les points de coordonnées $\{t(u_i), fco(u_i)\}$ ou $\{tu_c, fcou_c\}$ d'une loi réduite observée pour que celle-ci puisse être considérée comme normale.

On remarquera au passage, sur le diagramme type ci-contre, que l'on est en fait beaucoup plus exigeant ici en prenant le seuil 0,05 que le seuil 0,01, les écarts théoriques admissibles augmentant avec $(1-\alpha)$.

2.3.2. - Signification des coefficients de corrélation et de régression.

D'après la formule (10)

$$r_{xu}^2 = b_{xu} \cdot b_{ux}$$

tester si r_{xu} est significativement différent de zéro revient à tester si b_{xu} ou b_{ux} le sont aussi, c'est-à-dire à tester si les droites de régression sont ou non parallèles aux axes de coordonnées.

Deux méthodes peuvent être utilisées.

PREMIERE METHODE.

Considérons la régression $X = f(U)$ et la variance à $(n-2)$ degrés de liberté de la variable dépendante $X(U)$:

$$s_{X(U)}^2 = \sum_n \{x_i - x(u_i)\}^2 / n-2 \quad (44)$$

On démontre que la variance du coefficient de régression b_{xu} peut être estimée par :

$$s_{b_{xu}}^2 = s_{x(u)}^2 / (n-1) \cdot s_u^2 \quad (45)$$

En posant $b_{xu} = b$, $s_{bxu} = s_b$ et $r_{xu} = r$, et en se souvenant que :

$$x(u_i) = \bar{x} + b_{xu} \cdot (u_i - \bar{u})$$

$$b_{xu} = r_{xu} \cdot (s_x / s_u)$$

on peut donc estimer la variable $t = b/s_b$

$$t = b/s_b = (r/\sqrt{1-r^2}) \cdot \sqrt{n-2} \quad (46)$$

qui permet à l'aide de la table de STUDENT-FISHER de tester si b est significativement différent de zéro pour $(n-2)$ degrés de liberté, à différents seuils $(1 - \alpha)$ (cf. le paragraphe 221 et la table 3 de l'annexe 1 pour le test en question)

Comme r_{xu} ne privilégie ni X , ni U , on a de plus

$$b_{xu}/s_{bxu} = b_{ux}/s_{bux} \quad (47)$$

En d'autres termes, la formule (46) permet de tester à la fois si b_{xu} , b_{ux} et r_{xu} sont significativement différents de zéro.

Réalisation du test.

En réalité, on peut ne pas calculer t : une table des valeurs limites de $|r|$ aux seuils $\alpha = 0.05, 0.01$ et 0.001 a été dressée à partir de la formule suivante déduite de la formule (46) :

$$|r_\alpha| = |t_\alpha| / \sqrt{n-2+t_\alpha^2} \quad (48)$$

Conférez la table 7 à la fin de l'annexe 1.

Par exemple, pour 25 degrés de liberté, si l'on observe $|r| > 0,4869$ celui-ci n'a qu'une probabilité $\alpha = 0.01$ de réapparaître avec la même intensité (mais a, par contre, la probabilité $(1-\alpha) = 0.99$ d'être significativement différent de zéro).

Cependant, lorsque les moyens de calculs le permettent, on peut aussi déterminer directement la probabilité $P'(t)$ correspondant à la valeur observée t (formule 46) à l'aide des formules approchées indiquées à l'annexe 1 paragraphe 222.

Il est facile, ensuite de comparer $P'(t)$ à $\alpha = 0.05, 0.01, 0.001, \dots$ etc.

DEUXIEME METHODE

FISHER a montré que la variable

$$z = (1/2) \cdot L_N \cdot (1+r)/(1-r) \quad (49)$$

est distribuée presque normalement quelles que soient p , la valeur exacte de r et n l'effectif des couples (X,U) , et, par ailleurs, que la variance de z est

$$s_z^2 = 1/(n-3) \quad (50)$$

Par cette méthode on peut donc aisément déterminer l'intervalle de confiance de r ou comparer r à 0 ou à une valeur théorique $r \neq 0$, ou encore comparer deux ou plusieurs valeurs observées de r .

Intervalle de confiance de r .

Soient donc $z = (1/2) \cdot L_N \cdot (1+r)/(1-r)$, α le seuil de probabilité choisi, z_m et z_M les limites inférieure et supérieure de l'intervalle de confiance de z . Nous avons :

$$\begin{aligned} z_m &= z - t^\alpha \cdot s_z = z - t^\alpha / \sqrt{n-3} \\ z_M &= z + t^\alpha \cdot s_z = z + t^\alpha / \sqrt{n-3} \end{aligned} \quad (51)$$

et, en utilisant la fonction inverse,

$$r = thz = (e^{2z} - 1) / (e^{2z} + 1) \quad (52)$$

$$r_m = thz_m \quad (53)$$

$$r_M = thz_M$$

Du fait du caractère quasi-normal de la distribution de z , les limites z_m et z_M peuvent être calculées quel que soit l'effectif n : en effet t_α est égal : 1.96, 2.58 ou 3.29 pour les seuils classiques $\alpha = 0.05, 0.01$ ou 0.001 .

Comparaison de r à une valeur théorique r_0 .

Soient z et z_0 les transformées de r et r_0 .

Calculons

$$t = (z - z_0) / s_z = (z - z_0) \cdot \sqrt{n-3} \quad (54)$$

On peut aisément, à l'aide de la table des valeurs de la fonction de répartition (cf. le paragraphe 211 et la table 2 de l'annexe 1), ou grâce à la meilleure approximation de HASTINGS, déterminer la probabilité correspondant à la valeur observée de t ou comparer t observée à t_α .

Remarques.

1 - Lorsque le nombre de couples est faible, une meilleure estimation $\hat{\rho}$ de r peut être obtenue à partir de la formule suivante :

$$\hat{\rho} = r \{1 + (1-r^2)/2(n-3)\} \quad (55)$$

(travaux de OLKIN et PRATT, 1958, in DAGNELIE 1974)

2 - L'application de la formule (50) implique que l'on ait $n \geq 4$

2.3.3. Vérification de la linéarité des régressions.

Vérifier la linéarité des régressions consiste, plus précisément, en l'application de deux tests successifs de non-linéarité, le premier à la régression de X sur U , le second à la régression de U sur X . Ces tests se conduisent naturellement de la même façon dans les deux cas. Nous n'exposerons donc que la démarche à suivre pour tester la non-linéarité de $X = f(U)$, mais les deux tests doivent être réalisés, un lien linéaire testé au niveau de $X = f(U)$ n'excluant pas un lien non-linéaire au niveau $U = g(X)$.

2.3.3.1. Principe du test.

Le test repose sur la comparaison de deux variances de X estimées grâce à

une mise en classes préalable des données sur la variable agissante U, et représentant,

- la première, le carré moyen des écarts à la droite de régression des moyennes par classe de la variable dépendante X,
- la seconde, le carré moyen des écarts par rapport aux moyennes par classe des valeurs observées de X dans les classes.

2.3.3.2. Conduite du test.

Le test comportera les étapes suivantes :

1°/ - la mise en classes de U (k classes). Ce travail peut, en réalité, avoir été déjà entrepris pour le test de normalité de U (cf paragraphe 2.3.2.1) ;

2°/ - le calcul :

a) - des moyennes par classe de U. Soit \bar{u}_c la moyenne de la c^{eme} classe de U ;

b) - des moyennes correspondantes de X, observées \bar{x}_c et estimées \bar{x}'_c à partir de l'équation de la droite de régression

$$\bar{x}'_c = \bar{x} + b_{xu} \cdot (\bar{u}_c - \bar{u}) \quad (56)$$

c) - des variances s_D^2 et s_E^2 , carrés moyens des quantités Q_D et Q_E à (k-2) et (n-k) degrés de liberté, éléments de l'équation des sommes des carrés des écarts

$$Q = Q_C + Q_E = Q_R + Q_D + Q_E \quad (57)$$

établie sur l'identité

$$(x_i - \bar{x}) \equiv (\bar{x}_c - \bar{x}) + (x_{yc} - \bar{x}_c) \equiv (\bar{x}'_c - \bar{x}) + (\bar{x}_c - \bar{x}'_c) + (x_{yc} - \bar{x}_c) \quad (58)$$

identité dans laquelle x_i , qui est la i^{eme} observation de X, correspond à u_i qui est aussi la y^{eme} observation de U dans la c^{eme} classe de U. Il en résulte que :

$$u_i \equiv u_{yc} \longrightarrow x_i \equiv x_{yc}$$

Ces sommes des carrés des écarts sont donc données par :

$$Q = \sum_n (x_i - \bar{x})^2 = \sum_{c=1}^k \sum_{y=1}^{nuc} (x_{yc} - \bar{x})^2 \quad (59)$$

$$Q_C = \sum_{c=1}^k nu_c \cdot (\bar{x}_c - \bar{x})^2 \quad (60)$$

$$Q_E = \sum_{c=1}^k \sum_{y=1}^{nuc} (x_{yc} - \bar{x}_c)^2 \quad (61)$$

$$Q_R = \sum_{c=1}^k nu_c \cdot (\bar{x}'_c - \bar{x})^2 = b_{xu}^2 \cdot \sum_{c=1}^k nu_c \cdot (\bar{u}_c - u)^2 \quad (62)$$

$$Q_D = \sum_{c=1}^k nu_c \cdot (\bar{x}_c - \bar{x}'_c)^2 \quad (63)$$

Les carrés moyens s_X^2 , s_C^2 , s_E^2 , s_R^2 , et s_D^2 qui leur correspondent représentent différentes estimations de la variance σ^2 . Ainsi :

$s_X^2 = Q/(n-1)$ est, pour mémoire (cf paragraphe 22), l'estimation de la variance de la distribution marginale de X,

$s_C^2 = Q_C/(k-1)$, appelée "variance inter-classe", est une autre estimation de σ_X^2 , mais seulement si X n'est pas liée à U ;

$s_E^2 = Q_E/(n-k)$, appelée "variance intra-classe", est encore une autre estimation de σ_X^2 , mais indépendante, elle, de toute hypothèse sur l'influence de U sur X ;

$s_R^2 = Q_R/1$, qui représente la part de la variation de X imputable au "facteur systématique régression linéaire", est aussi une estimation de σ_X^2 si la régression linéaire n'est pas significative ;

$s_D^2 = Q_D/(k-2)$, enfin, est une dernière estimation de σ_X^2 si l'hypothèse de la linéarité est exacte ;

3°/ le calcul et la comparaison à F_{α} pour (k-2) et (n-k) degrés de liberté du rapport

$$F_{NI} = s_D^2 / s_E^2 \quad (64)$$

D'après ce qui précède, si $F_{NL} \geq F_{\alpha}$, l'hypothèse de la non-linéarité peut être acceptée au seuil α et la régression ne pourra donc pas être considérée comme linéaire (pour le test F voir le paragraphe 3 de l'annexe 1 et les tables 4 à la fin de la même annexe).

Remarque.

Une réalisation entièrement automatique du test de F_{NL} est possible à l'aide des formules qui permettent de déterminer la probabilité α correspondant à une valeur observée de F (pour ces formules cf le paragraphe 32 de l'annexe 1).

Si α observée < 0.05 , c'est-à-dire si $(1-\alpha) > 0.95$ on acceptera l'hypothèse de la non-linéarité.

3 - CORRELATION DANS LE CAS OU X ET U SONT DISTRIBUEES NORMALEMENT MAIS NE SONT PAS LIEES LINEAIREMENT.

Au cours des opérations qui viennent d'être décrites deux éventualités peuvent se produire :

- ou bien r_{xu} n'est pas significativement différent de zéro, mais on se demande cependant si les deux variables ne sont pas liées non-linéairement ;
- ou bien r_{xu} est significativement différent de zéro mais l'une des deux régressions, ou les deux, sont significativement non-linéaires.

Dans ces conditions, l'application aux résultats observés de la loi normale à deux variables pour le calcul de r_{xu} , b_{xu} , b_{ux} n'est plus possible. Deux nouveaux coefficients - les rapports de corrélation de X sur U et de U sur X - doivent être utilisés pour caractériser l'intensité de la liaison.

La conduite des calculs est la même dans les deux cas.

3.1 - Rapport de corrélation de X sur U.

Considérons donc la régression $X = f(U)$ et le classement de X sur U.

Nous avons vu au paragraphe précédent (parag.2332) que l'on pouvait décomposer Q, la somme des carrés des écarts de X, en deux composantes Q_C et Q_E telles que

$$Q = Q_C + Q_E$$

Nous avons vu également que $s_C^2 = Q_C/(k-1)$ et $s_E^2 = Q_E/(n-k)$ représentaient deux estimations indépendantes de σ_X^2 , la première si X n'était pas liée à U, la seconde indépendamment de toute hypothèse sur la nature de la liaison.

La comparaison de s_C^2 à s_E^2 permet ainsi de tester la signification de la liaison $X = f(U)$ mise en évidence par le classement, indépendamment de toute hypothèse sur la forme de la liaison.

Si

$$F = s_C^2 / s_E^2 \geq F_\alpha \quad (65)$$

pour $(k-1)$ et $(n-k)$ degrés de liberté, on conclura que les variables X et U sont liées significativement et que la liaison $X = f(U)$ explique

$$100 \cdot \frac{QC}{Q} = (\eta_{XU}^2)\% \quad (66)$$

des variations de X

La quantité $\eta_{XU}^2 = \frac{QC}{Q}$ est appelée "rapport de corrélation".

Elle est toujours positive et comprise entre 0 et 1.

Remarque.

On démontre qu'en théorie $\eta_{XU}^2 \geq r_{XU}^2$, l'égalité n'étant réalisée que lorsque la relation $X = f(U)$ est rigoureusement linéaire.

Par l'application de la méthode d'estimation de η_{XU}^2 qui vient d'être exposée, il arrive cependant que l'on observe des valeurs de η_{XU}^2 inférieures à celles des carrés des coefficients de corrélation correspondant. Ceci ne se produit, à notre connaissance, que lorsque n, l'effectif des couples (x_i, u_i) est faible ($n < 40$).

L'explication de cette situation paradoxale réside simplement dans le fait que l'estimation de η_{XU}^2 repose, en partie, sur des décisions empiriques : celles concernant la mise en classe des données (choix du nombre de classes, définition de l'étendue de référence), tandis que l'estimation de r_{XU}^2 ne dépend que des valeurs observées des couples (x_i, u_i) .

L'interprétation qu'il semble falloir donner à de telles situations est que η_{XU}^2 est voisin de r_{XU}^2 ou que les deux coefficients ne sont pas significativement

différents de zéro.

3.2. Rapport de corrélation de U sur X.

L'estimation de η_{UX}^2 , rapport de corrélation de U sur X doit toujours être faite après celle de η_{XU}^2 , que ce dernier soit ou non significatif.

En effet, non seulement la probabilité de mettre en évidence la liaison susceptible d'exister entre les deux variables est aussi élevée, a priori, à partir du classement de U sur X qu'à partir de celui de X sur U, mais surtout cette liaison peut très bien ne pas apparaître significative dans un cas et très hautement dans l'autre, les estimations des rapports de corrélation dépendant, aussi, de la forme réelle de la liaison, c'est-à-dire de celle du graphique des moyennes liées qui la traduit le mieux.

3.3. Linéarisation des relations $X = f(U)$ et $U = g(X)$

Montrer l'existence d'un lien entre les deux variables sans pouvoir traduire celui-ci par une fonction d'estimation particulière s'est évidemment pas très satisfaisant pour l'esprit, même si la forme générale de la liaison déduite de celle des graphiques des moyennes liées permet d'en décrire les traits essentiels et les conséquences éventuelles sur les variations de A.

Ceci est d'autant plus sensible que les régressions $X = f(U)$ et, ou $U = g(X)$ apparaissent nettement.

La linéarité des régressions est indispensable, par ailleurs, aux analyses de covariance.

S'agissant d'une façon générale des problèmes de la régression curvillinéaire, le lecteur aura intérêt à se reporter aux ouvrages traitant de la question, en particulier à ceux de DAGNELIE déjà cités (1969, paragraphique 2.10 p. 101 et suivantes, et 1970, paragraphe 18.7, p 297 et suivantes).

Le schéma général de la démarche, qui est exposé ci-dessous, n'est donné qu'à titre indicatif.

3.3.1. Principe d'une démarche générale.

Soient donc deux caractéristiques X et U distribuées normalement et liées significativement, mais de façon non-linéaire.

Soient, par ailleurs, $T(X)$ et $M(U)$ deux fonctions mathématiques particulières dont l'application à X et U permet de faire correspondre à tout couple (x_i, u_i) de valeurs observées, un couple (x'_i, u'_i) de valeurs transgénérées telles que :

$$\begin{aligned} x'_i &= T(x_i) \\ u'_i &= M(u_i), \end{aligned} \quad (67)$$

$T(X)$ et $M(U)$ étant choisies, — soit de façon empirique, soit sur des bases théoriques—, de telle sorte que les relations $X' = f(U')$ ou $U' = g(X')$ puissent apparaître linéaires.

Corrélons alors X' et U' en suivant les étapes décrites au paragraphe 2.

Si, à l'issue des tests, X' et U' apparaissent toujours liées significativement, mais linéairement cette fois-ci, et si d'autre part leurs distributions marginales sont encore normales, alors $T(X)$ et $M(U)$ seront considérées comme des fonctions de linéarisation de la régression.

Leur application aux valeurs observées de X et de U permettra d'estimer X à partir de U , et réciproquement, à l'aide des équations des droites de régressions des données transgénérées.

$$\begin{aligned} x'_i (u'_i) &= \bar{x}' + b_{x' u'} \cdot (u'_i - \bar{u}') \\ u'_i (x'_i) &= \bar{u}' + b_{u' x'} \cdot (x'_i - \bar{x}') \end{aligned} \quad (68)$$

et des fonctions inverses des fonctions de transgénération $T^{-1}(X)$ et $M^{-1}(U)$.

3.3.2. Fonctions de linéarisation utilisées couramment.

Pratiquement, un petit nombre seulement de fonctions sont utilisées pour linéariser. Ces fonctions sont les mêmes que celles mises en oeuvre pour répondre à un autre problème : celui de la normalisation des caractéristiques.

On se reportera donc au tableau ci-après présentant ces fonctions de transgénération .

4 - CORRELATION DANS LE CAS OU X ET U NE SONT PAS DISTRIBUEES NORMALEMENT.

Une autre éventualité peut se présenter au cours des opérations de corrélation

de X et U : la non-normalité de l'une des deux caractéristiques, ou même des deux.

Lorsqu'il en est ainsi, le problème se pose évidemment de savoir si l'on peut transformer la ou les variables pour que leurs distributions respectives deviennent normales (sans modifier le caractère linéaire de la liaison si celle-ci semblait déjà l'être, en ou linéarisant en même temps cette régression si elle semblait ne pas l'être).

Ce problème est en réalité étroitement lié au précédent ainsi qu'à celui de l'homogénéisation des variances (l'homogénéité des variances des caractéristiques homologues de deux ou plusieurs agrosystèmes est l'une des conditions requises pour leur comparaison comme on le verra au chapitre suivant).

4.1. Fonctions de transgénération classiques.

Le tableau ci-après regroupe les principales fonctions de transgénération susceptibles d'être utilisées, en précisant les conditions de leur application.

4.2 Règle empirique de choix d'une fonction de transgénération stabilisant les variances.

Lorsqu'aucune raison théorique ne permet de choisir une fonction de transgénération particulière et lorsque l'on dispose d'un nombre suffisant de séries de données d'une même caractéristique X, on peut, si une relation linéaire existe entre le logarithme de la moyenne et celui de la variance, utiliser comme fonction de transgénération

$$T(X) = X^{1-b/2} \quad (b \neq 2) \quad (69)$$

où b est la pente de la droite de régression établie sur les logarithmes de la variance et de la moyenne :

$$\log s_X^2 = b \cdot \log(\bar{X}) + \log k \quad (70)$$

Cas particuliers

1°/ si $b = 1$, $s_X^2 = k \cdot \bar{X}$ et $T(X) = \sqrt{X}$ (71)

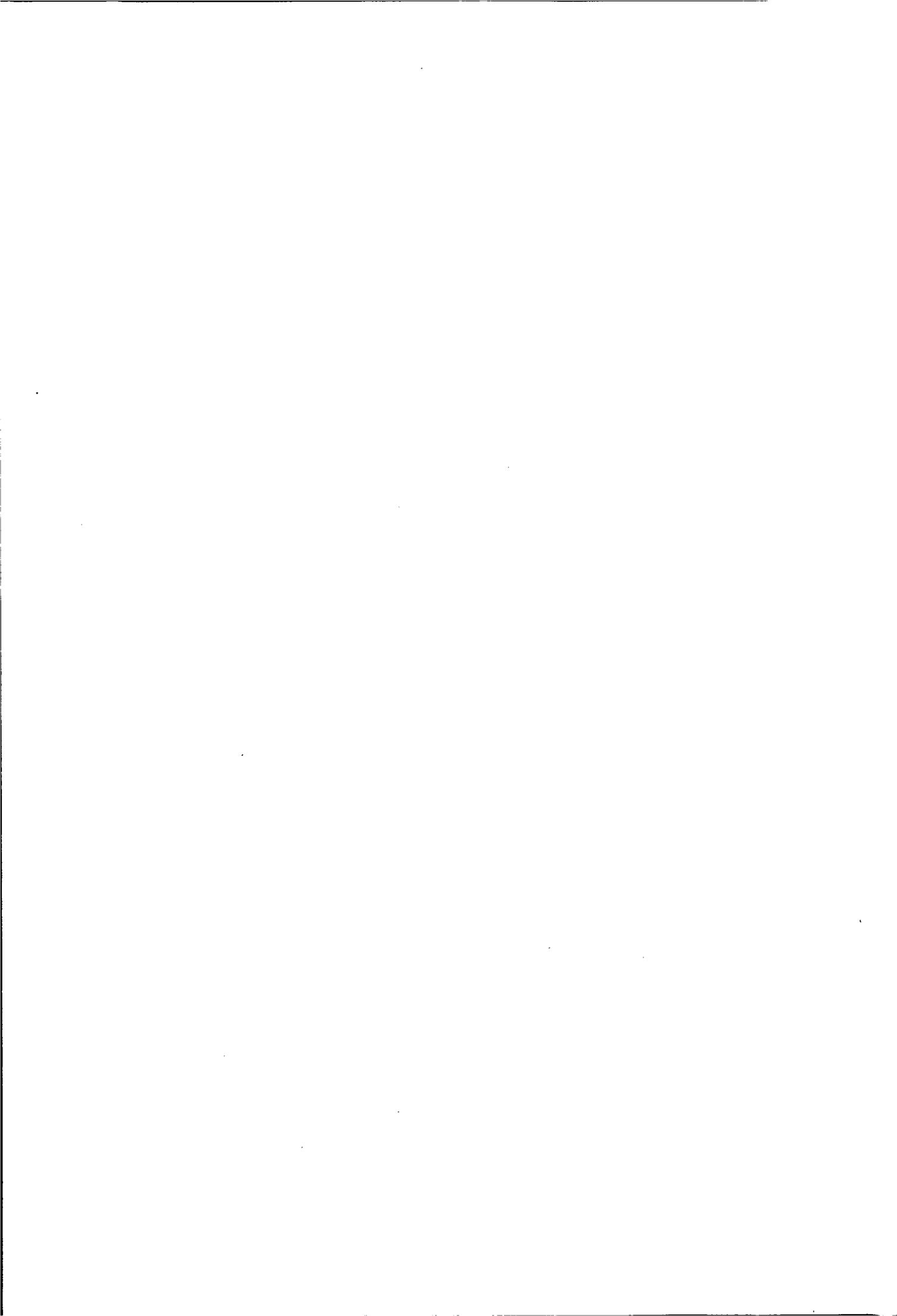
2°/ si $b = 2$, $s_X^2 = k \cdot (\bar{X})^2$ et $T(X) = \log_{10} X$ (ou $L_N X$) (72)

PRINCIPALES FONCTIONS DE TRANSGENERATION $X' = T(X)$ UTILISEES POUR LINEARISER LES REGRESSIONS, NORMALISER LES DISTRIBUTIONS OU STABILISER LES VARIANCES.

Fonctions de transgénération	Conditions d'application
$X' = \sqrt{X}$ <p>Cas particuliers :</p> $\left. \begin{aligned} X' &= \sqrt{X+1/2} \\ X' &= \sqrt{X+3/8} \end{aligned} \right\}$ $X' = \sqrt{X+a}$	<p>X est une variable aléatoire discontinue dont la distribution est proche de celle de l'une des lois de POISSON.</p> <p>Les valeurs observées de X, variable aléatoire discontinue, sont faibles.</p> <p>La valeur de la constante a est recherchée empiriquement</p>
$\left. \begin{aligned} X' &= L_N X \\ X' &= \log_{10} X \end{aligned} \right\}$ $X' = \log_{10} (X+1)$ $X' = \log_{10} (X+a)$	<p>X est distribuée selon une loi log-normale. Ou bien, X est une fonction exponentielle de U (U n'est pas transgénérée simultanément). Ou bien, X est une fonction puissance de U (U est aussi transgénérée par $L_N U$ ou $\log_{10} U$).</p> <p>Les valeurs observées de X sont faibles; $X' = 0$ si $X = 0$</p> <p>La valeur de la constante a est recherchée empiriquement.</p>
$X' = \arcsin \sqrt{x/n}$	<p>X est une variable binomiale. Ou bien, X est un rapport compris entre 0 et 1, qui peut être continu.</p>

Remarque :

Pour aucune de ces fonctions on n'observera $\bar{x}' = T(\bar{x})$



CORRELATION ENTRE X ET U

n couples de valeurs
 . observées (x_i, u_i)
 ou
 . transgénérées (x'_i, u'_i)
 $n \geq 4$

n ≥ 5 ?

n ≥ 20 ?

X et U normales ?

On applique à chaque variable non-normale une fonction de transgénération et on recommence l'analyse.

X et U normales ?

η^2_{xu} significatif ?

$r_{xu} \neq 0$?

$r_{xu} \neq 0$?

$r_{xu} \neq 0$?

η^2_{ux} significatif ?

η^2_{ux} significatif ?

U = g(X) linéaire ?

X = f(U) linéaire ?

U = g(X) linéaire ?

X et U ne sont pas liées significativement.

X et U sont liées significativement mais non linéairement. Il faut transgénérer X et, ou U et recommencer l'analyse si l'on désire estimer X en fonction de U et U en fonction de X. La transgénération portera :

- . plutôt sur U si seule $X = f(U)$ n'est pas linéaire,
- . ou si seul η^2_{xu} est significatif,
- . plutôt sur X dans les cas contraires,
- . sur X et U si les deux régressions ne sont pas linéaires ou si les deux η^2 sont significatifs.

X et U sont liées significativement et linéairement :

$$x(u)_i = x + b_{xu} (u_i - \bar{u})$$

$$u(x)_i = u + b_{ux} (x_i - \bar{x})$$

X et U sont liées significativement mais la linéarité, des régressions ne peut être testée.

X et U ne sont pas liées significativement de façon linéaire. Aucune autre étude n'est possible.

X et U sont liées significativement mais la linéarité de leurs régressions et la normalité de leurs distributions n'ont pu être testées

4.3. Remarques sur certains aspects de la non-linéarité.

1 - Du fait de l'influence des modalités de mise en classes des données sur les moyennes liées des caractéristiques et sur les résultats des tests de non-linéarité et de signification du rapport de corrélation (ainsi que sur la valeur de ce dernier), il semble qu'il serait utile de reprendre l'ensemble de ces calculs avec différents nombres de classes encadrant le nombre de classes donné par k_2 (formule 42).

En effet, l'application de fonctions de transgénération peut compliquer par la suite les problèmes d'estimation des valeurs d'une caractéristique en fonction de celles de l'autre.

2 - Lorsque le nombre de classes est suffisant ($n \geq 6$ en pratique), la mise en classe d'une variable sur l'autre peut faire apparaître une classe vide dans la partie centrale de la distribution (au niveau de la quatrième classe, par exemple, dans le cas où $n = 6$). Il y a lieu alors de se demander si le caractère bi-modale de la distribution ne correspondrait pas à un échantillonnage réalisé sur deux populations de la variable utilisée pour le classement.

3 - Dans le même ordre d'idée, l'établissement d'un diagramme des couples observées (x_i, u_i) peut toujours être utile. Une valeur "aberrante" de l'une des deux caractéristiques, voir la présence de deux populations peut apparaître de cette façon.

5. ORGANISATION GENERALE DES CALCULS.

Désirant savoir si un lien pourrait exister entre les deux caractéristiques observées X et U , et, dans l'affirmative, désirant le caractériser par son intensité et par une ou deux fonctions de régression (une si l'on peut considérer que l'une des deux caractéristiques est la variable agissante, deux s'il est impossible de les départager de ce point de vue) on peut résumer les opérations analytiques qui viennent d'être passées en revue par l'organigramme général ci-contre.

CONCLUSION.

Ce tableau permet de se rendre compte, finalement, que si l'on désire caractériser correctement la partie observée de A, le nombre de sites d'observations devra être au minimum de 20.

En réalité, il s'agit-là d'un effectif très faible, car de nombreuses caractéristiques "sol" présentent des coefficients de variation élevés dus à leur micro-hétérogénéité dans le sol. Les dimensions de cette micro-hétérogénéité du sol sont difficiles à cerner au demeurant.

Par ailleurs, les déterminations analytiques de certaines caractéristiques physiques ou chimiques manquent aussi, parfois, de précision.

Chaque fois que les moyens disponibles le permettront il sera donc préférable de prendre un nombre de sites d'observation très supérieur à 20, et, ou de faire faire en double les analyses chimiques ou physiques les plus imprécises.

Les situations où le nombre de sites d'observations est insuffisant pour caractériser de façon très complète la partie observée de A, sont cependant les plus nombreuses. Elles ne sont pas à rejeter pour autant : chaque fois qu'un ensemble cohérent et homogène d'observations a été constitué, il est toujours intéressant a priori d'en examiner le caractère systématique.