# THÈSES DE L'UNIVERSITÉ PARIS-SUD (1971-2012)

## **SYLVIA RICHARDSON**

Processus spatialement dépendants : convergence vers la normalité, tests d'association et applications, 1989

Thèse numérisée dans le cadre du programme de numérisation de la bibliothèque mathématique Jacques Hadamard - 2016

Mention de copyright :

Les fichiers des textes intégraux sont téléchargeables à titre individuel par l'utilisateur à des fins de recherche, d'étude ou de formation. Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale.

Toute copie ou impression de ce fichier doit contenir la présente page de garde.



63846

ORSAY n° d'ordre : 3559

## UNIVERSITE PARIS SUD

## Centre d'Orsay

## THESE

De Doctorat d'Etat Es Sciences Mathématiques

présentée pour obtenir le grade de

### **DOCTEUR ES-SCIENCES**

par

Sylvia RICHARDSON

<u>Sujet</u> : PROCESSUS SPATIALEMENT DEPENDANTS : CONVERGENCE VERS LA NORMALITE, TESTS D'ASSOCIATION ET APPLICATIONS.

Soutenue le 19 Octobre 1989 devant la Commission d'examen composée de :

MM. DACUNHA-CASTELLE BRETAGNOLLE HUBER LELLOUCH THOUVENOT Didier, Président Jean Catherine Joseph Jean-Paul

#### ABSTRACT

The first part of this thesis is concerned with the study of ergodic properties of stopping time transformations,  $T^v$  where T is an automorphism and v a transformation with values in Z. The demonstrations use cutting and stacking methods. Results when T is Bernoulli have been published in Z. Warhs. Verv. (vol 48, 1979, with K.M. Wilkinson) and density results when T is ergodic in the Proc. London Math. Soc (vol 43, 1981).

The second part of this thesis is composed of several papers, in probability and in statistics, about spatial processes and the testing of association. A first paper, published in J.A.P. (vol 18, 1981 with D. Hémon) derives the asymptotic variance of the correlation coefficient between two independent lattice processes. The speed of convergence in the Central Limit Theorem for weakly dependent processes in  $Z^d$  is the subject of the second paper (published in Z. Warhs. Verw (vol 66, 1984, with X. Guyon). The speed is dependent on the dimension d, on moments conditions and on mixing conditions.

The development of tests of association between spatial processes having a controlled Type I error form the statistical part of this thesis. These tests were proposed in a paper published in Biometrics (vol 45, 1989, with P. Clifford and D. Hémon). Their power and an extension to the testing of partial correlations was further studied (to appear in Statistics in Medecine). The applications discussed concern the field of geographical correlations in epidemiology (Int. J. Epidemiol vol 16, 1987, with I. Stücker and D. Hémon).

<u>Keywords</u>: Ergodic properties of stopping times, spatial processes, Central Limit Theorem, Speed of convergence, Correlation coefficient, Tests of association, Effective degrees of freedom, geographical epidemiology.

<u>A.M.S. code</u> : 28 D 05, 60 F 05, 60 G 60, 62 E 25, 62 H 20, 62 M 30, 62 P 10, 92 A 15.

# AVANT-PROPOS et REMERCIEMENTS

Mes recherches, commencées dans le département de Mathématiques de l'Université de Nottingham, ont d'abord porté sur un sujet de théorie ergodique suggéré par K.M. Wilkinson, mon directeur de Ph.D. J'ai étudié les propriétés ergodiques des temps d'arrêt auxquels ont été imposées des conditions telles que ces transformations soient des automorphismes d'un espace de Lebesgue continu. Ces temps d'arrêt peuvent aussi être considérés comme des temps d'arrêt probabilistes qui sont définis sur une suite de variables, mais la classe des temps d'arrêt probabilistes est strictement plus large.

Je me suis ensuite intéressée à d'autres types de processus, les processus spatialement dépendants et plus particulièrement à l'étude de tests d'association. Mes recherches ont été menées à la fois du point de vue probabiliste et du point de vue statistique. Elles ont aussi comporté un domaine d'application, les études de corrélations écologiques. Ces recherches ont été faites dans le Laboratoire de Statistiques Médicales de l'Université de Paris V et dans l'Unité 170 de l'INSERM en collaboration avec D. Hémon (U.170), X. Guyon (Université de Paris I) et P. Clifford (Université d'Oxford).

Je tiens tout d'abord à remercier Didier Dacunha-Castelle, Professeur, qui me fait l'honneur de présider le jury de cette thèse ainsi que Jean Bretagnolle, Professeur, qui a si gentiment accepté d'en assurer la direction.

Je dois remercier tout particulièrement mes amis et collègues avec qui j'ai signé des articles :

Keith Wilkinson qui m'a introduit à la théorie ergodique ; Xavier Guyon qui m'a fait comprendre les beautés et les difficultés des processus dans Z<sup>2</sup> ; Peter Clifford dont la rigueur et les intuitions statistiques m'ont beaucoup apporté ;

Denis Hémon qui, dès le départ, a su me montrer l'intérêt et les différents aspects de ce sujet et m'a soutenu activement tout au long de mes recherches.

Qu'ils trouvent ici l'expression de ma reconnaissance.

Je remercie chaleureusement Catherine Huber, Professeur, et Joseph Lellouch, Directeur de recherches, des encouragements qu'il m'ont toujours prodigués et d'avoir bien voulu faire partie du jury.

Je voudrais aussi remercier mes deux rapporteurs pour la commission des thèses, Monsieur J.P. Thouvenot, Directeur de recherches, et Monsieur E. Bolthausen, Professeur ; Monsieur J.P. Thouvenot a également eu la gentillesse de faire partie du jury.

Il me faut encore remercier les membres du Laboratoire de Statistiques Médicales de l'Université de Paris V pour le soutien qu'ils m'ont apporté ; ainsi que les chercheurs et techniciens de l'Unité 170 de l'INSERM qui ont beaucoup contribué à faire aboutir ce travail.

Ces remerciements ne seraient pas complets si je ne disais pas que c'est à Philippe Lazar, Directeur Général de l'INSERM, que je dois l'intérêt qui ne m'a pas quitté pour l'épidémiologie géographique. Qu'il trouve ici l'expression de ma gratitude.

## TABLE DES MATIERES

Première partie : Processus ergodiques : études des propriétés ergodiques des temps d'arrêt.

- 1. Temps d'arrêts et "tours"
- 2. Approximation des temps d'arrêts

<u>Deuxième partie</u> : Processus spatialement dépendants : convergence vers la normalité, tests d'association et applications.

- 1. Approche probabiliste
  - 1.1 Variance asymptotique de  $r_{XY}$
  - 1.2 Normalité asymptotique de rXY
  - 1.3 Vitesse de convergence du Théorème de la Limite Centrale
- 2. Tests d'association
  - 2.1 Approximation de la variance de  $r_{XY}$
  - 2.2 Tests d'association modifiés
  - 2.3 Puissance des tests modifiés
  - 2.4 Tests modifiés des corrélations partielles
- 3. Applications en épidémiologie à l'étude des corrélations écologiques
- 4. Conclusion

# <u>Première partie</u> : Processus ergodiques : études des propriétés ergodiques des temps d'arrêt.

Les définitions et les résultats classiques en théorie ergodique utilisés dans cette partie peuvent se trouver dans les ouvrages de Halmos (5) et Walters (9).

Les temps d'arrêt considérés sont de la forme  $\tau^v$  où  $\tau$  est un automorphisme (c'est-àdire une transformation bijective préservant la mesure) d'un espace de Lebesgue ( $\Omega, B, \mu$ ) et v une transformation mesurable à valeur dans N, telle que les ensembles  $\tau^v \{v = n\}$  forment une partition de  $\Omega$ . Dans ce cas la transformation  $\tau^v$  est-elle même un automorphisme. Cette classe de transformations inclut les temps de retour  $\tau^v_D$  à un ensemble D et la restriction de  $\tau^v_D$  à l'espace de Lebesgue (D,B<sub>D</sub>, $\mu_D$ ) coïncide avec  $\tau_D$  l'automorphisme induit.

Ces temps d'arrêt ont été introduits par Neveu (6) qui en a donné une décomposition en termes de produit de temps de retour associés à une suite décroissante d'ensemble  $\{D_i\}$ . Quand l'espérance de v, E(v), est finie, on peut montrer qu'elle est égale à un entier k qui coïncide avec le nombre d'ensembles  $\{D_i\}$ , i = 0..., k-1 dans la décomposition de Neveu. L'entropie de  $\tau^v$  est reliée à celle de  $\tau$  via l'espérance de v.

De nombreux auteurs se sont intéressés aux propriétés des transformations induites : Friedman (3), Friedman et Ornstein (4), Ornstein et Smorodinski (7). Nous avons cherché à voir comment leurs résultats pouvaient se généraliser aux temps d'arrêt ; une première différence étant que  $\tau^v$  n'est pas nécessairement ergodique quand  $\tau$  l'est, bien que  $\tau_D$  le soit également.

Mes recherches se sont portées essentiellement sur deux aspects. Dans un premier temps, j'ai étudié les temps d'arrêt liés à une transformation  $\tau$  linéaire par morceaux, définie sur [0,1] par une "tour". Dans un deuxième temps, la transformation  $\tau$  est seulement supposée ergodique.

#### 1. Temps d'arrêts et "tours"

Une tour, définie sur [0,1], est un ensemble de colonnes composées d'intervalles disjoints de même mesure, empilés les uns sur les autres. A l'intérieur des colonnes, une application partielle bijective est définie de manière à ce que chaque intervalle soit linéairement appliqué sur l'intervalle situé directement au-dessus. Au sommet de chaque colonne, la transformation est définie par récurrence en utilisant la méthode du "cutting and stacking". A l'étape n, la tour existant  $T_n$  est remplacée par une nouvelle tour  $T_{n+1}$  constituée par une copie réduite de moitié de  $T_n$ :  $1/2T_n$ , à la base (cutting) et par une copie de  $T_n$  au dessus de chaque colonne de  $1/2T_n$  (stacking). Ainsi à l'étape n, la transformation est étendue linéairement par morceaux à la moitié du sommet. Comme la mesure du sommet tend vers zéro quand n tend vers l'infini, la transformation est éventuellement définie sur l'ensemble des intervalles qui composent la tour. Elle préserve la mesure et elle est ergodique.

Cette classe de transformations, introduite par Chacon (1), a été largement étudiée (voir Friedman (2)). Friedman (3) et Shields (8) ont montré que les transformations définies par des tours sont markoviennes (Markov shifts). Il est facile de poser des conditions telles qu'elles soient, de plus, mélangeantes. Par isomorphisme, cette classe inclut donc les transformations de Bernoulli (Bernoulli shifts) d'entropie finie quelconque (Friedman et Ornstein (4)).

Quand  $\tau$  est une transformation définie par une tour T, on peut suivre explicitement l'action de  $\tau^{v}$  dans la tour quand les intervalles de T sont des sous-ensembles de la partition  $E_i = D_{i-1} - D_i$ ,  $0 \le i \le k$ , associée au temps d'arrêt v, E(v) = k. Pour certains intervalles (au plus k) dans chaque colonne, la transformation de  $\tau^{v}$  n'est pas définie. Ceci conduit à décomposer la tour T en une union de k tours  $\{T_i\}$  disjointes. A l'intérieur de chaque  $T_i$ , les intervalles sont transformés par  $\tau^{v}$ . L'action de  $\tau^{v}$  sur les sommets des colonnes peut être identifiée grâce à des fonctions associées à v, définies sur la partition des intervalles de la tour. Ces fonctions caractérisent de manière unique un ensemble de permutations qui indiquent quelle tour placer au dessus de quelle colonne dans la procédure de "cutting and stacking". Cette construction géométrique, que nous avons nommée : "tower blocks", peut être étudiée indépendamment des temps d'arrêts. La transformation qui lui est associée est markovienne et on peut donner des conditions explicites pour qu'elle soit ergodique ou de Bernoulli. On en déduit que si  $\tau$  est de Bernoulli d'entropie finie, l'ensemble des  $\tau^{v}$  qui sont de Bernoulli est dense (par rapport à la métrique uniforme) dans la famille des transformations { $\tau^{v}$ }. Ceci généralise un théorème de Friedman (3) pour les transformations induites.

Ces résultats faisaient partie de mon doctorat (Ph.D) que j'ai soutenu en 1978 à l'Université de Nottingham (Directeur K.M. Wilkinson); ils ont été publiés dans Z. Wahrscheinlichkeitstheorie verw.(liste des travaux (1)).

#### 2. Approximation des temps d'arrêt

Dans un deuxième temps, en supposant seulement  $\tau$  ergodique, j'ai montré comment on peut approcher une transformation  $\tau^{v}$  quelconque par une transformation  $\tau^{v'}$  de même entropie et ergodique sur l'ensemble {v'  $\neq$  0}. A nouveau, ces résultats sont fondés sur une construction explicite par induction de la transformation  $\tau^{v'}$ . L'étape inductive est fournie par le théorème de Rohlin, qui donne une représentation géométrique d'une transformation ergodique  $\tau$  quelconque à l'aide de  $\tau$ -colonnes. Le théorème de Rohlin a aussi été utilisé par Friedman et Ornstein (4) dans leur construction d'une transformation mélangeante. Les fonctions associées à v qui ont été introduites dans la construction précédente sont à nouveau sollicitées et le temps d'arrêt v' est défini comme la limite d'une suite de temps d'arrêt v<sup>(r)</sup>. A chaque étape le temps d'arrêt v<sup>(r)</sup> est légèrement modifié afin que la transformation limite soit ergodique. Une fois l'ergodicité obtenue, en utilisant des résultats connus sur les transformations induites (4,7), on peut en déduire des résultats analogues concernant la densité des  $\tau^{v}$  mélangeantes ou possédant la propriété de Kolmogorov.

Ces résultats font également partie de mon Ph.D et ont été publiés dans les Proceedings of the London Mathematical Society (liste des travaux (2)).

# <u>Deuxième partie</u> : Processus spatialement dépendants : convergence vers la normalité, tests d'association et applications.

Etudier ou tester l'association entre deux variables aléatoires est un problème classique en statistique. Pearson en 1896 introduisait déjà une mesure d'association : le coefficient de corrélation empirique. Dans le travail que je vais présenter, je me suis intéressée au cas où les variables sont observées en différents points d'un espace géographique et sont des réalisations de processus stochastiques pour lesquels une dépendance entre des variables indexées par des points géographiquement proches existe.

Cette dépendance ou autocorrélation spatiale peut avoir des origines diverses. Il peut s'agir d'une propriété intrinsèque au processus lui-même, comme l'existence d'interactions entre les différents sites, de phénomènes de diffusion ou de contagion. Cette autocorrélation peut également résulter indirectement de l'influence sur la variable étudiée d'un certain nombre d'autres facteurs aléatoires ou déterministes qui varient à une échelle spatialement plus large que celle du réseau de points considérés. C'est le cas de variables définies comme une moyenne sur une unité géographique de variables dichotomiques observées par individu. Nos exemples d'applications, choisis en épidémiologie, se situent dans ce cadre d'autocorrélation crée indirectement.

Cette dépendance spatiale a de multiples conséquences en particulier sur les conditions de convergence vers la normalité dans le Théorème de la Limite Centrale ainsi que sur la vitesse de cette convergence.

Du point de vue des tests d'association, les méthodes statistiques classiques fondées sur la corrélation ou la régression ne sont pas applicables directement au cas de variables spatialement dépendantes et il est nécessaire de les modifier. Les conséquences du fait de ne pas prendre en compte l'autocorrélation dans les modèles de régression ont été passées en revue par Johnston (13) pour les séries chronologiques et Cliff et Ord (5) pour les séries spatiales. Certains auteurs ont proposé des mesures d'association qui impliquent une transformation des données par filtrage pour réduire l'autocorrélation avant d'appliquer des techniques standards : ce sont les travaux de Haugh (10) développés dans le cadre des séries chronologiques. La généralisation de cette méthode à des données spatiales irrégulières soulève des difficultés.

Mes travaux ont considéré deux aspects : l'un probabiliste et l'autre statistique. J'ai aussi développé des applications en épidémiologie.

#### 1. Approche probabiliste

#### 1.1 Variance asymptotique du coefficient de corrélation

Il s'agit d'une étude sur la variance asymptotique du coefficient de corrélation entre deux processus mutuellement indépendants, chacun étant autocorrelé spatialement. Dans le cadre des séries chronologiques, la variance et les covariances des coefficients de corrélation croisés ont été étudiés par Bartlett (2,3).

La variance asymptotique est établie pour des processus gaussiens sur un lattice et une suite croissante  $\{D_n\}$  de domaines de Z<sup>2</sup> de forme quelconque, de cardinal  $|D_n|$ .

Supposons que {X(u)} et {Y(u)}, u  $\varepsilon Z^2$ , soient deux processus gaussiens stationnaires, d'espérance nulle et de variance finie,  $\sigma^2_X$  et  $\sigma^2_Y$ . Supposons de plus que les autocovariances de X et Y : C<sub>X</sub>(u,v) = E (X(u) X(v)) et C<sub>Y</sub>(u,v) = E (Y(u) Y(v)), u,v  $\varepsilon Z^2$ , soient à décroissance exponentielle.

Soit  $r_{XY}$ , le coefficient de corrélation empirique entre {X(u)} et {Y(u)}, u  $\varepsilon D_n$ . Alors :

 $\lim_{n \to \infty} |D_n| \operatorname{var} (r_{XY}) = \sum_{v \in \mathbb{Z}^2} C_X(0,v) C_Y(0,v) / \sigma_X^2 \sigma_Y^2$ 

où 0 est l'origine et la suite  $\{D_n\}$  peut être de forme quelconque du moment qu'asymptotiquement son bord  $\delta D_n$  soit négligeable par rapport à  $D_n$ . Cette condition sur la suite  $\{D_n\}$  est souvent retrouvée quand les processus spatiaux stationnaires sont étudiés. Dans le cas de processus markoviens conditionnels ou autorégressifs simultanés des plus proches voisins, une expression explicite de la variance en fonction des paramètres des processus est donnée. Ces calculs numériques illustrent bien l'influence considérable de l'autocorrélation sur la variance du coefficient de corrélation.

Ces travaux ont fait l'objet de deux publications dans le Journal of Applied Probability et la Revue de Statistique Appliquée en collaboration avec D. Hémon., (liste des travaux (3), (4)).

#### 1.2 Normalité asymptotique de rXY sous des hypothèses de mélange

Les résultats énoncés dans ce paragraphe ne sont pas publiés. On peut les considérer comme une annexe aux publications des paragraphes 1.1 et 1.3.

Dans le cadre de variables stationnaires dépendantes sur Z<sup>d</sup>,  $d \ge 1$ , les théorèmes de la limite centrale (T.L.C.) sont habituellement obtenus par l'imposition de conditions de mélange qui assurent en quelque sorte une faible dépendance pour des  $\sigma$ -algèbres engendrées par des variables éloignées, (Ibragimov et Linnick (12) et Hall et Heyde (8) pour d = 1, Bolthausen (4) pour d  $\ge$  1).

En particulier le champ {X(u)} est défini comme  $\alpha$ -mélangeant (mélange fort) si pour toutes parties A et B de Z<sup>d</sup>, sup  $|P(E \cap F) - P(E) |P(F)| \le \alpha (d (A,B))$ 

avec  $\alpha(m) \rightarrow 0$  quand  $m \rightarrow +\infty$ ,  $\mathcal{F}_A$  et  $\mathcal{F}_B$  étant les  $\sigma$ -algèbres engendrées par X sur A et B et d la distance du sup.

Pour les champs gaussiens stationnaires, cette propriété de mélange fort est liée à la décroissance des covariances qui peut être facilement vérifiée quand on connait la forme de la densité spectrale. Cette propriété peut aussi être vérifiée dans le cas de champs gaussiens markoviens non nécessairement stationnaires (Kunsch (14), Guyon (7)).

Une autre notion de mélange, introduite par Dobrushin (6), fait intervenir les parties A et B par l'intermédiaire de leurs cardinaux ou d'une fonction de ceux-ci. Elle est particulièrement bien adaptée aux champs de Gibbs (Kunsch (15), Guyon (7)) pour lesquels on ne sait pas vérifier le mélange fort. Bolthausen considère un coefficient de mélange  $\alpha_2(m)$  où  $|A| \le 2$ .

Une classe également intéressante de processus dépendants est celle des processus linéaires (ou moyennes mobiles infinies) qui ont souvent été étudiés pour les séries chronologiques (Anderson (1), Rosenblatt (17)). Pour ceux-ci, un T.L.C. peut être montré directement en imposant des conditions sur les coefficients de la moyenne mobile. Les processus linéaires ne sont pas nécessairement fortement mélangeants (Withers (23)).

Ces différentes notions peuvent être utilisées pour formuler des conditions entraînant la normalité asymptotique de  $r_{XY}$ . A titre d'exemple nous allons en formuler deux.

Soit {X(u)} et {Y(u)}, deux processus stationnaires sur Z<sup>2</sup> centrés et de variance finie,  $\sigma^2_X$  et  $\sigma^2_Y$ .

Soit {D<sub>n</sub>} une suite croissante de domaines de Z<sup>2</sup>, U D<sub>n</sub> = Z<sup>2</sup>  $\lim_{n \to \infty} |\delta D_n| / |D_n| = 0$ .

Les conditions alternatives suivantes sont considérées :

(i) {X(u)}, {Y(u)}  $\varepsilon$  L<sub>2+ $\delta$ </sub>, pour un  $\delta$ >0, sont mutuellement indépendants et mélangeants avec coefficients de mélange  $\alpha_2^X(m)$  et  $\alpha_2^Y(m)$  satisfaisant

$$\sum_{m=1}^{\infty} m \alpha_{2}^{X}(m) \frac{\delta/2+\delta}{\delta} < \infty \text{ et } \sum_{m=1}^{\infty} m \alpha_{2}^{Y}(m) \frac{\delta/2+\delta}{\delta} < \infty$$

(ii) {X(u)}, {Y(u)} sont linéaires : X(u) =  $\sum_{k \in \mathbb{Z}^2} g_k Z(u - k)$ , Y(u) =  $\sum_{l \in \mathbb{Z}^2} h_l W(u-l)$ 

avec Z(u) et W(u) i.i.d. de moyenne nulle et de variance finie,  $\sum |g_k| < \infty$ ,  $\sum |h_l| < \infty$  et les  $\sigma$ -algèbres engendrées par les {Z(u)} et les {W(u)} sont indépendantes.

Sous l'une ou l'autre de ces conditions, on montre alors que :  $\sum = \sum C_X(0,v) C_Y(0,v)$  $v \epsilon Z^2$ 

est absolument convergente et que si  $\Sigma > 0$ ,  $|D_n|^{1/2} r_{XY}$  converge en distribution quand n -->∞ vers une distribution normale, d'espérance nulle et de variance  $\Sigma / \sigma^2_X \sigma^2_Y$ .

La démonstration sous la condition (i) repose sur la vérification des hypothèses du T.L.C de Bolthausen (4). La convergence des variances empiriques vers  $\sigma^2_X$  et  $\sigma^2_Y$  est

à nouveau utilisé. La convergence des autres termes vers 0 étant assurée par les conditions de sommabilité des coefficients  $\{g_k\}$  ou  $\{h_l\}$ .

Pour d=1 Anderson (1) démontre un T.L.C. sous des conditions semblables à (ii). Ces conditions ont été affaiblies par la suite par Hannan (9) et Hall et Heyde (8) qui utilisent alors une approche de martingale. Cette approche ne se généralise pas naturellement à  $Z^d$ ,

#### 1.3 <u>Vitesse de convergence du théorème de la limite centrale pour des champs faiblement</u> <u>dépendants</u>

Il est aussi intéressant d'étudier la vitesse de convergence dans le T.L.C. pour des variables faiblement dépendantes indexées par Z<sup>d</sup>. Ceci a été fait dans le cadre général d'un champ {X(u)} centré, non nécessairement stationnaire, vérifiant sup  $||X(u)||_{2+\delta} < \infty$ , pour un  $\delta > 0$ , u $\mathbb{Z}^2$ 

et d'une suite strictement croissante de domaines à laquelle sont associés les suites  $(S_n), (\sigma_n^2)$ :

$$S_n = \sum X(u)$$
,  $\sigma_n^2 = Var S_n$ .  
 $u \epsilon D_n$ 

Ni la stationnarité, ni la forme des domaines  $D_n$  n'intervient plus explicitement contrairement à la conjecture de Prakasa Rao (16). Il faut seulement supposer que quand  $D_n$  croît, on apporte toujours du "nouveau" au sens de la variance, d'où une condition :

$$\lim \inf \sigma_n^2 |D_n|^{-1} = \alpha > 0.$$

Soit  $\Delta_n$ , la distance entre  $S_n$  renormalisée et une normale réduite de fonction de répartition  $\Phi$ :

$$\Delta_n = \sup_{x} |P(S_n / \sigma_n \le x) - \Phi(x)|$$

La technique utilisée par Tikhomirov (21) pour d = 1 a été généralisée. Elle consiste à évaluer la distance entre  $f_n(t)$ , la fonction caractéristique de  $S_n/\sigma_n$ , et exp (- $t^2/2$ ) et à utiliser ensuite l'inégalité de Berry-Esseen. Pour ce faire, d/dt  $f_n(t)$  est développée, conduisant à une équation différentielle qui fait intervenir - t  $f_n(t)$  et des termes de reste dont le contrôle est lié à la vitesse de convergence. Dans certains termes du reste, la dimension **d** apparaît explicitement dans les majorations qu'on effectue sur la covariance de sommes partielles,

majorations qui utilisent les inégalités classiques du mélange (cf (8)). La vitesse de convergence va donc dépendre à la fois de de la dimension **d**, de conditions sur les moments et de conditions sur le mélange.

<u>Cas d'un champ m-dépendant</u> (c.a.d.  $\alpha(k) = 0$  pour  $k \ge m$ )

a) si  $0 < \delta < 1$ ,  $\Delta_n = O(\sigma_n^{-\delta})$ .

b) si  $\delta \ge 1$ ,  $\Delta_n = O((\text{Log } \sigma_n)^{(d-1)/2} \sigma_n^{-1})$ 

Quand  $\delta < 1$ , la vitesse optimale est atteinte (Sergin (18)). Dans le cas de variables i.i.d., cela correspond à une vitesse en  $n^{-\delta/2}$ .

Quand  $\delta \ge 1$  et d = 1, on a également la vitesse optimale en  $\sigma_n^{-1}$ ; mais quand d > 1, il y a apparition d'un terme en (Log  $\sigma_n$ )<sup>(d-1)/2</sup> lié à l'ordre du développement de d/dt  $f_n(t)$ . La question reste ouverte de savoir si cette vitesse peut être améliorée sans conditions de moments plus fortes (Heinrich (11)).

<u>Cas d'un champ  $\alpha$ -mélangeant</u>,  $\alpha$  à décroissance exponentielle

a)	$\Delta_{\mathbf{n}} = \mathbf{O} \; ( \; (\text{Log } \sigma_{\mathbf{n}})^{d[1+\delta]} \; \sigma_{\mathbf{n}} \cdot \delta)$	,	$0 < \delta < 1$
	$\Delta_{\mathbf{n}} = \mathbf{O} \; ( \; (\text{Log } \sigma_{\mathbf{n}})^{2d} \; \sigma_{\mathbf{n}}^{-1})$	,	$\delta > 1$ , $d \ge 2$
	$\Delta_{\mathbf{n}} = \mathbf{O} \ ( \ (\text{Log } \sigma_{\mathbf{n}})^{5/2} \ \sigma_{\mathbf{n}}^{-1})$	,	$\delta > 1$ , $d = 1$
b) si c	le plus X est dans $L^{4+\delta}$ , $\delta > 0$ , a	lors $\Delta_i$	$n = O ((Log \sigma_n)^d \sigma_n^{-1})$

La vitesse a aussi été étudiée pour un mélange à décroissance puissance.

En utilisant une autre technique et un mélange de type Dobrushin, Takahata (19) obtient une vitesse optimale pour des champs m-dépendants sous une condition L<sup>8</sup>, tandis que pour des champs  $\alpha$ -mélangeants à décroissance exponentielle et sous une condition L<sup>8+\delta</sup>,  $\delta > 0$ , cette vitesse optimale est ralentie pour le facteur (Log  $\sigma_n$ )<sup>d</sup> comme dans le cas (b) qui impose une condition L<sup>4+\delta</sup>,  $\delta > 0$  mais un mélange plus fort.

De même que dans §1.2, les résultats formulés peuvent être appliqués en particulier à l'étude de r<sub>XY</sub> sous l'hypothèse nulle d'indépendance entre X et Y.

Ces résultats ont fait l'objet d'une publication dans Z. Wahrscheinlichkeitstheorie verw. en collaboration avec X. Guyon (liste des travaux (5)).

#### 2. Tests d'association

Plusieurs approches – qui se trouvent résumées dans les comptes rendus d'un colloque (liste des travaux (6)) – peuvent être envisagées pour tester l'association entre des variables autocorrelées. Je me suis intéressée au développement d'un test fondé sur le coefficient de corrélation empirique  $\mathbf{r}$  qui, bien que paramétrique, ne nécessite pas le choix et l'estimation de modèles paramétriques spécifiques (comme dans les modèles de régression avec paramétrisation spatiale de la matrice de variance-covariance des erreurs).

Les publications présentées dans cette partie constituent des étapes successives du développement de ce test : élaboration de la méthode, étude du risque de première espèce et de la puissance, extension au test du lien conditionnel.

Du point de vue statistique, la variance asymptotique du coefficient de corrélation **r** trouvée au §1.1 ne peut-être utilisée directement pour modifier une procédure de test car elle fait intervenir les autocorrélations pour tous les décalages et ne tient pas compte de l'estimation des moyennes et des variances.

#### 2.1 Approximation de la variance de rXY

Une bonne estimation de cette variance pour un domaine A de taille finie a été développée. Cette estimation est basée sur une approximation au premier ordre de la variance dans le cas de variables gaussiennes autocorrelées, X ~ N ( $\mu_X, \Sigma_X$ ), Y ~ N ( $\mu_Y, \Sigma_Y$ ) indépendantes, définies sur un ensemble A de N sites.

Dénotons par  $\Sigma_{\xi}$  et  $\Sigma_{\eta}$  les matrices de variance-covariance des vecteurs d'éléments  $\xi_{\alpha} = X_{\alpha} - \overline{X}, \eta_{\alpha} = Y_{\alpha} - \overline{Y}, \alpha \in A$  où  $\overline{X} = \sum_{\alpha \in A} X_{\alpha}/N$ , et  $\overline{Y} = \sum_{\alpha \in A} Y_{\alpha}/N$ .

Soit  $r_{XY}$ , le coefficient de corrélation empirique,  $r_{XY}=s_{XY}/s_X s_Y$ ,  $s_{XY} = N^{-1} \sum (X_{\alpha}-\overline{X})(Y_{\alpha}-\overline{Y})$ ,  $s_X^2 = N^{-1} \sum (X_{\alpha}-\overline{X})^2$  et  $\sigma_r^2$  sa variance.

Au premier ordre :

$$\sigma_{r}^{2} = \frac{\operatorname{tr} (\Sigma_{\xi} \Sigma_{\eta})}{\operatorname{tr} (\Sigma_{\xi}) \operatorname{tr} (\Sigma_{\eta})}$$
(1)

Dans le cas particulier où : (i)  $\Sigma_{\xi}$  et  $\Sigma_{\eta}$  commutent, (ii)  $\Sigma_{\eta}$  est proportionelle à une matrice idempotente, (iii) les valeurs propres de  $\Sigma_{\xi}$  sont nulles si les valeurs propres associées de  $\Sigma_{\eta}$  sont aussi nulles, l'approximation (1) est exacte. En effet dans ce cas,  $r_{XY}$  a pour densité  $f_{M}(\mathbf{r}) = (1-r^{2})^{1/2(M-4)} / B(1/2, 1/2(M-2)), |\mathbf{r}| \leq 1$  où B est la fonction beta et M=1+ rang ( $\Sigma_{\eta}$ ).

#### 2.2 Tests d'association modifiés

L'approximation (1) conduit à un estimateur de  $\sigma^2_r$  sous l'hypothèse qu'une partition  $\{S_k\}$  k=1,...K de AxA puisse être définie de telle sorte que la covariance de X et de Y soit constante dans chaque élément  $S_k$ . Quand les points de A sont répartis irrégulièrement dans le plan, les strates peuvent être indexées par une fonction discrétisée de la distance.

En définissant  $\hat{C}_{X}(k) = \sum_{S_{k}} (X_{\alpha} - \overline{X})(X_{\beta} - \overline{X})/N_{k}$ 

où  $N_k = |S_k|$ , on est conduit à l'estimateur de  $\sigma^2_r$ :

$$\hat{\sigma}_{r}^{2} = \frac{\sum_{k}^{N_{k}} \hat{C}_{X}(k) \hat{C}_{Y}(k)}{N^{2} s^{2}_{X} s^{2}_{Y}}$$
(2)

a) Tests modifiés

Des tests d'association modifiés pour tenir compte de l'autocorrélation peuvent alors être proposés en considérant soit une covariance renormalisée :

W = N s<sub>XY</sub>( $\sum N_k \hat{C}_X(k) \hat{C}_Y(k)$ )<sup>-1/2</sup>, que l'on teste suivant une N(0,1) en se basant sur un T.L.C, soit un test modifié de r<sub>XY</sub> : t<sub>M-2</sub>, fondé sur une modification,  $\hat{M}$ , des degrés de liberté du test sur la corrélation :

$$\hat{M} = \hat{\sigma}_{r}^{-2} + 1$$

Précisément r<sub>XY</sub> est testé suivant la densité  $f_{M}^{A}(r)$  définie au §2.1.

b) Risques de première espèce

Une première étude par simulations portant à la fois sur des processus à valeurs dans un lattice ou dans un réseau irrégulier montre que le risque d'erreur  $\alpha$  de ces tests modifiés reste alors proche de la valeur nominale. Dans cette étude, les performances des statistiques W et t<sub>M-2</sub> sont équivalentes.

Ces résultats ont fait l'objet d'un article écrit en collaboration avec P. Clifford et D. Hémon, paru dans Biometrics (liste des travaux (7)).

Comme les statistiques W et  $t_{M-2}$  ne sont pas fondées sur les mêmes approximations des distributions, les performances des tests W et  $t_{M-2}$  ont ensuite été comparées pour des échantillons de faible taille. On constate alors une supériorité du test  $t_{M-2}$  sur W.

Il a d'autre part été vérifié que les risques de première espèce restaient stables si différentes partitions isotropes  $\{S_k\}$  de AxA étaient définies. Ce risque ne s'accroît que si trop peu de strates sont considérées.

Dans cette étude par simulations il a aussi été montré que le risque de première espèce du test non paramétrique d'association proposé par Tjøstheim (22) n'était pas bien contrôlé et augmentait avec une autocorrélation positive.

#### 2.3 Puissance des test modifiés

La puissance des tests modifiés est étudiée sous une hypothèse alternative de modèle linéaire entre X et Y : H<sub>1</sub> : Y = aX + W , X ~ N( $\mu_X$ ,  $\Sigma_X$ ) , W ~ N( $\mu_W$ , $\Sigma_W$ ) et X et W indépendantes. Soit  $\Sigma_{\theta}$  la matrice de variance-covariance du vecteur d'éléments W<sub> $\alpha$ </sub> -  $\overline{W}$ ,  $\alpha \epsilon A$ .

Il est difficile de calculer la puissance des tests modifiés car l'on ne connaît pas exactement leurs lois sous H<sub>1</sub>. Cette puissance peut être évaluée par simulations. Par ailleurs il peut être intéressant de calculer la puissance  $\pi_T(s_{XY})$  d'un test de la covariance  $s_{XY}$  où l'estimateur N<sup>-2</sup>  $\sum N_k \hat{C}_X(k) \hat{C}_Y(k)$  de sa variance est remplacé par son évaluation théorique sous H1 :

N<sup>-2</sup> tr 
$$(\Sigma_{\xi} \Sigma_{\eta}) = N^{-2} [a^2 \operatorname{tr}(\Sigma_{\xi}^2) + \operatorname{tr}(\Sigma_{\xi} \Sigma_{\theta})].$$

Pour pouvoir comparer les puissances observées par simulations à une référence, la puissance  $\pi_{N*}(r)$  du test standard de la corrélation basé sur un nombre N\* de points compatible avec la variance empirique observée de  $r_{XY}$  a aussi été calculée. Des tableaux indiquant  $\pi_{T}(s_{XY})$ ,  $\pi_{N*}(r)$  et les puissances observées de W et t $M_{-2}$  par simulations sur un domaine irrégulier sont donnés dans (8).

On constate que dans l'ensemble toutes les puissances calculées ou observées sont proches. La différence ne devient notable que dans quelques cas de fortes autocorrélations pour l'un ou l'autre des processus.

Ainsi on peut dire que la puissance "théorique"  $\pi_T(s_{XY})$  donne une bonne indication dans la plupart des cas rencontrés. Les tests modifiés ne présentent pas de perte de puissance notable par rapport à un test classique fondé sur un nombre d'observations N\* "équivalent". *L'étude de la puissance par simulation et de certains aspects complémentaires ont fait l'objet d'une communication écrite (liste des travaux (8)).* 

#### 2.4 Tests modifiés des corrélations partielles

Il est souvent intéressant de pouvoir tester le lien conditionnel entre deux variables après conditionnement sur une ou plusieurs autres variables, ce qui revient à tester la corrélation partielle. Les méthodes développées précédemment s'étendent aisément au test de la corrélation partielle si l'on considère un cadre gaussien et les distributions conditionnelles appropriées.

En pratique les tests modifiés W et  $t_{M-2}$  seront calculés sur les résidus, obtenus par les moindres carrés, des régressions linéaires sur les variables de conditionnement. Une étude par simulations a permis de vérifier la bonne performance de cette méthode. *Ces résultats vont paraître dans Statistics in Medecine (liste des travaux (9)).* 

#### 3. Applications à l'étude des corrélations écologiques en épidémiologie

Les études où les variables considérées sont moyennées sur des groupes (par exemple les taux de mortalité par département) sont appelées des études de corrélations écologiques. Ces études sont sujettes à différentes sources de biais qu'il faut prendre en compte pour pouvoir en donner une interprétation valable.

Dans un premier temps, je me suis attachée à discuter ces biais quantitativement en comparant la forme des relations dose-effet individuelles ou agrégées et l'effet du niveau d'agrégation. Dans une deuxième partie, j'ai comparé sur des exemples les risques relatifs estimés par des études écologiques et des études individuelles. On peut montrer que les relations dose-effet données par les études écologiques et les études individuelles sont d'autant plus proches que le facteur de risque étudié est prépondérant dans l'ensemble des facteurs de risque relevants, que les fractions de populations concernées sont bien identifiées et que la relation dose-effet individuelle est proche de la linéarité.

Ces réflexions ont fait l'objet d'une publication dans l'International Journal of Epidemiology en collaboration avec I. Stücker et D. Hémon (liste des travaux (10)).

D'autre part, j'ai appliqué à différents problèmes épidémiologiques les méthodes proposées. Les résultats trouvés sont en bon accord avec certains facteurs de risque isolés dans des enquêtes individuelles et illustrent l'intérêt des méthodes statistiques proposées. Ces résultats sont présentés dans la dernière partie de l'article sur le test ajusté de la corrélation partielle (liste des travaux (9)).

#### 4. Conclusion

Dans cette deuxième partie, j'ai résumé un ensemble de travaux centrés sur l'étude de tests d'association entre des variables spatialement dépendantes. Ceci m'a conduit d'une part à m'intéresser à certaines techniques probabilistes, utiles dans le cadre de variables faiblement dépendantes pour étudier la consistance ou la convergence vers la normalité de certaines statistiques. D'autre part, j'ai été amenée à proposer un test d'association modifié fondé sur le coefficient de corrélation empirique. Le coefficient de corrélation est une mesure d'association qui, bien que paramétrique, ne nécessite pas l'identification et l'estimation de modèles particuliers. Le risque de première espèce et la puissance des tests modifiés ont été étudiés, ainsi que leur généralisation au test du lien conditionnel. J'ai montré l'intérêt d'utiliser ces test dans des applications épidémiologiques sur les corrélations écologiques.

Il serait intéressant de développer d'autres mesures d'association, en particulier non paramétriques. Un index d'association non paramétrique a été proposé par Tjøstheim sans que le risque d'erreur de la méthode proposée n'ait été bien contrôlé. Une approche fondée sur les permutations serait possible si on prenait soin de faire une inférence conditionnelle à la préservation d'une structure spatiale.

#### Références de la première partie

- Chacon, R.V. A geometric construction of measure preserving transformations, Proc. Fifth Berkeley Sym. Math. Stat. Prob., University of California Press, 1967, Vol II, part 2, 335-360.
- (2) Friedman, N.A. Introduction to Ergodic Theory. New York : van Nostrand 1970
- (3) Friedman, N.A. Bernoulli shifts induce Bernoulli shifts. Advances in Math. 10, 39-48 (1970)
- (4) Friedman, N.A., Ornstein, D.S. Ergodic transformations induce mixing transformations. Advances in Math. 10, 147-163 (1973)
- (5) Halmos, P.R. Lectures on Ergodic Theory. New-York - Chelsea 1956.
- (6) Neveu J. Temps d'arrêt d'un système dynamique.Z. Wahrscheinlichkeitstheorie verw. Gebiete 13, 81-94 (1969).
- (7) Ornstein, D.S., Smorodinski, M. Ergodic flows of positive entropy can be time changed to become K.flows. Israel J. Math. 26, 75-83 (1977).
- (8) Shields, P. Cutting and independent stacking of intervals. Math. Systems Theory 7, 1-4 (1973).
- (9) Walters, P. An Introduction to Ergodic Theory. Graduate Text in Mathematics. Vol 79 New-York, Springer Verlag 1981.

#### Références de la deuxième partie

- (1) Anderson, T.W. The statistical analysis of time series. New York : Wiley (1971).
- (2) Bartlett, M.S. Some aspects of the time-correlation problem in regard to tests of significance. Journal of the Royal Statistical Society, 98, 536-543 (1935).
- (3) Bartlett, M.S. An Introduction to Stochastic Processes, 2nd edition Cambridge University Press 1966.
- (4) Bolthausen, E. On the central limit theorem for stationary random fields. Ann. Probability 10, 1047-1050 (1982).
- (5) Cliff, A.D. and Ord, J.K. Model building and the analysis of spatial pattern in human geography. Journal of the Royal Statistical Society B, 37, 297-348 (1975).
- (6) Dobrushin, R.L. The description of a random field by its conditional distribution and its regularity condition. Theory Probability Appl. 13, 197-227 (1968).

- (7) Guyon, X. Estimation par pseudo-vraisemblance conditionnelle pour les champs : étude asymptotique et application aux champs markoviens. Proceedings of the 6th Franco-Belgian Meeting of Statisticians, November 1985. Edited by F. Droesbeke. Travaux et Recherches Vol 11, p. 15-62. Publications des Facultés universitaires St Louis, Bruxelles.
- (8) Hall, P., Heyde, C.C. Martingale limit theory and its application. New York : Academic Press 1980.
- (9) Hannan, E.J. Multiple Times Series. New York : Wiley 1970.
- (10) Haugh, L.D. and Box, G.E.P. Identification of dynamic regression (distributed lag) models connecting two time series.
   J. Amer. Statist. Assoc. 72, 121-130 (1977).
- (11) Heinrich, L. Asymptotic expansions in the Central Limit Theorem for a special class of m-dependent random fields. Mathematische Nachrichten 83-106 (1987).
- (12) Ibragimov, I.A., Linnik, Yu.V. Independent and stationary sequences of random variables. Groningen : Wolters-Nordhoff 1971.
- (13) Johnston, J. Econometric Methods, 2nd edition. New York : McGrawHill 1972.
- (14) Künsch, H.R. Reelwertige Zufallsfelder anf einem Gitter : Interpolation, variations sprinzip und Stat. Analyse, Thèse E.T.H. de Zurich, 1980.
- (15) Künsch, H.R. Decay of correlations under Dobrushin's uniqueness condition and its applications.
   Comm. Math. Phys. 84, 207-222 (1982).
- (16) Prakasa Rao, B.L. A non uniform estimates of the rate of convergence in the central limit theorem for m-dependent random fields.
   Z. Wahrscheinlichkeitstheorie verw. Gebiete 58, 247-256 (1981).
- (17) Rosenblatt, M. Stationary sequences and random fields. Boston : Birkhäuser 1985.
- (18) Shergin, V.V. An estimate of the remainder term in the central limit theorem for m-dependent random variables. Theory Probability Appl. 24, 782-796 (1979).
- (19) Takahata, H. On the rates in the central limit theorem for weakly dependent random fields. Z. Wahrscheinlichkeitstheorie verw. Gebiete 64, 445-456 (1983).
- (20) Tempelman, A.A. Ergodic theorems for dynamical systems. Trans. Moscow Math. Soc. 26, 94-132 (1972).
- (21) Tikhomirov, A.N. On the converge rate in the centra limit theorem for weakly dependent random variables. Theory Probability Appl. 25, 790-809 (1980).
- (22) Tjøstheim, D. A measure of association for spatial variables. Biometrika 65, 109-114 (1978).
- (23) Withers, C.S. Conditions for linear processes to be strong mixing.Z. Wahrscheinlichkeitstheorie verw. Gebiete 57, 477-480 (1981).

#### LISTE DES TRAVAUX

- S.T. Richardson, K.M. Wilkinson Stopping time transformations and towers.
   Z. Wahrscheinlichkeitstheorie verw. Gebiete 48. 259-284, (1979)
- (2) S.T. Richardson Stopping times for measure-preserving transformations : some approximation results. Proceedings of the London Mathematical Society (3) 43 : 273-294 (1981).
- (3) S.T. Richardson, D. Hémon
   On the variance of the sample correlation between two independent lattice processes.
   J. Appl. Prob. 18, 943-948 (1981).
- (4) S.T. Richardson, D. Hémon Autocorrélation spatiale : ses conséquences sur la corrélation empirique de deux processus spatiaux. Revue de Statistique Appliquée, 30 : 41-51 (1982).
- (5) X. Guyon, S.T. Richardson Vitesse de convergence du théorème de la limite centrale pour des champs faiblement dépendants.
  Z. Wahrscheinlichkeitstheorie verw. Gebiete, 66 : 297-314 (1984).
- (6) S.T. Richardson Testing the association between two spatial processes : a review of different approaches Spatial processes and spatial time series analysis. Proceedings of the 6th Franco-Belgian Meeting of Statisticians, November 1985. Edited by F. Droesbeke. Travaux et Recherches Vol 11, p. 193-200. Publications des Facultés universitaires St Louis, Bruxelles.
- (7) P. Clifford, S.T. Richardson, D. Hémon Assessing the significance of the correlation between two spatial processes. Biometrics, 45 : 123-134 (1989).
- (8) S.T. Richardson, P. Clifford Testing association between spatial processes. A paraître dans Proceedings of the A.M.S conference on Spatial Statistics and Imaging (June 1988).
- (9) S.T. Richardson
   A method for testing the significance of geographical correlations with application to industrial lung cancer in France.
   A paraître dans Statistics in Medecine.
- (10) S.T. Richardson, I. Stücker, D. Hémon Comparison of relative risks obtained in ecological and individual studies : some methodological considerations. International Journal of Epidemiology, Vol 16, n°1, 111-120 (1987).

# **PREMIERE PARTIE**

# PROCESSUS ERGODIQUES : ETUDES DES PROPRIETES ERGODIQUES DES TEMPS D'ARRET

# STOPPING TIME TRANSFORMATIONS AND TOWERS

Z. Wahrscheinlichkeitstheorie verw. Gebiete 48, 259–284 (1979)

## **Stopping Time Transformations and Towers**

S.T. Richardson<sup>1</sup> and K.M. Wilkinson<sup>2</sup>

<sup>1</sup> Department of Statistics, University of Warwick, Coventry CV47AL, England <sup>2</sup> Department of Mathematics, The University of Nottingham, University Park.

Nottingham NG72RD England

#### § 1. Introduction

In this paper we shall study properties of  $\tau^v$  where  $\tau$  is an automorphism of a Lebesgue space and v is a non-negative integer valued measurable function. We shall always assume that v satisfies the conditions, given by Neveu [4] and discussed in §4, which ensure that  $\tau^v$  is again an automorphism. Such a v is called a stopping time and  $\tau^v$  a stopping time transformation. The set of all stopping time transformations derived from a particular automorphism  $\tau$  is equivalent to the positive part of the full group of  $\tau$ .

We shall restrict attention to that class of automorphisms which are defined by means of a tower construction. This class contains many interesting transformations including Bernoulli shifts of any prescribed entropy. For an automorphism in this class, using a theorem of Neveu, we can follow the action of  $\tau^{v}$ in the tower and this leads to an analysis of  $\tau^{v}$  by means of a more general contruction which we call a tower block construction.

This new construction appears to be of interest in its own right and so because of this and also for reasons of clarity we shall commence by introducing the construction and discuss properties of the transformation so defined. In §§ 4 and 5 we show that in general  $\tau^v$ , where  $\tau$  is defined by means of a tower construction, can be represented by a tower block construction. We are then able to use the results of § 3 to obtain density results for  $\tau^v$ , our main results being firstly that the set of  $\tau^v$  which are ergodic and secondly the set of  $\tau^v$  which are Bernoulli are both dense in the set of stopping time transformations derived from the same tower transformation  $\tau$ . As a consequence of this last result we have that if  $\tau$  is Bernoulli the set of  $\tau^v$  which are also Bernoulli is dense, thus generalising a result of Friedman for the induced transformation [2]. Saleski [6] has also investigated properties of a stopping time transformation where  $\tau$  is Bernoulli, particularly in the case where v has expected value 2.

Throughout the paper we take  $\Omega = [0, 1]$ .  $\mathcal{B}$  the set of Borel subsets of  $\Omega$  and  $\mu$  Lebesgue measure on  $\mathcal{B}$ .

0044-3719/79 0048 0259/\$05.20
#### §2. Tower Block Construction

Our initial definitions follow Friedman ([1], 6 and [2]). A column C is an ordered set of disjoint left-closed, right-open subintervals of  $\Omega$ .

$$C = (L_i: 1 \leq i \leq h).$$

The sets  $L_i$ ,  $1 \le i \le h$ , are referred to as column levels and they are assumed to each have the same measure. C has base  $B(C) = L_1$ , top  $A(C) = L_h$ , width  $w(C) = \mu(L_1)$  and height h(C) = h. We let

$$C^* = \bigcup_{i=1}^h L_i$$

and define columns  $C_1, C_2$  to be disjoint if  $C_1^*$  and  $C_2^*$  are disjoint sets. We denote by  $\tau_C$  the one-to-one measure preserving mapping from  $C^* - L_h$  to  $C^* - L_1$  which maps  $L_i$  linearly onto  $L_{i+1}, 1 \le i \le h-1$ .

A tower T is an ordered set of disjoint columns.

$$T = (C_i: 1 \leq i \leq m).$$

T has base  $B(T) = \bigcup_{i=1}^{m} B(C_i)$ , top  $A(T) = \bigcup_{i=1}^{m} A(C_i)$  and width  $w(T) = \sum_{i=1}^{m} w(C_i)$ =  $\mu(B(T))$ . We let

$$T^* = \bigcup_{i=1}^m C_i^*$$

and define towers  $T_1, T_2$  to be disjoint if  $T_1^*$  and  $T_2^*$  are disjoint sets.  $\tau_T$  is the mapping from  $T^* - A(T)$  to  $T^* - B(T)$  consisting of the mappings  $\tau_{C_i}$  in each column  $C_i$ ,  $1 \le i \le m$ .

If  $C_1 = (L_{1i}: 1 \le i \le h_1)$  and  $C_2 = (L_{2i}: 1 \le i \le h_2)$  are two disjoint columns with  $w(C_1) = w(C_2)$  we define

$$C_1 * C_2 = (L_{11}, L_{12}, \dots, L_{1h_1}, L_{21}, L_{22}, \dots, L_{2h_2}),$$

so  $C_1 * C_2$  is the column obtained by placing  $C_2$  above  $C_1$ .

If  $T_1 = (C_{1i}: 1 \le i \le m_1)$  and  $T_2 = (C_{2i}: 1 \le i \le m_2)$  are disjoint towers we define

$$T_1 + T_2 = (C_{11}, C_{12}, \dots, C_{1m_1}, C_{21}, C_{22}, \dots, C_{2m_2}),$$

so  $T_1 + T_2$  is the tower obtained by placing  $T_2$  at the side of  $T_1$ .

Suppose  $\alpha_i \ge 0$ ,  $0 \le i \le r$ , and  $\sum_{i=0}^r \alpha_i = 1$ . If L = [a, b] we define

$$\alpha_{j}L = \left[a + (b-a)\sum_{i=0}^{j-1} \alpha_{i}, a + (b-a)\sum_{i=0}^{j} \alpha_{i}\right], \quad 0 \le j \le r.$$

Then for a column  $C = (L_i: 1 \le i \le h)$  we define

$$\alpha_j C = (\alpha_j L_i: 1 \le i \le h). \quad 0 \le j \le r.$$

and for a tower  $T = (C_i: 1 \leq i \leq m)$  we define

 $\alpha_i T = (\alpha_i C_i: 1 \le i \le m), \quad 0 \le j \le r.$ 

Note that  $w(\alpha_j C) = \alpha_j w(C)$  and  $w(\alpha_j T) = \alpha_j w(T)$ .

If C is a column and  $T = (C_i: 1 \le i \le m)$  is a tower such that C\* and T\* are disjoint and w(C) = w(T) we define a tower C \* T by

$$C * T = (\alpha_{i-1} C * C_i, 1 \leq i \leq m)$$

where  $\alpha_{i-1} = w(C_i)/w(C)$ ,  $1 \le i \le m$ . So C \* T is obtained by placing the tower T above C. If  $T_1 = (C_i: 1 \le i \le m)$  and  $T_2$  are disjoint towers with  $w(T_1) = w(T_2)$  we define a tower  $T_1 * T_2$  by

$$T_1 * T_2 = \sum_{j=1}^{m} C_j * \alpha_{j-1} T_2$$

where  $\alpha_{j-1} = w(C_j)/w(T_2)$ ,  $1 \le j \le m$ . So  $T_1 * T_2$  is obtained by placing a copy of  $T_2$  above each column in  $T_1$ .

If  $\alpha_0 = \alpha_1 = 1/2$  we adopt the notation

$$(\frac{1}{2}L)_i = \alpha_i L, \qquad (\frac{1}{2}C)_i = \alpha_i C, \qquad (\frac{1}{2}T)_i = \alpha_i T.$$

i=0, 1, where L, C, T are, respectively, a column level, a column and a tower. For any tower T we define a tower S(T) by

 $S(T) = (\frac{1}{2}T)_0 * (\frac{1}{2}T)_1.$ 

So S(T) consists of a copy of T in the base and a further copy of T above each column in the base copy. The base copy of T in S(T) will be referred to as the copy of rank 0 and, if T has m columns, there are m copies of T of rank 1 in S(T), one above each column in the rank 0 copy.

Letting  $S^0(T) = T$  we define  $S^u(T)$ ,  $u \ge 1$ , inductively by

$$S^{u+1}(T) = S(S^u(T)).$$

 $S^{u}(T)$ ,  $u \ge 0$ , consists of one rank 0 copy of T. m rank 1 copies of T.  $m^{2}$  rank 2 copies of T,..., and  $m^{2^{u}-1}$  copies of T of rank  $2^{u}-1$  and each copy of T of rank i,  $1 \le i \le 2^{u}-1$  is immediately above a column in a rank i-1 copy of T.  $S^{u}(T)$ ,  $u \ge 0$ , also contains  $m^{2^{u}}$  columns.

If the *i*th column in a tower T has the form

$$\alpha_{i_1}^1 c_1 * \alpha_{i_2}^2 c_2 * \ldots * \alpha_{i_l}^l c_l$$

where  $\alpha_s^j \ge 0$ ,  $\sum_{s=0}^{r_j} \alpha_s^j = 1$ ,  $0 \le i_j \le r_j$ ,  $1 \le j \le l$ , and  $c_j$ ,  $1 \le j \le l$ , are columns, it will sometimes be unimportant what the values of  $\alpha_{i_j}^j$ ,  $1 \le j \le l$ , are. In such a case we write

 $C_i(T) = [c_1, c_2, \dots, c_l]$ 

and take this to mean that the *i*th column of T consists of a copy of  $c_1$  with a copy of  $c_2$  above it. etc.

In particular, if  $T = (C_i, 1 \le i \le m)$ , then the  $m^2$  columns of S(T) are given by

$$C_{(i_0-1)m+i_1}(S(T)) = [C_{i_0}, C_{i_1}], \quad 1 \le i_j \le m, \ j = 0, 1,$$

and more generally, for  $u \ge 1$ , the  $m^{2^{u}}$  columns of  $S^{u}(T)$  are given by

$$C_{(i_0-1)m^{2^{u}-1}+(i_1-1)m^{2^{u}-2}+\dots+(i_{2^{u}-2}-1)m+i_{2^{u}-1}}(S^u(T))$$
  
=  $[C_{i_0}, C_{i_1}, \dots, C_{i_{2^{u}-1}}], \quad 1 \le i_j \le m, \quad 0 \le j \le 2^u - 1.$ 

Note that  $\tau_{S(T)}$  agrees with  $\tau_T$  on  $T^* - A(T)$  and that, in addition,  $\tau_{S(T)}$  is defined on half of A(T). In fact the set on which  $\tau_{S(T)}$  is undefined is A(S(T)) which has measure  $\frac{1}{2}w(T)$ . Similarly  $\tau_{S^u(T)}$ ,  $u \ge 2$ , agrees with  $\tau_{S^v(T)}$  on  $T^* - A(S^v(T))$ ,  $0 \le v \le u - 1$  and

$$\mu(A(S^{u}(T))) = \frac{1}{2^{u}} w(T).$$

Hence we may define a mapping  $\tau(T): T^* \to T^*$  by

$$\tau(T) = \lim_{u \to \infty} \tau_{S^u(T)}.$$

 $\tau(T)$  is invertible, preserves Lebesgue measure and is ergodic ([1]).

We now wish to define a more general construction for which the resulting transformation need not be ergodic. We shall use this construction to investigate properties of stopping time transformations. We start with an ordered set of towers

$$\mathbb{T} = (T_i: 1 \leq i \leq k)$$

each with *m* columns.

$$T_i = (C_{ij}: 1 \leq j \leq m)$$

and  $w(C_{i_1,j}) = w(C_{i_2,j}), 1 \le i_1, i_2 \le k, 1 \le j \le m$ . T is called a tower block and has base

$$B(\mathbf{T}) = \bigcup_{i=1}^{k} B(T_i), \quad \text{top .4}(\mathbf{T}) = \bigcup_{i=1}^{k} A(T_i)$$

and width  $w(\mathbf{T}) = \mu(B(\mathbf{T})) = k w(T_1)$ . We also put

$$\mathbf{T}^* = \bigcup_{i=1}^k T_i^*.$$

Let  $\Sigma = (\sigma_i: 1 \le i \le m)$  where each  $\sigma_i, 1 \le i \le m$  is a permutation of  $\{1, 2, \dots, k\}$  and define  $P(\Sigma, \mathbb{T})$  to be the tower block  $(P(\Sigma, T_i): 1 \le i \le k)$  where

$$P(\Sigma, T_i) = \sum_{j=1}^{m} \left(\frac{1}{2} C_{ij}\right)_0 * \alpha_{j-1} \left(\frac{1}{2} T_{\sigma_j(i)}\right)_1, \qquad 1 \le i \le k,$$

Stopping Time Transformations and Towers

and

$$\alpha_{i-1} = w(C_{ii})/w(T_i), \qquad 1 \le j \le m.$$

So  $P(\Sigma, T_i)$  has a copy of  $T_i$  in the base and then a copy of  $T_{\sigma_j(i)}$  above the *j*th column of  $T_i$ . In fact

$$C_{(i_0-1)m+i_1}(P(\Sigma, T_i)) = [C_{i,i_0}, C_{\sigma_{i_0}(i),i_1}], \quad 1 \le i_j \le m, \ j = 0, 1.$$

We shall let  $\tau_{\mathbb{T}}$  be the transformation defined on  $\mathbb{T}^* - A(\mathbb{T})$  which is  $\tau_{T_i}$  on each  $T_i^* - A(T_i)$ ,  $1 \leq i \leq k$ . Clearly  $\tau_{P(\Sigma, \mathbb{T})}$  agrees with  $\tau_{\mathbb{T}}$  on  $\mathbb{T}^* - A(\mathbb{T})$  and is also defined on half of  $\mathbb{T}^* - A(\mathbb{T})$ .

We now wish to extend the definition of  $\tau_{\pi}$  by defining an inductive tower block construction. Let  $\mathbb{T}$  and  $\Sigma$  be as above. For each  $u \ge 0$  we define an ordered collection of permutations of  $\{1, 2, \dots, k\}$ ,

$$\Sigma_u = \{ \sigma_i^u \colon 1 \leq i \leq m^{2^u} \},\$$

by  $\Sigma_0 = \Sigma$  (i.e.  $\sigma_i^0 = \sigma_i$ ,  $1 \le i \le m$ ) and for  $u \ge 0$ ,  $1 \le j \le m^{2^u}$ ,  $1 \le l \le m^{2^u}$ ,

(2.1) 
$$\sigma_{(j-1)m^{2^{u}}-l}^{u+1}(i) = \sigma_{l}^{u}(\sigma_{j}^{u}(i)),$$

for each  $1 \le i \le k$ . We then define  $P^0(T_i) = T_i$ .  $1 \le i \le k$ , and for each  $u \ge 0$  we inductively define

$$P^{u-1}(T_i) = P(\Sigma_u, P^u(T_i)), \quad 1 \le i \le k.$$

and define  $P^{\mu}(\mathbb{T})$ ,  $u \ge 0$ , to be the tower block

$$(P^{\boldsymbol{u}}(T_i): 1 \leq i \leq k).$$

Note that  $P^{u}(T_{i})$ ,  $1 \leq i \leq k$ ,  $u \geq 0$ , consists of a rank 0 copy of  $T_{i}$ , above each of the *m* columns of  $T_{i}$  there is a copy of one of the towers  $T_{i}$ ,  $1 \leq i \leq k$ , above each of the  $m^{2}$  columns in these rank 1 copies there is a copy of one of the towers  $T_{i}$ ,  $1 \leq i \leq k$ ... and at the top of  $P^{u}(T_{i})$  there are  $m^{2^{u}-1}$  copies of some of the towers  $T_{i}$ ,  $1 \leq i \leq k$ , of rank  $2^{u}-1$ , one above each column in the copies of rank  $2^{u}-2$ . Because of (2.1) for each  $u \geq 0$ ,  $1 \leq i \leq k$ .

$$(2.2) \quad C_{(i_0-1)m^{2^{u-1}}+(i_1-1)m^{2^{u-2}}+\dots+(i_2^{u}-2-1)m+i_2^{u}-1}P^u(T_i) \\ = [C_{ii_0}, C_{\sigma_{i_0}(i), i_1}, C_{\sigma_{i_1}(\sigma_{i_0}(i)), i_2, \dots}(\sigma_{i_2^{u}-3}(\dots(\sigma_{i_1}(\sigma_{i_0}(i)))\dots)), i_{2^{u}-1} \\ C_{\sigma_{i_2^{u}-2}}], \\ 1 \le i_i \le m, \ 0 \le j \le 2^u - 1.$$

So (2.1) ensures that above every copy of the column  $C_{ij}$  of rank  $0, 1, \ldots, 2^u - 2$  in  $P^u(\mathbb{T}), u \ge 1$ , we have a copy of the tower  $T_{\sigma_i(i)}$ . Note also that in  $P^u(\mathbb{T}), u \ge 0$ , there are  $m^l$  copies of each  $T_i, 1 \le i \le k$ , of rank  $l, 1 \le l \le 2^u - 1$  distributed among the various  $P^u(T_i), 1 \le i \le k$ .

As  $\tau_{P^u(\mathbb{T})}$ ,  $u \ge 1$ , agrees with each of  $\tau_{P^r(\mathbb{T})}$  on  $\mathbb{T}^* - \mathcal{A}(P^r(\mathbb{T}))$ ,  $0 \le v \le u - 1$ , and as

$$\mu(A(P^{u}(\mathbf{T}))) = \frac{1}{2^{u}} w(\mathbf{T}), \quad u \ge 0.$$

we can define a mapping  $\tau(\Sigma, \mathbf{I}): \mathbf{I}^* \to \mathbf{I}^*$  by

$$\tau(\Sigma, \mathbf{T}) = \lim_{u \to \infty} \tau_{P^u(\mathbf{T})}.$$

So  $\tau(\Sigma, \mathbb{T})$  maps a column level in  $\mathbb{T}$  which is a subset of  $\mathbb{T}^* - A(\mathbb{T})$  linearly onto the column level immediately above and maps  $A(C_{ij}), 1 \leq j \leq k, 1 \leq j \leq m$ , onto  $B(T_{\sigma_j(i)})$ .

The above construction in the case k = 1 yields

 $P^u(\mathbf{T}) = S^u(T_1), \qquad u \ge 0,$ 

where  $S^{u}(T_{1})$ ,  $u \ge 0$ , is as defined in Friedman's construction. Note that  $\tau(\Sigma, \mathbb{T})$  is invertible and preserves Lebesgue measure. As an example of  $\tau(\Sigma, \mathbb{T})$  which is not ergodic take  $k \ge 2$  and each  $\sigma_{i}$ ,  $1 \le i \le m$ , to be the identity permutation. In this case

$$P^u(T_i) = S^u(T_i)$$

for each  $u \ge 0$ ,  $1 \le i \le k$  and so  $\tau(\Sigma, \mathbb{T})$  leaves each set  $T_i^*$ ,  $1 \le i \le k$ , invariant.

# § 3. Mixing Properties of $\tau(\Sigma, \mathbb{T})$

Throughout this section we take  $\tau = \tau(\Sigma, \mathbb{T})$  where  $\Sigma$  and  $\mathbb{T}$  are as defined in §2. We shall first show that  $\tau$  is a Markov shift, thus generalising a result of Friedman [2] and Shields [7].

Let  $\mathscr{L}$  be the partition of  $\mathbb{T}^*$  consisting of the column levels in  $\mathbb{T}$  and let  $\mathscr{L}_n = \bigvee_{i=0}^{n-1} \tau^i \mathscr{L}$ . If  $L \in \mathscr{L}$ ,  $n \ge 1$ , we shall let

 $\mathscr{L}_n(L) = \{ M \in \mathscr{L}_n \colon M \subset L \}.$ 

A partition  $\mathcal{L}$  is a Markov partition for a transformation  $\tau$  with invariant measure  $\mu$  if we can find  $p(L, M) \ge 0$ . L.  $M \in \mathcal{L}$ , satisfying

(3.1) For each 
$$M \in \mathscr{L}$$
,  $\mu(M) = \sum_{L \in \mathscr{L}} \mu(L) p(L, M)$ .

(3.2) For each 
$$L \in \mathscr{L}$$
,  $\sum_{M \in \mathscr{L}} p(L, M) = 1$ , and

(3.3) For each  $D \in \mathscr{L}_n(L)$ ,  $L \in \mathscr{L}$  and  $M \in \mathscr{L}$ .

 $\mu(\tau D \cap M) = p(L, M) \,\mu(D).$ 

We first show that an atom of  $\mathcal{L}_n$ ,  $n \ge 1$ , can be approximated by the union of certain column levels in  $P^{\mu}(\mathbf{T})$ .  $u \ge 0$ .

**Lemma 3.1.** Let  $L \in \mathscr{L}$  and  $M \in \mathscr{L}_n(L)$ ,  $n \ge 1$ . Then for each  $u \ge 0$  we can write

$$M = D_u \cup E_u$$

where  $D_u$  is a union of copies of L in  $P^u(\mathbb{T})$  and  $E_u = \bigcup_{j=0}^{n-1} \tau^j(B(P^u(\mathbb{T}))).$ 

*Proof.* Let  $\mathscr{L} = (L_i: 1 \le i \le s)$ . Each column level in  $P^{\mu}(\mathbb{T})$  is a subset of an atom of  $\mathscr{L}$  and if F is a column level of  $P^{\mu}(\mathbb{T})$  we write l(F) = i if  $F \subset L_i$ ,  $1 \le i \le s$ , i.e. l(F) = i if and only if F is a copy of  $L_i$ . Now for each i,  $1 \le i \le s$ , there are  $2^{\mu}m^{2^{\mu}-1}$  column levels F of  $P^{\mu}(\mathbb{T})$  with l(F) = i. Let

$$M = L_{i_1} \cap \tau L_{i_2} \cap \ldots \cap \tau^{n-1} L_{i_n}$$

and let  $F_j$ ,  $1 \le j \le 2^u m^{2^u - 1}$ , be the copies of  $L_{i_1}$  in  $P^u(\mathbb{T})$ . Define  $J_1$  to be the set of those j,  $1 \le j \le 2^u m^{2^u - 1}$ , for which the column of  $P^u(\mathbb{T})$  containing  $F_j$  has no more than n-2 column levels beneath  $F_j$ . For  $j \in J_1^c$  we let  $F_{j,-1} = F_j$  and for  $1 \le p \le n-1$  we let  $F_{j,-p+1}$  denote the column level immediately below  $F_{j,-p}$ . Define

$$J_2 = \{ j \in J_1^c : \ l(F_{j, p}) = i_p, \ 1 \le p \le n \}$$

and

$$J_3 = \{ j \in J_1^c : l(F_{j, p}) \neq i_p \text{ for at least one } p, 1 \leq p \leq n \}$$

and let

$$D_u = \bigcup_{j \in J_2} F_j, \qquad E_u = M - D_u.$$

As  $E_u \subset \bigcup_{j \in J_1} F_j \subset \bigcup_{i=0}^{n-1} \tau^i B(P^u(\mathbf{T}))$ , we have the result.

**Corollary 3.2.** Given  $\varepsilon > 0$  we can find  $U = U(n, \varepsilon)$  such that  $\mu(E_u) < \varepsilon$  for  $u \ge U$ .

Proof. 
$$\mu\left(\bigcup_{i=0}^{n-1}\tau^{i}B(P^{u}(\mathbb{T}))\right)=\frac{(n-1)}{2^{u}}w(\mathbb{T}).$$

**Lemma 3.3.**  $\mathcal{L}$  is a Markov partition for  $\tau$ .

*Proof.* Define for  $L, M \in \mathscr{L}$ 

$$p(L, M) = \begin{cases} 1 & L \notin A(\mathbb{T}), \ M = \tau L, \\ 0 & L \notin A(\mathbb{T}), \ M \neq \tau L, \\ \frac{\mu(M)}{w(T_1)} & L = A(C_{ij}), \ M \subset B(T_{\sigma_j(i)}), \ 1 \leq i \leq k, \ 1 \leq j \leq m, \\ 0 & L = A(C_{ij}), \ M \notin B(T_{\sigma_j(i)}), \ 1 \leq i \leq k, \ 1 \leq j \leq m. \end{cases}$$

We must check that this definition satisfies (3.1), (3.2), (3.3).

Let  $M \in \mathscr{L}$ ,  $M \notin B(\mathbb{T})$  and let  $L_1 = \tau^{-1} M$ . Then

$$p(L, M) = \begin{cases} 1 & L = L_1 \\ 0 & L = L_1 \end{cases}$$

and so

$$\sum_{L \in \mathscr{L}} \mu(L) p(L, M) = \mu(L_1) = \mu(M)$$

and (3.1) is true in this case. If now  $M \subset B(T_i)$ ,  $1 \leq i \leq k$ , then

$$p(L, M) = \begin{cases} \frac{\mu(M)}{w(T_1)}, & L = A(C_{lj}) & \text{with } \sigma_j(l) = i \\ 0 & \text{otherwise} \end{cases}$$

and so

$$\sum_{L \in \mathscr{L}} \mu(L) p(L, M) = \sum_{\{(j, l): \sigma_j(l) = i\}} \mu(A(C_{lj})) \frac{\mu(M)}{w(T_1)} = \mu(M)$$

as

$$\sum_{\{(j,\,l):\,\sigma_j(l)\,=\,i\}}\mu(A(C_{lj}))=\sum_{j\,=\,1}^m\mu(A(C_{\sigma_j^{-1}(i),\,j}))=w(T_1).$$

So (3.1) holds.

Now if  $L \not\subset A(\mathbf{T})$ .

$$p(L, M) = \begin{cases} 1 & M = \tau L \\ 0 & M \neq \tau L \end{cases}$$

and so

$$\sum_{M \in \mathscr{L}} p(L, M) = 1$$

and if  $L = 4(C_{ij}), 1 \leq i \leq k, 1 \leq j \leq m$ .

$$p(L, M) = \begin{cases} \frac{\mu(M)}{w(T_1)}, & M \subset B(T_{\sigma_j(i)}), \\ 0, & M \notin B(T_{\sigma_j(i)}). \end{cases}$$

and so

$$\sum_{M \in \mathscr{L}} p(L, M) = \sum_{M \in B(T_{\sigma_1}(\mu))} \frac{\mu(M)}{w(T_1)} = 1.$$

So (3.2) holds.

Now suppose  $D \in \mathscr{L}_n(L)$ .  $L \not \subset \mathcal{A}(\mathbb{T})$  and let  $M_1 = \neg L$ . Clearly

$$\mu(\tau D \cap M) = 0$$

for  $M \neq M_1$ . Recall that by Corollary 3.2 we can write for  $u \ge U(n, \varepsilon)$ 

$$D = D_u \cup E_u$$

where  $D_u$  is a union of copies of L in  $P^u(\mathbb{T})$  and  $\mu(E_u) < \varepsilon$ . As  $L \neq A(\mathbb{T})$  no copy of L in  $P^u(\mathbb{T})$  is a subset of  $A(P^u(\mathbb{T}))$ . Hence the column level in  $P^u(\mathbb{T})$  immediately above each copy of L where  $L \subset D_u$  will be a copy of  $M_1$ . Hence

$$\mu(\tau D_u \cap M_1) = \mu(D_u)$$

Stopping Time Transformations and Towers

and so

$$\mu(D) - \varepsilon < \mu(\tau D_{\mu} \cap M_{\mu}) \leq \mu(\tau D \cap M_{\mu}) \leq \mu(D).$$

As  $\varepsilon$  is arbitrary we have

 $\mu(\tau D \cap M_1) = \mu(D)$ 

as required for this choice of L since  $p(L, M_1) = 1$ .

Now suppose  $D \in \mathcal{L}_n(L)$  and  $L = A(C_{ij})$ ,  $1 \le i \le k$ ,  $1 \le j \le m$ . As  $\tau$  maps every copy of  $A(C_{ij})$  linearly to the base of a copy of  $T_{\sigma_j(i)}$  we have

$$\mu(\tau D \cap M) = 0$$

for  $M \not\subset B(T_{\sigma_i(i)})$ . We write, as before,

$$D = D_u \cup E_u$$

and let

$$D_u = D_u^1 \cup F_u$$

where  $D_u^1$  is a union of copies of L of rank less than  $2^u - 1$  and  $F_u$  is a union of copies of L of rank 2<sup>u</sup>. If u is large enough we have, for  $\varepsilon > 0$ ,

$$\mu(E_u \cup F_u) < \varepsilon.$$

Now, if  $M \subset B(T_{\sigma_i(i)})$ , we have

$$\mu(\tau D_u^1 \cap M) = \frac{\mu(M)}{w(T_1)} \mu(D_u^1)$$

and so

$$(\mu(D) - \varepsilon) \frac{\mu(M)}{w(T_1)} < \mu(\tau D_u^1 \cap M) \leq \mu(\tau D \cap M).$$

But

$$\mu(\tau D \cap M) = \mu(\tau D_u^1 \cap M) + \mu(\tau(E_u \cup F_u) \cap M) < \frac{\mu(M)}{w(T_1)} + \varepsilon.$$

As  $\varepsilon$  is arbitrary we have

$$\mu(\tau D \cap M) = \frac{\mu(M)}{w(T_1)} \mu(D)$$

as required for this choice of L. M as  $p(L, M) = \frac{\mu(M)}{w(T_1)}$ . Hence (3.3) holds and the proof of the lemma is complete.

For  $u \ge 1$ , let  $\mathscr{L}^{u}$  denote the partition of  $\mathbb{T}^*$  consisting of the column levels in  $P^u(\mathbb{T})$ . Analogous to Lemma 3.3 we have

**Lemma 3.4.** For each  $u \ge 1$ ,  $\mathcal{L}^u$  is a Markov partition for  $\tau$ .

*Proof.* Define for L,  $M \in \mathscr{L}^{"}$ 

$$p_{u}(L, M) = \begin{cases} 1 & L \notin A(P^{u}(\mathbf{T})), \ M = \tau L \\ 0 & L \notin A(P^{u}(\mathbf{T})), \ M \neq \tau L \\ \frac{\mu(M)}{w(P^{u}(T_{1}))} & L = A(C_{il}^{u}), \ M \in B(T_{\sigma_{i}^{u}(i)}) \\ 0 & L = A(C_{il}^{u}), \ M \notin B(T_{\sigma_{i}^{u}(i)}) \end{cases}$$

where  $C_{il}^{u}$  is the *l*th column of  $P^{u}(T_{i})$ . That this definition satisfies (3.1), (3.2), (3.3) may be checked in a way analogous to the proof of Lemma 3.3.

Let  $\mathscr{A}_0$  be the  $\sigma$ -algebra  $\bigvee_{-\infty}^{\infty} \tau^i \mathscr{L}$  and for  $u \ge 1$  let  $\mathscr{A}_u$  be the  $\sigma$ -algebra  $\bigvee_{-\infty}^{\infty} \tau^i \mathscr{L}^u$ .

**Corollary 3.5.**  $\tau$  is a Markov shift on each  $\mathcal{A}_u$ ,  $u \ge 0$ .

 $\tau$  is ergodic on  $\mathscr{A}_0$  if for each L,  $M \in \mathscr{L}$  we can find  $h \in \mathbb{Z}$  such that

$$(3.4) \quad \mu(\tau^h L \cap M) > 0.$$

Now if L is a column level in  $C_{ij}$  and M is a column level in  $C_{pq}$ ,  $1 \le i, p \le k$ ,  $1 \le j, q \le m$  and (i,j) = (p,q) then (3.4) holds with h the difference in the heights of L and M in the column  $C_{ij}$ . If  $(i,j) \ne (p,q)$  then (3.4) holds if and only if for some  $u \ge 1$  there is a column in  $P^{u}(\mathbb{T})$  having a copy of  $C_{ij}$  in the base and a copy of  $C_{pq}$  in the same column or vice versa. Now by (2.2) a typical column of  $P^{u}(T_{i})$  with a copy of  $C_{ij}$  in the case has the form

$$[C_{ij}, C_{\sigma_j(i), i_1}, C_{\sigma_{i_1}(\sigma_j(i)), i_2}, \dots, C_{\sigma_{i_2}^{u} - 2}^{(\sigma_{i_2}^{u} - 3)} (\dots (\sigma_{i_1}(\sigma_j(i))) \dots), i_2^{u} - 1]]$$

where  $1 \leq i_t \leq m$ ,  $1 \leq t \leq 2^{"} - 1$ . Hence we have

**Lemma 3.6.**  $\tau$  is ergodic on  $\mathscr{A}_0$  if and only if for each  $1 \leq i, p \leq k, 1 \leq j, q \leq m$ , at least one of the following holds

- (i) (i, j) = (p, q),
- (ii)  $\sigma_i(i) = p$  or  $\sigma_a(p) = i$ .

(iii) for some  $u \ge 2$  we can find  $1 \le s \le 2^u - 2$  and  $1 \le i_t \le m$ ,  $1 \le t \le s$  such that

$$\sigma_{i_s}(\sigma_{i_{s-1}}(\ldots(\sigma_{i_1}(\sigma_j(i)))\ldots)) = p$$

or

 $\sigma_{i_s}(\sigma_{i_{s-1}}(\ldots(\sigma_{i_1}(\sigma_q(p)))\ldots))=i.$ 

**Corollary 3.7.** If for some  $j_1, 1 \leq j_1 \leq m, \sigma_{j_1}$  is the permutation

 $\sigma_{i_1}(i) = i + 1 \pmod{k}$ 

then  $\tau$  is ergodic on  $\mathscr{A}_0$ .

*Proof.* If  $(i, j) \neq (p, q)$  and  $\sigma_i(i) \neq p$ , suppose that

 $p = \sigma_i(i) + s \pmod{k},$ 

then  $\sigma_{j_1}^{(s)}(\sigma_j(i)) = p$ , where  $\sigma_{j_1}^{(s)}$  denotes the sth power of  $\sigma_{j_1}$ . In a similar way we can prove

**Lemma 3.8.** Suppose  $u \ge 0$ . If for some  $j_1$ ,  $1 \le j_1 \le m^{2^u}$ ,  $\sigma_{j_1}^u$  is the permutation

 $\sigma_{i_1}^u(i) = i + 1 \pmod{k}$ 

then  $\tau$  is ergodic on  $\mathscr{A}_{u}$ .

**Theorem 3.9.**  $\tau$  is ergodic on  $\mathscr{B}(\mathbb{T}^*)$ , the Borel subsets of  $\mathbb{T}^*$ , if for some  $1 \leq j_1$ ,  $j_2 \leq m, \sigma_{j_1}, \sigma_{j_2}$  are defined by

 $\sigma_{j_1}(i) = i$  $\sigma_{j_2}(i) = i + 1 \pmod{k}.$ 

*Proof.* Note that for each  $u \ge 0$ ,  $l \in \{1, 2, ..., k\}$ 

$$\sigma_{(j_1-1)m^{2^{u-1}}-(j_1-1)m^{2^{u-2}}+\dots+(j_1-1)m+j_1}^{(l)}(l) = \sigma_{j_1}^{(2^u)}(l) = l$$

and

$$\sigma_{(j_1-1)m^{2^{n}-1}-(j_1-1)m^{2^{n}-2}-\dots-(j_1-1)m-j_2}^{u}(l) = \sigma_{j_2}(\sigma_{j_1}^{(2^n-1)}(l)) = l+1 \pmod{k},$$

where  $\sigma_{j_1}^{(2^u)}$  and  $\sigma_{j_1}^{(2^u-1)}$  denote the 2"th and  $(2^u-1)$ th powers of  $\sigma_{j_1}$  respectively. Hence, by Lemma 3.8,  $\tau$  is ergodic on each  $\mathcal{A}_u$ ,  $u \ge 0$ . As  $\mathscr{B}(\mathbb{T}^*)$  is the  $\sigma$ -algebra generated by  $\{\mathcal{A}_u: u \ge 0\}$ ,  $\tau$  is ergodic on  $\mathscr{B}(\mathbb{T}^*)$ .

We now investigate the conditions under which  $\tau$  is mixing. As  $\tau$  is a Markov Shift on each  $\mathcal{A}_u$ ,  $u \ge 0$ ,  $\tau$  will be mixing on  $\mathcal{A}_u$  if it is ergodic and aperiodic. In addition to conditions on the permutations we have to impose a condition on heights of certain columns of  $\mathbb{T}$ , as was the case for Friedman's construction, to obtain the aperiodicity of  $\tau$ .

**Lemma 3.10.** If for some  $j_1, j_2, j_3, 1 \le j_i \le m, i = 1, 2, 3$ , and some  $i, 1 \le i \le k$ , we have

(i) 
$$\sigma_{i}(l) = l+1 \pmod{k}, \quad l \in \{1, 2, \dots, k\},\$$

(ii) 
$$\sigma_{i,i}(l) = \sigma_{i,i}(l) = l, \quad l \in \{1, 2, \dots, k\},\$$

(iii) if h is the height of  $C_{ii}$ , then  $C_{ii}$  has height h+1.

then  $\tau$  is mixing on  $\mathscr{A}_0$ .

*Proof.*  $\tau$  is ergodic on  $\mathscr{A}_0$  by Lemma 3.7. Moreover

$$C_{(j_2-1)m^3+(j_2-1)m^2+(j_2-1)m+p}(P^2(T_i))$$
  
=  $[C_{ij_2}, C_{ij_2}, C_{ij_2}, C_{ip}], \quad 1 \le p \le m$ 

and

S.T. Richardson and K.M. Wilkinson

$$C_{(j_2-1)m^3+(j_3-1)m^2+(j_2-1)m+q}(P^2(T_i)) = [C_{ij}, C_{ij}, C_{ij}, C_{ia}], \quad 1 \le q \le m.$$

and so for any column level L in  $C_{ij_2}$  we have

$$\mu(\tau^{2h}L \cap L) > 0$$
 and  $\mu(\tau^{2h+1}L \cap L) > 0$ 

which implies that  $\tau$  is aperiodic and therefore mixing.

In a similar way we can show

**Lemma 3.11.** Suppose  $u \ge 0$ . If for some  $j_1, j_2, j_3, 1 \le j_i \le m^{2^u}$ , i = 1, 2, 3, and some *i*,  $1 \le i \le k$ .

- (i)  $\sigma_{j_1}^u(l) = l+1 \pmod{k}, \quad l \in \{1, 2, \dots, k\},$
- (ii)  $\sigma_{j_2}^u(l) = \sigma_{j_3}^u(l) = l, \quad l \in \{1, 2, \dots, k\},\$
- (iii) the heights of the  $j_2$ th and  $j_3$ th columns of  $P^u(T_i)$  differ by 1,

then  $\tau$  is mixing on  $\mathscr{A}_{\mu}$ .

However it turns out that if the conditions of Lemma 3.10 are satisfied by  $\mathbb{T}$  and  $\Sigma$  then the conditions of Lemma 3.11 are satisfied by  $P^{u}(\mathbb{T})$  and  $\Sigma_{u}$  for each  $u \ge 0$ .

**Theorem 3.12.** If **T** and  $\Sigma$  satisfy the conditions of Lemma 3.10 then  $\tau$  is mixing on each  $\mathscr{A}_u$ ,  $u \ge 0$ .

Proof. Note that

$$C_{(j_2-1)m^{2^{u}-1}+(j_2-1)m^{2^{u}-2}+\dots+(j_2-1)m+j_2}(P^u(T_i)) = [C_{ij_2}, C_{ij_2}, \dots, C_{ij_2}, C_{ij_2}]$$

and

$$C_{(j_2-1)m^{2^{n-1}}+(j_2-1)m^{2^{n-2}}+\dots+(j_2-1)m+j_3}(P^u(T_i)) = [C_{ij_2}, C_{ij_2}, \dots, C_{ij_2}, C_{ij_3}]$$

have heights differing by 1 and also for  $l \in \{1, 2, ..., k\}$ 

$$\sigma^{u}_{(j_2-1)m^{2^{u}-1}+(j_2-1)m^{2^{u}-2}+\dots+(j_2-1)m+j_2}(l) = \sigma^{(2^{u})}_{j_2}(l) = l$$

and

$$\sigma_{j_2-1}^{u} = \sigma_{j_3}(\sigma_{j_2}^{(2^u-1)}(l)) = l,$$

where  $\sigma_{j_2}^{(2^n)}$  and  $\sigma_{j_2}^{(2^n-1)}$  represent the 2<sup>*u*</sup>th and  $(2^n-1)$ th powers of  $\sigma_{j_2}$  respectively. Moreover for  $1 \leq l \leq k$ 

$$\sigma_{(j_2-1)m^{2^{n-1}}+(j_2-1)m^{2^{n-1}}+\dots+(j_2-1)m+j_1}^n(l) = \sigma_{j_1}(\sigma_{j_2}^{(2^n-1)}(l)) = l+1 \pmod{k}.$$

Hence the conditions of Lemma 3.11 hold and so  $\tau$  is mixing on each  $\mathscr{A}_u, u \ge 0$ . **Corollary 3.13.** If  $\mathbb{T}$  and  $\Sigma$  satisfy the conditions of Lemma 3.10,  $\tau$  is a Bernoulli shift on  $\mathscr{A}_u, u \ge 0$ . Stopping Time Transformations and Towers

*Proof.* In [3] it is shown that a mixing Markov shift is weak Bernoulli and hence a Bernoulli shift.

**Theorem 3.14.** If  $\mathbb{T}$  and  $\Sigma$  satisfy the conditions of Lemma 3.10 then  $\tau$  is Bernoulli on  $\mathscr{B}(\mathbb{T}^*)$ .

*Proof.* By Corollary 3.13,  $\tau$  is Bernoulli on each of the  $\sigma$ -algebras  $\mathscr{A}_u$ ,  $u \ge 0$ , which form an increasing sequence. Hence by a result of Ornstein [5]  $\tau$  is Bernoulli on  $\mathscr{B}(\mathbb{T}^*)$ , the  $\sigma$ -algebra generated by  $\mathscr{A}_u$ ,  $u \ge 0$ .

# §4. Stopping Time Transformations

Let  $\tau$  be an automorphism of  $(\Omega, \mathcal{B}, \mu)$ . If v is a measurable mapping from  $(\Omega, \mathcal{B})$  to the natural numbers  $\mathbb{N}$  then the transformation  $\tau^{v}$  defined by

 $\tau^{v}\omega = \tau^{n}\omega, \quad \omega \in \{v=n\}, \ n \in \mathbb{N}.$ 

is an automorphism of  $(\Omega, \mathcal{B}, \mu)$  if and only if the sets  $\tau^n \{v = n\}$ ,  $n \in \mathbb{N}$ , form a partition of  $\Omega$ . If this is the case v is called a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  and  $\tau^v$  is called a stopping time transformation. Neveu [4] established

**Theorem 4.1.** (i) If  $\Omega \supset D_0 \supset D_1 \dots$  is a sequence of measurable sets such that either  $D_{k+1} \neq \emptyset$  and  $D_k = \emptyset$  for some  $k \ge 1$  or  $\bigcap_{n=0}^{\infty} D_n = \emptyset$  then there is a stopping time v of  $(\Omega, \mathscr{B}, \mu, \tau)$  such that

(4.1) 
$$\tau^{v} \omega = \begin{cases} \omega & \omega \in E_{0} = D_{0}^{c} \\ \tau_{D_{0}}(\tau_{D_{1}}(\dots(\tau_{D_{n-1}}(\omega))\dots)) & \omega \in E_{n}, n \ge 1 \end{cases}$$

where  $E_n = D_{n-1} - D_n$ ,  $n \ge 1$ , and for any  $D \in \mathcal{B}$ ,  $\tau_D$  represents the induced transformation of  $\tau$  on D.

(ii) Corresponding to any stopping time v of  $(\Omega, \mathcal{B}, \mu, \tau)$  there is a sequence of measurable sets  $\Omega \supset D_0 \supset D_1 \dots$  as in (i) such that  $\tau^v$  has the form (4.1).

(iii) If v is a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  and  $\tau$  is ergodic then  $\int v d\mu$  is either a non-negative integer or  $+\infty$ . If  $\int v d\mu = k < \infty$  then the sets  $D_n$ .  $n \ge 0$ . defined in (ii) satisfy  $D_{k-1} \neq \emptyset$ .  $D_k = \emptyset$  whereas if  $\int v d\mu = +\infty$  we have  $D_n \neq \emptyset$ .  $n \ge 0$  and  $\bigcap_{n=0}^{\infty} D_n = \emptyset$ .

In the sequel we shall assume that  $\int v d\mu = k < \infty$ , noting that if  $\int v d\mu = +\infty$ and  $\Omega \supset D_0 \supset D_1$ ... are the sets defined by Theorem 4.1(ii) the transformation

$$\tau^{v_k} = \begin{cases} \omega & \omega \in D_k \\ \tau^v \omega & \omega \in D_k^c \end{cases}$$

differs from  $\tau^{\nu}$  on a set of small measure for large k. As  $\tau^{\nu}$  is the identity transformation on  $\{\nu = 0\}$ , for our study of mixing properties of  $\tau^{\nu}$  we shall restrict attention to the behaviour of  $\tau^{\nu}$  on  $D_0 = \bigcup_{n=1}^{k} E_n$ . Note that the sets  $E_n$ .

 $0 \leq n \leq k$ , form a partition of  $\Omega$ . We shall only consider those transformations  $\tau$  for which there is a tower T satisfying  $T^* = \Omega$  and  $\tau = \tau(T)$  and we shall assume that our stopping time is chosen so that each of the sets  $E_n$ ,  $0 \leq n \leq k$ , defined in Theorem 4.1 is a union of column levels in T. If a column level L in T is a subset of  $E_n$ ,  $0 \leq n \leq k$ , we say that L has name n and write n(L) = n. We shall also assume that in each column of T there is at least one column level with name k. Note that if for some tower T and stopping time v these conditions are not satisfied then, as the towers  $S^u(T)$ ,  $u \geq 1$ , have smaller and smaller column widths, we can modify v on as small a subset of  $\Omega$  as we please to have these conditions satisfied by  $S^u(T)$  for some  $u \geq 1$  and a slightly modified stopping time.

As we restrict attention to the action of  $\tau^{v}$  on  $\bigcup_{n=1}^{k} E_{n}$  we delete from T any column level with name 0 and denote by  $\hat{T}$  the resulting tower. Note that the deleted set is invariant under  $\tau^{v}$  and that  $\widehat{S^{u}(T)} = S^{u}(\hat{T}), u \ge 1$ .

Recall that in the definition of  $\tau$ ,  $\tau$  maps each column level in T except tops of columns, linearly onto another column level within the same column of T. If  $C = (L_i: 1 \le i \le h)$  is a column in  $\hat{T}$ , in order to define  $\tau^v$  on a level  $L_i$  with name n, by Theorem 4.1, we first have to define  $\tau_{D_{n-1}}$  on  $L_i$ , then, if  $n \ge 2$ , define  $\tau_{D_{n-2}}$ on  $\tau_{D_{n-1}}(L_i)$ , etc. As  $D_{n-1} = \bigcup_{l=n}^{k} E_l$ ,  $1 \le n \le k$ ,  $\tau_{D_{n-1}}(L_i)$  is another column level within C if and only if

$$\{j: i+1 \leq j \leq h, n(L_i) \geq n\}$$

is non-empty, in which case  $\tau_{D_{n-1}}$  maps  $L_i$  linearly onto  $L_{i_1}$  where

 $i_1 = \min \{j: i+1 \leq j \leq h, n(L_j) \geq n\}.$ 

Then, if  $n \ge 2$ ,  $\tau_{D_{n-1}}(L_{i_1})$  is another column level within C if and only if

$${j: i_1 + 1 \le j \le h. n(L_j) \ge n - 1}$$

is non-empty, in which case  $\tau_{D_{n-2}}$  maps  $L_{i_1}$  linearly onto  $L_{i_2}$  where

$$i_2 = \min \{j: i_1 + 1 \le j \le h, n(L_j) \ge n - 1\},\$$

Clearly there may be column levels other than A(C), i.e. some  $L_i$ ,  $1 \le i \le h-1$ , for which, if  $n(L_i) = n$ , some or all of the sets

$$\tau_{D_{n-1}}(\tau_{D_{n-1}-1}(\dots(\tau_{D_{n-1}}(L_i))\dots)), \quad 1 \le l \le n.$$

are not column levels within C. In our tower block construction for  $z^v$  we shall be increased in the number of these sets which are not column levels within C. For a column level  $L_i$ ,  $1 \le i \le h$ , with name n we shall define  $u(L_i) = 0$  if each of

$$\tau_{D_{n-l}}(\tau_{D_{n-l-1}}(\dots(\tau_{D_{n-1}}(L_{i}))\dots)), \quad 1 \leq l \leq n.$$

Stopping Time Transformations and Towers

is a column level within C and  $u(L_i) = u \ge 1$  if each of

 $\tau_{D_{n-l}}(\tau_{D_{n-l+1}}(\dots(\tau_{D_{n-1}}(L_{i}))\dots)), \quad 1 \leq l \leq n-u,$ 

is a column level within C but

 $\tau_{D_{\mu-1}}(\tau_{D_{\mu}}(\ldots(\tau_{D_{n-1}}(L_i))\ldots)).$ 

is not a column level in  $C(u(L_i) = n \text{ meaning } \tau_{D_{n-1}}(L_i)$  is not a column level in C). So  $u(L_i)$  counts the number of steps in the iterative definition of  $\tau^{\nu}$  on  $L_i$ given by Theorem 4.1 which are left undefined within the column containing  $L_i$ .

Equivalently, if  $n(L_i) = n$ ,  $u(L_i) = 0$  if the set

 $\{i_1, i_2, \dots, i_n: i < i_1 < i_2 < \dots < i_n \leq h, n(L_{i_i}) \geq n - j + 1, 1 \leq j \leq n\}$ 

is non-empty,  $u(L_i) = u$ , with  $1 \le u \le n-1$  if the sets

$$\{i_1, i_2, \dots, i_s: i < i_1 < i_2 < \dots < i_s \leq h, n(L_i) \geq n - j + 1, 1 \leq j \leq s\}$$

are non-empty for  $s \leq n-u$  and empty for  $s \geq n-u+1$  and  $u(L_i) = n$  if the set

 $\{j: i < j \leq h, n(L_j) \geq n\}$ 

is empty.

Note that for any column level L in  $\hat{T}$  with n(L) = n we have  $0 \leq u(L) \leq n$ . Any column level L in  $\hat{T}$  with  $u(L) \ge 1$  is called a v-top. Note that a column level L is not a v-top if and only if there is another column level L in the same column as L and above L such that  $\tau^{v}(L) = L$ . If C is a column in  $\hat{T}$  we let  $A_1(C), A_2(C), \dots, A_l(C)$  denote the v-tops in C where we number from the top of C, so  $A_1(C) = A(C)$ .  $A_2(C)$  is the first v-top in C below  $A_1(C)$ , etc. We now study properties of the v-tops in  $\hat{T}$  and, in particular, we shall show that each column in  $\hat{T}$  may be assumed to contain exactly k v-tops.

**Lemma 4.2.** For each column C in  $\hat{T}$ ,  $n(A_i(C)) \ge i$ ,  $1 \le i \le l$ .

*Proof.* Clearly  $n(A_1(C)) \ge 1$ . Suppose  $n(A_i(C)) \ge i$ ,  $1 \le i \le j$ . Then if  $n(A_{j+1}(C)) \le j$ we have  $u(A_{j+1}(C)) = 0$  and so  $A_{j+1}(C)$  is not a v-top.

**Corollary 4.3.** There are at most k x-tops in each column of  $\hat{T}$ .

*Proof.* For any column level L in  $\hat{T}$  we have  $n(L) \leq k$ .

**Lemma 4.4.** If C is a column in  $\hat{T}$  and L is a column level of C lying between  $A_i(C)$  and  $A_{i+1}(C)$ ,  $1 \le i \le l-1$ , then  $n(L) \le i$ .

*Proof.* The result is clearly true for i=1 since  $A_2(C)$  must be the first column level in C (counting from the top of C) with name greater than or equal to 2. Now suppose the result is true for  $1 \le i \le j$ . If  $j \le l-2$ , let L be the first column level below  $A_{j+1}(C)$  with  $n(L) = n \ge j+2$ . We show that L is a v-top and so L  $=A_{j+2}(C)$  and the result is true for  $1 \leq i \leq j+1$ .

Now  $\tau_{D_{n-1}}(L)$  is a column level in  $\overline{C}$  if and only if there is a column level L'above L with  $n(L) \ge n$  and so by the inductive hypothesis and the definition of L. L' can only be one of  $A_i(C)$ ,  $1 \le i \le j+1$ . Similarly each of the sets

$$\tau_{D_{n-t}}(\tau_{D_{n-t+1}}(\dots(\tau_{D_{n-1}}(L))\dots)), \quad 1 \le t \le n,$$

can be a column level in C if and only if it is equal to one of  $A_i(C)$ ,  $1 \le i \le j+1$ . Thus we see that  $u(L) \ge 1$  and so L is a v-top.

**Corollary 4.5.** If  $n(A_i(C)) = n$  and  $u(A_i(C)) = u$ , then each of

$$\tau_{D_{n-i}}(\tau_{D_{n-i+1}}(\dots(\tau_{D_{n-1}}(A_j(C))\dots))), \quad 1 \le i \le n-u,$$

is a v-top.

**Corollary 4.6.** Suppose  $\hat{T} = (C_j: 1 \le j \le m)$  and  $S(\hat{T}) = (C'_j: 1 \le j \le m^2)$ . If  $l_j$  is the number of v-tops in  $C_j$ ,  $1 \le j \le m$ , and  $l'_j$  is the number of v-tops in  $C'_j$ ,  $1 \le j \le m^2$ , then if  $l_j \le k - 1$ .

 $l'_{(i-1)m+j} \ge l_j + 1$ 

for each  $1 \leq i \leq m$ .

*Proof.* Recall that for  $1 \leq i \leq m$ 

$$C_{(i-1)m+i}(S(\hat{T})) = [C_i, C_i].$$

The copies of  $A_n(C_j)$ ,  $1 \le n \le l_j$ , will again be v-tops in  $C_{(i-1)m+j}$ . Moreover if L is the highest column level in  $C_i$  with  $n(L) \ge l_j + 1$ , which exists since we have assumed that each column in T contains a column level with name k, the copy of L in the base of  $C_{(i-1)m+j}$  will be a v-top in  $C_{(i-1)m+j}$ .

**Corollary 4.7.** For any tower T each column in  $S^{k-1}(\hat{T})$  contains exactly k v-tops.

Because of this last corollary we shall assume in the sequel that T is chosen so that each column in  $\hat{T}$  contains exactly k v-tops. If  $\hat{T} = (C_j: 1 \le j \le m)$  we shall let

$$A^{\mathbf{v}}(\hat{T}) = \bigcup_{j=1}^{m} \bigcup_{n=1}^{k} A_n(C_j).$$

**Lemma 4.8.** If C is a column in  $\hat{T}$  with v-tops  $A_i(C)$ ,  $1 \leq i \leq k$ , then  $i \neq j$  implies

$$u(A_i(C)) \neq u(A_i(C)).$$

*Proof.* Suppose  $i \neq j$  and  $u(A_i(C)) = u(A_j(C)) = u$ , say, and let  $n(A_i(C)) = n_i$ ,  $n(A_j(C)) = n_j$ . If

$$L_{i} = \begin{cases} \tau_{D_{u}}(\tau_{D_{u-1}}(\dots(\tau_{D_{n_{i}-1}}(A_{i}(C)))\dots)), & u \leq n_{i}-1 \\ A_{i}(C), & u = n_{i} \end{cases}$$

and

$$L_{j} = \begin{cases} \tau_{D_{u}}(\tau_{D_{u-1}}(\dots,\tau_{D_{n_{j}-1}}(A_{j}(C)))\dots)), & u \leq n_{j}-1 \\ A_{j}(C), & u = n_{j} \end{cases}$$

then  $L_i \neq L_j$  and we may as well assume that  $L_j$  is above  $L_i$  in C. We shall see that whatever the values of u,  $n_i$ ,  $n_j$ , then  $\tau_{D_{u-1}}(L_i)$  is a column level in C, contradicting the assumption that  $u(A_i(C)) = u$ . If  $u = n_j$  the result follows from the fact that  $n(L_j) = u$ , whereas if  $u \leq n_j - 1$  the result follows from the fact that  $n(L_i) \geq u + 1$ .

# § 5. The Tower Block Construction for $\tau^{v}$

In this section se shall show that if  $\tau = \tau(T)$  and  $\nu$  is a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  with  $\int v d\mu = k$  then we can find a tower block  $\mathbb{T}$  and a set of permutations  $\Sigma$  such that  $\tau^{\nu} = \tau(\Sigma, \mathbb{T})$ .

Let C be a column in  $\hat{T}$  and let  $L_n(C)$  be the first column level (counting from the base) which satisfies  $n(L_n(C)) \ge n$ . Then define  $B_1(C) = L_1(C)$  and for  $n \ge 2$ 

$$B_n(C) = \tau_{D_0}(\tau_{D_1}(\dots(\tau_{D_{n-2}}(L_n(C)))\dots)),$$

where we assume that T is chosen so that each  $B_n(C)$ .  $1 \le n \le k$ , is a column level in C. (If this is not the case for T it will be for  $S^u(T)$  for some  $u \ge 1$ .)

Now for each  $B_n(C)$ ,  $1 \le n \le k$ , there is  $i_n \ge 0$  such that  $(\tau^v)^{i_n} B_n(C)$  is a v-top in C and in the column

 $D_n(C) = \{(\tau^v)^j B_n(C) : 0 \le j \le i_n\}$ 

 $\tau^{v}$  maps each column level except the v-top  $(\tau^{v})^{i_{n}}B_{n}(C)$  linearly onto the column level immediately above it. Since  $\tau^{v}$  is invertible we have

$$\bigcup_{n=1}^{k} D_n(C)^* = C^*$$

Now if  $\hat{T} = (C_j: 1 \le j \le m)$  we define for  $1 \le n \le k$ 

$$T_n = (C_{ni}: 1 \leq j \leq m)$$

where  $C_{ni} = D_n(C_i)$  and we let

$$\mathbb{T} = (T_n: 1 \leq n \leq k).$$

Note that  $\mathbf{T}^* = \hat{T}^*$ . We also clearly have

Lemma 5.1. (i)  $A(\mathbf{T}) = A^{v}(\hat{T})$ ,

(ii)  $\tau_{\mathbb{T}}(\omega) = \tau^{v}(\omega), \quad \omega \in \mathbb{T}^{*} - A(\mathbb{T}).$ 

We shall define  $\Sigma = \{\sigma_j : 1 \le j \le m\}$  by  $\sigma_j(n) = u(A(C_{nj}))$ . Because of Lemma 4.8, each  $\sigma_j$ ,  $1 \le j \le m$ , is a permutation of  $\{1, 2, \dots, k\}$ .

**Lemma 5.2.** For each  $u \ge 0$ .

(i) 
$$A(P^{u}(\mathbf{T})) = A^{v}(S^{u}(\hat{T})).$$

(ii)  $\tau_{P^u(\mathbb{T})}(\omega) = \tau^v(\omega), \quad \omega \in \mathbb{T}^* - \mathcal{A}(P^u(\mathbb{T})).$ 

*Proof.* The result is true for u=0 by Lemma 5.1. Suppose the result is true for  $0 \leq u \leq v$  and let

$$P^{v}(T_{i}) = (C_{ij}: 1 \leq j \leq m^{2^{v}}),$$
  
$$S^{v}(\hat{T}) = (C_{j}: 1 \leq j \leq m^{2^{v}}).$$

Then

$$A(P^{\nu+1}(\mathbf{T})) = \bigcup_{n=1}^{k} A(P^{\nu+1}(T_n))$$
  
=  $\bigcup_{j=1}^{m^{2^{\nu}}} \bigcup_{n=1}^{k} A(\alpha_{j-1}(\frac{1}{2}(P^{\nu}(T_{\sigma_j^{\nu}(n)}))_1))$ 

where 
$$\alpha_{j-1} = w(C_{ij})/w(P^{v}(T_{i}))$$
  

$$= \bigcup_{j=1}^{m^{2^{v}}} \bigcup_{n=1}^{k} A(\alpha_{j-1}(\frac{1}{2}(P^{v}(T_{n}))_{1}))$$

$$= \bigcup_{n=1}^{k} A(\frac{1}{2}(P^{v}(T_{n}))_{1})$$

$$= A^{v}(\frac{1}{2}(S^{v}(\hat{T}))_{1}) = A^{v}(S^{v-1}(\hat{T})).$$

So (i) holds for u = v + 1. To check that (ii) holds for u = v + 1, it is sufficient to check that  $\tau^{v}$  agrees with  $\tau_{P^{v-1}(\mathbb{T})}$  on  $\mathcal{A}(P^{v}(\mathbb{T})) - \mathcal{A}(P^{v-1}(\mathbb{T}))$ . Now for each  $1 \leq i \leq k$ ,  $1 \leq j \leq m^{2^{v}}$ .  $\mathcal{A}(C_{ij})$  is a v-top in  $C_{j}$  in  $S^{v}(\hat{T})$ . Suppose

that

$$u(A(C_{ij})) = u, \quad n(A(C_{ij})) = n.$$

So each of

$$\tau_{D_{n-l}}(\tau_{D_{n-l+1}}(\dots(\tau_{D_{n-1}}(A(C_{ij})))\dots))), \quad 1 \leq l \leq n-u$$

is a column level in  $S^{v}(\hat{T})$  above  $A(C_{ij})$  but

 $\tau_{D_{u-1}}(\tau_{D_{u}}(\dots(\tau_{D_{n-1}}(A(C_{ij})))\dots))$ 

is not a column level in  $S^{r}(\hat{T})$ . Now looking at  $S^{r+1}(\hat{T})$  we have

$$\tau_{D_{u-1}}(\tau_{D_{u}}(\dots(\tau_{D_{n-1}}((\frac{1}{2}A(C_{ij}))_{0}))\dots)) = \bigcup_{l=1}^{m^{2^{v}}} \alpha_{j-1}(\frac{1}{2}L_{u}(C_{l}))_{1}$$

where  $\alpha_{j-1} = w(C_j)/w(S^{r}(\hat{T}))$ , and so

$$\tau^{v}((\frac{1}{2}A(C_{ij}))_{0}) = \bigcup_{l=1}^{m^{2^{v}}} \alpha_{j-1}(\frac{1}{2}B_{u}(C_{l}))_{1}$$
$$= \bigcup_{l=1}^{m^{2^{v}}} \alpha_{j-1}(\frac{1}{2}B_{\tau_{j}'(i)}(C_{l}))_{1}$$
$$= B(\alpha_{j-1}(\frac{1}{2}P^{v}(T_{\tau_{j}'(i)}))_{1})$$

and so  $\tau^{v}$  agrees with  $\tau_{P^{v}(\mathbb{T})}$  on each  $A(C_{ij}), 1 \leq i \leq k, 1 \leq j \leq m^{2^{v}}$  and so

 $\tau_{P^{\nu+1}(\mathbf{T})}(\omega) = \tau^{\nu}(\omega)$ 

for  $\omega \in \mathbb{T}^* - A(P^{\nu+1}(\mathbb{T}))$  and so, by induction, the lemma is proved. Hence we have

**Theorem 5.3.**  $\tau^{\nu} = \tau(\Sigma, \mathbb{T}).$ 

And the following theorem follows from the results of  $\S 3$ .

**Theorem 5.4.** (i)  $\tau^{v}$  is ergodic on  $\{v \neq 0\}$  if for some  $1 \leq j_1, j_2 \leq m, \sigma_{j_1}, \sigma_{j_2}$  are defined by

 $\sigma_{i_1}(i) = i, \quad \sigma_{i_2}(i) = i+1 \pmod{k}, \ i \in \{1, 2, \dots, k\}.$ 

(ii)  $\tau^{v}$  is Bernoulli on  $\{v \neq 0\}$  if for some  $j_1, j_2, j_3, 1 \leq j_i \leq m, i = 1, 2, 3$ , and some  $i, 1 \leq i \leq k$  we have

(a) 
$$\sigma_{i_1}(l) = l + 1 \pmod{k}, \quad l \in \{1, 2, ..., k\},\$$

(b)  $\sigma_{j_2}(l) = \sigma_{j_3}(l) = l, \quad l \in \{1, 2, ..., k\},\$ 

(c) If  $T_i = (C_{ij}, 1 \le j \le m)$  then the heights of  $C_{ij}$ , and  $C_{ij}$ , differ by one.

### §6. Density Results

In order to investigate density results we need to investigate to what extent a stopping time needs to be modified in order to satisfy the conditions for  $\tau^v$  to be ergodic or Bernoulli as outlined in Theorem 5.4. These conditions are mainly concerned with the permutations  $\Sigma$  which are in turn determined by the names of column levels in T. So in this section we shall assign names between 0 and  $k \ge 1$  to column levels in T and take v to be the corresponding stopping time of  $(\Omega, \mathcal{B}, \mu, \tau(T))$  as defined by Theorem 4.1(i). It will be convenient when assigning names to column levels to insist that if v is the stopping time which arises then the kv-tops in that column are the top k levels in that column.

**Lemma 6.1.** Let  $A_i$ ,  $1 \le i \le k$ , be the top k levels (numbered from the top) in a column C of  $\hat{T}$ . A necessary and sufficient condition for each  $A_i$ ,  $1 \le i \le k$ , to be a v-top is

$$n(A_i) \in \{i, i+1, \dots, k\}, \qquad 1 \leq i \leq k.$$

*Proof.* That this condition is necessary is shown in Lemma 4.2. If  $n(A_i) = n_i \ge i$ , then not all of

$$\tau_{D_{n_i-j}}(\tau_{D_{n_i-j+1}}(\dots(\tau_{D_{n_i-1}}(A_i))\dots)), \quad 0 \leq j \leq n_i$$

can be column levels in C as there are only i-1 column levels above  $A_i$  in C.

We shall define an *n*-sequence of order k to be a collection  $\{n_i: 1 \le i \le k\}$  satisfying  $n_i \in \{i, i+1, ..., k\}, 1 \le i \le k$ . Thus if we give the *i*th level (counting from

the top) of a column in  $\hat{T}$  the name  $n_i$ ,  $1 \le i \le k$ , where  $\{n_i: 1 \le i \le k\}$  is an *n*-sequence of order k, then these k column levels are v-tops. Note that the number of different *n*-sequences of order k is k!. If  $\{n_i: 1\le i\le k\}$  is an *n*-sequence of order k then each of

$$\{l, n_1 + 1, n_2 + 1, \dots, n_k + 1\}, \quad 1 \le l \le k + 1,$$

is an *n*-sequence of order k+1 and in fact we can obtain each of the (k+1)! *n*-sequences of order k+1 from the *n*-sequences of order k by this method.

A collection  $\{u_i: 1 \le i \le k\}$  is called a *u*-sequence of order *k* if it is a permutation of  $\{1, 2, ..., k\}$ . Corresponding to each *n*-sequence of order *k* we define a *u*-sequence of order *k* by  $u_1 = n_1$  and for  $j \ge 2$ ,  $u_j = u$ ,  $1 \le u \le n_j - 1$ , if the sets

$$\{i_1, i_2, \dots, i_s: 1 \le i_s < i_{s-1} < \dots < i_1 < j, \ n_{i_t} \ge n_j - t + 1, \ 1 \le t \le s\}$$

are non-empty for  $s \leq n_j - u$  and empty for  $s \geq n_j - u + 1$  and  $u_j = n_j$  if the set

$$\{i: 1 \leq i < j, n_i \geq n_i\}$$

is empty.

Equivalently, if  $A_i$ ,  $1 \le i \le k$ , are the top k levels (counting from the top) of a column in  $\hat{T}$  with  $n(A_i) = n_i$ , then  $u_i = u(A_i)$ . That the collection  $\{u_i: 1 \le i \le k\}$  is a u-sequence of order k follows from Lemma 4.8.

**Lemma 6.2.** Let  $\{n_i: 1 \le i \le k\}$  be a n-sequence of order k and  $\{u_i: 1 \le i \le k\}$  the corresponding u-sequence. If  $\{n'_i: 1 \le i \le k+1\}$  is the n-sequence of order k+1 defined for some  $l, 1 \le l \le k+1$ , by

$$n'_{1} = l$$
  
$$n'_{j+1} = n_{j} + 1 \qquad 1 \le j \le k$$

then the corresponding u-sequence of order k + 1,  $\{u'_i: 1 \leq i \leq k + 1\}$  is defined by

$$u'_1 = l$$

and for  $1 \leq j \leq k$ .

$$u_{j+1}' \begin{cases} u_j & 1 \leq u_j \leq l-1 \\ u_j+1 & l \leq u_j \leq k. \end{cases}$$

*Proof.*  $u'_1 = l$  by definition. Now for  $2 \le j \le k+1$ ,  $1 \le s \le n'_j - 1$ .

$$\{i_1, i_2, \dots, i_s \colon 2 \leq i_s < i_{s-1} < \dots < i_1 < j, \ n'_{i_t} \geq n'_j - t + 1, \ 1 \leq t \leq s\}$$
  
=  $\{j_1 + 1, j_2 + 1, \dots, j_s + 1 \colon 1 \leq j_s < j_{s-1} < \dots < j_1 < j - 1,$   
 $n_{j_t} \geq n_{j-1} - t + 1, \ 1 \leq t \leq s\}.$ 

The set on the right hand side is non-empty for  $s \le n_{j-1} - u_{j-1}$  and empty for  $s \ge n_{j-1} - u_{j-1} + 1$ . Since  $n'_1 = l$ , if  $1 \le u_{j-1} \le l-1$  and  $s = n_{j-1} - u_{j-1}$ .

Stopping Time Transformations and Towers

$$\{i_1, i_2, \dots, i_s, i_{s-1} \colon 1 \leq i_{s-1} < i_s < \dots < i_1 < j, \ n'_{i_t} \geq n'_j - t + 1, \ 1 \leq t \leq s + 1\}$$
  
=  $\{1\} \cup \{i_1, i_2, \dots, i_s \colon 2 \leq i_s < \dots < i_1 < j, \ n'_{i_t} \geq n'_j - t + 1, \ 1 \leq t \leq s\}$ 

and so  $u'_i = u_{j-1}$  in this case, whereas if  $l \leq u_{j-1} \leq k$  and  $s = n_{j-1} - u_{j-1}$ 

$$\{i_1, i_2, \dots, i_s, i_{s-1} \colon 1 \leq i_{s-1} < i_s < \dots < i_1 < j, \ n'_{i_t} \geq n'_j - t + 1, \ 1 \leq t \leq s + 1 \}$$
  
=  $\{i_1, i_2, \dots, i_s \colon 2 \leq i_s < \dots < i_1 < j, \ n'_{i_t} \geq n'_j - t + 1, \ 1 \leq t \leq s \}$ 

and so  $u'_{i-1} = u_{i-1} + 1$ .

**Theorem 6.3.** The mapping which takes an n-sequence of order k to the corresponding u-sequence of order k is one-to-one.

*Proof.* Note that there are k! *n*-sequences of order k and k! *u*-sequences of order k. Hence it is sufficient to check that each *n*-sequence of order k gives rise to a unique *u*-sequence of order k. We proceed by induction on k. Clearly the result is true for k = 1. Suppose to each *n*-sequence of order k there is associated a unique *u*-sequence of order k. Lemma 6.2 shows that the procedure we have used to proceed from *n*-sequences of order k to *n*-sequences of order k+1 will preserve the uniqueness of definition of *u*-sequences.

As a result of this theorem there is a unique *n*-sequence of order k which gives rise to any particular *u*-sequence of order k. As  $\Sigma$  depends on the value of *u* on *v*-tops we want to see how we must change the names in columns of  $\hat{T}$  in order to obtain the permutations required for the resulting transformations to be ergodic or Bernoulli. Let  $\hat{T} = (C_j; 1 \le j \le m)$  and for each  $1 \le j \le m$  let  $A_i(C_j)$ ,  $1 \le i \le k$ , represent the top k levels (counting from the top) of the column  $C_j$ .

**Theorem 6.4.** Let  $\rho$  be a permutation of  $\{1, 2, ..., k\}$ . We can assign names to  $A_i(C_i), 1 \leq i \leq k$ , so that  $\sigma_i = \rho$ .

*Proof.* Let  $B_n(C_j)$ .  $1 \le n \le k$ , be as defined in § 5 and first of all assign the name k to each  $A_i(C_j)$ .  $1 \le i \le k$ . This ensures that each  $A_i(C_j)$ .  $1 \le i \le k$ , is a v-top. Now for each n,  $1 \le n \le k$ , we can find  $t(n) \in \{1, 2, ..., k\}$  and  $i_n \ge 0$  such that

 $(\tau^{v})^{i_{n}} B_{n}(C_{j}) = A_{t(n)}(C_{j}).$ 

Note that t is a permutation of  $\{1, 2, \dots, k\}$ .

Now we shall assign a new name  $n_i$  to  $A_i(C_j)$ ,  $1 \le i \le k$ , with  $\{n_i: 1 \le i \le k\}$  a particular *n*-sequence of order k and we denote by v' the stopping time associated with the partition  $E'_i$ ,  $1 \le i \le k$ , by Theorem 4.1, where  $E'_i$  is the union of column levels now having name i. Note that changing the names in this way does not affect the permutation t. i.e., we shall still have for each  $1 \le n \le k$ 

$$(\tau^{\mathbf{v}'})^{i_n} B_n(C_j) = A_{\tau(n)}(C_j),$$

for, if  $i_n = 0$ , this is clearly true and if  $i_n \ge 1$  and

$$n((\tau^{v})^{i_{n}-1} B_{n}(C_{j})) = l$$

for some l with  $1 \leq l \leq k$ , then for some h,  $1 \leq h \leq l$ ,

$$\tau_{D_{h}}(\tau_{D_{h+1}}(\ldots(\tau_{D_{l-1}}((\tau^{v'})^{i_{n-1}}B_{n}(C_{j})))\ldots)) \neq \bigcup_{i=1}^{k} A_{i}(C_{j})$$

and

$$\tau_{D_{h-1}}(\tau_{D_{h}}(\ldots(\tau_{D_{l-1}}((\tau^{v'})^{i_{n-1}}B_{n}(C_{j})))\ldots)) \subset \bigcup_{i=1}^{k} A_{i}(C_{j}).$$

But because  $n(A_k(C_j)) = k$  we have

$$\tau_{D_{h-1}}(\tau_{D_h}(\ldots(\tau_{D_{l-1}}((\tau^{v'})^{i_n-1}B_n(C_j)))\ldots)) = A_k(C_j)$$

and so

$$\tau_{D_0}(\tau_{D_1}(\ldots(\tau_{D_{l-1}}((\tau^{\nu'})^{i_n-1}B_n(C_j)))\ldots)) = A_{k-h+1}(C_j)$$

irrespective of the particular values of  $n(A_i(C_j)) \in \{i, i+1, ..., k\}, 1 \leq i \leq k$ .

In order to have  $\sigma_j = \rho$  the particular *n*-sequence of order k we choose is the one corresponding to the *u*-sequence of order k defined by

$$u_i = \rho(t^{-1}(i)), \qquad 1 \leq i \leq k.$$

Then, for  $1 \le i \le k$ ,  $\sigma_i(i) = u_{t(i)} = \rho(i)$ , as required.

In discussing the density of stopping time transformations with specified properties we shall work with the uniform metric

 $d(\tau_1, \tau_2) = \mu\{\omega: \tau_1(\omega) \neq \tau_2(\omega)\}$ 

where  $\tau_1, \tau_2$  are both automorphisms of  $(\Omega, \mathcal{B}, \mu)$ . It will be of use to relate the proximity of two stopping time transformations with the proximity of the associated partitions which are defined by Neveu's Theorem.

**Lemma 6.5.** Let v and v' be stopping times for  $(\Omega, \mathcal{B}, \mu, \tau)$  with  $\int v d\mu = \int v' d\mu = k \ge 1$ and let  $E_n$ ,  $0 \le n \le k$ , and  $E'_n$ ,  $0 \le n \le k$  be as defined for v and v' respectively by Theorem 4.1. For any  $\varepsilon > 0$  there is  $\eta(\varepsilon, k)$  such that

$$\mu(E_n \, \exists \, E'_n) < \eta(\varepsilon, k), \qquad 0 \leq n \leq k.$$

implies

$$d(\tau^{v},\tau^{v'}) < \varepsilon$$

*Proof.* The result follows from the fact that, for  $1 \le n \le k$ ,

$$\mu(\{\omega \in E_n \cap E'_n : \tau^{v}(\omega) \neq \tau^{v'}(\omega)\}) \leq \sum_{l=0}^{n-1} \mu(D_l \varDelta D'_l)$$

where

$$D_l = \bigcup_{m=l+1}^k E_m$$
 and  $D'_l = \bigcup_{m=l+1}^k E'_m$ .

For the case n = 1 we note that

$$\{\omega \in E_1 \cap E'_1 : \tau^{v}(\omega) = \tau^{v'}(\omega)\} = \{\omega \in E_1 \cap E'_1 : \tau_{D_0}(\omega) = \tau_{D'_0}(\omega)\}$$
$$= (E_1 \cap E'_1) \cap \tau_{D_0}^{-1}(D_0 \cap D'_0) \cap \tau_{D'_0}^{-1}(D_0 \cap D'_0).$$

Hence

$$\{\omega \in E_1 \cap E_1' : \tau^{v}(\omega) \neq \tau^{v'}(\omega)\} = (E_1 \cap E_1') \cap [\tau_{D_0}^{-1}(D_0 \setminus D_0') \cup \tau_{D_0'}^{-1}(D_0' \setminus D_0)]$$

and so

$$\mu(\{\omega \in E_1 \cap E'_1 \colon \tau^{\nu}(\omega) \neq \tau^{\nu'}(\omega)\}) \leq \mu(D_0 \setminus D'_0) + \mu(D'_0 \setminus D_0) = \mu(D_0 \varDelta D'_0).$$

For the case n=2 we note that if

$$\omega \in E_2 \cap E_2' \cap \tau_{D_1}^{-1} [(D_1 \cap D_1') \cap \tau_{D_0}^{-1} (D_0 \cap D_0')] \cap \tau_{D_1'}^{-1} [(D_1 \cap D_1') \cap \tau_{D_0'}^{-1} (D_0 \cap D_0')]$$

then

$$\tau_{D_{1}}(\omega) = \tau_{D_{1}'}(\omega) \text{ and } \tau_{D_{0}}(\tau_{D_{1}}(\omega)) = \tau_{D_{0}'}(\tau_{D_{1}'}(\omega)), \text{ so}$$

$$\{\omega \in E_{2} \cap E_{2}': \tau^{v}(\omega) = \tau^{v'}(\omega)\}$$

$$\cong E_{2} \cap E_{2}' \cap \tau_{D_{1}}^{-1} [(D_{1} \cap D_{1}') \cap \tau_{D_{0}}^{-1} (D_{0} \cap D_{0}')]$$

$$\cap \tau_{D_{1}'}^{-1} [(D_{1} \cap D_{1}') \cap \tau_{D_{0}'}^{-1} (D_{0} \cap D_{0}')]$$

or

$$\{ \boldsymbol{\omega} \in \boldsymbol{E}_{2} \cap \boldsymbol{E}_{2}^{\prime} : \tau^{\boldsymbol{v}}(\boldsymbol{\omega}) \neq \tau^{\boldsymbol{v}^{\prime}}(\boldsymbol{\omega}) \}$$
$$\subseteq \boldsymbol{E}_{2} \cap \boldsymbol{E}_{2}^{\prime} \cap (\tau_{D_{1}}^{-1} [(\boldsymbol{D}_{1} \cup \boldsymbol{D}_{1}^{\prime}) \cup \tau_{D_{0}}^{-1} (\boldsymbol{D}_{0} \cup \boldsymbol{D}_{0}^{\prime})]$$
$$\cup \tau_{D_{1}^{\prime}}^{-1} [(\boldsymbol{D}_{1}^{\prime} \setminus \boldsymbol{D}_{1}) \cup \tau_{D_{0}}^{-1} (\boldsymbol{D}_{0}^{\prime} \setminus \boldsymbol{D}_{0})]).$$

Hence

$$\mu(\{\omega \in E_2 \cap E'_2 \colon \tau^{\mathbf{v}}(\omega) \neq \tau^{\mathbf{v}'}(\omega)\})$$

$$\leq \mu(D_1 \setminus D'_1) + \mu(D_0 \setminus D'_0) + \mu(D'_1 \setminus D_1) + \mu(D'_0 \setminus D_0)$$

$$= \mu(D_0 \sqcup D'_0) + \mu(D_1 \sqcup D'_1).$$

The cases  $n \ge 3$  follow in an analogous way.

Now

$$d(\tau^{\mathbf{v}}, \tau^{\mathbf{v}}) = \mu(\{\omega : \tau^{\mathbf{v}}(\omega) \neq \tau^{\mathbf{v}'}(\omega)\})$$

$$\leq \sum_{n=0}^{k} \mu(E_n \setminus E'_n) + \sum_{n=1}^{k} \mu(\{\omega \in E_n \cap E'_n : \tau^{\mathbf{v}}(\omega) \neq \tau^{\mathbf{v}'}(\omega)\})$$

$$\leq \sum_{n=0}^{k} \mu(E_n \sqcup E'_n) + \sum_{n=1}^{k} \sum_{l=0}^{n-1} \mu(D_l \sqcup D'_l)$$

$$\leq \sum_{n=0}^{k} \mu(E'_n \sqcup E'_n) + \sum_{n=1}^{k} \sum_{l=0}^{n-1} \sum_{m=l+1}^{k} \mu(E_m \sqcup E'_m)$$

which establishes the result.

**Theorem 6.6.** Let  $\tau = \tau(T)$ . If v is a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  with  $\int v d\mu = k \ge 1$ then for any  $\varepsilon > 0$  there is a stopping time v' with  $\int v' d\mu = k$  such that  $\tau^{v'}$  is ergodic on  $\{v' \ne 0\}$  and

 $d(\tau^{v},\tau^{v'}) < \varepsilon.$ 

*Proof.* Let  $E_n$ ,  $0 \le n \le k$ , be as defined for v by Theorem 4.1. As the column levels in  $S^v(T)$ ,  $v \ge 0$ , generate  $\mathscr{B}$  we can find  $v \ge 0$  such that

(i) There are disjoint sets  $E'_n$ ,  $0 \le n \le k$ , which are unions of column levels in  $S^v(T)$  with  $\mu(E_n \Delta E'_n) < \eta(\varepsilon, k)/2$ .

(ii)  $\mu(B(S^{\nu}(T))) < \eta(\varepsilon, k)/2k$ .

In  $C_1$  give each column level name *n* if it is a subset of  $E'_n$ ,  $0 \le n \le k$ , and then modify the top k column levels so that

$$\sigma_1(i) = i + 1 \pmod{k}, \quad i \in \{1, 2, \dots, k\}.$$

In  $C_2$  give each column level name *n* if it is a subset of  $E'_n$ ,  $0 \le n \le k$ , and then modify the top k column levels so that

 $\sigma_2(i) = i, \quad i \in \{1, 2, \dots, k\}.$ 

For  $3 \le i \le m$ , give each column level in  $C_i$  name *n* if it is a subset of  $E'_n$  and then assign the name *k* to each of the top *k* column levels. Now let  $F_n$  be the union of those column levels now having name *n*,  $0 \le n \le k$  and define *v'* so that

$$\tau^{v'}\omega = \begin{cases} \omega & \omega \in F_0 \\ \tau_{D_0}(\tau_{D_1}(\dots(\tau_{D_{n-1}}(\omega))\dots)), & \omega \in F_n, \ 1 \le n \le k. \end{cases}$$

where  $D_{n-1} = \bigcup_{l=n}^{k} F_l$ ,  $1 \le n \le k$ . Note that because of Theorem 5.4(i)

 $\tau^{v'}$  is ergodic on  $\{v' \neq 0\}$ .

Also

$$\mu(E'_n \Delta F_n) < k \, \mu(A(S^{\mathsf{v}}(T))) < \mu(\varepsilon, k)/2$$

and so

$$\mu(E_n \, \varDelta F_n) < \mu(E_n \, \varDelta E'_n) + \mu(E'_n \, \varDelta F_n) < \eta(\varepsilon, k).$$

Hence, by Lemma 6.5, we have

$$d(\tau^{v},\tau^{v'}) < \varepsilon.$$

**Theorem 6.7.** Let  $\tau = \tau(T)$ . If v is a stopping time for  $(\Omega, \mathscr{B}, \mu, \tau)$  with  $\int v d\mu = k \ge 1$ then for any  $\varepsilon > 0$  there is a stopping time v with  $\int v' d\mu = k$  such that  $\tau^v$  is Bernoulli on  $\{v' \neq 0\}$  and

$$d(\tau^{v},\tau^{v'}) < \varepsilon.$$

*Proof.* Let  $E_n$ ,  $0 \le n \le k$ , be as defined for v by Theorem 4.1. As the column levels in  $S^v(T)$ ,  $v \ge 0$ , generate  $\mathscr{B}$  we can find  $v \ge 0$  such that

(i) there is a partition  $F_n$ ,  $0 \le n \le k$ , of  $\Omega$ , with each  $F_n$  a union of column levels in  $S^{v}(T)$  and

 $\mu(E_n \varDelta F_n) < \eta(\varepsilon/2, k).$ 

(ii) If  $S^{\mathfrak{v}}(T) = (C_i: 1 \leq i \leq m)$  then  $\mu(C_i^*) < \eta(\varepsilon/2, k)/4, 1 \leq i \leq m$ .

(iii)  $\mu(B(S^{v}(T))) < \eta(\varepsilon/2, k)/2k$ ,

(iv)  $C_2$  has height greater than k.

Let  $\theta$  be the stopping time associated with the partition  $F_n$ ,  $0 \le n \le k$ , by Theorem 4.1. In  $C_1$  give each column level name *n* if it is a subset of  $F_n$  and then modify the top *k* column levels so that

$$\sigma_1(i) = i + 1 \pmod{k}, \quad i \in \{1, 2, \dots, k\}.$$

In  $C_2$  give each column level except the top k + 1 levels name 0, give the (k + 1)st column level (counting from the top) name 1 and allocate names to the top k levels so as to give

$$\sigma_{2}(i) = i, \quad i \in \{1, 2, \dots, k\}.$$

In  $C_3$  give each column level except the top k levels name 0 and give the top k levels names which give

 $\sigma_3(i) = i, \quad i \in \{1, 2, \dots, k\}.$ 

For  $4 \le j \le m$  give the top k column levels in  $C_j$  name k and any other column level name n if it is a subset of  $F_n$ . Now let  $E'_n$  be the union of column levels now having name n,  $0 \le n \le k$ , and define v' so that

$$\tau^{v'}\omega = \begin{cases} \omega & \omega \in E'_{0} \\ \tau_{D_{0}}(\tau_{D_{1}}(\dots(\tau_{D_{n-1}}(\omega))\dots)), & \omega \in E'_{n}, \ 1 \leq n \leq k, \end{cases}$$

where  $D_{n-1} = \bigcup_{l=n}^{k} E'_{l}$ ,  $1 \le n \le k$ . Note that the second and third columns of  $T_1$ , as defined in §5 have heights 2 and 1 respectively and so, by Theorem 5.4(ii)  $\tau^{v}$  is Bernoulli on  $\{v' \ne 0\}$ . Moreover, for  $0 \le n \le k$ .

$$\mu(F_n \, \, \exists \, E'_n) < \mu(C_2^*) + \mu(C_3^*) + k \, \mu(A(S^{\rm v}(T))) < \eta(\varepsilon/2, k)$$

which implies that

$$d(\tau^{\theta},\tau^{v'}) < \varepsilon/2.$$

Hence

 $d(\tau^{v},\tau^{v'}) < \varepsilon.$ 

**Corollary 6.8.** If  $\tau$  is a Bernoulli shift of finite entropy then the set of stopping time transformations  $\tau^v$  for which  $\tau^v$  is Bernoulli is dense in the set of stopping time transformations.

*Proof.* Any Bernoulli shift of finite entropy is isomorphic to  $\tau = \tau(T)$  for some T [2].

# References

- 1. Friedman. N.A.: Introduction to Ergodic Theory. New York: van Nostrand 1970
- 2. Friedman, N.A.: Bernoulli shifts induce Bernoulli shifts. Advances in Math. 10, 39-48 (1970)
- 3. Friedman. N.A., Orstein. D.S.: On isomorphism of weak Bernoulli transformations. Advances in Math. 5. 365-394 (1970)
- 4. Neveu, J.: Temps d'arrêt d'un système dynamique. Z. Wahrscheinlichkeitstheorie verw. Gebiete 13. 81-94 (1969)
- 5. Ornstein, D.S.: Bernoulli shifts with infinite entropy are isomorphic. Advances in Math. 5, 339–348 (1970)
- 6. Saleski, A.: Stopping times for Bernoulli automorphisms. Pacific J. Math. 52, 547-551 (1974)
- 7. Shields. P.: Cutting and Independent Stacking of Intervals. Math. Systems Theory 7, 1-4 (1973)

Received Janurary 15, 1977; in revised form July 18, 1978

# STOPPING TIMES FOR MEASURE-PRESERVING TRANSFORMATIONS : SOME APPROXIMATION RESULTS

# STOPPING TIMES FOR MEASURE-PRESERVING TRANSFORMATIONS: SOME APPROXIMATION RESULTS

# By S. T. RICHARDSON<sup>†</sup>

[Received 24 July 1979]

# 1. Introduction

In this paper we prove some approximation results within the class  $\mathscr{S}(\tau)$  of stopping-time transformations derived from an ergodic automorphism  $\tau$  of a Lebesgue space.

If  $\nu$  is a non-negative integer-valued measurable function then  $\tau^{\nu}$  belongs to  $\mathscr{S}(\tau)$  if and only if  $\tau^{\nu}$  is again an automorphism, and  $\nu$  is then called a stopping time. Stopping-time transformations were first discussed by Neveu in [6].  $\mathscr{S}(\tau)$  includes induced transformations and coincides with the positive part of the full group of  $\tau$ .

It is known, see [1] and [2], that  $\mathscr{S}(\tau)$  is separable and complete with respect to the uniform metric and that the entropy of  $\tau^{\nu}$  is related to that of  $\tau$  via the expected value of  $\nu$ . In §5 we approximate any stopping-time transformation  $\tau^{\nu}$  by a stopping-time transformation  $\tau^{\nu_1}$ , which has the same entropy and is ergodic on the set  $\{\nu_1 \neq 0\}$ . We then deduce density results for the set of mixing or Kolmogorov stopping-time transformations using results of [4] and [7]. These approximations are based upon a construction which is given in §4. In a previous paper [8] we proved that when  $\tau$  is a Bernoulli shift the set of  $\tau^{\nu}$  which are also Bernoulli shifts is dense. Throughout the paper we take  $\Omega = [0, 1]$ ,  $\mathscr{B}$  to be the set of Borel subsets of  $\Omega$ .  $\mu$  to be Lebesgue measure on  $\mathscr{B}$ , and  $\tau$  to be an ergodic automorphism of  $(\Omega, \mathscr{B}, \mu)$ .

# 2. Preliminaries

Let  $\nu$  be a measurable mapping form  $(\Omega, \mathscr{B})$  to the natural numbers  $\mathbb{N}$ . Then the transformation  $\tau^{\nu}$  defined by

$$\tau^{\mathbf{v}}(\omega) = \tau^{\mathbf{n}}(\omega), \quad \omega \in \{\nu = n\}, \quad n \in \mathbb{N},$$

<sup>+</sup> The author wishes to thank her supervisor Dr K. M. Wilkinson for much help and encouragement in the preparation of this paper which is part of her Ph.D. thesis at the University of Nottingham.

Proc. London Math. Soc. (3) 43 (1981) 273-294

5388.3.43

is an automorphism of  $(\Omega, \mathcal{B}, \mu)$  if and only if the sets  $\{\tau^n \{\nu = n\}: n \in \mathbb{N}\}\$  form a partition of  $\Omega$ . If this is the case, following Neveu, we call  $\nu$  a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  and  $\tau^{\nu}$  a stopping-time transformation. Neveu [6] established the following decomposition theorem:

THEOREM 1. (i) If  $\Omega \supset D_0 \supset D_1 \dots$  is a sequence of measurable sets such that either  $D_{k-1} \neq \emptyset$  and  $D_k = \emptyset$  for some  $k \ge 1$  or  $\bigcap_{n=0}^{\infty} D_n = \emptyset$  then there is a stopping time  $\nu$  of  $(\Omega, \mathcal{B}, \mu, \tau)$  such that

$$\tau^{\nu}\omega = \begin{cases} \omega, & \text{if } \omega \in E_0 = D_0^c, \\ \tau_{D_0}(\tau_{D_1}(\dots(\tau_{D_{n-1}}(\omega))\dots)), & \text{if } \omega \in E_n \text{ for } n \ge 1, \end{cases}$$
(2.1)

where  $E_n = D_{n-1} \setminus D_n$  for  $n \ge 1$ , and for any  $D \in \mathcal{B}$ ,  $\tau_D$  represents the induced transformation of  $\tau$  on D.

(ii) Corresponding to any stopping time  $\nu$  of  $(\Omega, \mathcal{B}, \mu, \tau)$  there is a sequence of measurable sets  $\Omega \supset D_0 \supset D_1 \dots$  as in (i) such that  $\tau^{\nu}$  has the form (2.1).

(iii) If v is a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  and  $\tau$  is ergodic then  $\int v d\mu$  is either a non-negative integer or  $+\infty$ . If  $\int v d\mu = k < \infty$  then t've sets  $D_n$ with  $n \ge 0$ , defined in (ii) satisfy  $D_{k-1} \ne \emptyset$ ,  $D_k = \emptyset$ , whereas if  $\int v d\mu = +\infty$  we have  $D_n \ne \emptyset$ , for  $n \ge 0$ , and  $\bigcap_{n=0}^{\infty} D_n = \emptyset$ .

If we let  $E_0 = D_0^c$  and  $E_i = D_{i-1} \setminus D_i$ , for  $i \ge 1$ , then, for a stopping time  $\nu$  with  $\int \nu d\mu = k < \infty$ , the sets  $\{E_i\}_{i=0}^k$  form a partition of  $\Omega$  that we call the partition associated with  $\nu$ . Conversely, given a partition  $\{E_i\}_{i=0}^k$  of  $\Omega$ , we say that the stopping time  $\nu$  defined as in (2.1) with  $D_0^c = E_0$  and  $D_i = \bigcup_{j=i+1}^k E_j$ , for  $1 \le i \le k-1$ , is the stopping time corresponding to the partition  $\{E_i\}_{i=0}^k$ .

Note that any stopping-time transformation  $\tau^{\nu}$  with  $\int \nu d\mu = +\infty$  can be approximated by a stopping-time transformation  $\tau^{\nu_k}$  with  $\int \nu_k d\mu = k < \infty$ . Indeed if  $\{D_i\}_{i=1}^{\infty}$  is the sequence of sets. decreasing to  $\emptyset$ , associated with  $\nu$  by Theorem 1 (ii), we define

$$\tau^{\nu_k}(\omega) = \begin{cases} \omega, & \text{if } \omega \in D_k, \\ \tau^{\nu}(\omega). & \text{if } \omega \in D_k^{c}, \end{cases}$$

and we choose k large enough so that  $\tau^{\nu}$  and  $\tau^{\nu_k}$  differ on a set of arbitrary small measure.

The entropy of a stopping-time transformation is related to the number of sets in the decomposition as follows: if  $\int \nu d\mu = k < \infty$ , then

$$h(\tau^{\nu}) = kh(\tau). \tag{2.2}$$

The entropy of  $\tau^{\nu}$  was first studied by Belinskaya in [2], where (2.2) is proved under the additional hypothesis of ergodicity of  $\tau^{\nu}$ . By

decomposition of  $\tau^{\nu}$  into its ergodic component and use of a result of Rohlin's [9] relating the entropy of an automorphism to the entropy of its ergodic components, Belinskaya's result easily leads to (2.2).

The construction in §4 involves the geometric representation of an ergodic transformation via Rohlin's theorem. Let us recall some terminology, introduced principally in [3, 4, and 8].

Let  $\{L_i: 0 \le i < h\}$  be a collection of disjoint measurable sets. We say that  $C = (L_i: 0 \le i < h)$  is a  $\tau$ -column if  $\tau(L_i) = L_{i+1}$ , for  $0 \le i < h-1$ . The sets  $L_i$ , with  $0 \le i < h$ , are the column levels of C; C has base  $B(C) = L_0$ , top  $A(C) = L_{h-1}$ , height h(C) = h, and width  $w(C) = \mu(L_0)$ . We let  $C^* = \bigcup_{i=0}^{h-1} L_i$ . Note that, even though  $\tau$  is an automorphism in  $\Omega$ , in a  $\tau$ -column C,  $\tau$  is not defined within C on A(C). If  $F \subset L_0$ , the column ( $\tau^i F: 0 \le i < h$ ) is called a subcolumn (or copy) of C with base F.

A  $\tau$ -tower T is a finite set of disjoint columns  $T = (C_j: 1 \le j \le m); T$ has base  $B(T) = \bigcup_{j=1}^{m} B(C_j)$  and top  $A(T) = \bigcup_{j=1}^{m} A(C_j)$ . We let  $T^* = \bigcup_{j=1}^{m} C_j^*$  and denote by  $\mathscr{L}(T)$  the set of column levels of T. In the sequel, when we simply say column (respectively tower) we mean  $\tau$ -column (respectively  $\tau$ -tower).

We say that a column C is *pure* with respect to a partition  $P = \{p_i\}_{i=1}^n$  if each level of C is a subset of an atom of P. Given any column C. we *purify* C with respect to P if we divide C into subcolumns that are pure with respect to P. This is done by taking the bases of these subcolumns equal to

$$B(C) \cap \bigcap_{j=0}^{h(C)-1} \tau^{-j} p_{ij},$$

for all choices of  $i_j \in \{1, ..., n\}$ ,  $0 \leq j < h(C)$ . We purify a tower with respect to P if we purify each column in the tower.

Let us now consider a fixed stopping time  $\nu$ , with associated partition  $\{E_i\}_{i=0}^k$ , and a  $\tau$ -column C, pure with respect to  $\{E_i\}_{i=0}^k$ . We can then label each column level of C with respect to  $\{E_i\}_{i=0}^k$ , that is, we define a function n on  $\mathscr{L}(C)$ , taking values in  $\{0, \ldots, k\}$  by

$$n(L) = n \quad \text{if } L \subset E_n, \text{ for } 0 \leq n \leq k.$$

and say that L has name n. Levels with name 0 are invariant under  $\tau^{\nu}$ and will not be involved in the discussion that follows. The action of  $\tau^{\nu}$ within the column C can now be described by an appropriate sequence of names. Suppose that a column level  $L_i$  has name  $n \ge 1$ . Then  $\tau^{\nu}(L_i)$  is another column level in C (we also say  $\tau^{\nu}(L_i)$  is defined within C) if and only if there is a sequence of indices

$$i < i_1 < i_2 < \ldots < i_n < h$$

## 8. T. RICHARDSON

with  $n(L_{i_t}) \ge n-t+1$  and  $1 \le t \le n$ . Further,  $\tau^{\nu}(L_i) = L_{i_n}$  if each index  $i_t$  satisfies  $i_t = \min\{j: i_{t-1}+1 \le j < h, n(L_j) > n-t+1\}$  for  $1 \le t \le n$ .

Clearly there may be column levels other than A(C), that is, some  $L_i$ , with  $0 \leq i \leq h-1$ , for which  $\tau^{\nu}(L_i)$  is not defined within C. For a column level  $L_i$  with  $\varkappa(L) = n$ , we shall define  $\varkappa(L_i) = u$ , where  $1 \leq u \leq n$ , if the set

$$\{i_1, \dots, i_s: i < i_1 < \dots < i_s < h, \ n(L_{i_t}) \ge n-t+1, 1 \le t \le s\}$$

is non-empty for  $s \leq n-u$  and empty for  $s \geq n-u+1$ , and  $\omega(L_i) = 0$  if  $\tau^{\mathsf{v}}(L_i)$  is defined within the column. So  $\omega(L_i)$  counts the number of steps in the iterative definition of  $\tau^{\mathsf{v}}$  on  $L_i$  given by Theorem 1 which are left undefined within the column containing  $L_i$ .

In the sequel, we refer to [8] for properties of the functions n and u that we quote without proof.

Any column level L in C with  $\omega(L) \ge 1$  is called a *v*-top in C and it can be shown that there are at most k *v*-tops in any column. Furthermore, any level L in C with  $\omega(L) = 0$  is mapped by some power of  $\tau^{\nu}$  onto a unique *v*-top.

If we now have two disjoint  $\tau$ -columns  $C_0$  and  $C_1$  with  $\tau(\mathcal{A}(C_0)) = B(C_1)$ , we denote by  $C_0 * C_1$  the column defined by

$$(\tau^{i}B(C_{0}): 0 \leq i < h(C_{0}) + h(C_{1})).$$

Similarly, given two disjoint towers  $T_1 = (C_{1j}: 1 \le j \le m_1)$  and  $T_2 = (C_{2l}: 1 \le l \le m_2)$  with  $\tau(A(T_1)) = B(T_2)$ , we then define a tower with base  $B(T_1)$  and top  $A(T_2)$  denoted by  $T_1 * T_2$ . The bases of the columns of  $T_1 * T_2$  are given by

$$\{B(C_{1i}) \cap \tau^{-h(C_{1i})}B(C_{2l}): 1 \leq j \leq m_1, 1 \leq l \leq m_2\},\$$

so  $T_1 * T_2$  has at most  $m_1 m_2$  columns. If  $C_0$  and  $C_1$  are both pure with respect to  $\{E_i\}_{i=0}^k$  and if L is a  $\nu$ -top for  $C_0$ , then  $\tau^{\nu}(L)$  might be defined within  $C_0 * C_1$ . Precisely, we define at most k levels in  $C_1$ :  $\{B_i(C_1): 1 \leq i \leq k\}$ , called the  $\nu$ -bases of  $C_1$ , by

$$\begin{cases} B_1(C_1) = L_1(C_1), \\ B_i(C_1) = \tau_{D_0}(\tau_{D_1}(\dots(\tau_{D_{i-2}}(L_i(C_1)))\dots)), & \text{for } 2 \leq i \leq k, \end{cases}$$

if  $\tau_{D_0}(\tau_{D_1}(...(\tau_{D_{i-2}}(L_i(C_1)))...))$  is a column level within  $C_1$  and where  $L_i(C_1)$  is the first column level of  $C_1$  (counting from the base) which satisfies  $n(L_i(C_1)) \ge i$ . It follows that when L is a  $\nu$ -top for  $C_0$  with  $\alpha(L) = u \ge 1$  and when  $B_u(C_1)$  exists. we have

$$\tau^{\nu}(L) = B_{u}(C_{1}). \tag{2.3}$$

Note that, for a column  $C = (L_i: 0 \le i < h)$ , pure with respect to  $\{E_i\}_{i=0}^k$ . if there exists a sequence of indices  $\{i_j\}_{j=0}^{k-1}$  with  $0 \le i_0 < i_1 < \ldots < i_{k-1} \le h-1$  and  $\varkappa(L_{ij}) \ge k-j$ , for  $0 \le j \le k-1$ . then C possesses exactly k v-bases. However, any level in C is the image under some power of  $\tau^{\mathsf{v}}$  of a v-basis.

Thus the name function characterizes the action of  $\tau^{\nu}$  within a column  $C_0$ ; further, if the column  $C_0$  is extended by considering  $C_0 * C_1$ , then  $\tau^{\nu}$  maps  $\nu$ -tops of  $C_0$  to  $\nu$ -bases of  $C_1$  and this action is characterized by the function  $\omega$ .

# 3. The *u*-blocks

The definitions concerning  $n, \alpha, \nu$ -tops, and  $\nu$ -bases given in the previous paragraph depend on  $\nu$  only through the partition  $\{E_i\}_{i=0}^k$  associated with  $\nu$ . We could as well have stated these definitions with respect to a partition and taken  $\nu$  to be the corresponding stopping time. Since  $\nu$  and  $\{E_i\}_{i=0}^k$  are interchangeable in this way, we simply refer to n.  $\alpha$  functions as defined with respect to a name partition  $\{E_i\}_{i=0}^k$ , and similarly for  $\nu$ -tops and  $\nu$ -bases.

In what follows it will be convenient to impose conditions on the name structure such that the  $k \nu$ -tops in a column C. pure with respect to a name partition  $\{E_i\}_{i=0}^k$ , are the top k levels. We denote these top k levels by  $D(C) = (A_j(C): 1 \le j \le k)$ , numbering from the top so that  $A_1(C) = A(C)$ ; if  $T = (C_j: 1 \le j \le k)$  then we let  $D(T) = (D(C_j): 1 \le j \le k)$ .

A necessary and sufficient condition for each  $A_i(C)$  to be a  $\nu$ -top with respect to a name partition  $\{E_i\}_{i=0}^k$  is

$$\varkappa(A_i(C)) \in \{i, i+1, \dots, k\}, \quad \text{for } 1 \le i \le k.$$

$$(3.1)$$

In this case  $\{ \omega(A_i(C)): 1 \leq i \leq k \}$  is a permutation of  $\{1, 2, ..., k\}$ . Further, we note that if two name partitions  $\{E_i\}_{i=0}^k$  and  $\{E'_i\}_{i=0}^k$ , with associated stopping times  $\nu$  and  $\nu'$  respectively, both satisfy (3.1) and coincide on  $C \setminus D(C)$ , then the sets of  $\nu$  and  $\nu'$ -tops coincide. If L is a level in C with  $\omega(C) = 0$  and  $\tau^{\nu}(L) = A_j(C)$  for some j, where  $1 \leq j \leq k$ , then  $\omega'(L) = 0$  and  $\tau^{\nu'}(L) = A_j(C)$ .

We now define a specific structure of names, called a  $\omega$ -block. which we shall use later to control the action of  $\tau^{\nu}$ .

A tower  $T = (C_j: 1 \le j \le k)$  is a  $\alpha$ -block with respect to a name partition  $\{E_i\}_{i=0}^k$  if

(i)  $h(C_j) = k$ , for  $1 \le j \le k$  and  $w(C_j) = w(C_l)$ , for  $1 \le j, l \le k$ .

(ii) the columns of T are pure with respect to  $\{E_i\}_{i=0}^k$ ,

(iii)  $\alpha(A_i(C_i)) = i + j - 1 \pmod{k}$ , for  $1 \le i, j \le k$ .

### S. T. RICHARDSON

A  $\alpha$ -block with k = 3 is illustrated below in Fig. 1. Above each level we have listed its name and the value of the  $\alpha$ -function.

Now suppose that there is no underlying name partition and that C is a  $\tau$ -column of height k. Then we can define a name partition so as to convert C into a  $\omega$ -block. This is done by first dividing  $A_1(C)$  into k disjoint subsets of equal measure, which we denote by  $\{A_1(C)(j): 1 \leq j \leq k\}$ , and then defining

$$C_{j} = (\tau^{i} A_{1}(C)(j)) : -(k-1) \leq i \leq 0).$$

Thus  $T = (C_j: 1 \le j \le 1)$  satisfies condition (i) of the definition of a  $\alpha$ -block. Secondly we assign names to column levels in each column  $C_j$  so that conditions (ii) and (iii) are satisfied: for  $1 \le j \le k$ , we let

$$\varkappa(A_i(C_j)) = \begin{cases} (j-1)+i, & \text{if } 1 \leq i \leq k-j+1, \\ i, & \text{if } k-j+1 < i \leq k. \end{cases}$$

In general, we start with an underlying name partition  $\{E_i\}_{i=0}^k$  and a column C, with  $h(C) \ge k$ , pure with respect to  $\{E_i\}_{i=0}^k$ . We then, as above, subdivide C into k subcolumns of equal width,  $\{C_j: 1 \le j \le k\}$  and modify the names on the top k levels of each  $C_j$ , and only on these, to yield a  $\alpha$ -block. We call the resulting tower U(C). Thus the new name partition  $\{E_i\}_{i=0}^k$  defined is such that

(i) the tower  $\{D(C_j): 1 \leq j \leq k\}$  is a  $\alpha$ -block with respect to  $\{E'_i\}_{i=0}^k$ , (ii)  $E_i \cap (\Omega \setminus D(C)^*) = E'_i \cap (\Omega \setminus D(C)^*)$ , for  $1 \leq i \leq k$ .

We refer to  $\{E'_i\}_{i=0}^k$  as the name partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of U(C).

If  $T = (C_l: 1 \le l \le m)$  is a tower with columns of equal height h, where  $h \ge k$ , which are pure with respect to a name partition  $\{E_i\}_{i=0}^k$ , then we let U(T) denote the tower resulting after modification of each  $C_l$ to  $U(C_l)$ , with  $1 \le l \le m$ . Thus U(T) has km columns and we refer to the new name partition as the partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of U(T).

To outline the fundamental property of  $\omega$ -blocks, we consider the following situation. We have a column  $C = C_0 * C_1$ , where  $h(C_0) \ge k$ .  $h(C_1) \ge k$ , w(C) > 0, and both  $C_0$  and  $C_1$  are pure with respect to a

name partition  $\{E_i\}_{i=0}^k$ . We suppose that the k  $\nu$ -bases of  $C_1$ .  $\{B_l(C_1): 1 \leq l \leq k\}$ , exist. We form  $U(C_0)$  and let  $\{E'_i\}_{i=0}^k$  be the partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of  $U(C_0)$ . Note that  $C_1$  is still pure with respect to  $\{E'_i\}_{i=0}^k$ . We adopt the following notation: if L is a column level in C of height h, then define

$$L(j) = \tau^{h-h(C_0)} A_1(C_0)(j), \quad \text{for } 1 \le j \le k.$$

Thus  $L = \bigcup_{j=1}^{k} L(j)$ . Let  $\nu'$  be the stopping time associated with  $\{E'_i\}_{i=0}^k$ . Note that  $\{B_l(C_1): 1 \leq l \leq k\}$  are also the  $\nu'$ -bases of  $C_1$ . We now show that the existence of a  $\alpha$ -block at the top of  $C_0$  controls 'uniformly' the intersection between images under  $\tau^{\nu'}$  of levels in  $C_0$  and levels in  $C_1$ .

PROPOSITION 1. Let  $C = C_0 * C_1$ ,  $\{E'_i\}_{i=0}^k$ , and  $\nu'$  be as above. Let L be any level in  $C_0$ , and let L' be any level in  $C_1$ , both included in  $E'_0^\circ$ . Then there exists a unique t > 0 such that

$$(\tau^{\mathbf{v}'_{\cdot}})^{t}L \cap L' \neq \emptyset,$$

and, for this value of t,

$$\mu((\tau^{\mathbf{v}'_{i}})^{t}L \cap L') = w(C)/k.$$

*Proof.* We first show that for each *i*, with  $1 \le i \le k$ . *k* subsets of  $A_i(C_0)$  are mapped by  $\tau^{\mathbf{v}'}$  to different  $\mathbf{v}'$ -bases. Precisely, by definition of a  $\alpha$ -block.  $\alpha(A_i(C_0)(j)) = i+j-1 \pmod{k}$  and thus, by (2.3).

$$\tau^{\mathbf{v}'}(A_i(C_0)(j)) = B_{i+j-1}(C_1)(j),$$

where from now on, all indices are taken modulo k. For  $1 \le i, l \le k$  fixed, there exists a unique j = j(i, l) such that

$$i+j-1 = l \pmod{k},$$

and, for  $j' \neq j$ ,  $\tau^{\mathbf{v}'}(A_i(C_0)(j')) \cap B_l(C_1) = \emptyset$ . Hence, for any  $1 \leq i, l \leq k$ .  $\mu(\tau^{\mathbf{v}'}(A_i(C_0)) \cap B_l(C_1)) = \mu(B_l(C_1)(j)) = w(C)/k.$ (3.2)

Now for a fixed j, with  $1 \leq j \leq k$ , L(j) belongs to a subcolumn of  $C_0$ , pure with respect to  $\{E'_i\}_{i=0}^k$ , with  $\{A_i(C_0)(j): 1 \leq i \leq k\}$  as its  $\nu'$ -tops. Hence there exists a unique  $t_0 \in \mathbb{N}$ , such that

$$(\tau^{\nu'})^{t_0} L(j) = A_i(C_0)(j), \tag{3.3}$$

for some *i*. with  $1 \le i \le k$ . Moreover, as remarked at the beginning of the paragraph, *i* is independent of *j* because, by definition of  $U(C_0)$ , for any *j*, with  $1 \le j \le k$ , the names of  $\{A_i(C_0)(j): 1 \le i \le k\}$  satisfy (3.1).

S. T. RICHARDSON

Hence (3.3) holds for any j, with  $1 \le j \le k$ , and

$$(\tau^{\nu'})^{t_0}L = \bigcup_{j=1}^k (\tau^{\nu'})^{t_0}L(j) = \bigcup_{j=1}^k A_i(C_0)(j) = A_i(C_0).$$
(3.4)

Finally, as  $C_1$  is pure with respect to  $\{E'_i\}_{i=0}^k$ , there exists a unique l, with  $1 \leq l \leq k$ , and a unique  $t_1 \in \mathbb{N}$  such that

$$(\tau^{\nu'_{l}})^{t_{1}}B_{l}(C_{1}) = L'. \tag{3.5}$$

Letting  $t = t_0 + t_1 + 1$ , we have, by (3.3), (3.4), and (3.5),

$$(\tau^{\nu'})^{t}L \cap L' = (\tau^{\nu'})^{t_{1}}[(\tau^{\nu'})A_{i}(C_{0}) \cap B_{l}(C_{1})],$$

for unique  $1 \leq i, l \leq k$ , and the result follows from (3.2).

Note that Proposition 1 still holds if, instead of supposing that  $C_1$  is pure with respect to  $\{E_i\}_{i=0}^k$ , we only suppose that the first  $h(C_1) - k$ levels of  $C_1$  are pure with respect to  $\{E_i\}_{i=0}^k$  whilst  $D(C_1)$  is partitioned by  $\{E_i\}_{i=0}^k$  into any number of  $\alpha$ -blocks.

If we now have *n* disjoint  $\tau$ -columns, satisfying  $\tau(A(C_j)) = B(C_{j+1})$ , where  $0 \leq j < n$ , we denote by  $C = C_0 * C_1 * \ldots * C_n$  the column defined by  $(\tau^i B(C_0): 0 \leq i < \sum_{j=0}^n h(C_j))$ . If  $h(C_j) \geq k$  and each  $C_j$ , with  $0 \leq j \leq n$  is pure with respect to a name partition  $\{E_i\}_{i=0}^k$ , we construct several  $\alpha$ -blocks by forming successively

$$T_{0} = U(C_{0}).$$

$$T_{1} = U(T_{0} * C_{1}),$$

$$\vdots$$

$$T_{l} = U(T_{l-1} * C_{l}), \text{ for } 1 \leq l \leq n-1,$$

$$\vdots$$

$$T_{n} = T_{n-1} * C_{n}.$$

By abuse of notation we sometimes denote  $T_n$  by  $U(C_0 * C_1 * ... * C_n)$ and, when there is no possible confusion about the determination of the  $C_j$   $(0 \le j \le n)$ , we simply say that we have formed independent  $\omega$ blocks within C. In this process, at each stage, we define a new name partition  $\{E_i^{(l)}\}_{i=0}^k$ , for  $1 \le l \le n$ , with  $\{E_i^{(1)}\}_{i=0}^k$  the name partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of  $U(C_0)$ , and  $\{E_i^{(l+1)}\}_{i=0}^k$  the name partition corresponding to  $\{E_i^{(l)}\}_{i=0}^k$  after formation of  $U(T_{l-1} * C_l)$ , for  $1 \le l \le n-1$ . Note that the names of levels in  $D(C_n)$ have not been modified. We refer to  $\{E_i^{(n)}\}_{i=0}^k$  as the partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of  $U(C_0 * ... * C_n)$ . We can

describe this construction as follows. For  $0 \le l \le n-1$ ,  $T^l$  consists of  $k^{l+1}$  columns, pure with respect to  $\{E_i^{(l)}\}_{i=0}^k$ , and  $D(T_l)$  is a set of  $k^l$   $\alpha$ -blocks.  $T_n$  consists of  $k^n$  columns, pure with respect to  $\{E_i^{(n)}\}_{i=0}^k$ , and  $\{E_i^{(n)}\}_{i=0}^k$  coincides with  $\{E_i\}_{i=0}^k$  on  $\Omega \setminus (\bigcup_{l=0}^{n-1} D(C_l)^*)$ . For an illustration of the formation of independent  $\alpha$ -blocks in the case where k = 3 and n = 3, see Fig. 2.

We now state a generalization of Proposition 1 which can be proved by induction.

PROPOSITION 2. Let  $C = C_0 * C_1 * \ldots * C_n$  be a column, pure with respect to a name partition  $\{E_i\}_{i=0}^k$ , with w(C) > 0 and  $h(C_j) \ge k$ , for  $0 \le j \le h$ . Suppose also that  $C_n$  possesses k v-bases with respect to  $\{E_i\}_{i=0}^k$ . Let  $\{E'_i\}_{i=0}^k$ be the partition corresponding to  $\{E_i\}_{i=0}^k$  after formation of  $U(C_0 * \ldots * C_n)$ , and let v' be the stopping time associated with  $\{E'_i\}_{i=0}^k$ . Then, for any level L in  $C_{i_0}$  and level L' in  $C_{i_1}$ , where  $0 \le i_0 < i_1 \le n$  and both levels are included in  $E'_0^c$ , there exist  $t_j$ , with  $1 \le j \le k^{i_1-i_0-1}$ , such that

$$\mu([\bigcup_{j} (\tau^{v'})^{t_j} L] \cap L') = \frac{w(C)}{k},$$

the union being taken over j ranging from 1 to  $k^{i_1-i_0-1}$ .

In the sequel, we adopt the notation

$$C = [C_0, C_1, \dots, C_n]$$

if  $C = c_0 * c_1 * \ldots * c_n$ , where  $c_i$  is a copy of  $C_i$  for each *i*. Forming independent  $\omega$ -blocks within *C* is then taken to mean forming  $U(c_0 * \ldots * c_n)$ . We conclude this section by quoting two results that we shall use in the next section. Suppose that  $P = \{p_i\}_{i \in I}$  is a partition of  $\Omega$ . with *I* a finite set of indices, and that  $\sigma$  is an ergodic automorphism. The following theorem is referred to as the strong Rohlin theorem.

THEOREM 2. Given n. a positive integer, and  $\varepsilon > 0$ , there exists a  $\sigma$ -column C of height n and base B such that

- (i)  $\mu(C^*) > 1 \varepsilon$ .
- (ii) d(P/B) = d(P), where

$$d(P) = \{\mu(p_i)\}_{i \in I} \quad and \quad d(P/B) = \left\{\frac{\mu(p_i \cap B)}{\mu(B)}\right\}_i$$

For a set  $A \in \bigvee_{i=0}^{m-1} \sigma^{-i} P$ , with  $m \ge 1$  and

$$A = p_{i_0} \cap \sigma^{-1} p_{i_1} \cap \ldots \cap \sigma^{-(m-1)} p_{i_{m-1}},$$

εI








FIG. 2. Forming independent  $\alpha$ -blocks in  $C = C_0 * C_1 * C_2 * C_3$ .

we call the sequence  $(i_0, \ldots, i_{m-1})$  the *P*-*m*-name of *A* with respect to  $\sigma$ . As a consequence of the ergodic theorem, the following holds: for any  $\varepsilon > 0$ , there exists N > 0 such that for  $m \ge N$ , there is a collection  $\mathscr{E}_m$  of sets in  $\bigvee_{i=0}^{m-1} \sigma^{-i}P$  of total measure at least  $1 - \varepsilon$  such that for all  $p_i \in P$  and  $A \in \mathscr{E}_m$ ,

$$\left|f_{\mathcal{A}}(i,m)-\mu(p_i)\right|<\varepsilon,$$

where  $f_A(i, m)$  is the relative frequency of the index *i* in the *P*-*m*-name of *A* with respect to  $\sigma$ .

#### 4. Construction

Throughout this section  $\nu$  is a given stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  with  $\int \nu d\mu = k \ge 1$  and associated partition  $\{E_i\}_{i=0}^k$ .

We construct a stopping time  $\nu'$ , with  $\int \nu' d\mu = k$ , for which the automorphism  $\tau^{\nu'}$ , whilst remaining close to  $\tau^{\nu}$ , is ergodic on the set  $\{\nu' \neq 0\}$ . This involves constructing a sequence of  $\tau$ -towers  $\{T_r\}_{r\geq 1}$  and of stopping times  $\{\nu^{(r)}\}_{r\geq 1}$ . As in [4] where Friedman and Ornstein induced a mixing transformation, the recurrence step is provided by a Rohlin's theorem, though here we need to use the strong form of Rohlin's theorem (Theorem 2). The stopping time  $\nu'$  is defined as the limit of the sequence  $\{\nu^{(r)}\}_{r\geq 1}$ , where at the *r*th stage we modify the existing name partition by introducing independent  $\alpha$ -blocks and take  $\nu^{(r)}$  to be the stopping time associated with this modified partition.

We shall use the following notation.  $\{S_r\}_{r\geq 0}$  is a decreasing sequence of sets formed inductively in the construction, and  $\tau_r = \tau_{S_r}$ , the induced transformation on  $S_r$ . For each  $r \in \mathbb{N}$ , the dyadic intervals of order rinduce on  $S_r$  a partition denoted by  $Q_r$ , that is

$$Q_r = \{S_r \cap I_{r,j} : 1 \leq j \leq 2^r\},\$$

where  $I_{r,j} = [(j-1)2^{-r}, j2^{-r})$ , for  $1 \le j \le 2^r$  and  $r \ge 1$ . We shall purify the towers  $T_r$ , where  $S_r = T_r^*$ , with respect to  $Q_r$  in order that the set of levels of the successive towers generates the Borel  $\sigma$ -algebra.

The construction involves a decreasing sequence of positive real numbers  $\{\eta_r\}_{r\geq 1}$ , where  $\eta_1 < 1$  is given and  $\{\eta_r\}_{r\geq 2}$  are determined inductively. Later, in §5, further conditions will be imposed on the sequence  $\{\eta_r\}_{r\geq 1}$  to prove approximation results.

The first step of the induction process constructs a  $\tau$ -tower  $T_1$ , pure with respect to  $\{E_i\}_{i=0}^k$  and such that the top k levels in each column have, after modification if necessary, name k.

Let  $S_0 = \Omega$ : Choose a positive integer t(1) satisfying

$$t(1) > 2k/\eta, \tag{4.1}$$

and apply Rohlin's theorem to the ergodic transformation  $\tau_0 = \tau$  to find  $F_1 \subset S_0$  which has t(1) disjoint images under  $\tau_0$  and is such that

$$\mu \left( \bigcup_{i=0}^{i(1)-1} \tau_0^{i} F_1 \right) > \mu(S_0) - \eta_1.$$
(4.2)

Next we purify the  $\tau_0$ -column ( $\tau_0^i F_i: 0 \leq i < t(1)$ ) with respect to the partition  $\{E_i\}_{i=0}^k$ , and so obtain a tower  $T_1$  where each column level  $L \in \mathscr{L}(T_1)$  is given a unique name by

$$n(L) = i$$
 if and only if  $L \subset E_i$ , for  $0 \leq i \leq k$ .

Suppose that  $T_1$  has  $m_1$  columns, say,  $T_1 = (C_{1j}; 1 \le j \le m_1)$ , with bases  $B_{1j} = B(C_{1j})$ . Let  $w_1$  be the width of the smallest column of  $T_1$ and let  $H_1 = h(C_{1j}) = t(1)$ , where  $1 \le j \le m_1$ . Note that  $H_1 \ge k$ . We now change to k the names of all levels in  $D(T_1)$ . This defines a new name function  $m_1$  on  $\mathcal{L}(T_1)$ . Precisely, for  $L \in T_1 \setminus D(T_1)$ , we have  $m_1(L) = m(L)$ , whilst for  $L \in D(T_1)$ .  $m_1(L) = k$ . Letting  $S_1 = T_1^*$ , we define a modified name partition  $G^{(1)} = \{G_i^{(1)}\}_{i=0}^k$  by

$$G_i^{(1)} = \bigcup \{ L \in \mathscr{L}(T_1) : \, \varkappa_1(L) = i \} \cup (E_i \cap (\Omega \setminus S_1)), \quad \text{for } 0 \leq i \leq k.$$

Note that

$$\mu(E_i \Delta G_i^{(1)}) < 2k\mu(F_1). \quad \text{for } 0 \le i \le k,$$

and so. by (4.1),

$$\mu(E_i \Delta G_i^{(1)}) < \eta_1, \quad \text{for } 0 \le i \le k.$$
(4.3)

We finally define the first modified stopping time  $\nu^{(1)}$  to be the stopping time corresponding to the partition  $\{E_i^{(1)}\}_{i=0}^k$  where

$$\begin{split} E_0^{(1)} &= G_0^{(1)} \cup (\Omega \backslash S_1), \\ E_i^{(1)} &= G_i^{(1)} \cap S_1, \quad \text{for } 1 \leq i \leq k. \end{split}$$

In other words, when considering  $\nu^{(1)}$ , we reduce our attention to  $S_1$ .

We now describe the general induction step which differs from the construction of  $T_1$ .

Suppose that at stage r. where  $r \ge 1$ ,  $T_r$  is a tower with  $m_r$  columns, say  $T_r = (C_{rj}: 1 \le j \le m_r)$  with bases  $B_{rj} = B(C_{rj})$ , and  $F_r = \bigcup_{j=1}^{m_r} B_{rj}$ . Let  $S_r = T_r^*$ , let  $w_r$  be the width of the smallest level in  $T_r$ , and let  $H_r$  be the common height of all columns of  $T_r$ . Suppose also that the  $m_r$ -name of levels in  $D(T_r)$  is k.

Define  $P_r$  to be the partition of  $F_r$  given by  $P_r = \{B_{rj}\}_{j=1}^{m_r}$ , and let  $\sigma_r = \tau_{F_r}$ . Choose

$$\eta_{r+1} < w_r / 2^{r+1} H_r. \tag{4.4}$$

284

Since  $\sigma_r$  is ergodic there exists  $N_r > 0$  such that for  $n \ge N_r$ , there is a collection  $\xi_{r,n}$  of atoms of  $\bigvee_{i=0}^{n-1} \sigma_r^{-i} P_r$  of total measure at least  $\mu(F_r) - \frac{1}{2}\eta_{r+1}$  and such that, for all  $B_{rj} \in P_r$  and  $A \in \xi_{r,n}$ ,

$$\left| f_A(B_{rj}, n) - \mu(B_{rj}) \right| < \frac{1}{2} \eta_{r+1}.$$
 (4.5)

Recall that  $f_A(B_{rj}, n)$  is a relative frequency of occurrence of the index jin the  $P_r$ -n-name of A with respect to  $\sigma_r$ . Inequalities (4.4) and (4.5) and the definition of  $w_r$  imply that, for any  $n \ge \max(N_r, 2/\eta_{r+1})$  and any set  $A \in \xi_{r,n}, f_A(B_{rj}, n) \ge 2/n$  for any j, where  $1 \le j \le m_r$ .

Now we choose t(r+1) a positive integer such that

$$t(r+1) > \max(N_r, 4k/\eta_{r+1}),$$
 (4.6)

and apply the strong Rohlin theorem on  $F_r$  to find  $F_{r+1} \subset F_r$ , which has t(r+1) disjoint images under  $\sigma_r$  and is such that

$$d\left(\bigvee_{i=0}^{i(r+1)-1} \sigma_{r}^{-i} P_{r} / F_{r+1}\right) = d\left(\bigvee_{i=0}^{i(r+1)-1} \sigma_{r}^{-i} P_{r}\right)$$
(4.7)

and

$$\mu\left(F_r\right) \bigvee_{i=0}^{t(r+1)-1} \sigma_r^i F_{r+1} \right) < \frac{\eta_{r+1}}{2t(r)}.$$

$$(4.8)$$

Let  $T_{r,1}$  be the  $\tau_r$ -tower with base  $F_{r+1}$  and columns pure with respect to  $\mathscr{L}(T_r)$ . For  $0 \leq i \leq t(r+1)-1$  and  $1 \leq j \leq m_r$ . (4.7) implies that  $\mu(\sigma_r^{-i}B_{rj} \cap F_r) = \mu(B_{rj})\mu(F_r) > 0$ . In other words, for any  $0 \leq i \leq t(r+1)-1$ ,  $\sigma_r^i F_{r+1}$  intersects the base of any column of  $T_r$ . Let  $T_r^{(i)}$  be the tower with columns consisting of the subcolumns of  $C_{rj}$  with base  $\sigma_r^i F_{r+1} \cap B_{rj}$ , for  $1 \leq j \leq m_r$ . If we follow the successive images of  $F_{r+1}$  under  $\tau_r = \tau_{Sr}$ , we see that first  $F_{r+1}$  climbs through the subtower  $T_r^{(0)}$ . Further, note that if  $\omega \in A(T_r^{(0)})$ , then  $\tau_r(\omega) \subset F_r$ , and thus  $\tau_r(A(T_r^{(0)})) = \tau_{Fr}(F_{r+1}) = \sigma_r(F_{r+1})$ . Hence next, under  $\tau_r$ ,  $F_{r+1}$  climbs through the levels of  $T_r^{(1)}$ , and so on. Consequently

$$T_{r,1} = T_r^{(0)} * T_r^{(1)} * \dots * T_r^{(t(r+1)-1)}.$$

So each column of  $T_{r,1}$  is a product of t(r+1) columns which are copies of the columns of  $T_r$ . A typical column C of  $T_{r,1}$  can thus be written as

$$C = [C_{ri_0}, \dots, C_{ri_{(r+1)-1}}].$$
(4.9)

where  $1 \leq i_l \leq m_r$  and  $0 \leq l \leq t(r+1)-1$ , and its base is then

$$F_{r+1} \cap B_{ri_0} \cap \sigma_r^{-1} B_{ri_1} \cap \ldots \cap \sigma_r^{-(t(r+1)-1)} B_{ri_{t(r+1)-1}}$$

We call a column of  $T_{r,1}$  regular if there appear at least two distinct

copies of each of the  $C_{rj}$ , for  $1 \leq j \leq m_r$ , in (4.9). We now delete from  $T_{r,1}$  all columns which are not regular and call the resulting tower  $T_{r,2}$ .

LEMMA 1. 
$$\mu(T_{r,1}^*) - \mu(T_{r,2}^*) < \frac{1}{2}\eta_{r+1}$$
.

*Proof.* Let  $R_r$  be the union of all bases of the columns of  $T_{r,1}$  which are not regular. Then the choice of t(r+1) in (4.6) implies that if a base

$$B = F_{r+1} \cap B_{ri_0} \cap \ldots \cap (\sigma_r)^{-(t(r+1)-1)} B_{ri_t(r+1)-1}$$

of  $T_{r,1}$  belongs to  $R_r$  then

$$B_{ri_0} \cap \ldots \cap (\sigma_r)^{-(t(r+1)-1)} B_{ri_{t(r+1)-1}}$$

belongs to  $\xi_{r,t(r+1)}^{c}$ . Thus

$$R_r \subset \xi^{c}_{r,t(r+1)} \cap F_{r+1},$$

and since  $\xi_{r,t(r+1)}^c$  is a union of atoms of  $\bigvee_{i=0}^{t(r+1)-1} (\sigma_r)^{-i} P_r$ . (4.7) then implies that

$$\mu(R_r) \leq \mu(\xi_{r,i(r+1)}^c) \mu(F_{r+1}) < \frac{1}{2} \eta_{r+1} \mu(F_{r+1})$$

so. finally.

$$\mu(T^*_{r,1}) - \mu(T^*_{r,2}) < \frac{1}{2}\eta_{r+1}.$$

We then purify  $T_{r,2}$  with respect to  $Q_r$ , the dyadic intervals of order r. and obtain a tower  $T_{r,3}$ . The columns of  $T_{r,3}$  are subcolumns of the columns of  $T_{r,2}$  and so, when considered as products of copies of columns of  $T_r$ , are regular.

We now form independent  $\alpha$ -blocks, the  $r \cdot \alpha$ -blocks, within each column of  $T_{r,3}$ . This defines the tower  $T_{r+1}$ . So  $T_{r+1}^* = T_{r,3}^*$  and  $D(T_{r+1})^* = D(T_{r,3})^*$ . We let  $G^{(r+1)} = \{G_i^{(r+1)}\}_{i=0}^k$  be the partition corresponding to  $G^{(r)}$  after formation of independent  $\alpha$ -blocks within each column of  $T_{r,3}$ , and we let  $n_{r+1}$  be the new name function associated with it. The tower  $T_{r+1}$  is thus pure with respect to  $G^{(r+1)}$ .

The set  $J_r$  of the levels of  $T_{r,3}$  whose  $n_r$ -name has been altered by the r- $\omega$ -blocks, is a subset of  $\{D(C_{r,j}): 1 \leq j \leq m_r\}$  and so is included in  $D(T_r)^*$ . Since all levels in  $D(T_r)$  have an  $n_r$ -name equal to k, and since in forming  $\omega$ -blocks we never give a name equal to 0, it follows that

$$G_0^{(r+1)} = G_0^{(r)},$$
  

$$G_i^{(r+1)} \supset G_i^{(r)}, \text{ for } 1 \le i \le k-1.$$
  

$$G_k^{(r+1)} \subset G_k^{(r)}.$$

286

Note that the set  $J_r$  is disjoint from  $D(T_{r,3})^*$  and that since  $D(T_{r,3}) \subset D(T_r)$ , the  $\varkappa_{r+1}$ -name of the levels in  $D(T_{r,3})$  is k. It is also clear that the formation of the r- $\varkappa$ -blocks has not altered the names of the levels of  $T_{r,3}$  included in any of the *i*- $\varkappa$ -blocks, for  $i \leq r-1$ , as they are disjoint from  $D(T_{i,3})^* = D(T_{i+1})^*$  and thus disjoint from  $D(T_r)$ . Now

$$\mu(G_i^{(r)} \Delta G_i^{(r+1)}) < 2k\mu(F_r),$$

and so, by our choice of t(r),

$$\mu(G_i^{(r)} \Delta G_i^{(r+1)}) < \eta_r.$$
(4.10)

We finally define  $S_{r+1} = T_{r+1}^*$  and let  $\nu^{(r+1)}$  be the stopping time corresponding to the partition  $\{E_i^{(r+1)}\}_{i=0}^k$  where

$$E_0^{(r+1)} = G^{(r+1)} \cup (\Omega \setminus S_{r+1}),$$
  

$$E_i^{(r+1)} = G_i^{(r+1)} \cap S_{r+1}, \text{ for } 1 \le i \le k.$$

This completes the induction step.

We are now ready to define the limiting stopping time  $\nu'$ . First note that the monotonic properties of the sequence  $\{G^{(r)}\}_{r\geq 1}$  imply that if we set  $G_i = \bigcup_{r=1}^{\infty} G_i^{(r)}$ . for  $0 \leq i \leq k-1$ . and  $G_k = \bigcap_{r=1}^{\infty} G_k^{(r)}$ , then  $G = \{G_i\}_{i=0}^k$  is a partition of  $\Omega$ . Letting  $S = \bigcap_{r=1}^{\infty} S_r$ , we can then define a partition  $\{E'_i\}_{i=0}^k$  of  $\Omega$  by

$$\begin{split} E'_{0} &= G_{0} \cup (\Omega \setminus S), \\ E'_{i} &= G_{i} \cap S, \quad \text{for } 1 \leq i \leq k, \end{split}$$

and we take  $\nu'$  to be the stopping time corresponding to the partition  $\{E'_i\}_{i=0}^k$ .

Note that as  $\{S_r\}_{r\geq 0}$  is a decreasing sequence of sets,  $S_0 = \Omega$ . Also

$$\mu(S) > 1 - \sum_{r=1}^{\infty} \eta_r$$
 (4.11)

since, by (4.2),  $\mu(S_0 \setminus S_1) < \eta_1$ , and, by (4.8) and Lemma 1.  $\mu(S_r \setminus S_{r+1}) \leq (\mu(S_r) - \mu(T^*_{r,1})) + (\mu(T^*_{r,1}) - \mu(T^*_{r,2}))$   $\leq \frac{1}{2}\eta_{r+1} + \frac{1}{2}\eta_{r+1}.$ 

LEMMA 2.

$$\mu(E_0 \Delta E'_0) < 2 \sum_{r=1}^{\infty} \eta_r,$$
$$\mu(E_i \Delta E'_i) < 3 \sum_{r=1}^{\infty} \eta_r, \quad for \ 1 \le i \le k.$$

Proof.

$$\mu(E_0 \Delta E'_0) \leq \mu(E_0 \Delta G_0^{(1)}) + \mu(\Omega \setminus S) < \eta_1 + \sum_{r=1}^{\infty} \eta_r$$

by (4.3) and (4.11). For  $1 \leq i \leq k-1$ .

$$E'_i = G_i \cap S$$
 and  $G_i = G_i^{(1)} \cup \bigcup_{r=1}^{\infty} (G_i^{(r+1)} \setminus G_i^{(r)});$ 

thus

$$\mu(E_i \Delta E'_i) \leq \mu(E_i \Delta G_i^{(1)}) + \sum_{r=1}^{\infty} \mu(G_i^{(r+1)} \setminus G_i^{(r)}) + \mu(\Omega \setminus S)$$
$$< \eta_1 + 2 \sum_{r=1}^{\infty} \eta_r,$$

by (4.3), (4.10), and (4.11). A similar result can be proved for  $\mu(E_k \Delta E'_k)$ . Note that the lemma implies. in particular, that

$$\mu(E'_{k}) > \mu(E_{k}) - 3 \sum_{r=1}^{\infty} \eta_{r}.$$
(4.12)

#### 5. Approximation results

The metric chosen here will be the uniform metric d defined by

 $d(\tau,\tau') = \mu \big\{ \omega \colon \tau(\omega) \neq \tau'(\omega) \big\}.$ 

where  $\tau$  and  $\tau'$  are both automorphisms of  $(\Omega, \mathcal{B}, \mu)$ .

We shall use the following fact:

LEMMA 3 [8]. Let  $\nu$  and  $\nu'$  be two stopping times for  $(\Omega, \mathcal{B}, \mu, \tau)$ , let  $\int \nu d\mu = \int \nu' d\mu = k$ . with  $\{E_i\}_{i=0}^k$  and  $\{E'_i\}_{i=0}^k$  their respective associated partitions. For any  $\varepsilon > 0$ , if  $\mu(E_i \Delta E'_i) < \varepsilon/k^2(k+2)$ , then  $d(\tau^{\nu}, \tau^{\nu'}) < \varepsilon$ .

We now give a preparatory result based on the following lemma which is due to Knopp.

LEMMA 4. Let  $(X, \mathcal{F}, m)$  be a probability space where  $\mathcal{F}$  is the  $\sigma$ -algebra generated by a ring  $\mathcal{R}$  of subsets of X. Suppose that there is a subset  $\mathcal{I}$  of  $\mathcal{R}$ such that every  $R \in \mathcal{R}$  is a finite disjoint union of sets in  $\mathcal{I}$ . Then if  $F \in \mathcal{F}$ is such that

$$m(F \cap I) \ge cm(I),$$

for any  $I \in \mathcal{I}$ , where c > 0 is a constant independent of I, it follows that m(F) = 1.

We use the following notation:  $\mu'$  represents the normalized Lebesgue measure on the set  $\{\nu' \neq 0\}$ , and  $\mathscr{L}_r^+ = \{L \in \mathscr{L}(T_r): L \subset G_0^{(r)c}\}$ .

 $\mathbf{288}$ 

STOPPING-TIME TRANSFORMATIONS

COROLLARY 1. If for any  $r \ge 1$  and  $L_1, L_2 \in \mathscr{L}_r^+$ , we have

$$\mu'\left(\left(\bigcup_{i=-\infty}^{\infty} (\tau^{\nu'})^i L_1\right) \cap L_2\right) \ge c\mu'(L_2),\tag{5.1}$$

where c is independent of  $L_2$ , then, for any  $L \in \mathscr{L}_r^+$ ,

$$\mu'\left(\bigcup_{i=-\infty}^{\infty}(\tau^{\nu'_{\cdot}})^{i}L\right)=1.$$

*Proof.* Fix  $r \ge 1$  and let  $L \in \mathscr{L}_r^+$ . Note that

$$\{\nu' \neq 0\} = G_0^{\,\mathfrak{c}} \cap S = (G_0^{(r)})^{\mathfrak{c}} \cap S.$$

If we let  $\mathscr{R}$  be the ring of finite disjoint unions of elements of  $\mathscr{L}_{p}^{+}$ , where  $p \geq r$ , then it follows by construction that  $\mathscr{R} \cap S$  generates the Borel  $\sigma$ -algebra restricted to  $\{\nu' \neq 0\}$ . Replacing, if necessary, L by one of its copies in  $T_{p}$  and using (5.1), we see that for any  $L' \in \mathscr{L}_{p}^{+}$ , with  $p \geq r$ .

$$\mu'\left(\left(\bigcup_{i=-\infty}^{\infty} (\tau^{\nu'})^i L\right) \cap L'\right) \ge c\mu'(L').$$
(5.2)

The result now follows from Lemma 4.

We now prove our main result which shows that for a suitable choice of  $\{\eta_r\}_{r\geq 1}$ , the stopping-time transformation  $\tau^{\nu'}$  is ergodic on  $\{\nu'\neq 0\}$ , has the same entropy as  $\tau^{\nu}$ , and is arbitrarily close to  $\tau^{\nu}$ .

THEOREM 3. Let  $\tau$  be an ergodic automorphism of  $(\Omega, \mathcal{B}, \mu)$ , and let v be a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$  with  $\int v d\mu = k$ , where k is a positive integer. Then for any  $\varepsilon > 0$ , there exists a stopping time v' for  $(\Omega, \mathcal{B}, \mu, \tau)$  such that the following statements hold:

- (i)  $\tau^{\mathbf{v}'}$  is ergodic in  $\{\mathbf{v}' \neq 0\}$ ;
- (ii)  $h(\tau^{\nu'}) = h(\tau^{\nu});$
- (iii)  $d(\tau^{\nu}, \tau^{\nu'}) < \varepsilon$ .

*Proof.* Let  $\{E_i\}_{i=0}^k$  be the partition associated with  $\nu$ . By Theorem 1 (iii).  $E_k$  is a non-empty set and so there exists  $\delta > 0$  such that  $\mu(E_k) > \delta$ . We now choose a decreasing sequence of positive numbers  $\{\eta'_r\}_{r\geq 1}$  satisfying

$$\sum_{r=1}^{\infty} \eta'_r < \min(\delta/6, \varepsilon/3k^2(k+2)).$$
(5.3)

Next we construct  $T_1$  as in §4 for  $\eta_1 = \eta'_1$ . We then choose  $\eta_2 < \eta'_2$  satisfying (4.4) and construct  $T_2$ , and so on. That is, in general, we 5388.3.43

289

choose  $\eta_{r+1}$  such that

$$\eta_{r+1} < \min(\eta'_{r+1}, w_r/2^{r+1}H_r), \text{ where } r \ge 1,$$
 (5.4)

and define the stopping time  $\nu'$  to be as constructed in §4 for that choice of  $\{\eta_r\}_{r\geq 1}$ .

(i) To prove that  $\tau^{\nu'}$  is ergodic in  $\{\nu' \neq 0\}$ , we first show that for any  $L \in \mathscr{L}_r^+$ , with  $r \ge 1$ ,

$$\mu'\left(\bigcup_{i=-\infty}^{\infty} (\tau^{\nu'})^i L\right) = 1.$$
(5.5)

By Corollary 1, it will be sufficient to show that for any  $L, L' \in \mathscr{L}_r^+$ ,

$$\mu'\left(\left(\bigcup_{i=-\infty}^{\infty}(\tau^{\nu'})^{i}L\right)\cap L'\right) \ge c\mu'(L'),$$

with a constant independent of L'.

So suppose that  $L \in C_{rl_0}$  and  $L' \in C_{rl_1}$ , where  $1 \leq l_0, l_1 \leq m_r$ . Now L' is distributed as distinct copies among the columns of  $T_{r,3}$ ; precisely, we write

$$L' \cap S_{r+1} = \bigcup_{i \in I} l'_i,$$

where  $l'_i$  is a copy of L' in a column of  $T_{r,3}$  and I is a finite index set. By construction, for any  $i \in I$ , in the column of  $T_{r,3}$  containing  $l'_i$ , there is a copy of  $C_{rl_0}$ , that we denote by  $l_i$ , above or below the copy of  $C_{rl_1}$  containing  $l'_i$ .

Recall that we have formed independent  $\alpha$ -blocks within each column of  $T_{r,3}$  and that  $\nu^{(r+1)}$  is the associated stopping time. Therefore we can apply Proposition 2 to each  $l_i$  and  $l'_i$ . For the sake of simplification, we shall include more powers of  $\tau^{\nu(r+1)}$  than those given by Proposition 2 and hence we have the following inequality:

$$\mu\left(\left(\bigcup_{j=-H_{r+1}}^{H_{r+1}}(\tau^{\nu(r+1)})^{j}l_{i}\right)\cap l_{i}'\right) \geq \mu(l_{i}')/k.$$

Now regrouping over  $i \in I$ , we obtain

$$\mu\left(\left(\bigcup_{j=-H_{r+1}}^{H_{r+1}} (\tau^{\nu(r+1)})^{j} L \cap S_{r+1}\right) \cap (L' \cap S_{r+1})\right) \ge \mu(L' \cap S_{r+1})/k. \quad (5.6)$$

Now for  $r \ge 1$  and by the choice of  $\eta_{r+1}$  in (5.4),

$$\mu(S_{r+1} \setminus S) < \sum_{i=r+1}^{\infty} \eta_{i+1} < \sum_{i=r+1}^{\infty} \frac{w_i}{H_i} \frac{1}{2^{i+1}}.$$

 $\mathbf{290}$ 

and so

$$\mu(S_{r+1} \setminus S) < w_{r+1} / 4H_{r+1}. \tag{5.7}$$

Hence writing  $L \cap S_{r+1}$  as  $L \cap S_{r+1} = (L \cap S) \cup (L \cap (S_{r+1} \setminus S))$  and similarly for  $L' \cap S_{r+1}$ , we deduce. from (5.6) and (5.7), that

$$\mu \left( \left( \bigcup_{j=-H_{r+1}}^{H_{r+1}} (\tau^{\nu(r+1)})^{j} (L \cap S) \right) \cap (L' \cap S) \right) \\
\geqslant (\mu(L' \cap S_{r+1})/k) - (2H_{r+1} + 3)\mu(S_{r+1} \setminus S) \\
\geqslant (\mu(L' \cap S_{r+1})/k) - w_{r+1}.$$
(5.8)

We now notice that, by (4.6),

$$w_{r+1} < \frac{1}{t(r+1)} < \frac{\eta_{r+1}}{4k} < \frac{w_r}{4k}.$$
 (5.9)

and that as L' is a level in  $T_r$ ,

$$\mu(L' \cap S_{r+1}) > \mu(L') - \eta_{r+1}$$
  
>  $\mu(L') - \mu(L')/2^{r+1}.$  (5.10)

Combining (5.9). (5.10). and (5.8). we then have

$$\mu\left(\left(\bigcup_{j=-H_{r+1}}^{H_{r+1}}(\tau^{\nu(r+1)})^{j}(L\cap S)\right)\cap (L'\cap S)\right) \ge \frac{\mu(L')}{k} - \frac{\mu(L')}{2^{r+1}k} - \frac{\mu(L')}{4k}$$

and therefore in any case

$$\mu\left(\left(\bigcup_{j=-H_{r+1}}^{H_{r+1}} (\tau^{v(r+1)})^j (L \cap S)\right) \cap (L' \cap S)\right) \ge \mu(L')/2k$$
$$\ge \mu(L' \cap S)/2k. \quad (5.11)$$

Note that the iterates of  $L \cap S$  under  $\tau^{v(r+1)}$  which make (5.11) hold are all *within* the columns of  $T_{r+1}$  and so the formation of  $\omega$ -blocks of order  $i \ge r+1$  will not affect them. So using (5.11), we finally obtain

$$\mu\left(\left(\bigcup_{j=-\infty}^{\infty} (\tau^{\mathbf{v}'_{j}})^{j}(L \cap S)\right) \cap (L' \cap S)\right)$$
  
$$\geq \mu\left(\left(\bigcup_{j=-H_{r+1}}^{H_{r+1}} (\tau^{\mathbf{v}'_{j}})^{j}(L \cap S)\right) \cap (L' \cap S)\right)$$
  
$$\geq \mu(L' \cap S)/2k.$$

or equivalently, as  $(\tau^{v'})S = S$  and normalizing, we have

$$\mu'\left(\left(\bigcup_{j=-\infty}^{\infty}(\tau^{\mathbf{v}'_{j}})^{j}L\right)\cap L'\right) \geqslant \mu'(L')/2k.$$

Thus (5.5) holds. Now, since the sets in  $\{(\bigcup_{r=1}^{\infty} \mathscr{L}_r^+) \cap S\}$  generate the Borel  $\sigma$ -algebra restricted to  $\{\nu' \neq 0\}$ , the ergodicity of  $\tau^{\nu'}$  on  $\{\nu' \neq 0\}$  follows by a standard approximation argument (see [3, §6], for instance).

(ii) Recall that  $\mu(E_k) > \delta$  and, by (5.3),  $\sum_{r=1}^{\infty} \eta'_r < \frac{1}{6}\delta$ . Using (4.12), we see that it then follows that the set  $E'_k$  is non-empty and hence that  $\int \nu' d\mu = k$ . As a consequence of (2.2), we have

$$h(\tau^{\mathbf{v}'}) = kh(\tau) = h(\tau^{\mathbf{v}}).$$

(iii) If we use Lemma 2 and equation (5.3), we see that

$$\mu(E_i \Delta E'_i) < 3 \sum_{r=1}^{\infty} \eta_r < \frac{\varepsilon}{k^2(k+2)}, \quad \text{for } 0 \le i \le k$$

Lemma 3 then implies that  $d(\tau^{\nu}, \tau^{\nu'}) < \epsilon$ .

COROLLARY 2. The conclusions (i) and (iii) of Theorem 3 still hold when  $\int \nu d\mu = +\infty$ .

*Proof.* We remarked earlier that we can approximate a stopping-time transformation  $\tau^{\nu}$  with  $\int \nu d\mu = +\infty$  by a stopping-time transformation  $\tau^{\nu_k}$  with  $\int \nu_k d\mu = k < \infty$ .

We are now in a position to study indirectly other ergodic properties of stopping-time transformations by using Theorem 3 and known density results for induced transformations.

THEOREM 4. Let  $\tau$  be an ergodic automorphism of  $(\Omega, \mathcal{B}, \mu)$ , and let  $\nu$  be a stopping time for  $(\Omega, \mathcal{B}, \mu, \tau)$ .

(i) For any  $\varepsilon > 0$ , there exists a stopping-time transformation  $\tau^{\nu_1}$  which is strong mixing on  $\{\nu_1 \neq 0\}$  and such that  $d(\tau^{\nu}, \tau^{\nu_1}) < \varepsilon$ .

(ii) Suppose that  $\tau$  has positive entropy. For any  $\varepsilon > 0$ , there exists a stopping-time transformation  $\tau^{\nu_2}$  which is a Kolmogorov automorphism on  $\{\nu_2 \neq 0\}$  and such that  $d(\tau^{\nu}, \tau^{\nu_2}) < \varepsilon$ .

*Proof.* (i) In [4] Friedman and Ornstein proved that given any ergodic automorphism  $\tau$ . the sets A such that  $\tau_A$  is strong mixing are dense.

So given  $\varepsilon > 0$ , we first use Corollary 2 to obtain a stopping time  $\nu'$  such that  $\tau^{\nu'}$  is ergodic on  $\{\nu' \neq 0\}$  and  $d(\tau^{\nu}, \tau^{\nu'}) < \frac{1}{2}\varepsilon$ . Let  $\tau^{\nu'}|_{\{\nu' \neq 0\}}$  denote the restriction of  $\tau^{\nu'}$  to  $\{\nu' \neq 0\}$ . Then there exists a set  $A \subset \{\nu' \neq 0\}$  such

that

(a)  $\mu(A)/\mu(\{\nu' \neq 0\}) > 1 - \frac{1}{6}\varepsilon$ . and (b)  $(\tau^{\nu'}|_{\{\nu' \neq 0\}})_A$  is strong mixing.

Let  $\tau^{\nu_1}$  be the transformation defined by

 $\tau$ 

$$\tau^{\nu_1}(\omega) = \begin{cases} (\tau^{\nu'})_A(\omega), & \text{if } \omega \in A, \\ \omega, & \text{if } \omega \in A^c. \end{cases}$$
(5.12)

Then clearly  $\tau^{\nu_1}$  is a stopping-time transformation. Moreover, since  $\tau^{\nu'_1}|_{\{\nu'\neq 0\}} = (\tau^{\nu'_1})_{\{\nu'\neq 0\}}$ , we have

$$\begin{aligned} {}^{\mathbf{v}_{1}} |_{\{\mathbf{v}_{1} \neq 0\}} &= (\tau^{\mathbf{v}'})_{A} \\ &= (\tau^{\mathbf{v}'})_{\{\mathbf{v}' \neq 0\} \cap A} \\ &= (\tau^{\mathbf{v}'} |_{\{\mathbf{v}' \neq 0\}})_{A}, \end{aligned}$$

which is strong mixing by (b). Also

$$d(\tau^{\mathbf{v}'_{\cdot}},\tau^{\mathbf{v}_{1}}) < 2\mu(\{\mathbf{v}'\neq 0\}\setminus A) + \mu(A^{\mathbf{c}}\setminus\{\mathbf{v}'=0\}) < \frac{1}{2}\varepsilon,$$

which in turn implies that  $d(\tau^{\nu}, \tau^{\nu_1}) < \varepsilon$ .

(ii) In [7], Ornstein and Smorodinsky proved that given any ergodic automorphism  $\tau$  of positive entropy. the sets A such that  $\tau_A$  is a K-automorphism are dense.

As in (i) we use Corollary 2 to obtain a stopping time  $\nu'$  with finite expectation such that  $\tau^{\nu'_1}|_{\{\nu'\neq 0\}}$  is ergodic and  $d(\tau^{\nu}, \tau^{\nu'_1}) < \frac{1}{2}\varepsilon$ . We now check that the entropy of  $\tau^{\nu'_1}|_{\{\nu'\neq 0\}}$  is strictly positive. We have

$$h(\tau^{\nu'}|_{\{\nu'\neq 0\}}) = h((\tau^{\nu'})_{\{\nu'\neq 0\}})$$
  
=  $kh(\tau) \cdot \mu(\{\nu'\neq 0\}),$ 

where  $\int \nu' d\mu = k < \infty$ .  $\mu(\{\nu' \neq 0\}) \ge \mu(E'_k) > 0$ , and we have used (2.2) to deduce the last equality. As  $h(\tau) > 0$  by hypothesis, it then follows that  $h(\tau^{\nu'}|_{\{\nu'\neq 0\}}) > 0$ , and we can then choose a set A, with  $\mu(A)/\mu(\{\nu' \neq 0\}) > 1 - \frac{1}{6}\varepsilon$ . such that  $(\tau^{\nu'}|_{\{\nu'\neq 0\}})_A$  is a K-automorphism. We then define  $\tau^{\nu_2}$  similarly to  $\tau^{\nu_1}$  in (5.12).

#### REFERENCES

- 1. R. M. BELINSKAYA. 'Partitions of Lebesgue space in trajectories defined by ergodic automorphisms'. Funkcional. Anal. i Priložen 2 (1968) 4-16; Functional Anal. Appl. 2 (1968) 190-99.
- 2. Generalised degrees of automorphism and entropy', Sibirsk. Mat. Z. 11 (1970) 739–49: Siberian Math. J. 11 (1970) 559–66.
- 3. N. A. FRIEDMAN. Introduction to ergodic theory (van Nostrand, New York, 1970).
- 4. and D. S. ORNSTEIN, 'Ergodic transformations induce mixing transformations , Adv. in Math. 10 (1973) 147-63.

- 5. S. KAKUTANI. Induced measure preserving transformations', Proc. Imp. Acad. Sci. Tokyo 19 (1943) 635-41.
- 6. J. NEVEU, 'Temps d'arrêt d'un système dynamique', Z. Wahrsch. Verw. Gebiete 13 (1969) 81-94.
- 7. D. S. ORNSTEIN and M. SMORODINSKY, 'Ergodic flows of positive entropy can be time changed to become K-flows', Israel J. Math. 26 (1977) 75-83.
- 8. S. T. RICHARDSON and K. M. WILKINSON, 'Stopping time transformations and towers', Z. Wahrsch. Verw. Gebiete 48 (1979) 259-84.
- 9. V. A. ROHLIN, Lectures on the entropy theory of measure preserving transformations', Uspehi Mat. Nauk 22 (1967) 3-56; Russian Math. Surveys 22 (1967) 1-52.

University of Warwick

Present address: U.170 INSERM 16 bis avenue Paul-Vaillant-Couturier 94800 Villejuif, France

## **DEUXIEME PARTIE**

## PROCESSUS SPATIALEMENT DEPENDANTS : CONVERGENCE VERS LA NORMALITE, TESTS D'ASSOCIATION ET APPLICATIONS

### ON THE VARIANCE OF THE SAMPLE CORRELATION BETWEEN TWO INDEPENDENT LATTICE PROCESSES

### ON THE VARIANCE OF THE SAMPLE CORRELATION BETWEEN TWO INDEPENDENT LATTICE PROCESSES

SYLVIA RICHARDSON\* AND

DENIS HEMON.\* Institut National de la Santé et de la Recherche Médicale

#### Abstract

Consider two stochastically independent, stationary Gaussian lattice processes with zero means,  $\{X(u), u \in \mathbb{Z}^2\}$  and  $\{Y(u), u \in \mathbb{Z}^2\}$ . An asymptotic expression for the variance of the sample correlation between  $\{X(u)\}$  and  $\{Y(u)\}$  over a finite square is derived. This expression also holds for a wide class of domains in  $\mathbb{Z}^2$ . As an illustration, the asymptotic variance of the correlation between two first-order autonormal schemes is evaluated.

AUTONORMAL SCHEMES; LATTICE PROCESSES; SAMPLE CROSS-CORRELATION

#### 1. Introduction

We shall consider throughout an infinite square lattice,  $\mathbb{Z} \times \mathbb{Z}$ , and random variables  $X(u) = X_{(u_1,u_2)}$ ,  $Y(v) = Y_{(v_1,v_2)}$  associated with each point  $(u_1, u_2)$  or  $(v_1, v_2)$  of the lattice. The processes  $\{X(u)\}$  and  $\{Y(v)\}$  are assumed to be mutually independent and both stationary Gaussian with zero means and finite variance  $\sigma_X^2$  and  $\sigma_Y^2$  respectively. Let  $C_n, n \ge 1$ , be any square region in  $\mathbb{Z}^2$ containing  $n^2$  points,  $C_n = \{(i, j); s \le i \le s + n - 1; t \le j \le t + n - 1\}$  for some integers s, t. We consider in this paper the variance of the sample correlation coefficient  $r_{XY}$  between the  $n^2$  observations from  $\{X(u)\}, \{Y(u)\}$  over  $C_n$ :

$$r_{XY} = \frac{\sum_{u \in C_n} X(u) Y(u)}{\left[\sum_{u \in C_n} X(u)^2\right]^{\frac{1}{2}} \left[\sum_{u \in C_n} Y(u)^2\right]^{\frac{1}{2}}}$$

Interest in lattice processes has arisen recently in many contexts, see Besag (1974) and Tjøstheim (1978) for further references. For time series, the asympto-

Received 7 July 1980; revision received 2 October 1980.

<sup>\*</sup> Postal address: U170 Statistiques, 16 bis avenue Paul Vaillant Couturier, 94800 Villejuif, France.

tic variance and covariances of the cross correlation were derived by Bartlett (1935), (1966), and their asymptotic distribution by Hannan (1970). These results are commonly used in the analysis of relationships between two time series (Haugh and Box (1977), Pierce (1979)). Similarly, spatial cross-correlation coefficients have been proposed for studying relationships between spatial series (Bennett (1980)).

#### 2. The asymptotic variance of $r_{XY}$

Variances and covariances will be denoted as follows:

$$C_{XX}(u, v) = E[X(u)X(v)], \qquad C_{YY}(u, v) = E[Y(u)Y(v)]$$
$$C_{XY}(u, v) = E[X(u)Y(v)],$$

in particular  $\operatorname{Var}[X(u)] = C_{xx}(0,0) = \sigma_x^2$ , where in this case 0 is the point in  $\mathbb{Z}^2$  with both coordinates equal to 0. It will be supposed that the autocovariances of  $\{X(u)\}$  and  $\{Y(u)\}$  are dominated by exponentially decreasing functions of their coordinates, that is there exists  $\theta_1, \theta_2, 0 < \theta_1, \theta_2 < 1$ , such that

(1)  
$$\begin{aligned} |C_{XX}(u,v)| &\leq c_1 \theta_1^{|u_1-v_1|+|u_2-v_2|} \\ |C_{YY}(u,v)| &\leq c_2 \theta_2^{|u_1-v_1|+|u_2-v_2|} \end{aligned}$$

where  $c_1, c_2$  are positive constants. Note that Condition (1) follows from Cauchy's inequalities whenever the autocovariance generating functions of  $\{X(u)\}$  and  $\{Y(u)\}$  have a Laurent expansion in the neighbourhood of the unit torus, as is the case for the autonormal schemes of Besag (1974) or the simultaneous autoregressive schemes of Whittle (1954). It will be shown that

(2) 
$$\lim_{n \to \infty} n^2 \operatorname{Var}(r_{XY}) = \frac{1}{\sigma_X^2 \sigma_Y^2} \sum_{v \in \mathbb{Z}^2} C_{XX}(0, v) C_{YY}(0, v),$$

where, due to (1), the series on the right-hand side is absolutely convergent.

To obtain this result, we first derive  $Var[\sum_{u \in C_n} X(u) Y(u)]$ . Since  $\{X(u)\}$  and  $\{Y(u)\}$  are independent and both have mean 0,

$$\operatorname{Var}\left[\sum_{u\in C_n} X(u)Y(u)\right] = \sum_{u\in C_n} \sum_{w\in C_n} C_{XX}(u,w)C_{YY}(u,w).$$

For  $0 < |v_1|, |v_2| \le n - 1$ , we let  $D(v) = D_{(v_1, v_2)} = \{(u, w) \in C_n \times C_n; u_1 - w_1 = v_1, u_2 - w_2 = v_2\}$ . With this notation,

$$\sum_{u \in C_n} \sum_{w \in C_n} C_{XX}(u, w) C_{YY}(u, w) = \sum_{\substack{|v_1| \le n-1 \\ |v_2| \le n-1}} \sum_{\substack{(u, w) \in D(v) \\ |v_2| \le n-1}} C_{XX}(u, w) C_{YY}(u, w).$$

Note that the number of couples (u, w) in D(v) is equal to the number of points in the intersection of  $C_n$  and  $T_{-v}(C_n)$ , the translation of  $C_n$  by -v, that is  $(n - |v_1|) \cdot (n - |v_2|)$ . Because  $\{X(u)\}$  and  $\{Y(u)\}$  are stationary  $C_{xx}(u, w) = C_{xx}(0, v)$  for  $(u, w) \in D(v)$ , so that

$$\operatorname{Var}\left[\sum_{u \in C_n} X(u) Y(u)\right] = \sum_{\substack{|v_1| \leq n-1 \\ |v_2| \leq n-1}} (n - |v_1|) \cdot (n - |v_2|) C_{XX}(0, v) C_{YY}(0, v).$$

Hence, using Condition (1),

(3) 
$$\frac{1}{n^2} \operatorname{Var} \left[ \sum_{u \in C_n} X(u) Y(u) \right] = \sum_{\substack{|v_1| \leq n-1 \\ |v_2| \leq n-1}} C_{XX}(0, v) C_{YY}(0, v) + O\left(\frac{1}{n}\right).$$

Next, we expand  $Var(r_{XY})$  in powers of  $n^{-2}$  using standard techniques. Let

$$A_n = \left[\frac{1}{n^2}\sum_{u \in C_n} X(u) Y(u)\right]^2, \qquad B_n = \left[\frac{1}{n^2}\sum_{u \in C_n} X(u)^2\right] \cdot \left[\frac{1}{n^2}\sum_{u \in C_n} Y(u)^2\right],$$

so that  $E(r_{XY}^2) = E(A_n/B_n)$ . The first term in the expansion of  $E(r_{XY}^2)$  in powers of  $n^{-2}$  is given by

$$\frac{E(A_n)}{E(B_n)} = \frac{1}{n^4 \sigma_N^2 \sigma_Y^2} E\left[\left(\sum_{u \in C_n} X(u) Y(u)\right)^2\right].$$

To show that  $E(r_{XY}^2)$  equals  $E(A_n)/E(B_n)$  plus terms of order  $n^{-4}$  we must consider the variances and covariances of  $A_n$  and  $B_n$ . Note that relation (3) implies that  $E(A_n)$  is of order  $n^{-2}$ . Now

$$E(A_{n}^{2}) = \frac{1}{n^{2}} \sum_{u,v,w,x \in C_{n}} E[X(u)X(v)X(w)X(x)]E[Y(u)Y(v)Y(w)Y(x)],$$

but since  $\{X(u)\}$  is Gaussian,

$$E[X(u)X(v)X(w)X(x)] = C_{xx}(u,v)C_{xx}(w,x) + C_{xx}(u,w)C_{xx}(v,x) + C_{xx}(u,x)C_{xx}(v,w),$$

and consequently  $E(A_n^2)$  can be expressed as the sum of nine terms, three symmetric in X, Y and six cross-product terms. Each symmetric term is equal to  $E(A_n)^2$  and hence of order  $n^{-4}$ . All cross-product terms are of similar form and so we consider only one. If  $\theta = \max(\theta_1, \theta_2)$ , and  $||u - v|| = |u_1 - v_1| + |u_2 - v_2|$ , we have:

$$\frac{1}{n^{s}} \sum_{u,v,w,x \in C_{n}} C_{XX}(u,v) C_{XX}(w,x) C_{YY}(u,w) C_{YY}(v,x)$$

$$\leq \frac{C_{1}C_{2}}{n^{s}} \sum_{u,v,w,x \in C_{n}} \theta^{||u-v||+||w-x||+||v-x||+||u-w||}$$

$$\leq \frac{C_{1}C_{2}}{n^{4}} \left(\sum_{v \in \mathbb{Z}^{2}} \theta^{||v||}\right)^{2}.$$

Consequently  $\operatorname{Var}(A_n)$ , and similarly  $\operatorname{Var}(B_n)$ ,  $\operatorname{Cov}(A_n, B_n)$  are of order  $n^{-4}$ . Similarly  $E(r_{XY}) = 0$  up to order  $n^{-2}$  so that the first term in the expansion of  $\operatorname{Var}(r_{XY})$  is given by

(4) 
$$\operatorname{Var}(r_{XY}) = \frac{1}{n^4 \sigma_X^2 \sigma_Y^2} \operatorname{Var}\left[\sum_{u \in C_n} X(u) Y(u)\right] + O\left(\frac{1}{n^4}\right).$$

Using (3) and (4), (2) follows immediately.

Expression (2) can be extended to sequences of domains in  $\mathbb{Z}^2$  which are not necessarily square-shaped. Consider an infinite sequence of domains  $\{D_n\}_{n\geq 1}$ with  $\lim_{n\to\infty} (b_n/d_n) = 0$ , where  $d_n$  is the number of lattice points in  $D_n$  and  $b_n$  the number of lattice points on the boundary of  $D_n$ . Comparing the number of points in  $D_n$  and in  $D_n \cap T_{-v}(D_n)$ , for v fixed in  $\mathbb{Z}^2$ , it can easily be shown that (2) still holds in this case. On the other hand, if we now consider an infinite sequence of rectangles of fixed width m and increasing length n, an expression other than (2) can be obtained for the asymptotic variance of  $r_{XY}$ , namely

$$\lim_{n \to \infty} nm \operatorname{Var}(r_{XY}) = \frac{1}{\sigma_X^2 \sigma_Y^2} \sum_{\substack{v_1 \in \mathbf{Z} \\ 0 \le |v_2| \le m-1}} C_{XX}(v, 0) C_{YY}(v, 0) - \frac{1}{m} \sum_{\substack{v_1 \in \mathbf{Z} \\ 0 \le |v_2| \le m-1}} |v_2| C_{XX}(v, 0) C_{YY}(v, 0),$$

an expression which is Bartlett's result for time series when m = 1. Result (2) is therefore specifically related to sequences of domains such that boundary effects can be asymptotically neglected.

#### 3. Application to symmetric first-order autonormal schemes

As an illustration of (2), suppose that  $\{X(u)\}\$  and  $\{Y(u)\}\$  are both stationary symmetric first-order autonormal schemes with autocovariance generating functions given by

$$F_{X}(z_{1}, z_{2}) = \frac{1}{1 - a(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})},$$
  

$$F_{Y}(z_{1}, z_{2}) = \frac{1}{1 - b(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})},$$
  

$$1 - \delta < |z_{1}|, |z_{2}| < 1 + \delta, |a| < \frac{1}{4}, |b| < \frac{1}{4}$$

The existence of such processes is discussed in Rozanov (1967) and Moran (1973b).

As noted by Besag (1972), the autocovariances of many spatial processes are not readily expressed in an analytically tractable form. However, using a remark of Quenouille (1947) for time series,  $\sum_{v \in \mathbb{Z}^2} C_{XX}(0, v) C_{YY}(0, v)$  can be calculated as the constant term in the Laurent expansion of  $F_X(z_1, z_2)F_Y(z_1, z_2)$ . Denoting this constant term by  $P_1(a, b)$ , we have

$$P_1(a,b) = \sum_{\substack{n=0\\m+n \text{ even}}}^{\infty} \sum_{\substack{m=0\\m+n \text{ even}}}^{\infty} a^n b^m \left(\frac{m+n}{\frac{1}{2}(m+n)}\right)^2,$$

so that

(5) 
$$P_{1}(a,b) = \begin{cases} \frac{1}{a-b} \left[ a \sum_{t=0}^{\infty} {\binom{2t}{t}}^{2} a^{2t} - b \sum_{t=0}^{\infty} {\binom{2t}{t}}^{2} b^{2t} \right], & a \neq b \\ \sum_{t=0}^{\infty} (2t+1) {\binom{2t}{t}}^{2} a^{2t}, & a = b. \end{cases}$$

Moran (1973a) noted that

$$\sum_{t=0}^{\infty} {\binom{2t}{t}}^2 a^{2t} = \frac{2}{\pi} K(16a^2)$$

where  $K(16a^2) = \int_0^{\pi/2} (1 - 16a^2 \sin^2 \theta)^{-\frac{1}{2}} d\theta$ , is the complete elliptic integral of the first kind. Similarly, using standard results (Abramowitz and Stegun (1965), 15.1.1., 17.3.10) we have

$$\sum_{t=0}^{\infty} (2t+1) \binom{2t}{t} a^{2t} = \frac{2}{\pi} \frac{1}{1-16a^2} E(16a^2),$$

where  $E(16a^2) = \int_0^{\pi/2} (1 - 16a^2 \sin^2 \theta)^{\frac{1}{2}} d\theta$ , is the complete elliptic integral of the second kind.

Consequently

(6) 
$$P_{1}(a,b) = \begin{cases} \frac{2}{\pi} \frac{1}{a-b} \left[ aK(16a^{2}) - bK(16b^{2}) \right], & a \neq b \\ \frac{2}{\pi} \frac{1}{1-16a^{2}} E(16a^{2}), & a = b. \end{cases}$$

Moran (1973a) also showed that  $\sigma_X^2 = (2/\pi)K(16a^2)$ ,  $\sigma_Y^2 = (2/\pi)K(16b^2)$ . We conclude from (2) and (6) that for  $a \neq b$ 

$$\lim_{n \to \infty} n^2 \operatorname{Var}(r_{XY}) = \frac{\pi}{2} \frac{1}{a-b} \left[ \frac{a}{K(16b^2)} - \frac{b}{K(16a^2)} \right]$$

and for a = b

$$\lim_{n \to \infty} n^2 \operatorname{Var}(r_{XY}) = \frac{\pi}{2} \frac{1}{1 - 16a^2} \frac{E(16a^2)}{[K(16a^2)]^2} \, .$$

Numerical evaluation of these last two expressions shows that  $n^2 \operatorname{Var}(r_{NY})$  increases more than exponentially with *a* and *b*.

#### Acknowledgements

We should like to thank Drs J. K. Ord and C. M. Triggs for helpful discussions and Mrs M. Guerrois for typing the manuscript.

#### References

ABRAMOWITZ, M. AND STEGUN, I. A. (1965) Handbook of Mathematical Functions. Dover, New York.

BARTLETT, M. S. (1935) Some aspects of the time-correlation problem in regard to tests of significance. J. R. Statist. Soc. 98, 536-543.

BARTLETT, M. S. (1966) An Introduction to Stochastic Processes, 2nd edn. Cambridge University Press.

BENNETT, R. J. (1980) Spatial Time Series. Pion, London.

BESAG, J. E. (1972) On the correlation structure of some two-dimensional stationary processes. Biometrika 59, 43-48.

BESAG, J. E. (1974) Spatial interaction and the statistical analysis of lattice systems (with discussion). J. R. Statist. Soc. B 36, 192-236.

HANNAN, E. J. (1970) Multiple Time Series. Wiley, New York.

HAUGH, L. D. AND BOX, G. E. P. (1977) Identification of dynamic regression (distributed lag) models connecting two time series. J. Amer. Statist. Assoc. 72, 121-130.

MORAN, P. A. P. (1973a) A Gaussian Markovian process on a square lattice. J. Appl. Prob. 10, 54-62.

MORAN, P. A. P. (1973b) Necessary conditions for Markovian processes on a lattice. J. Appl. Prob. 10, 605-612.

PIERCE, D. A. (1979) R<sup>2</sup> measures for time series J. Amer. Statist. Assoc. 74, 901-910.

QUENOUILLE, M. H. (1947) Notes on the calculation of autocorrelations of linear autoregressive schemes. *Biometrika* 34, 365-367.

ROZANOV. YU. A. (1967) On the Gaussian homogeneous fields with given conditional distributions. *Theory Prob. Appl.* 12, 381-391.

TJØSTHEIM, D. (1978) Statistical spatial series modelling. Adv. Appl. Prob. 10, 130-154. WHITTLE, P. (1954) On stationary processes in the plane. Biometrika 41, 434-449.

### AUTOCORRELATION SPATIALE : SES CONSEQUENCES SUR LA CORRELATION EMPIRIQUE DE DEUX PROCESSUS SPATIAUX

### AUTOCORRELATION SPATIALE : SES CONSEQUENCES SUR LA CORRELATION EMPIRIQUE DE DEUX PROCESSUS SPATIAUX

Sylvia RICHARDSON, Denis HEMON (\*)

#### RESUME

L'existence d'autocorrélations pose des problèmes quant aux tests d'indépendance stochastique de deux processus spatiaux.

Pour apprécier de façon quantitative ces problèmes, on étudie la variance du coefficient de corrélation empirique des deux processus. Cette variance est évaluée pour différentes valeurs des autocorrélations et pour plusieurs types de processus. Il ressort clairement de ces évaluations que l'influence des autocorrélations ne peut être négligée. Ceci conduit à discuter les approches proposées pour tester l'indépendance de deux processus spatiaux.

#### 1. INTRODUCTION

Les méthodes classiques d'analyse statistique visant à tester l'indépendance de deux variables aléatoires X et Y supposent que l'on dispose d'un échantillon de réalisations indépendantes du couple (X, Y).

Cependant lorsque ces variables sont liées à un espace géographique, elles présentent le plus souvent un certain degré d'autocorrélation spatiale. Ceci pose des problèmes du point de vue de l'analyse statistique, comme le soulignent LEBART [15], CLIFF et ORD [8] et UNWIN et HEPPLE [26]. Ces auteurs remarquent en particulier que si l'on néglige l'autocorrélation, on est souvent conduit à sous-estimer le risque de première espèce des tests employés. Pour apprécier de façon quantitative ces problèmes, nous étudions ici la variance du coefficient de corrélation empirique.

Comme dans le cas de séries temporelles, l'existence d'une corrélation entre deux processus spatiaux peut résulter de l'existence de deux tendances.. L'estimation et l'interprétation d'une tendance spatiale a fait l'objet de nombreux travaux, dont on trouvera une bibliographie dans UNWIN et HEPPLE [26]. Aussi ne consi-

(\*) INSERM U.170, 16 Bis av. P.V. Couturier, 94800 VILLEJUIF.

Revue de Statistique Appliquée, 1982, vol. XXX, n° 1 41

Mots-Clés : Processus spatiaux, autocorrélation, corrélation croisée.

dérons-nous que des processus stationnaires définis en chaque point d'un quadrillage régulier, processus étudiés en particulier par BESAG [6] et TJØSHEIM [25].

Soient  $\{X(u)\}$  et  $\{Y(u)\}$  deux processus stationnaires, gaussiens, centrés et de variances finies  $\sigma_X^2$  et  $\sigma_Y^2$  respectivement, définis en tout point  $u = (u_1, u_2)$  de Z x Z.Les autocovariances de  $\{X(u)\}$  seront notées :

$$C_{XX}(u, v) = E[X(u) X(v)], C_{YY}(u, v) = E[Y(u) Y(v)],$$

en particulier Var  $[X(u)] = C_{XX}(0, 0) = \sigma_X^2$ , 0 représentant le point de  $Z^2$ dont les deux cordonnées sont égales à zéro. Pour tout sous-ensemble  $D_n$  de  $Z^2$ contenant  $d_n$  points, le coefficient de corrélation empirique  $r_{XY}$  entre les observations de  $\{X(u)\}, \{Y(u)\}$  sur  $D_n$  est défini par :

$$r_{XY} = \frac{\sum_{u \in D_n} X(u) Y(u)}{\left[\sum_{u \in D_n} X(u)^2\right]^{1/2} \left[\sum_{u \in D_n} Y(u)^2\right]^{1/2}}$$

Sous l'hypothèse d'indépendance mutuelle entre  $\{X(u)\}$  et  $\{Y(u)\}$  nous avons démontré précédemment ([23]) que la variance asymptotique de  $r_{XY}$  est donnée par :

$$\lim_{n \to \infty} d_n \operatorname{Var}(\mathbf{r}_{XY}) = \frac{1}{\sigma_X^2 \sigma_Y^2} \sum_{\mathbf{v} \in \mathbf{Z}^2} C_{XX}(0, \mathbf{v}) C_{YY}(0, \mathbf{v}), \quad (1)$$

pour toute suite de domaines  $D_n$  de  $Z^2$  où les effets de bords s'estompent à l'infini, (c'est-à-dire tels que  $\lim_{n\to\infty} \frac{b_n}{d_n} = 0$ ,  $b_n$  étant le nombre de points de  $Z^2$  sur la frontière de  $D_n$ ). Notons que la série intervenant dans (1) converge absolument pour une large classe de processus, en particulier ceux dont les autocovariances  $C_{XX}(0, v)$  et  $C_{YY}(0, v)$  sont dominées par une fonction exponentielle décroissante du type  $c\theta^{|v|}$ , avec  $|\theta| < 1$  et c > 0, [23].

L'expression (1) est analogue à celle de BARTLETT dans le cas de séries temporelles, [2], [3]. Remarquons qu'il suffit que l'un des processus  $\{X(u)\}, \{Y(u)\}$ soit non autocorrelé pour que la variance asymptotique de  $r_{XY}$  ait la même valeur que dans le cas d'un échantillon de  $d_n$  couples (X,Y) indépendants c'est-à-dire  $1/d_n$ . Ceci est par exemple le cas en expérimentation agronomique où des parcelles disposées selon un lattice reçoivent un traitement tiré au sort. Par contre si les deux processus  $\{X(u)\}$  et  $\{Y(u)\}$  présentent tous deux une autocorrélation positive la relation (1) montre que la variance asymptotique de  $r_{XY}$  est supérieure à  $1/d_n$ .

Au § 2 nous évaluerons numériquement cette variance pour différents types de processus spatiaux. Au § 3 nous discuterons les différentes approches qui ont été proposées pour tester l'indépendance stochastique de deux processus spatiaux.

#### 2. EVALUATION DE LA VARIANCE ASYMPTOTIQUE DE $r_{xY}$ POUR PLUSIEURS TYPES DE PROCESSUS

#### 2.1. PROCESSUS DONT L'AUTOCOVARIANCE DECROIT GEOMETRIQUE-MENT

Considérons tout d'abord le cas où  $\{X(t), t \in Z\}$  et  $\{Y(t), t \in Z\}$  sont deux processus autoregressifs d'ordre 1 sur Z. Dans ce cas  $\{X(t), t \in Z\}$  et  $\{Y(t), t \in Z\}$  sont également des processus de Markov, leur autocovariance est donnée par :

$$C_{XX}(0, k) = E[X(0)X(k)] = \sigma_X^2 \lambda^{|k|}, k \in \mathbb{Z}, |\lambda_1| < 1$$
  

$$C_{YY}(0, k) = E[Y(0)Y(k)] = \sigma_Y^2 \lambda^{|k|}, k \in \mathbb{Z}, |\lambda_1| < 1.$$
 (2)

La variance asymptotique de  $r_{XY}$ , défini sur un intervalle de n points, est égale à :

$$\lim_{n \to \infty} n \operatorname{Var} (\mathbf{r}_{XY}) = \frac{1 + \lambda_1 \lambda_2}{1 - \lambda_1 \lambda_2}, \qquad (3)$$

formule obtenue par BARTLETT [2].

Cherchant à étendre la forme d'autocovariance donnée par [2] à un processus dans le plan, MARTIN [16] considère les processus "doublement géométriques" dont la structure stochastique est celle du produit de deux processus unidimensionnels satisfaisants (2); pour ces processus :

$$C_{XX}(0, v) = \sigma_X^2 \lambda_1^{|v_1| + |v_2|}, |\lambda_1| < 1, v \in Z^2.$$

Contrairement au cas temporel {X(u),  $u \in Z^2$ } n'est plus un processus markovien et l'équation "autorégressive" qui le définit n'a pas d'interprétation naturelle. Si l'on suppose que {X(u)} et {Y(u)} sont tous deux des processus doublement géométriques, de paramètre  $\lambda_1$  et  $\lambda_2$  respectivement, l'expression (1) devient alors l'extension directe à deux dimensions de la formule (3):

$$\lim_{n \to \infty} d_n \operatorname{Var} [r_{XY}] = \left(\frac{1 + \lambda_1 \lambda_2}{1 - \lambda_1 \lambda_2}\right)^2$$
(4)

#### 2.2. PROCESSUS MARKOVIENS ET "AUTOREGRESSIFS" DES PLUS PROCHES VOISINS

L'autocovariance de beaucoup de processus spatiaux n'est pas exprimable sous une forme analytique simple, BESAG [5]. Cependant l'expression  $\sum_{v \in \mathbb{Z}^2} C_{XX}(0, v) C_{YY}(0, v)$  peut être calculée comme terme constant du produit des développements de Laurent des fonctions génératrices d'autocovariance des processus {X(u)} et {Y(u)}.

#### a) Processus markoviens

Nous supposons que  $\{X(u)\}$  et  $\{Y(u)\}$  sont tous deux des processus stationnaires gaussiens, markoviens des plus proches voisins, définis par les fonctions génératrices des autocovariances suivantes :

$$F_{X}(z_{1}, z_{2}) = [1 - a(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})]^{-1}$$
  

$$F_{Y}(z_{1}, z_{2}) = [1 - b(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})]^{-1}$$
(5)

 $1 - \delta < |z_1|, |z_2| < 1 + \delta$ , |a| < 1/4, |b| < 1/4. L'existence de tels processus a été discutée par ROZANOV [24] et MORAN [19].

Ces processus satisfont aux relations suivantes :

$$X_{u_1,u_2} = a(X_{u_1-1,u_2} + X_{u_1+1,u_2} + X_{u_1,u_2-1} + X_{u_1,u_2+1}) + \epsilon_{u_1,u_2}, (u_1, u_2) Z \times Z.$$
(6)

Si  $\partial X(u) = \{X_{u_1-1,u_2}, X_{u_1+1,u_2}, X_{u_1,u_2-1}, X_{u_1,u_2+1}\}$  est l'ensemble des plus proches voisins de X(u), chaque variable aléatoire  $\epsilon(u)$  intervenant dans (6) est gaussienne ; conditionellement à  $\partial X(u)$  elle est indépendante de  $\{X(u)\}$  avec  $|v_1 - u_1| + |v_2 - u_2| > 1$  et centrée.

L'expression de la variance asymptotique de  $r_{XY}$  que l'on calcule à partir des fonctions génératirces des autocovariances (5) est donnée dans l'appendice. Cette expression, qui fait intervenir les intégrales elliptiques des paramètres a et b, a été tabulée par intégration numérique avec trois décimales de précision (Table 1), les valeurs des paramètres a et b sont choisies de telle sorte que les autocorrélations d'ordre 1\*,  $\rho_X$  et  $\rho_Y$ , varient de 0,0 à 0,7 par pas de 0,1.

TABLE 1 Valeur asymptotique de  $d_n Var(r_{XY})$  pour deux processus Markoviens

$\rho_{\rm X}$	0,00	0,100	0,200	0,300	0,400	0,500	0,600	0,700	b
0,000	1,00								0,000
0,100	1,00	1,042							0,095
0,200	1,00	1,088	1,195						0,168
0,300	1,00	1,138	1,321	1,558					0,2123
0,400	1,00	1,192	1,469	1,866	2,444				0,2355
0,500	1,00	1,248	1,637	2,251	3,270	4,989			0,24565
0,600	1,00	1,307	1,824	2,720	4,419	7,917	15,903		0,24918
0,700	1,00	1,368	2,025	3.256	5,883	12,432	33,523	140,603	0,24994
a	0,00	0,095	0,168	0,2123	0,2355	0,24565	0,24918	0,249949	

 $\rho_X$  et  $\rho_Y$  sont les autocorrélations d'ordre 1 de  $\{X(u)\}$  et  $\{Y(u)\}$  respectivement. Les valeurs correspondantes de a et b peuvent être lues respectivement en bas et à droite. Pour a >b, les valeurs sont obtenues par symétrie.

44

Revue de Statistique Appliquée, 1982, vol. XXX, n° 1

#### b) Processus autorégressifs

Supposons maintenant que  $\{X(u)\}$  et  $\{Y(u)\}$  soient des processus stationnaires gaussiens associés aux fonctions génératrices d'autocovariance suivantes :

$$F_{X}(z_{1}, z_{2}) = [1 - a(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})]^{-2}$$
  
$$F_{Y}(z_{1}, z_{2}) = [1 - b(z_{1} + z_{1}^{-1} + z_{2} + z_{2}^{-1})]^{-2}$$

 $1 - \delta < |z_1|, |z_2| < 1 + \delta$ , |a| < 1/4, |b| < 1/4. Ces processus, qui furent considérés pour la première fois par WHITTLE [27], correspondent aux équations autorégressives :

$$X_{u_1, u_2} = a(X_{u_1 - 1, u_2} + X_{u_1 + 1, u_2} + X_{u_1, u_2 - 1} + X_{u_1, u_2 + 1}) + \epsilon_{u_1, u_2}$$
(7)

dans lesquelles la suite d'innovations {E (u)},  $u = (u_1, u_2) Z \times Z$ }, est une suite de variables aléatoires gaussiennes indépendantes, d'espérance nulle. Notons que pour le processus markovien (6),  $\epsilon(u)$  et  $\epsilon(v)$  sont corrélées quand  $|u_1 - v_1| + |u_2 - v_2| \le 2$ . L'expression de la variance asymptotique de  $r_{XY}$ en fonction des paramètres des processus est donnée dans l'appendice et a été également tabulée (Table 2).

\*L'autocorrélation d'ordre 1,  $\rho_X$ , est définie comme le coefficient de corrélation entre X(u) et l'un quelconque de ses plus proches voisins.

ρ <sub>Y</sub>	0,00	0,100	0,200	0,300	0,400	0.500	0,600	0,700	0,800	0,900	Ъ
0.000	1,00										0,0000
0,100	1,00	1,041									0,0493
0,200	1,00	1,085	1,178								0.0945
0.300	1,00	1,130	1,282	1,458							0,1335
0,400	1,00	1,178	1,394	1,655	1,960						0,165
0.500	1,00	1,227	1,515	1,878	2,320	2,864					0,1902
0.600	1,00	1,279	1,647	2,130	2,743	3,530	4,540				0,2099
0.700	1,00	1,332	1,788	2,411	3,237	4,347	5,844	7,894			0,225
0,800	1,00	1,387	1,940	2,729	3,823	5,369	7,581	10,826	15,893		0,2364
0,900	1,00	1,444	2,094	3,096	4,541	6,710	10,045	15.403	24,838	44,846	0,2447
a	0,000	0,0493	0,0945	0,1335	0,165	0,1902	0,2099	0,225	0,2364	0,2447	

TABLE 2

Valeur asymptotique de  $d_n Var(r_{XY})$  pour deux processus autoregressifs

 $\rho_X$  et  $\rho_Y$  sont les autocorrélations d'ordre 1 de  $\{X(u)\}$  et  $\{Y(u)\}$  respectivement. les valeurs correspondantes de a et b peuvent être lues respectivement en bas et à droite. Pour a > b, les valeurs sont obtenues par symétrie.

Revue de Statistique Appliquée, 1982, vol. XXX, n° 1

En comparant les résultats obtenus pour les différents processus envisagés (formules (3), (4), Tables 1,2), on remarque tout d'abord que l'effet de l'autocorrélation intrinsèque des processus  $\{X(u)\}$  et  $\{Y(u)\}$  sur la variance du coefficient de corrélation est numériquement bien plus important dans le cas de processus spatiaux que dans celui des séries temporelles. Par exemple pour  $\rho_X = \rho_Y = 0.3$ , la "taille corrigée" d'un N-échantillon est de N/1,20 pour une série temporelle satisfaisant (2), tandis qu'elle est de N/1,56 pour le processus markovien dans  $Z^2$  défini par (6).

Par ailleurs les valeurs données pour les processus markoviens (Table 1) sont plus élevées que celles correspondant aux processus autorégressifs (Table 2) et aux processus doublement géométrique. Néanmoins quand  $\rho_X$  est plus petit que 0,3, les valeurs données pour les trois processus sont voisines.

#### 3. DISCUSSION

L'analyse de la corrélation de deux séries temporelles est classiquement conduite en étudiant les corrélations croisées ou le spectre croisé. La plupart des concepts utilisés dans ce type d'analyse peuvent être étendus au cas des séries spatiales. Avant d'envisager de façon détaillée les différentes méthodes proposées, il est bon de rappeler les difficultés intrinsèques à l'étude des processus spatiaux.

La première de ces difficultés tient à l'absence d'un ordre naturel sur  $Z^2$ . Ceci exclue en particulier la possibilité d'utiliser la distinction entre le passé et le futur pour interpréter les corrélations, PIERCE et HAUGH [22]. De ce point de vue, il est clair que l'analyse de séries à la fois spatiales et temporelles est particulièrement informative, GRANGER [11]. Une seconde difficulté tient à l'arbitraire du choix de la taille des unités géographiques considérées : on sait en effet que les résultats des analyses statistiques dépendent du degré d'agglomération de ces unités. Enfin, les observations sont souvent disposées de manière irrégulière dans l'espace ce qui interdit une analyse spectrale classique. Dans le cas d'un quadrillage irrégulier, on peut toutefois donner un sens à une décomposition hiérarchique en variations locales, départementales et régionales, comme le proposent par exemple CLIFF et ORD [9] et CURRY [10].

Les coefficients de corrélation croisés sont couramment employés dans l'analyse de la liaison de deux séries temporelles, HAUGH et BOX [14], PIERCE [21]. HANNAN [13] a établi la normalité asymptotique de ces coefficients sous certaines hypothèses. BENNETT [4] a suggéré d'entendre cette approche à l'étude de la liaison entre deux processus spatiaux. Si l'on se place dans le cas de processus stationnaires sur  $Z^2$  et d'un domaine d'échantillonnage regulier assez grand, on peut estimer les autocorrélations  $\frac{C_{XX}(0, v)}{\sigma_X^2}$ ,  $\frac{C_{YY}(0, v)}{\sigma_Y^2}$  d'un grand nombre de processus de manière consistante, GUYON et PRUM [12]. On peut alors estimer de façon consistante une approximation de la variance de  $r_{XY}$  en utilisant une version tronquée de la formule (1).

Plusieurs auteurs ont étudié des modifications du modèle de régression linéaire. LEBART [15] fixe la structure de la matrice de variance-covariance des résidus en la décomposant suivant des "niveaux de contiguités" décroissants. Cela

46

lui permet d'estimer les paramètres du modèle par la méthode des moindres carrés généralisés et d'apprécier la perte d'information due à la dépendance spatiale. Cette décomposition de la matrice de variance-covariance peut s'appliquer au cas où les variables sont mesurées aux sommets d'un graphe quelconque ; elle suppose que la covariance de deux variables dépende uniquement de leur niveau de contiguité sur le graphe et qu'elle soit nulle au delà d'une certaine distance. Généralisant les modèles de régression des séries chronologiques, ORD [20] étudie différents modèles mixtes régressifs-autorégressifs dans lesquels interviennent une matrice de voisinage supposée connue à un facteur de proportionalité près,  $\rho$ . Il discute la qualité des estimateurs de régressions et de  $\rho$  et compare l'estimateur du maximum de vraisemblance à plusieurs alternatives. Comme dans le cas de Lebart, ces modèles ne supposent pas une structure régulière de l'espace et la matrice de voisinage est adaptée au cas particulier considéré. Une approche empirique visant à prendre en compte l'effet des autocorrélations dans le modèle linéaire a été proposée par MARTIN [17]. Elle consiste à transformer  $\{X(u)\}$  et  $\{Y(u)\}$  en leur différence d'ordre 1, c'est-à-dire à supposer essentiellement  $\rho = 1$  avant d'effectuer une régression linéaire. Martin étudie l'effet de cette procédure par des méthodes de Monte Carlo et montre qu'en supposant  $\rho = 1$  les estimateurs des paramètres de régression sont en général meilleurs qu'en négligeant les autocorrélations ( $\rho = 0$ ).

Plusieurs auteurs se sont préoccupés de la définition d'une mesure de corrélation sur des unités que l'on pouvait en principe subdiviser ou regrouper. Les méthodes proposées sont toutes fondées sur l'idée d'une décomposition hiérarchique des variations de chaque processus. CURRY [10] estime les paramètres d'un modèle dans lequel la valeur prise par le processus Y en un point u est une combinaison linéaire des différentes composantes hiérarchiques du processus X. CLIFF et ORD [9] utilisent une analyse de variance emboîtée en postulant que la corrélation entre  $\{X(u)\}$  et  $\{Y(u)\}$  est dûe à un facteur commun à chaque niveau hiérarchique. BESAG [7] propose un test non paramétrique qui fait intervenir toutes les permutations d'unités à l'intérieur d'un même bloc hiérarchique.

Les méthodes décrites ci-dessus apportent des solutions partielles au problème général de l'analyse de la corrélation entre deux processus spatiaux. Certaines présentent un caractère empirique impliquant la connaissance de la structure de la matrice de variance-covariance. D'autres adoptent des méthodes développées pour les séries chronologiques, dans ce cas la validité de certains résultats asymptotiques utilisés reste à établir. Les problèmes d'inférence statistique concernant l'indépendance de deux processus spatiaux restent donc encore ouverts, l'importance de l'effet des autocorrélations justifie que des travaux complémentaires leur soient consacrés.

#### REMERCIEMENTS

Nous remercions vivement Micheline IMBERT pour la préparation de ce manuscrit.

#### APPENDICE

# a) LA VARIANCE ASYMPTOTIQUE DE $r_{XY}$ DANS LE CAS DE DEUX PROCESSUS DE MARKOV

Le terme constant  $P_1(a, b)$  du produit  $F_X(z_1, z_2)$ .  $F_Y(z_1, z_2)$  est égal à

$$P_{1}(a,b) = \begin{cases} \frac{1}{a-b} \left[ a \sum_{t=0}^{\infty} {\binom{2t}{t}}^{2} a^{2t} - b \sum_{t=0}^{\infty} {\binom{2t}{t}}^{2} b^{2t} \right], a \neq b \\ \sum_{t=0}^{\infty} (2t+1) {\binom{2t}{t}}^{2} a^{2t} , a = b \end{cases}$$
(8)

Les séries constituant (i) peuvent être exprimées au moyen des séries hypergéométriques de Gauss, ABRAWOWITZ et STEGUN [1] 15.1.1 :

$$F(k, \ell, m, z) = \sum_{t=0}^{\infty} \frac{(k)_t (\ell)_t}{(m)_t} \frac{z^t}{t!}$$

 $(k)_t = k(k + 1)... (k + t - 1).$  On obtient

$$\sum_{t=0}^{\infty} {\binom{2t}{t}^2}_{a^{2t}} = F(1/2, 1/2, 1, 16a^2)$$
$$\sum_{t=0}^{\infty} {(2t+1)} {\binom{2t}{t}^2}_{a^{2t}} = F(3/2, 1/2, 1, 16a^2)$$

On peut utiliser les relations qui lient les séries hypergéométriques aux intégrales elliptiques de la  $1^{ere}$  et  $2^e$  espèce quand on veut les évaluer numériquement :

F (1/2 . 1/2 . 1 , z) = 
$$\frac{2}{\Pi}$$
 K (z)  
F (3/2 , 1/2 . 1 , z) =  $\frac{2}{\Pi} \frac{1}{1 - z}$  E (z)

avec

K(z) = 
$$\int_0^{\pi/2} (1 - z \sin^2 \theta)^{-\frac{1}{2}} d\theta$$
 et E(z) =  $\int_0^{\pi/2} (1 - z \sin^2 \theta)^{\frac{1}{2}} d\theta$ ,

ABRAMOVITZ et STEGUN [1] 17.39 et 17.3.10.

En conséquence

$$P_{1}(a,b) = \begin{cases} \frac{2}{\Pi} \frac{1}{a-b} [aK(16a^{2}) - bK(16b^{2})], a \neq b \\ \frac{2}{\Pi} \frac{1}{1-16a^{2}} E(16a^{2}), a = b. \end{cases}$$

48

Revue de Statistique Appliquee, 1982, voi. XXX, n° 1

MORAN [18] a remarqué que :

$$\sigma_{\rm X}^2 = \frac{2}{\Pi} \ {\rm K} (16a^2) , \sigma_{\rm Y}^2 = \frac{2}{\Pi} \ {\rm K} (16b^2) .$$

On obtient donc finalement pour  $a \neq b$ :

$$\lim_{n \to \infty} d_n \operatorname{Var} [r_{XY}] = \frac{\Pi}{2} \frac{1}{a-b} \left[ \frac{a}{K(16b^2)} - \frac{b}{K(16a^2)} \right],$$

tandis que pour a = b:

$$\lim_{n \to \infty} d_n \operatorname{Var}[r_{XY}] = \frac{\Pi}{2} \frac{1}{1 - 16a^2} \frac{\mathsf{E}(16a^2)}{[\mathsf{K}(16a^2)]^2}$$

# b) LA VARIANCE ASYMPTOTIQUE DE $\mathbf{r}_{X\,Y}$ DANS LE CAS DE DEUX PROCESSUS AUTOREGRESSIFS

Le terme constant P<sub>2</sub> (a, b) du produit F<sub>X</sub> (z<sub>1</sub>, z<sub>2</sub>) . F<sub>Y</sub>(z<sub>1</sub>, z<sub>2</sub>) est égal à  

$$\begin{pmatrix}
\frac{2}{\Pi} \frac{1}{(a-b)^2} \left[ \frac{b^2}{1-16b^2} E(16b^2) + \frac{a^2}{1-16a^2} E(16a^2) \right] \\
+ \frac{2}{\Pi} \frac{2ab}{(a-b)^3} \left[ bK(16b^2) - aK(16a^2) \right] , a \neq b \\
F(5/2, 1/2, 1, 16a^2) + 20a^2 F(7/2, 3/2, 2, 16a^2), a = b.$$

On peut exprimer F(5/2, 1/2, 1, z) et F(7/2, 3/2, 2, z) en fonction des intégrales elliptiques K(z) et E(z) à l'aide des relations de récurrence suivantes :

$$5 F\left(\frac{7}{2}, \frac{3}{2}, 2, z\right) = 2 F\left(\frac{5}{2}, \frac{3}{2}, 1, z\right) + 3 F\left(\frac{5}{2}, \frac{3}{2}, 2, z\right)$$

$$3 F\left(\frac{5}{2}, \frac{3}{2}, 2, z\right) = 2 \frac{(1-z)}{z} \left[F\left(\frac{5}{2}, \frac{3}{2}, 1, z\right) - F\left(\frac{5}{2}, \frac{1}{2}, 1, z\right)\right]$$

$$2 F\left(\frac{5}{2}, \frac{3}{2}, 1, z\right) = \frac{2}{1-z} \left[4 F\left(\frac{5}{2}, \frac{1}{2}, 1, z\right) - 3 F\left(\frac{3}{2}, \frac{1}{2}, 1, z\right)\right]$$

$$F\left(\frac{5}{2}, \frac{1}{2}, 1, z\right) = \frac{1}{3(1-z)} \left[2(2-z) F\left(\frac{3}{2}, \frac{1}{2}, 1, z\right) - F\left(\frac{1}{2}, \frac{1}{2}, 1, z\right)\right]$$

Finalement remarquant que dans ce cas  $\sigma_X^2 = P_1(a, a)$  et  $\sigma_Y^2 = P_1(b, b)$  est défini par (8), on obtient :

$$\lim_{n \to \infty} d_n \text{ Var } [r_{XY}] = \frac{P_2(a, b)}{P_1(a, a) P_1(b, b)}$$

Revue de Statistique Appliquée, 1982, vol. XXX, n° 1

#### REFERENCES

- [1] M. ABRAMOWITZ and I.A. STEGUN. Handbook of Mathematical Functions. Dover, New-York, 1965.
- [2] M.S. BARTLETT. Some aspects of the time-correlation problem in regard to tests of significance. J.R. Statist. Soc., 98, 536-543, 1935.
- [3] M.S. BARTLETT. An Introduction to Stochastic Processes, 2nd edn. Cambridge University Press, Cambridge, 1966.
- [4] R.J. BENNETT. Spatial Time Series. Pion, London, 1980.
- [5] J.E. BESAG. On the correlation structure of some two-dimensional stationary processes. *Biometrika*, 59, 43-48, 1972.
- [6] J.E. BESAG. Spatial interaction and the statistical analysis of lattice system (with discussion) J.R. Statist, Soc., B 36, 192-236, 1974.
- [7] J.E. BESAG and P.J. DIGGLE. Simple Monte Carlo tests for spatial pattern. J.R. Statist. Soc., C 26, 327-333, 1977.
- [8] A.D. CLIFF and J.K. ORD. Model building and the analysis of spatial pattern in human geography (with discussion). J.R. Statist. Soc., B 37, 297-348, 1975.
- [9] A.D. CLIFF and J.K. ORD. The effects of spatial autocorrelation on geographical modelling. University of Warwick Technical Report, 1978.
- [10] L. CURRY. A bivariate spatial regression operator. Canadian Geographer, 16, 1-14, 1972.
- [11] C.W.J. GRANGER. Spatial data and time series analysis. In Studies in Regional Sciences (A.J. Scott, ed), 1-24. London : Pion, 1969.
- [12] X. GUYON et B. PRUM. Estimations et tests relatifs aux processus spatiaux réguliers du second ordre, Publication n° 201, Université d'Orsay, 1976.
- [13] E.J. HANNAN. Multiple Time Series. John Wiley and Sons, New-York, 1970.
- [14] L.D. HAUGH and G.E.P. BOX. Identification of dynamic regression (distributed lag) models connecting two time series. J. Amer. Statist. Assoc., 72, 121-130, 1977.
- [15] L. LEBART. Analyse statistique de la continguité. Pub. Inst. Stat. Univ. Paris, 18, 81-112, 1969.
- [16] R.J. MARTIN. A subclass of lattice processes applied to a problem in planar sampling. *Biometrika*, 66, 209-17, 1979.
- [17] R.L. MARTIN. On spatial dependence, bias, and the use of first spatial differences in regression analysis. Area, 6, 185-94, 1974.
- [18] P.A.P. MORAN. A Gaussian Markovian process on a square lattice. J. Appl. Prob., 10, 54-62, 1973.
- [19] P.A.P. MORAN. Necessary conditions for Markovian processes on a lattice. J. Appl. Prob., 10, 605-612, 1973.
- [20] J.K. ORD. Estimation methods for models of spatial interaction. J. Am. Statist. Assoc., 70, 120-6, 1975.

- [21] D.A. PIERCE. R<sup>2</sup> Measures for time series. J. Amer. Statist. Assoc., 74, 901-910, 1979.
- [22] D.A. PIERCE and L.D. HAUGH. Causality in temoral systems, characterizations and a survey. J. Econometrics, 5, 265-293, 1977.
- [23] S.T. RICHARSON and D. HEMON. On the variance of the sample correlation between two independant lattice processes. J. Appl. Prob., 18, 943-8, 1981.
- [24] Yu. A. ROZANOV. On the Gaussian homogeneous fields with given conditional distributions. *Theor. Prob. Appl.*, **12**, 381-391, 1967.
- [25] D. TJOSTHEIM. Statistical spatial series modelling. Adv. Appl. Probl., 10, 130-154, 1978.
- [26] D.J. UNWIN and L.W. HEPPLE. The statistical analysis of spatial series. The Statistician, 23, 211-27, 1974.
- [27] P. WHITTLE. On stationary processes in the plane. Biometrika, 41, 434-449, 1954.
## VITESSE DE CONVERGENCE DU THEOREME DE LA LIMITE CENTRALE POUR DES CHAMPS FAIBLEMENT DEPENDANTS

Z. Wahrscheinlichkeitstheorie verw. Gebiete 66, 297-314 (1984)

Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete © Springer-Verlag 1984

### Vitesse de convergence du théorème de la limite centrale pour des champs faiblement dépendants

Xavier Guyon<sup>1</sup> et Sylvia Richardson<sup>2</sup>

 <sup>1</sup> Laboratoire de Statistiques. Bâtiment 425, Université de Paris XI, F-91405 Orsay, France
 <sup>2</sup> Laboratoire de Statistiques. Université de Paris V, 45, rue des Saints Pères. F-75006 Paris, France

**Résumé.** Nous étudions la vitesse de convergence du théorème de la limite centrale pour des champs de  $\mathbb{Z}^d$ , faiblement dépendants: *m*-dépendant ou  $\alpha$ -fortement mélangeant. Dès que le champ est dans  $L^{2+\delta}$ ,  $\delta > 0$ . la vitesse de convergence obtenue est  $\sigma_n^{-(\delta \wedge 1)}$  avec un facteur  $(\log \sigma_n)^a$  qui intervient quand  $\alpha$  est à décroissance exponentielle et dans le cas *m*-dépendant quand  $\delta \ge 1$ . Le cas où  $\alpha$  est à décroissance puissance est aussi étudié. Ces résultats ne font intervenir ni la stationarité, ni la géométrie des domaines sur lesquels le T.L.C. est étudié<sup>1</sup>.

#### §1. Introduction

Nous nous intéressons dans ce travail à la vitesse de convergence dans le théorème de la limite centrale (T.L.C.) pour  $X = \{X_j, j \in \mathbb{Z}^d\}$  un champ defini sur  $\mathbb{Z}^d$ , centré, faiblement dépendant, vérifiant:

$$\sup_{j} \|X_{j}\|_{2+\delta} = \|X\|_{2+\delta} < \infty, \quad \delta > 0.$$
 (1-1)

Deux types de faible dépendance sont étudiés classiquement: la *m*-dépendance et les concepts de mélange. Parmi ceux-ci nous nous limiterons ici à une notion de  $\alpha$ -mélange fort.

La litterature dans le cas d=1 est abondante, traitant également le cas de processus accroissements de martingale (voir par exemple Ibragimov [10], Hall et Heyde [11] et Bolthausen [2]).

Pour le cas m-dépendant sans hypothèse de stationnarité. Petrov [17],

<sup>&</sup>lt;sup>1</sup> Au moment où est soumis cet article, les auteurs on pris connaissance de l'article de H. Takahata [23] qui étudie exactement le même problème. La condition de mélange (Dobrushin) est moins forte que la nôtre et les résultats sont obtenus en utilisant la technique de Stein. Pour des champs *m*-dépendant et sous une condition  $L^8$ , la vitesse obtenue est optimale tandis que pour des champs  $\alpha$ -fortement mélangeant à décroissance exponentielle et sous une condition  $L^{8+\delta}$ ,  $\delta > 0$ , cette vitesse optimale est ralentie par un facteur (log  $\sigma_n$ )<sup>d</sup>.

Egorov [8], Sergin [20] et Maejina [13] utilisent la technique de Bernstein (formation de paquets) et améliorent successivement les vitesses obtenues. Avec l'article de Sergin [21], la vitesse optimale est obtenue. Sur la base d'une autre technique d'évaluation d'espérence conditionnelle et sous une condition  $L^8$ , Stein [22] obtient également cette vitesse optimale dans le cas stationnaire et une vitesse presque optimale pour un processus  $\alpha$ -fortement mélangeant,  $\alpha$ étant à décroissance exponentielle. Bolthausen [4] obtient aussi une vitesse optimale dans le cas de chaines de Markov récurrentes et fortement mélangeantes.

S'inspirant de l'article de Stein, Tikhomirov [24], sur la base d'une évaluation directe de la dérivée de la fonction caractéristique  $f_n(t)$  de la somme étudiée, obtient des résultats analogues pour un processus stationnaire sous la seule condition  $L^{2+\delta}$ ,  $\delta > 0$ . L'hypothèse de stationnairté est levée par Schneider [19]. Comme on va le voir dans le cas spatial, ni la stationnarité, ni la géométrie des domaines sur lesquels on somme X, ne jouent un rôle particulier dans le problème de vitesse de convergence. Dans ce travail, nous suivrons, en la généralisant à  $\mathbb{Z}^d$ . la technique de Tikhomirov.

Dans l'étude des théorèmes limites pour les champs  $(d \ge 2)$  deux notions de mélange apparaissent. L'une où seule la distance entre les deux sous-ensembles sur lesquels on examine les tribus intervient (Deo [5], Gorodetskii [9]). L'autre, introduite par Dobrushin [7] dans le contexte des champs de Gibbs, fait intervenir également les sous-ensembles eux-mêmes, parfois uniquement par l'intermédiaire de leur cardinal (Nahapetian [14], Neaderhauser [15, 16], Bolthausen [3]). C'est la première notion, plus forte, que nous considérons.

Dans le cas *m*-dépendant, la vitesse de convergence du T.L.C. a été étudiée par Leonenko [12] qui obtient le même résultat que Petrov dans le cas  $d \ge 1$  et par Prakasa-Rao [18]. C'est la conjecture annoncée par ce dernier qui a motivé notre travail: sa conjecture est inexacte et la forme du domaine n'intervient pas. Le cas de champs mélangeants est étudié par Neaderhauser [16] et Takahata [23] (voir note (1)). Berkes et Morrow [1] donnent une évaluation de l'approximation des sommes partielles multidimensionnelles par un drap brownien.

#### 2. Notations et resultats

Soit  $X = (X_j)_{j \in \mathbb{Z}^d}$  un champ centré vérifiant (1-1); soit  $(A_n)$  une suite strictement croissante de domaines de  $\mathbb{Z}^d$  à laquelle sont associées les suites  $(S_n)$ ,  $(\sigma_n)$ :

$$S_n = \sum_{A_n} X_j, \qquad \sigma_n^2 = \operatorname{Var} S_n. \tag{2-1}$$

La seule condition que devra vérifier la suite  $(A_n)$ , relative au processus X, est:

$$\liminf \sigma_n^2 \cdot |A_n|^{-1} = \alpha > 0$$
 (2-2)

(|A| est le cardinal de A).

Théorème de la limite centrale pour des champs

Soit  $\phi$  la fonction de répartition de la normale réduite. La distance entre  $S_n$  renormalisée et la normale réduite est évaluée par:

$$\Delta_n = \sup_{x} |P(S_n / \sigma_n \le x) - \phi(x)|.$$
(2-3)

On considérera par la suite la notion de mélange fort suivante: soit A. B deux parties de  $\mathbb{Z}^d$ , d(A, B) leur distance, où la distance entre deux points x. y de  $\mathbb{Z}^d$  est celle du sup:  $d(x, y) = \sup_{1 \le i \le d} |x_i - y_i|$ , et  $\mathcal{F}_A$ ,  $\mathcal{F}_B$  les  $\sigma$ -algèbres engendrées par X sur A. B.

On dira que le champ X est  $\alpha$ -mélangeant si, pour tout couple de parties A, B, on a:

$$\sup_{E\in\mathcal{F}_A, F\in\mathcal{F}_B} |P(E\cap F) - P(E)P(F)| \leq \alpha(d(A, B))$$
(2-4)

où  $\alpha(d) \rightarrow 0$  quand  $d \rightarrow +\infty$ .

On dira que la champ est *m*-dependant si  $\alpha(d)=0$  pour  $d \ge m$ . On a les résultats suivants:

**Théorème 1.** Soit X un champ centré et m-dépendant de  $\mathbb{Z}^d$ .  $(A_n)$  une suite strictement croissante de domaines de  $\mathbb{Z}^d$ . vérifiants (1-1), (2-2). Alors:

a) Si 
$$0 < \delta < 1$$
,  $\Delta_n = O(\sigma_n^{-\delta})$ .  
b) Si  $\delta \ge 1$ ,  $\Delta_n = O((\log \sigma_n)^{\frac{d-1}{2}} \cdot \sigma_n^{-1})$ .

Notant  $a \wedge b$  l'inf de a et de b.

**Théorème 2.** Supposons que X soit  $\alpha$ -mélangeant,  $\alpha$  à décroissance exponentielle; alors:

$$\exists_n = O\left[ (\operatorname{Log} \sigma_n)^{d\left[(1-\delta) \wedge 2\right]} \cdot \sigma_n^{-(\delta \wedge 1)} \right]$$

Si de plus X est dans  $L^{4+\delta}$ ,  $\delta > 0$ , alors:

$$\Delta_n = O\left[(\log \sigma_n)^d \cdot \sigma_n^{-1}\right]$$

**Théorème 3.** Supposons que X soit  $\alpha$ -mélangeant, où  $\alpha$  est à décroissance puissance :  $\alpha(m) = O(m^{-\alpha})$ . Posons :

$$a = \beta \cdot \frac{2d(2+\delta)\left[(1+\delta) \wedge 2\right]}{\delta \cdot (\delta \wedge 1)}.$$

Alors, dès que  $\beta > 1$ , on a le contrôle suivant de  $\Delta_n$ :

$$\Delta_n = O(\sigma_n^{-E(\beta,\delta)}), \quad o\dot{u}$$
$$E(\beta,\delta) = (\delta \wedge 1) \cdot \frac{2(\beta-1)}{2\beta + (\delta \wedge 1)}.$$

#### 3. La décomposition de Tikhomirov

Les résultats énoncés au §2 sont obtenus en évaluant la distance entre  $f_n(t)$ , la fonction caractéristique de  $S_n$ , et exp $\left(-\frac{t^2}{2}\right)$  puis en appliquant l'inégalité de

Esseen (voir [10]). A cette fin, Tikhomirov utilise un développement de  $f'_n(t)$ . Nous employerons les notations suivantes:

$$S_{j}^{(l)} = \frac{1}{\sigma_{n}} \sum_{A_{n} \cap \overline{B(j, lm)}} X_{j}, \ l \ge 1, \qquad S_{j}^{(0)} = \frac{1}{\sigma_{n}} \sum_{j \in A_{n}} X_{j}$$
  
$$\zeta_{j}^{(l)} = e^{it(S_{j}^{(l-1)} - S_{j}^{(l)})} - 1, \ l \ge 1,$$

où B(j, lm) représente la boule de centre j et de rayon lm, et  $\overline{B(j, lm)}$  son complément.

Pour k > 2, on a par induction ([24]):

$$f'_{n}(t) = \frac{i}{\sigma_{n}} \bigg[ \sum_{j \in A_{n}} E(X_{j} e^{itS_{j}^{(1)}}) + \sum_{r=2}^{k} \sum_{j \in A_{n}} E\left(X_{j} \prod_{l=1}^{r-1} \zeta_{j}^{(l)} e^{itS_{j}^{(r)}}\right) \\ + \sum_{i \in A_{n}} E\left(X_{j} \prod_{l=1}^{k} \zeta_{j}^{(l)} e^{itS_{j}^{(k)}}\right) \bigg].$$

En écrivant:

$$E\left(X_{j}\prod_{l=1}^{r-1}\xi_{j}^{(l)}e^{itS_{j}^{(r)}}\right) = E\left(X_{j}\prod_{l=1}^{r-1}\xi_{j}^{(l)}\left\{e^{itS_{j}^{(r)}} - E(e^{itS_{j}^{(r)}})\right\}\right) + E\left(X_{j}\prod_{l=1}^{r-1}\xi_{j}^{(l)}\right)E(e^{itS_{j}^{(r)}})$$

et en ajoutant et soustrayant  $f_n(t)$  de  $E(e^{itS_j^{(r)}})$  dans le 2ème membre, on obtient:

$$f'_{n}(t) = T_{0} + \sum_{r=2}^{k} (T_{1}(r) + T_{2}(r) + T_{3}(r)) + R_{k}$$
(3-1)

avec

$$T_0 = \frac{i}{\sigma_n} \sum_{j \in A_n} E(X_j e^{itS_j^{(1)}})$$

et pour  $2 \leq r \leq k$ , en notant  $a_j(r) = E\left(X_j \prod_{l=1}^{r-1} \xi_j^{(l)}\right)$ 

$$\begin{split} T_{1}(r) &= \frac{i}{\sigma_{n}} \sum_{j \in A_{n}} a_{j}(r) f_{n}(t), \qquad T_{2}(r) = \frac{i}{\sigma_{n}} \sum_{j \in A_{n}} a_{j}(r) \{ E(e^{itS_{j}^{(r)}}) - f_{n}(t) \} \\ T_{3}(r) &= \frac{i}{\sigma_{n}} \sum_{j \in A_{n}} E\left( X_{j} \prod_{l=1}^{r-1} \xi_{j}^{(l)} \{ e^{itS_{j}^{(r)}} - E(e^{itS_{j}^{(r)}}) \} \right) \\ R_{k} &= \frac{i}{\sigma_{n}} \sum_{j \in A_{n}} E\left( X_{j} \prod_{l=1}^{k} \xi_{j}^{(l)} e^{itS_{j}^{(k)}} \right). \end{split}$$

Dans le cas d'un processus *m*-dépendant.  $T_0$  et  $T_3(r)$  sont nuls et la décomposition (3-1) se réduit à:

$$f_n'(t) = \sum_{r=2}^k (T_1(r) + T_2(r)) + R_k.$$
(3-2)

Dans la décomposition (3-1) ou (3-2) de  $f'_n(t)$ ,  $T_1(2)$  contribuera un terme en  $-tf_n(t)$  à l'équation différentielle. La vitesse de convergence sera obtenue en évaluant l'ordre des termes restants. En écrivant l'équation différentielle sous la forme:

$$f'_{n}(t) - tf_{n}(t) + a_{n}(t)f_{n}(t) + h_{n}(t) \text{ avec}$$

$$a_{n}(t) = a_{n,0}\theta_{n,0}(t) + a_{n,1}\theta_{n,1}(t)|t| + a_{n,2}\theta_{n,2}(t)|t|^{(1+\delta)^{2}}, \ \delta > 0 \qquad (3-3)$$

$$h_{n}(t) = b_{n,0}\theta_{n,0}(t) + b_{n,1}\theta_{n,1}(t)|t| + \sum_{\Delta} b_{n,s}\theta_{n,s}(t)|t|^{s}$$

où  $\Delta$  est une partie finie de ]1,  $+\infty$  [;  $a_{n,i}, b_{n,j}$  réels positifs, on utilisera:

**Lemme 3.1.** Supposons que  $f_n(t)$  vérifie (3-3), les fonctions  $\theta$ ,  $\hat{\theta}$  étant continues, uniformément bornées en n, j, t pour  $|t| \leq T_0(n)$ ,  $a_n = \max\{a_{n0}, a_{n1}, a_{n2}\}$  tendant vers 0.

Soit

$$\theta_{n,2} = \sup_{|t| \le T_0(n)} |\theta_{n,2}(t)|, \quad T_1(n) = \inf\left(T_0(n), \frac{1}{4a_{n,2}\theta_{n,2}}\right),$$

alors pour n grand:

$$\Delta_n \leq \inf_{T \leq T_1(n)} \left( \max\left\{ \varepsilon_n(T), \frac{1}{T} \right\} \right)$$

оù

$$\varepsilon_n(T) = \int_{-T}^{T} \left| \frac{f_n(t) - e^{-t^2/2}}{t} \right| dt \leq Ba_n + C \left[ b_{n,0} + b_{n,1} \log T + \sum_{J} b_{n,s} T^{s-1} \right]$$

B. C étant des constantes finies.

#### 4. Vitesse de convergence pour un processus *m*-dependant

#### 4.1. Lemmes fondamentaux

Dans ce qui suit, nous utiliserons à plusieurs reprises les résultats suivants:

$$|e^{ix} - 1| \leq x$$

$$e^{ix} - 1 - ix = \theta_j(x) |x|^{(1+\delta) \wedge 2} \quad \text{où} \quad \theta_j(x) \quad \text{est bornée sur } \mathbb{R} \text{ par } \theta_j.$$

$$(4-1)$$

$$\frac{d-1}{d-1}$$

Lemme 4.1. 
$$|a_j(r)| \le c_d ||x||_2 \left(\frac{|t|b_m(r-1)^{-2}}{\sigma_n}\right)^{r-1} \times (r-1)^{\frac{d-1}{4}} = a(r), r \ge 2.$$

оù

$$c_d = (2\sqrt{2\pi})^d$$
 et  $b_m = (2m)^d \sqrt{d} ||X||_2 e^{\frac{-(d-1)^2}{2}}$ 

Démonstration.  $|a_j(r)| \leq E \left| X_j \prod_{\substack{l=1\\l \text{ pair}}}^{r-1} \zeta_j^{(l)} \prod_{\substack{l=1\\l \text{ impair}}}^{r-1} \zeta_j^{(l)} \right|.$ 

Comme les  $\{\xi_j^{(l)}, l \text{ pair (resp. impair)}\}$  sont indépendants entre eux, en appliquant l'inégalité de Hölder, on obtient:

$$|a_{j}(r)| \leq ||X||_{2} \prod_{l=1}^{r-1} ||\xi_{j}^{(l)}||_{2}.$$

Soit  $C_{n,l}$  l'intersection de  $A_n$  et de la couronne de centre j et de rayons [(l - 1)m, lm].

En utilisant (4-1):

$$\|\xi_{j}\|_{2}^{2} < \frac{t^{2}}{\sigma_{n}^{2}} E((\sum_{C_{n,l}} X_{j})^{2}),$$
$$E(\sum_{C_{n,l}} X_{j})^{2} \leq (2m)^{d} \times 2dm(2lm)^{d-1} \|X\|_{2}^{2}$$

à cause de la m-dépendance. On obtient donc:

$$|a_{j}(r)| \leq ||X||_{2} \left( \frac{|t|}{\sigma_{n}} (2m)^{d} \sqrt{d} ||X||_{2} \right)^{r-1} \left( (r-1)! \right)^{\frac{d-1}{2}}$$

et la démonstration est achevée en utilisant l'inégalité de Sterling. Lemme 4.2.

$$|T_{2}(r)| \leq a(r) \frac{|t|}{\sigma_{n}^{2}} \{ |A_{n}|^{\frac{1}{2}} r^{d} B_{1}(m, ||X||_{2}) + |A_{n}| r^{\frac{d}{2}} B_{2}(m, ||X||_{2}) |f_{n}(t)| \}$$

avec

$$B_1(m, \|X\|_2) = 3^{d/2} (2m)^{3d/2} \|X\|_2$$

et

$$B_2(m, ||X||_2) = (2m)^d ||X||_2.$$

Démonstration. Posons:

$$\begin{split} \eta_{j}(r) &= e^{it(S_{j}^{(r)} - S_{j}^{(0)})} - 1, \quad r \ge 2 \\ |T_{2}(r)| &\leq \frac{1}{\sigma_{n}} \left\{ E |\sum_{j \in A_{n}} a_{j}(r) [\eta_{j}(r) - E(\eta_{j}(r))]| + |f_{n}(t)| |\sum_{j \in A_{n}} a_{j}(r) E(\eta_{j}(r))| \right\} \\ E |\sum_{j \in A_{n}} a_{j}(r) [\eta_{j}(r) - E(\eta_{j}(r))]| &\leq a(r) \cdot \sum_{A_{n} \ge A_{n}} \cos(\eta_{j}(r), \eta_{p}(r))^{1/2} \end{split}$$

 $\operatorname{cov}(\eta_j^{(r)}, \eta_p^{(r)}) = 0$  dès que d(j, p) > 3rm. Dans le cas contraire, en utilisant (4-1) et la *m*-dépendance, on obtient:

$$\operatorname{cov}(\eta_{j}^{(r)},\eta_{p}^{(r)}) \leq \frac{t^{2}}{\sigma_{n}^{2}}(2m)^{d}(2rm)^{d} ||X||_{2}^{2}.$$

et de manière analogue:

$$|E(\eta_{j}(r))| \leq \frac{|t|}{\sigma_{n}} (2m)^{d/2} (2rm)^{d/2} ||X||_{2}$$

ce qui permet de conclure la démonstration.

Théorème de la limite centrale pour des champs

#### Lemme 4.3.

$$T_1(2) = f_n(t)(-t + |t|^{(1+\delta)/2} \frac{|A_n|}{\sigma_n^{(2+\delta)/3}} \theta_{\delta}(t)(2m)^{d[(1+\delta)/2]} \|X\|_{2+\delta}^{(2+\delta)/3}$$

où  $\theta_{\delta}(t)$  est une fonction bornée sur **R**.

Démonstration. Utilisant (4-1):

$$E(X_{j}\xi_{j}^{(1)}) = E[X_{j}\{it(S_{j}^{(0)} - S_{j}^{(1)}) + \theta_{\delta}(t(S_{j}^{(0)} - S_{j}^{(1)}))|t|^{(1+\delta)^{2}}|S_{j}^{(0)} - S_{j}^{(1)}|^{(1+\delta)^{2}}\}]$$

Notons:

$$\sum_{j \in A_n} E(X_j S_j^{(0)}) = \sigma_n,$$
  
$$\sum_{j \in A_n} E(X_j S_j^{(1)}) = 0 \quad \text{à cause de la } m\text{-dépendance.}$$

La démonstration est conclue en évaluant:

$$|E(X_{j}|S_{j}^{(0)} - S_{j}^{(1)}|^{(1+\delta)\wedge 2})| \leq |X|_{2+\delta} E[|S_{j}^{(0)} - S_{j}^{(1)}|^{[(1+\delta)\wedge 2](\frac{2+\delta}{1+\delta})]^{\frac{1+\delta}{2+\delta}} \\ \leq |X|_{2+\delta} \sigma_{n}^{-[(1+\delta)\wedge 2]}(2m)^{d[(1+\delta)\wedge 2]} ||X||_{2+\delta}^{(1+\delta)\wedge 2}.$$

#### 4.2. La décomposition de $f_n^{(i)}(t)$

On utilisera la décomposition (3-2) et l'on se place sous l'hypothèse (2-2).

Examen de  $T_1(2)$ . D'après le lemme 4-3.  $T_1(2)$  contribuera une expression de la forme:

$$-tf_{n}(t) + \frac{(2m)^{d[(1+\delta)\wedge 2]}}{\sigma_{n}^{\delta\wedge 1}} \theta_{n,2}^{*}|t|^{(1+\delta)\wedge 2} f_{n}(t)$$

à l'équation différentielle. La vitesse de convergence sera donc limitée (Lemme 3.1) par  $a_{n,2}^* = \sigma_n^{-(\delta \wedge 1)}$  avec une constante en *m* égale à  $(2m)^{d[(1-\delta) \wedge 2]}$ . *Examen de*  $R_k$ . Soit  $\varepsilon > 0$ .  $k_n$  le premier entier supérieur à  $\varepsilon \log \sigma_n$  et  $T_0(n)$  $= \frac{\sigma_n}{eb_m(\varepsilon \log \sigma_n)^{(d-1)/2}}$ . Alors pour  $k = k_n - 1$ .  $|t| \leq T_0(n)$ , utilisant le Lemme 4.1.:

$$|R_k| \leq \frac{|A_n|}{\sigma_n} c_d ||X||_2 \left(\frac{1}{e}\right)^k k^{\frac{d-1}{4}}.$$

Choisissons maintenant  $\varepsilon$  assez grand pour que:

$$\left(\frac{1}{\varepsilon}\right)^{\varepsilon \log \sigma_n - 1} \times (\varepsilon \log \sigma_n)^{\frac{d-1}{4}} < \sigma_n^{-2}$$

alors  $R_k = \frac{1}{\sigma_n} \hat{\theta}_{n,0}(t)$  avec  $|\hat{\theta}_{n,0}(t)| < c_d ||X||_2$ .

Examen des termes  $\sum_{3}^{k} T_{1}(r)$ ,  $k = k_{n} - 1$ . En appliquant à nouveau le Lemme 4.1 et en faisant sortir du crochet une puissance  $(1 + \delta) \wedge 2$ , on obtient pour  $|t| \leq T_{0}(n)$ :

$$\begin{split} \left| \sum_{r=3}^{k} \frac{i}{\sigma_{n}} \sum_{j \in A_{n}} E\left( X_{j} \prod_{l=1}^{r-1} \xi_{j}^{(l)} \right) \right| &\leq \frac{|A_{n}|}{\sigma_{n}} \left( \frac{|t| b_{m}}{\sigma_{n}} \right)^{(1+\delta) \wedge 2} \\ & \cdot \left( \sum_{r=3}^{\infty} \left( \frac{1}{e} \right)^{r - [(2+\delta) \wedge 3]} \cdot (r-1)^{\left(\frac{d-1}{4}\right)[(3+2\delta) \wedge 5]} \right) \\ & \Rightarrow \sum_{3}^{r} T_{1}(r) = \frac{|t|^{(1+\delta) \wedge 2}}{\sigma_{n}^{\delta \wedge 1}} \theta_{n,2}(t) f_{n}(t), \\ avec |\theta_{n,2}(t)| \leq b_{m}^{(1+\delta) \wedge 2} \cdot e^{-3} \cdot \Gamma\left( \frac{5d-1}{4} \right). \end{split}$$

Examen des termes  $\sum_{2}^{k} T_2(r)$ .  $k = k_n - 1$ . D'après le Lemme 4.2, ces termes font intervenir deux expressions:

(a) 
$$L_1 = \frac{|t|}{\sigma_n^2} |A_n|^{1/2} B_1(m, ||X||_2) \sum_{r=2}^k a(r) r^d$$
  
(b)  $L_2 = \frac{|t|}{\sigma_n^2} |A_n| B_2(m, ||X||_2) |f_n(t)| \sum_{r=2}^k a(r) r^d$ 

L'évaluation de  $L_1$  où  $L_2$  se fait de manière semblable, en appliquant le Lemme 4.1 et en faisant sortit du crochet une puissance  $(\delta \wedge 1)$ .

On obtient pour  $|t| \leq T_0(n)$ :

$$L_1 \leq \frac{|t|^{(1-\delta) \wedge 2}}{\sigma_n^{(1-\delta) \wedge 2}} \theta_2(t)$$

avec  $|\theta_{2}(t)| \leq B_{1}(m, ||X||_{2}) c_{d} ||X||_{2} b_{m}^{\delta \wedge 1} e^{-3} \Gamma\left(\frac{7d+1}{4}\right)$ 

$$L_2 \leq \frac{|t|^{(1+\delta) \wedge 2}}{\sigma_n^{\delta \wedge 1}}$$

avec  $|\theta_{3}(t)| \leq B_{2}(m, ||X||_{2}) \times c_{d} ||X||_{2} b_{m}^{\delta \wedge 1} e^{-3} \Gamma\left(\frac{5d+1}{4}\right).$ 

Par conséquence, pour  $|t| \leq T_0(n)$ :

$$\sum_{r=2}^{k} T_2(r) = \frac{|t|^{(1+\delta)/2}}{\sigma_n^{\delta/1}} \theta'_{n,2}(t) f_n(t) + \frac{|t|^{(1+\delta)/2}}{\sigma_n^{(1+\delta)/2}} \hat{\theta}_{n,2}(t)$$

avec  $|\theta'_{n,2}(t)| \leq |\theta_3(t)|$  et  $|\hat{\theta}_{n,2}(t)| \leq \theta_2(t)$ .

Théorème de la limite centrale pour des champs

#### 4.3. La vitesse de convergence

On déduit du §4.1 que pour  $|t| \leq T_0(n)$ ,  $f'_n(t)$  satisfait à l'équation suivante:

$$f'_{n}(t) = -tf_{n}(t) + a_{n,2}|t|^{(1+\delta)\wedge 2} \theta_{n,2}(t)f_{n}(t) + b_{n,0}\hat{\theta}_{n,0}(t) + b_{n,2}(t)|t|^{(1+\delta)\wedge 2}$$

où  $a_{n,2} = \sigma_n^{-(\delta \wedge 1)}, b_{n,0} \le \sigma_n^{-1}$  et  $b_{n,2} = \sigma_n^{-[(1+\delta) \wedge 2]}$ . Nous devons maintenant distinguer deux cas:

(a) 
$$\delta < 1: T_1(n) = \inf (T_0(n), a_{n,2}^{-1}) = a_{n,2}^{-1}$$
  
 $\varepsilon_n(T) \le \frac{B}{\sigma_n^{\delta}} + C \left[ \frac{1}{\sigma_n} + \frac{T^{\delta}}{\sigma_n^{(1+\delta) \wedge 2}} \right], \quad 0 \le T \le T_1(n)$ 

et la vitesse de convergence est donc en  $\sigma_n^{-\delta}$ .

(b) 
$$\delta \ge 1$$
:  $T_1(n) = \inf (T_0(n), a_{n,2}^{-1}) = T_0(n) = \frac{\sigma_n}{e b_m (\varepsilon \log \sigma_n)^{(d-1)/2}}$   
 $\varepsilon_n(T) \le \frac{B}{\sigma_n} + C \left[ \frac{1}{\sigma_n} + \frac{T}{\sigma_n^2} \right], \quad 0 \le T \le T_1(n).$ 

La vitesse de convergence est donc limitée par  $\frac{1}{T_1(n)}$  ce qui conclut la démonstration du Théorème 1 (b).

Remarque. Pour  $\delta < 1$ . la vitesse optimale est atteinte. Le seul terme restrictif est le coefficient de  $|t|^{1+\delta}$  dans  $T_1(2)$ . Dans ce cas il n'est pas nécessaire de pousser le développement de  $f'_n(t)$  jusqu'à un ordre  $k = k_n - 1$  non borné. Il suffirait de choisir  $k = \left[\frac{1}{1-\delta}\right] + 1$  pour obtenir cette vitesse.

#### 5. Cas d'un champ x-mélangeant: démonstration des théorèmes 2 et 3

Si X est  $\alpha$ -mélangeant, nous avons les inégalités suivantes (cf. Hall et Heyde [11]).

$$|Y_1| \leq C_1, |Y_2| \leq C_2: |E(Y_1|Y_2) - E(Y_1)E(Y_2)| \leq 4C_1 C_2 \alpha(Y_1, Y_2)$$
(\alpha)

 $(\alpha(Y_1, Y_2)$  est à prendre en un sens naturel).

Cett inégalité se généralise sans difficultés à:

$$|E(Y_1 | Y_2 ... Y_r)| \leq \prod_{1}^{r} |E(Y_k)| + 4(r-1) C_1 ... C_r \alpha(m)$$

dès que  $|Y_k| \leq C_k$ , k = 1, r, les Y étant deux à deux "distantes" de m.

$$|X| \leq C, \ Y \in L^{p}, \ p > 1:$$
  
$$|E(XY) - E(X) E(Y)| \leq 6 C ||Y||_{p} \alpha(X, Y)^{1 - \frac{1}{p}}$$
(\beta)

$$X \in L^{p}, Y \in L^{q}, \frac{1}{p} + \frac{1}{q} < 1:$$

$$|E(XY) - E(X)E(Y)| \leq 8 ||X||_{p} ||Y||_{q} \alpha(X, Y)^{1 - \frac{1}{p} - \frac{1}{q}}$$
(7)

Soit la condition:

$$C^{2} = \sum_{0}^{\infty} v^{d-1} \alpha(v)^{\frac{\delta}{2+\delta}} < \infty$$
 (C0)

On a le résultat préliminaire:

Lemme 5.0. Sous (C0).  $E(\sum_{A} X_{j})^{2} \leq 16 \cdot C^{2} \cdot |A| ||X||_{2+\delta}^{2}$ .

Démonstration. Utilisant (?) avec  $p = q = 2 + \delta$ , on obtient:

$$E(\sum_{A} X_{j})^{2} \leq 8 \|X\|_{2+\delta}^{2} |A| \cdot \sum_{A} \alpha(j)^{\frac{\delta}{2+\delta}}.$$

Le résultat s'obtient en constatant que le nombre de j à une distance v de l'origine est majoré par  $2dv^{d-1}$ .

#### 5.1. Les trois lemmes fondamentaux

Soit *m* un entier positif.

Lemme 5.1. Supposons que t vérifie:  $S = 4d \frac{|t|}{\sigma_n} (2m)^d C ||X||_{2+\delta} \le 1$ . Alors:

$$|a_{j}(r)| \leq ||X||_{2+\delta} \left\{ 6r \cdot 2^{r} \alpha(m)^{\frac{\delta}{2+\delta}} + C_{1}(r-1)^{\frac{5(d-1)}{4}} \\ \cdot \left[ C_{2} \frac{|t|}{\sigma_{n}} m^{\frac{d}{2}} (r-1)^{\frac{d-1}{2}} \right] \left( cr \right) \right\} = a(r)$$

avec: e(r) = 1 si r = 2: e(r) = 2p si r = 2p + 1 ou 2p + 2.  $p \ge 1$ 

$$C_1 = (2\sqrt{2\pi})^{\frac{d-1}{2}}, \quad C_2 = 4d \, 2^{\frac{d}{2}} \cdot C \cdot e^{\frac{1-d}{2}} \|X\|_{2+\delta}.$$

Démonstration. Cas. r-1=1: Il suffit d'utiliser l'inégalité de Hölder habituelle pour le produit  $X_j \cdot \xi_j^{(1)}$  et le Lemme 5.1.1.

Cas,  $r-1 \ge 2$ : Des  $\xi_j^{(l)}$  apparaissent pour *l* pair et *l* impair. Notons  $\Pi'$  le produit relatif aux *l* impair.  $\Pi''$  relatif aux *l* pair. Par Hölder;

$$|a_{j}(r)| \leq \|\Pi' \xi_{2}^{(\beta)}\|_{2} \|X_{j}\Pi'' \xi_{j}^{(\beta)}\|_{2}.$$
(5-1)  
Utilisant ( $\alpha$ ) pour  $\Pi'$ . ( $\beta$ ) avec  $p = \frac{2+\delta}{2}$  pour  $X_{j}\Pi''$ . on a:

Théorème de la limite centrale pour des champs

$$E|\Pi' \xi_j^{(l)}|^2 \leq r \cdot 2^{r+1} \alpha(m) + \Pi' E(\xi_j^{(l)})^2$$
$$E|X_j \Pi'' \xi_j^{(l)}|^2 \leq 6 ||X||_{2+\delta}^2 \cdot 2^{r-1} \alpha(m)^{\frac{\delta}{2+\delta}} + ||X||_{2+\delta} E(\Pi'' \xi_j^{(l)})^2.$$

L'inégalité ( $\alpha$ ) peut à nouveau être utilisée pour évaluer  $E(\Pi'')^2$ . Réunissant ces majorations dans (5-1), on obtient:

$$|a_{j}(r)| \leq ||X||_{2+\delta}(u+v)$$
  

$$u = 6r 2^{r} \alpha(m)^{\delta/2+\delta}, \quad v' = \Pi' E(\xi_{j}^{(l)})^{2}, \quad v'' = \Pi'' E(\xi_{j}^{(l)})^{2}$$
(5-2)  
et  $v = \max\{v', v''\}.$ 

Le leme 5.0. permet d'évaluer v', v''

$$v' \leq (\Pi' l)^{\frac{d-1}{2}} S^2 [\frac{r}{2}], \quad v'' \leq (\Pi'' l)^{\frac{d-1}{2}} S^2 [\frac{r-1}{2}]$$

([a] est la partie entière de a); utilisant la majoration de Sterling d'un factoriel et  $|S| \leq 1$ , on obtient:

$$v \leq (r-1)^{\frac{3(d-1)}{4}} \cdot C_1 \cdot \left(\frac{r-1}{e}\right)^{\frac{(r-1)(d-1)}{2}} S^{2\left[\frac{r-1}{2}\right]}.$$

L'examen r=2p+1. r=2p+2 pour  $p \ge 1$  conduit directement au résultat anoncé.

**Lemme 5.2.** Majoration du terme  $T_2(r)$ . Dès que  $m \ge 2^{d-2} \cdot C^2$ , on a:

$$|T_{2}(r)| \leq a(r) \frac{|t|}{\sigma_{n}^{2}} \{ |f_{n}(t)| C_{3} |A_{n}| (rm)^{\frac{d}{2}} + C_{4} |A_{n}|^{\frac{1}{2}} (rm)^{\frac{3d}{2}} \}$$

avec:  $C_3 = 4.2^{d-2} \sqrt{d} \|X\|_{2+\delta} \cdot C$ :  $C_4 = 2^{2(d-1)} \|X\|_{2+\delta}$ .

Démonstration. La démonstration suit la même démarche que celle du Lemme 4.2. dont nous reprenons les notations.

Dans un premier temps, il faut majorer  $\sum_{j,p} \operatorname{cov}(\eta_j(r), \eta_p(r))$ . Utilisant l'inégalité (7) avec  $p = q = 2 + \delta$ , on obtient:

$$\left|\sum_{A_n \land A_n} \operatorname{cov}(\eta_j(r), \eta_p(r))\right| \leq 8|A_n| \|\eta(r)\|_{2+\delta}^2 \sum_{A_n} \alpha(d(B(0, r m), B(j, r m)))^{\frac{\delta}{2+\delta}}$$

où  $\|\eta(r)\|_{2+\delta}$  majore les  $\|\eta_j(r)\|_{2+\delta}$ .

- Si  $d(0, j) \leq 2rm$ , la distance entre les deux boules est nulle. Il y a  $(4rm)^d$  tels *j*. - Si d(0, j) = 2rm + v, v > 0, la distance entre les deux boules est v. Il y a  $(4rm + 2v)^{d-1}$  tels *j*. Une suite de majorations élémentaires permet d'obtenir, sous

(CO):  

$$\sum_{j \in A_n} \alpha(d(B(0, r m), B(j, r m))) \leq 2^{2d+1} (r m)^d$$

dès que  $m \ge 2^{d-2} \cdot C^2$ . D'autre part:

$$\|\eta(r)\|_{2+\delta} \leq \frac{|t|}{\sigma_n} (2rm)^d \|X\|_{2+\delta}.$$

Ceci conduit à:

$$E|\sum_{A_n} a_j(r) (\eta_j(r) - E \eta_j(r))| \leq 2^{2(d+1)} |A_n|^{\frac{1}{2}} \frac{|t|}{\sigma_n} a(r) (rm)^{\frac{3d}{2}} ||X||_{2+\delta}$$

D'autre part, on a:

$$\sum_{A_n} a_j(r) E(\eta_j(r)) | \leq a(r) |A_n| ||\eta(r)||_2.$$

Le Lemme 5.0. permet alors de conclure.

**Lemme 5.3.** Evaluation de  $T_1(2)$ 

$$T_{1}(2) = \left(-t + \theta_{1}(t) |A_{n}| \alpha(m)^{\frac{\delta}{2-\delta}} \cdot t + \theta_{2}(t) \frac{m^{d[(1+\delta) \wedge 2]}}{\sigma_{n}^{\delta}} |t|^{(1+\delta) \wedge 2}\right) f_{n}(t)$$

оù

$$\begin{aligned} \|\theta_1(t)\| &\leq 8 \, \alpha \, \|X\|_{2+\delta} \\ \|\theta_2(t)\| &\leq 2^{d[(1+\delta)\wedge 2]} \, \alpha \, \theta_\delta \, \|X\|_{(2+\delta)\wedge 3}^{(2+\delta)\wedge 3} \quad (\alpha, \theta_\delta \text{ définis en } (2.2) \text{ et } (4.1)). \end{aligned}$$

*Démonstration*. La seule modification par rapport au Lemme 4.3 est l'évaluation de  $E(X_j S_j^{(1)})$ . Utilisant (7) avec  $p=q=2+\delta$  et l'inégalité de Minkowski pour  $S_j^{(1)}|_{2+\delta}$ , on obtient le lemme.

Par la suite les constantes ou les majorants de certaines fonctions de t ne seront pas explicités: de telles constantes notées  $\theta$ ... seront universelles une fois fixées  $\delta$ , d, X et  $(A_n)$ .

#### 5.2. La décomposition de $f'_n(t)$

On part de la décomposition (3.1) de  $f'_n(t)$ .

Majoration de 
$$T_0 + \sum_{2}^{k} T_3(r)$$
: par application de  $(\beta)$   
 $|T_0| \le 6 \propto \sigma_r ||X|_{2-\delta} \propto (m)^{\frac{\delta+1}{\delta+2}}$ 

$$|T_0| \leq 6 \alpha \cdot \sigma_n ||X||_{2+\delta} \alpha(m)^{\delta+2}$$
  
$$|T_3(r)| \leq 6 \alpha \cdot \sigma_n \cdot 2^r ||X||_{2+\delta} \alpha(m)^{\delta+2}$$

et donc:

$$\left|T_{0} + \sum_{2}^{k} T_{3}(r)\right| \leq \theta(t) \,\alpha(m)^{\frac{\delta+1}{\delta+2}} \cdot 2^{k}$$
(5-3)

Examen de  $T_1(2)$ : (Lemme 5.3)

$$T_{1}(2) = f_{n}(t) \left\{ -t + \theta_{1}(t) \tilde{a}_{n1} t + \theta_{2}(t) a_{n2} |t|^{(1+\delta) \wedge 2} \right\}$$
  
où:  $\tilde{a}_{n,1} = |A_{n}| \alpha(m)^{\frac{\delta}{\delta - 2}}, \quad a_{n2} = m^{d[(1+\delta) \wedge 2]} \sigma_{n}^{\delta \wedge 1}$  (5-4)

Théorème de la limite centrale pour des champs

Examen de  $R_k$ : Soient  $\varepsilon > 0$ , K > 0. Choisissons:

$$k = k_n = [\varepsilon \log \sigma_n] + 1$$
  
$$|t| \le T_0(n) = \sigma_n / (e^K C_2 m^{\frac{d}{2}} (\varepsilon \log \sigma_n)^{\frac{d-1}{2}})$$

Si:

$$0 \le k \le k_n, \left[ C_2 \frac{|t|}{\sigma_n} m^{\frac{d}{2}} (k-1)^{\frac{d-1}{2}} \right] \le e^{-K}$$

En

$$k = k_n : k_n^{\frac{5(d-1)}{4}} \left[ C_2 \frac{|t|}{\sigma_n} m^{\frac{d}{2}} (k_n - 1)^{\frac{d-1}{2}} \right]^{k_n - 1} \leq (\varepsilon \log \sigma_n)^{\frac{5(d-1)}{4}} \sigma_n^{-K\varepsilon}$$

A  $\varepsilon > 0$  fixé, choisir  $K(\varepsilon)$  tel que ce second membre soit inférieur à  $\sigma_n^{-2}$ . On obtient par application du Lemme 5.1:

$$R_{k_n} = \theta_3(t) k_n 2^{k_n} \sigma_n \alpha(m)^{\frac{\delta}{\delta+2}} + \theta_4(t) \sigma_n^{-1}$$
Examen de  $\sum_{3}^{k_n} T_1(r)$ :  $\left| \sum_{3}^{k_n} T_1(r) \right| \le \theta \cdot \sigma_n f_n(t) \sum_{3}^{k_n} a(r).$ 
(5-5)

Dans la majoration de a(r) deux termes apparaissent: le terme en  $\alpha(m)$  et le terme complémentaire: de ce second terme, sortons du crochet une puissance  $(1 + \delta) \wedge 2$  (possible puisque  $r \ge 3$ ). En constatant que la série:

$$\sum_{r} r^{a} e^{-K(\varepsilon) e(r)} < \infty$$

on en déduit:

$$\sum_{3}^{k_{n}} T_{1}(r) = f_{n}(t) \left\{ \theta_{5}(t) \sigma_{n} k_{n} 2^{k_{n}} \alpha(m)^{\frac{\delta}{2+\delta}} + \theta_{6}(t) \frac{m^{\frac{a}{2}[(1+\delta)\wedge2]}}{\sigma_{n}^{\delta\wedge1}} |t|^{(1+\delta)\wedge2} \right\}$$
(5-6)

Examen de  $\sum_{2}^{\kappa_n} T_2(r)$ : Utilisant le Lemme 5.2, sortant cette fois-ci du crochet une puissance  $\delta \wedge 1$  (possible puisque  $r \ge 2$ ), on obtient de façon analogue, une première quantité facteur de  $f_n(t)$ , une autre non facteur de  $f_n(t)$ :

$$f_{n}(t)\left\{\theta_{7}(t) m^{\frac{d}{2}} \alpha(m)^{\frac{\delta}{\delta+2}} k_{n}^{\frac{d}{2}+1} 2^{k_{n}}|t| + \theta_{8}(t) \frac{m^{\frac{d}{2}[(1+\delta)\wedge2]}}{\sigma_{n}^{\delta\wedge1}}|t|^{(1+\delta)\wedge2}\right\}$$
(5-7)

$$\frac{|t|}{\sigma_n} m^{\frac{3d}{2}} \{\theta_9(t) + \theta_{10}(t) k_n^{1 + \frac{3d}{2}} 2^{k_n} \alpha(m)^{\frac{\delta}{\delta + 2}}\}$$
(5-8)

Reportant les évaluations (5-3) à (5-8) dans le Lemme 3.1 d'intégration, on obtient:

$$a_{n0} = k_n 2^{k_n} \alpha(m)^{\frac{\delta}{2+\delta}},$$

$$a_{n1} = \max\{\tilde{a}_{n1}, m^{\frac{d}{2}} k_n^{\frac{d}{2}+1} 2^{k_n} \alpha(m)^{\frac{\delta}{2+\delta}}\}, \quad \tilde{a}_{n1} = |A_n| \alpha(m)^{\frac{\delta}{2-\delta}},$$

$$a_{n2} = m^{d[(1+\delta) \wedge 2]} / \sigma_n^{\delta \wedge 1},$$

$$b_{n0} = \max\{\sigma_n^{-1}, \sigma_n k_n 2^{k_n} \alpha(m)^{\frac{\delta}{2+\delta}}\},$$

$$b_{n1} = \frac{m^{\frac{3d}{2}}}{\sigma_n} \max\{1, k_n^{\frac{3d}{2}+1} 2^{k_n} \alpha(m)^{\frac{\delta}{\delta+2}}\}.$$
(5-9)

On va choisir  $\varepsilon$  et m = m(n) tels que:

 $\tilde{a}_{n1}$  est de l'ordre de  $a_{n2}$ ,  $a_{n1}$  est de l'ordre de  $\tilde{a}_{n1}$ ,  $a_{n0}$  et  $b_{n0}$  sont d'un ordre inférieur à  $a_{n2}$ ,  $b_{n1}$  est de l'ordre de  $m^{\frac{3d}{2}} \cdot \sigma_n^{-1}$ .

Alors:

$$\Delta_n = O\left( \max\left\{ a_{n2}, T_1(n)^{-1}, \frac{\log T_1(n)}{\sigma_n} m^{\frac{3d}{2}} \right\} \right)$$

оù

$$T_1(n) = \inf \{ T_0(n), a_{n2}^{-1} \}$$

Dès que  $a_{n2}$  tend vers 0 les cinq conditions ci-dessus se réduisent aux deux premières, soit, en notant  $\tilde{\varepsilon} = \varepsilon \log 2$ :

$$\sigma_n^2 \chi(m)^{\frac{\delta}{2+\delta}} \leq \frac{m^{d[(1+\delta)\wedge 2]}}{\sigma_n^{\delta\wedge 1}}$$
(C1)

$$m^{\frac{d}{2}} (\log \sigma_n)^{1+\frac{d}{2}} \sigma_n^{\widetilde{\epsilon}} \chi(m)^{\frac{\delta}{\delta+2}} \leq a_{n2}$$
(C2)

#### 5.3. Cas d'un mélange exponentiel

Le choix  $m(n) = a \log \sigma_n$ , avec a assez grand et  $\tilde{\epsilon} < 1$  assure (CO) et (C1)-(C2) pour tout  $m \ge m(n)$ . Pour ce choix:

$$T_0(n) = C \cdot \frac{\sigma_n}{(\log \sigma_n)^{d-\frac{1}{2}}}, \quad T_1(n) = a_{n2}^{-1}$$

La vitesse est en  $(\log \sigma_n)^{d[(1+\delta) \wedge 2]} / \sigma_n^{\delta \wedge 1}$ .

*Remarques.* i) Dans le cas d=1, X stationnaire,  $A_n = [1, n]$  et  $\delta > 1$ . le résultat peut être amélioré. Le Théorème 1 de [25] est applicable ( $\alpha$  est à décroissance exponentielle)

$$\|\sum_{A} X_{j}\|_{r} \le K |A|^{1/2} \quad \text{si } 2 \le r < 2 + \delta.$$
(5-10)

L'évaluation de  $a_{n2}$  dans le Lemme 5.3 est ainsi modifiée:

$$E(|X_{j}|) \sum_{A} |X_{k}|^{(1+\delta) \wedge 2} \leq ||X||_{2+\delta} ||\sum_{A} ||X_{k}||_{r})^{(1+\delta) \wedge 2}$$

avec  $r = [(1+\delta) \wedge 2] \frac{2+\delta}{1+\delta}$ . Si  $\delta > 1$ , alors  $2 < r < 2+\delta$ . On en déduit que  $a_{n2} = m \cdot \sigma_n^{-1}$ . Donc si X est dans  $L^{3+\gamma}$ ,  $\gamma > 0$ , la vitesse est en  $(\log \sigma_n)/\sigma_n$ .

ii) Dans le cas d=1, et sans stationarité, mais sous une condition  $L^{4+\gamma}\gamma > 0$ , on obtient la même vitesse en utilisant l'inégalité (5-10) démontrée en [6] pour les entiers pairs  $r \ge 4$ . Quand  $d \ge 2$  et si X est dans  $L^{4+\gamma}, \gamma > 0$ , la vitesse de convergence est en  $(\log \sigma_n)^d / \sigma_n$ . En effet, l'inégalité (5-10) est une conséquence du lemme suivant, qui présente un intérêt propre:

**Lemme 5.4.** Soit q un entier supérieur ou égal à 2, X dans  $L^{q+\delta}$ ,  $\delta > 0$ ,  $\alpha$ -mélangeant, satisfaisant:

$$C = \sum_{k \ge 0} k^{d(q - \left\lfloor \frac{q}{2} \right\rfloor) - 1} \alpha(k)^{\frac{\delta}{q + \delta}} < \infty.$$

Alors, pour toute partie A de  $\mathbb{Z}^d$ , on a:

$$E(\sum_{A} X_{j})^{q} \leq K \cdot |A|^{\left\lfloor \frac{q}{2} \right\rfloor}$$

où K est une constante finie ne dépendant que de q.d.X.

Démonstration.  $E(\sum_{A} X_{j})^{q} = \sum_{A^{q}} E(X_{j1}, X_{j2}, \dots, X_{jq}).$ 

Notons  $J = \{j_1, j_2, \dots, j_q\}$ ,  $D = D(J) = \max_{i=1,q} \{d(\{j_i\}, J \setminus \{j_i\})\}$ . X étant centré, on a, utilisant (7):

$$|E(X_{j_1}X_{j_2}\ldots X_{j_q})| \leq 8\,\alpha(k)^{\overline{q+\delta}} \, \|X\|_{q+\delta}^q$$

Reste à dénombrer l'ensemble des J de  $A^q$  tels que D(J) = k. Ce dénombrement résultera de:

**Lemme 5.5.** Il existe un recouvrement de J constitué d'au plus  $\begin{bmatrix} \frac{q}{2} \end{bmatrix}$  boules, centrées en des points de J. de rayon  $3^q D(J)$ , l'un au moins des points de J étant exactement à distance D(J) d'un des centres.

Admettons provisoirement ce lemme. Le nombre de J tels que D(J) = k est donc majoré par:  $q! |A|^{[q/2]} (2^d d \cdot k^{d-1}) (2 \cdot 3^q k)^{d(q-[q/2]-1)}$  et donc:

$$|E(\sum_{A} X_{j})^{q}| \leq 8q! ||X||_{q+\delta}^{q} \cdot 2^{d} \cdot d \cdot (2 \cdot 3^{q})^{d(q - \left\lfloor\frac{q}{2}\right\rfloor - 1)} \sum_{k \geq 0} k^{d(q - \left\lfloor\frac{q}{2}\right\rfloor) - 1} \alpha(k)^{\frac{\delta}{q+\delta}}.$$

Démontrons le lemme auxiliaire. Il est facile de le vérifier pour q=2,3. Supposons que  $j_1$  réalise D(J),  $j_2$  étant le point le plus proche de  $j_1$ . Deux cas se présentent: *ler cas:*  $B(j_2, D(J)) \cap \{j_3, \dots, j_q\} = \emptyset$ . Alors  $D(j_3, j_4, \dots, j_q)$  est inférieur ou égal à D(J), et donc avec  $\left[\frac{q-2}{2}\right]$  boules, de rayon  $3^{q-2}D(J)$ , centrées en certains points de  $j_3, \dots, j_q$ , on recouvre  $j_3, \dots, j_q$ . Adjoindre alors la boule  $B(j_2, D(J))$  et augmenter les rayons donne le résultat annoncé à l'ordre q.

2ème cas: Supposons que  $j_3$  appartienne à  $B(j_2, D(J))$ . Montrons alors que  $D(j_2, j_3, ..., j_q) \leq 2D(J)$ . Si cela n'est pas vrai, alors  $D' = D(j_2, ..., j_q) > 2D(J)$ . Soit j' le j qui réalise D'. Alors j'  $\pm j_2$  puisque  $d(j_2, j_3) \leq D(J)$ . Donc:

$$d(j', j_1) \ge D' - D(J) > D(J)$$

et donc:

$$d(\{j'\}, J \setminus \{j'\}) > D(J)$$

ce qui contredit l'hypothése que D(J) est réalisé en  $j_1$ .

On peut donc recouvrir  $j_2, ..., j_q$  avec  $\left[\frac{q-1}{2}\right]$  boules centrées en certains de ces j, de rayon  $3^{q-1} \times 2D(J)$ . Il suffit d'augmenter de D(J) le rayon de la boule contenant  $j_2$  pour englober  $j_1$ . Donc  $\left[\frac{q-1}{2}\right]$  boules de rayon  $3^{q-1} \times 2D(J)$  + D(J) recouvre J. D'où le résultat.

Pour obtenir la vitesse annoncée, il suffit alors, dans l'inégalité (5-10) de majorer la norme r par la norme 4. On en déduit, comme dans la remarque:

$$a_{n2} = (\log \sigma_n)^d / \sigma_n.$$

5.4. Cas d'un mélange à décroissance puissance

$$\alpha(m) = O(m^{-a}), \qquad a > 0$$

(C0) est équivalente à:  $a \delta > d(2+\delta)$ ,

(C1) est vérifiée si 
$$m = m(n) = \sigma_n^h, h = \frac{(2+\delta)\left[(2+\delta) \land 3\right]}{a\,\delta + d(2+\delta)\left[(1+\delta) \land 2\right]}.$$

Choisissons alors m = m(n):

Le terme  $a_{n2}$  tend vers 0 dès que:

$$a = \beta \frac{2d(2+\delta)\left[(1+\delta) \land 2\right]}{\delta \cdot (\delta \land 1)} \quad \text{avec } \beta > 1 \tag{A}$$

La condition (A) implique (C0). D'autre part, on vérifie facilement que si  $\tilde{\epsilon} < 1$ , et si (A) est vérifiée. (C2) est vraie. Pour ce choix m(n):

$$T_0(n) = \frac{\sigma_n}{(\log \sigma_n)^{\frac{d-1}{2}} \sigma_n^{\frac{d}{2} \cdot h}} \ge a_{n2}^{-1}$$
$$\Box_n = 0 \left( \max\left\{a_{n2}, \frac{\log T_0(n)}{\sigma_n} m(n)^{\frac{3d}{2}}\right\} \right)$$

et donc:

Sous (A),  $a_{n2}$  réalise ce max:

$$\Delta_n = O(a_{n2}), \quad \text{avec:}$$

$$a_{n2} = \sigma_n^{-E(\beta, \,\delta, \,d)}$$

$$E(\beta, \delta, d) = \delta \wedge 1 - d\left[(1+\delta) \wedge 2\right] \frac{(2+\delta)\left[(2+\delta) \wedge 3\right]}{a\,\delta + d(2+\delta)\left[(1+\delta) \wedge 2\right]}$$

$$= (\delta \wedge 1) \cdot \frac{2(\beta-1)}{2\beta + (\delta \wedge 1)} = E(\beta, \delta).$$

La condition (A) porte explicitement sur  $\delta, d, \beta$ . L'exposant dans la vitesse ne dépend de d que par l'intermédiaire de (A).

#### References

- 1. Berkes. I., Morrow, G.J.: Strong invariance principales for mixing random fields. Z. Wahrscheinlichkeitstheorie verw. Gebiete 57, 13-37 (1981)
- 2. Bolthausen. E.: Exact convergence rates in some martingale central limit theorem. Ann. Probability 10. 672-688 (1982)
- 3. Bolthausen, E.: On the central limit theorem for stationary random fields. Ann. Probability 10, 1047-1050 (1982)
- 4. Bolthausen, E.: The Berry-Esseen theorem for strongly mixing Harris recurrent Markov chains. Z. Wahrscheinlichkeitstheorie verw. Gebiete 60, 283-289 (1982)
- 5. Deo, C.M.: A functional central limit theorem for stationary random fields. Ann. Probability 3, 708-715 (1975)
- 6. Doukhan, P., Portal, F.: Principes d'invariance faibles, avec vitesse dans un cadre mélangeant. Prépublication Orsay n° 83T28 (1983)
- 7. Dobrushin, R.L.: The description of a random field by its conditional distribution and its regularity condition. Theory Probability Appl. 13, 197-227 (1968)
- 8. Egorov, V.A.: Certain limit theorems for *m*-dependent random variables. Liet. Mat. Rink 10, 51-59 (1970) (en Russe)
- 9. Gorodetskii, V.V.: The invariance principle for stationary random fields satisfying the strong mixing condition. Theory Probability Appl. 27, 380-385 (1982)
- 10. Ibragimov, I.A., Linnik, Yu.V.: Independent and stationary sequences of random variables. Groningen: Wolters-Nordhoff 1971
- 11. Hall, P., Heyde, C.C.: Martingale limit theory and its application. New York: Academic Press 1980
- 12. Leonenko, N.N.: Estimate of the rate of convergence in the central limit theorem for mdependent random fields. Mathematical Notes 17, 76-78 (1975)
- 13. Maejima: A non uniform estimate in the central limit theorem for *m*-dependent random variables. Keio Engineering Reports 31, 15-20 (1978)
- 14. Nahapetian, B.S.: The central limit theorem for random fields. In Multicomponent Random Systems, pp. 531-542, ed. by R.L. Dobrushin and Ya.G. Sinai. New York: Marcel Dekker 1980
- 15. Neaderhauser, C.C.: Limit theorems for multiply indexed mixing random variables, with application to Gibbs random fields. Ann. Probability 6, 207-215 (1978)
- 16. Neaderhauser, C.C.: Some limit theorems for random fields. Comm. Math. Phys. 61, 293-305 (1978)
- 17. Petrov, V.V.: On the central limit theorem for *m*-dependent variables. Selected Transl. Math. Statist. Probability 9. 83-88 (1970)
- 18. Prakasa Rao, B.L.: A non uniform estimate of the rate of convergence in the central limit theorem for *m*-dependent random fields. Z. Wahrscheinlichkeitstheorie verw. Gebiete 58, 247-256 (1981)

- 19. Schneider, E.: On the speed of convergence in the random central limit theorem of  $\phi$ -mixing processes. Z. Wahrscheinlichkeitstheorie verw. Gebiete 58, 125-138 (1981)
- 20. Shergin, V.V.: An estimate of the remainder term in the central limit theorem for *m*-dependent random variables. Lith. Math. J. vol 16, 637-641 (1977) (trad. de Litov Math. Sbornik, vol. 15, 245-250 (1976)
- 21. Shergin, V.V.: On the convergence rate in the central limit theorem for *m*-dependent random variables. Theory Probability Appl. 24, 782-796 (1979)
- 22. Stein, C.: A Bound for the error in the normal approximation of a sum of dependent random variables. Proc. Berkeley Sympos. Math. Statist Probability 2, 583-603 Univ. Calif. (1972)
- 23. Takahata, H.: On the rates in the central limit theorem for weakly dependent random fields. Z. Wahrscheinlichkeitstheorie verw. Gebiete 64, 445-456 (1983)
- 24. Tikhominov, A.N.: On the converge rate in the central limit theorem for weakly dependent random variables. Theory Probability Appl. 25. 790-809 (1980)
- 25. Yokoyama, R.: Moment bounds for stationary mixing sequences. Z. Wahrscheinlichkeitstheorie verw. Gebiete 52, 45-57 (1980)

Reçu le 22 Avril 1982; en forme revisée le 5 Décembre 1983

# TESTING THE ASSOCIATION BETWEEN TWO SPATIAL PROCESSES : A REVIEW OF DIFFERENT APPROACHES

### TESTING THE ASSOCIATION BETWEEN TWO SPATIAL PROCESSES A REVIEW OF DIFFERENT APPROACHES

bу

Sylvia RICHARDSON

Laboratoire de Statistiques Médicales, Université de Paris V

et

Institut National de la Santé et de la Recherche Médicale, Villejuif

France

#### Abstract

This paper presents a short review of the statistical problems involved in testing association between two autocorrelated variables.

#### Résumé

Cet article contient une brève synthèse des problèmes statistiques rencontrés lors d'un test d'association entre deux variables auto-corrélées.

#### 1. INTRODUCTION

A question which is frequently asked in statistics is whether the apparent association between two variables is due to something other than chance. Of particular interest are cases where the variables are observed at a variety of spatial locations.

Typical examples are found in the fields of geography and regional science, for instance when relating consumption of agricultural output and road accessibility (Cliff and Ord [12]). Examples in sociology and political science are described by Doreian [17], with reference for instance to voting behaviour and socio-economic or political factors. In epidemiology, etiological clues to environmental risk-factors are sometimes sought through their joint analysis with disease incidence or mortality maps (Doll [16]).

The data in these examples consist of a set A of N locations, and a pair of variables indexed by their locations. These variables will typically exhibit some degree of spatial autocorrelation. The first step in studying the relationship between these variables is often the testing of statistical independence between pairs of variables,  $\{(X_{\alpha}, Y_{\alpha}), \alpha \in A\}$ . This problem cannot be solved by the usual methods based on correlation and linear regression which require that both  $\{X_{\alpha}, \alpha \in A\}$  and  $\{Y_{\alpha}, x \in A\}$  represent independent samples.

The consequences of neglecting the autocorrelation in the X and Y processes were first pointed out in the case of stationary times series by Bartlett [3]. Bartlett calculated the asymptotic variance of the empirical correlation coefficient  $r_{XY}$ . When X and Y are mutually independent and normally distributed, this variance depends only upon the autocorrelations of both processes. This formula was extended to stationary spatial process in  $\mathbb{Z}^2$  by Richardson and Hémon [32]. Numerical calculations given in [33] show how, when X and Y are both positively autocorrelated, first-order nearest neighbour isotropic autogressive or Markovian processes, the asymptotic variance of  $r_{XY}$  increases steeply with the first-order autocorrelation of each process. Moreover, for two dependent processes, Bivand [7], in a Monte Carlo study, made similar observations about the increase in the standard deviation of Fisher's Z transformation of  $r_{XY}$ .

In regression analysis of Y on X the effect of autocorrelation in the dependent variables and in the error term has been analysed by several authors. In the case of positive autocorrelation, Johnston [26] for time series, Cliff and Ord [12] and Bennett [4] for spatial series, show that the variance of the slope is inflated and that the least-square estimator of the variance of the errors is biased downwards, leading to over significant t and F statistics. An illustration of the bias between the true variance of the slope and the OLS variance estimate is given by Martin [29] in a Monte Carlo study. Moran's I statistic for testing autocorrelation has been extended by Cliff and Ord [12] to the testing of autocorrelation among

regression residuals. Further work in this area was done by Brandsma and Ketellapper [10] who compared the power of different tests for autocorrelation among regression residuals and concluded that overall Moran's test with ordinary least squares residuals seems to perform best.

The problem of dealing with spatial autocorrelation in testing for association may be tackled in various ways. Time series methods such as the prewhitening of each series before analysis of the cross correlations could be extended. Alternatively, standard measures of association such as the correlation coefficient can be adapted to take account of autocorrelation. This approach will be developed in a subsequent paper and is outlined in Clifford and Richardson [13]. Moreover classical regression analysis can be extended to take specifically into account the spatial structure of the data. Finally new indices of association can be proposed.

#### 2. REVIEW

In time series, Haugh [23] has proposed that univariate time series model should be fitted to each series. The cross correlation coefficients between the innovation processes are then calculated. This set of cross correlation coefficients is found to be asymptotically independent and normally distributed and can be used to perform a  $arepsilon^2$  test of independence of the two series. He also uses prewhitening as a first step in identifying (dynamic) regression models relating the two series. Using this approach, Pierce [31] found only weak evidence of a relationship between pairs of economic time series which were traditionally considered as related. In a subsequent paper, Geweke [18] compares several tests of independence between stationary time series, in particular Haugh's test and an F test on the regression parameters of a mixed regressive-autoregressive model between the two series. He concludes that in many cases the rate of type II error of Haugh's test is larger than that F test of the regression coefficients. For spatial series, no comparable study has been done to date although the prewhitening of series has been discussed by Griffith [20]. It is likely that a similar conclusion to Geweke would hold and that tests of independence or estimates of regression coefficients based on the original series will be more efficient than those based on their residuals after prewhitening.

As for time series, spatial structure has been traditionally accounted for in regression models in two ways, either by considering a mixed simultaneous regressive autoregressive model :

$$Y = L(Y) + X\beta + \varepsilon , \qquad (1)$$

where z is independent of Y and L(Y), represents some linear function of the value of Y in locations z,  $y \neq \alpha$  or by considering a regression model with

autocorrelated error term :

 $Y = X\beta + U$ ,

where the spatial autocorrelation is reflected in the structure of  $\Sigma$ , the variance covariance matrix of the error U (see, for example, Cliff and Ord [12], Upton and Fingleton [37]).

The interpretation of  $\beta$  is essentially different for the regressions (1) and (2), hence these equations should not be treated as alternatives but the choice between them should be motivated by the problem considered (Cox [15]). In (1) the spatial autocorrelation is concentrated on Y and represents the effect of X when the influence of neighbouring values of Y has already been taken into account, whereas for (2) the autocorrelation is generated solely by U.

There are cases when equation (1) does not provide a plausible explanatory model, for instance when relating the incidence of a non contagious disease to environmental factors, since the neighbouring disease rates have no direct causal influence. On the other hand, for some examples, a propagation effect of Y is conceivable like in models of regional income diffusion considered by Haining [21], or those on the Huk rebellion considered by Doreian [17].

The estimation of models (1) and (2) was first discussed by Ord [30]. The autoregressive term in (1) is usually modelled as :  $L(Y) = \rho WY$ , where W is a matrix of weights representing contiguity introduced by Cliff and Ord [12], and  $\rho$  is a parameter summarising the overall level of autocorrelation which has to be estimated. If the error term is multivariate normal with uncorrelated components, the likelihood of Y can be written down and the parameters can be estimated, usually by direct search. The least-square estimates are not consistent for the simultaneous representation in (1) but Cliff and Ord [12] and Besag [5] have pointed out that this is not the case in the equivalent Markovian scheme.

Mardia and Marshall [28] have studied the asymptotic properties of the maximum likelihood estimators for model (2). Supposing that Y is a Gaussian process and that its matrix of variance-covariance  $\Sigma$  can be parametrised in term of a finite number of parameters  $\theta$ , they give conditions which ensure the consistency and the asymptotic normality of  $(\hat{\beta}, \hat{\theta})$ , the maximum likelihood estimates of  $(\beta, \theta)$ .

Models for  $\Sigma$  have been basically of two kinds, either a specific autoregressive model for U is supposed to hold :

$$U = \rho WU + \varepsilon , \varepsilon \sim N(0, \sigma^2 I) , \qquad (3)$$

which implies that  $\Sigma = \sigma^2 (I - \rho W)^{-1} (I - \rho W)^{-1}$ , or the elements of  $\Sigma$  are assumed

to depend on distance with a functional form suggested by the variogram of the ordinary least-square regression residuals.

Equation (3) has been fitted by Bodson and Peeters [9], Cliff and Ord [12] and Bivand [8]. Equations (2) and (3) can be combined to give :

$$Y = \Im WY + XB - \wp WX B + \varepsilon, \tag{4}$$

and so giving a particular case of a larger class of regression models in which Y is also regressed on WX. Tests of the adequacy of (4) within this larger class are discussed by Burridge [11]. A Bayesian analysis of model (3) has been given by Hepple [24]. Using a uniform diffuse prior for  $\beta = (\beta_1, \beta_2)$ ,  $\phi$  and  $\log \sigma$ , Hepple shows how the conditional posterior distribution of  $\beta_2$  is sensitive to value of  $\phi$  and he derives the bivariate posterior distribution for  $\beta_2$  and  $\phi$ , with mode corresponding to the maximum likelihood estimates of  $\beta_2$  and  $\phi$ .

The arbitrary nature of the definition of M is a common criticism of model (3). Arora and Brown's 52 alternative approach either needs the availability of spatial data at different time intervals or assumes that the spatial autocorrelation is entirely created through an unknown common error term. Their approach does not apply to data on which spatial autocorrelation decreases with distance.

The choice of parametrisation of [] by plotting the variogram of the least square residuals has been discussed by Ripley [34], Cook and Pocock [14] and Agterberg [1]. Care has to be taken when interpretation the variogram since it is sensitive to the number of pairs of data points used to estimate the empirical covariance at a particular distance.

The number of pairs will vary with the distance, typically increasing at first. Riplay advises the use of cross-validation by successive deleting of data points to assess the fit of the model chosen for the covariance function. Cook and Pocock fitted an exponential function decreasing with distance, whilst Agterberg uses a quadratic one.

Once a particular parametrisation is chosen, the parameters may be estimated by maximum likelihood, either by grid search or by iteration. Mardia and Marshall [28] recommend updating the parameters by scoring, Cliff and Ord [12] suggest finding a conditional maximum for  $\beta$ ,  $\beta$  being fixed, then estimating  $\beta$  by generalised least squares and iterating until convergence; a convergence however which is not necessarily guaranteed (see Haining [22]).

Other ways to model 10, for instance by a spectral decomposition (Streitberg .351) or with a finite number of contiguity matrices of increasing order (Lebart 27 have been proposed but are infrequently used.

Clearly, the regression methods described above are computationally involved and entail a certain degree of arbitrariness. Bivand [8] discusses model selection and introduces the use of Akaike's criterion. In a trend surface analysis, Haining [22] compares three different models for  $\Sigma$  in (2).

A non parametric index of association based on the locations of the ranks of two spatially defined variables  $\{(X_{\alpha}, Y_{\alpha}), \alpha \in A\}$  has been introduced by Tjøstheim [36] and generalised by Hubert and Golledge [25]. This index is defined as the sum of the distances between the locations of similar ranks, with a suitably chosen distance function.

In a simulation study, Glick [19] illustrates the performance of this test in comparison to Spearman's rank correlation. When the X process has no autocorrelation and the Y process corresponds to a simple spatial shift of the X process, Tjøstheim's index  $\Lambda$  detects an association whilst Sperman's rank correlation is non significant. The situation is reversed when the dependence between the X and the Y process has no spatial component. The randomisation model used by Tjøstheim and, Hubert and Golledge to generate a reference distribution against which observed values of  $\Lambda$  are compared under the null hypothesis is that of a spatial redistribution of one of the variables, the other staying fixed. This full set of permutation clearly does not preserve the autocorrelation of the variable and hence the significance of  $\Lambda$  is not correctly assessed in this case. One would need to consider a smaller set of permutations which would respect in some way the spatial autocorrelation of the variable. This would also be the case for the Monte Carlo method proposed by Besag and Dingle [6] for the detection of spatial association between two sets of continuous quadrat counts over a simple region.

#### 3. CONCLUSION

In this review we have tried to outline different approaches to the problem of testing association between spatial variables. Further work on the modification of existing measures, in particular with the aim of taking into account the spatial scale of the processes would be of interest and of practical use. Studies comparing the power of these methods would give useful insight into this complex problem.

#### REFERENCES

- [1] Agterberg, F.P. (1984), "Trend surface analysis", in Spatial Statistics and Models, edited by E.L. Gaile and C.J. Willmott. Dordrecht, Reidel.
- [2] Arora, S.S. and M. Brown (1977), "Alternative approaches to spatial autocorrelation : an improvement over current practice". International Regional Science Review, 2(1).
- [3] Bartlett, M.S. (1935), "Some aspects of time-correlation problem in regard to test of significance". Journal of the Royal Statistical Society, 98, 536-543.
- [4] Bennett, R.J. (1979), Spatial Time Series. London, Pion.
- [5] Besag, J.E. (1975), "Statistical analysis of non-lattice data". The Statistician, 24(3), 179-195.
- [6] Besag, J.E. and P.J. Diggle (1977), "Simple Monte Carlo tests for spatial pattern". Applied Statistics, 26(3), 327-333.
- [7] Bivand, R.S. (1980), "A Monte Carlo study of correlation coefficient estimation with spatially autocorrelated observations". *Quaestiones Geographicae*, 6, 5-10.
- [8] Bivand, R.S. (1984), "Regression modeling with spatial dependence : an application of some class selection and estimation methods". Geographical Analysis, 16(1), 25-37.
- 193 Bodson, P. and D. Peeters (1975), "Estimation of the coefficients of a linear regression in the presence of spatial autocorrelation. An application to a Belgian labour-demand function". Explorement and Planning A, 7, 456-472.
- [10] Brandsma, A.S. and R.H. Ketellapper (1979), "Further evidence on alternative procedures for testing of spatial autocorrelation among regression disturbances", in Exploratory and Explanatory Statistical Analysis of Spatial Data, edited by Bartels and Ketellapper.
- [11] Burridge, P. (1981), "Testing for a common factor in a spatial autogresssion model". Environment and Planning A, 13, 795-800.
- [12] Cliff, A.D. and J.K. Ord (1981), Spatial Processes. Models and Applications. London, Pion.
- [13] Clifford, P. and S. Richardson (1985), "Testing the association between two spatial processes". *Statistics and Decisions*, Suppl. issue, 2, 155-160.
- [14] Cook, D.G. and S.J. Pocock (1983), "Multiple regression in geographic mortality studies with allowance for spartially correlated errors". *Biometrics*, 39, 361-371.
- [15] Cox, D.R. (1981), "Statistical analysis of time series : some recent developments". Scandinavian Journal of Statistics, 8, 93-115.
- [16] Doll, R. (1980), "The epidemiology of cancer". Cancer, 45, 2475-2485.
- [17] Doreian, P. (1981), "Estimating linear models with spatially distributed data", in Scotalogical Methodology, edited by Leinhardt and Samuel. San Francisco, Jossey-Bass Publishers, 359-388.
- [18] Geweke, J. (1981), "A comparison of tests of the independence of two covariancestationary time series". Journal of the American Statistical Association, 76(374), 363-373.

- [19] Glick, B.J. (1982), "A spatial rank order correlation measure". Geographical Analysis, 14(2), 177-181.
- [20] Griffith, D.A. (1980), "Towards a theory of spatial statistics". Jeographical Analysis, 12(4), 325-339.
- [21] Haining, R. (1985), "Income diffusion and regional economic structure". Presented to the NATO advanced studies institute, Hanstholm, Denmark.
- [22] Haining, R. (1985), "Trend surface models with global and local scales of variation". Presented at a conference on Applied Spatial Statistics, Cambridge College of Arts and Technology, April 18, 1985.
- [23] Haugh, L.D. (1976), "Checking the independence of two covariance-stationary time series : a univariate residual crosscorrelation approach". Journal of the American Statistical Association, 71, 378-385.
- [24] Hepple, L.W. (1984), "Bayesian analysis of the linear model with spatial dependence", in *Spatial Statistics and Models*, edited by Gaile and Reidel. Dordrecht, Reidel.
- [25] Hubert, L.J. and R.G. Golledge (1982), "Measuring association between spatially defined variables : Tjøstheim's index and some extensions". Jeographical Analysis, 14(3), 273-278.
- [26] Johnston, J. (1972), Econometric Methods, 2nd edition. New York, McGraw Hill.
- [27] Lebart, L. (1969), "Analyse statistique de la contiguité". Eubl. Int. Stat. Endo. Pario, 18, 81-112.
- [23] Mardia, K.V. and R.J. Marshall (1984), "Maximum likelihood estimation of models for residual covariance in spatial regression". Sigmetrika, 71(1), 135-146.
- [29] Martin, R.L. (1974), "On spatial dependence, bias, and the use of first spatial differences in regression analysis". Appl. 6, 185-194.
- [30] Ord, K. (1975), "Estimation methods for models of spatial interaction". Journal of the American Crastilation American, 70, 120-126.
- [31] Pierce, D.A. (1977), "Relationships and the lack thereof between economic time series, with special reference to money and interest rates". Internal of the American Statistical Association, 72(357), 11-26.
- [32] Richardson, S. and D. Hémon (1981), "On the variance of the same correlation between two independent lattice processes". Journal of Applied Probability, 18, 943-948.
- [33] Richardson, S. and D. Hémon (1982), "Autocorrélation spatiale : ses conséquences sur la corrélation empirique de deux processus spatiaux". Revue de Stationie Appliquée, 30(1), 41-51.
- [34] Ripley, B.D. (1981), Sparial Statistics. New York, Wiley.
- [35] Streitberg, B. (1979), "Multivariate models of dependent spatial data", in Exploratory and Exploratory Statistical Analysis of Spatial Lata, edited by Bartels and Ketellapper.
- [36] Tjøstheim, D. (1978), "A measure of association for spatial variables". Biometrika, 65(1), 109-114.
- [37] Upton, G. and B. Fingleton (1983), Sparing Sata Analysis by Example. New York, Wiley.

# ASSESSING THE SIGNIFICANCE OF THE CORRELATION BETWEEN TWO SPATIAL PROCESSES

### Assessing the Significance of the Correlation Between Two Spatial Processes

Peter Clifford,<sup>1</sup> Sylvia Richardson,<sup>2</sup> and Denis Hémon<sup>3</sup>

<sup>1</sup> Mathematical Institute, University of Oxford, 24–29 St Giles, Oxford OX1 3LB, England

<sup>2</sup> INSERM U170, 16 Avenue Paul Vaillant Couturier, 94807 Villejuif Cedex, France, and Laboratoire de Statistiques Médicales, Université de Paris V, 45, rue des Saints Pères, 75006 Paris, France

<sup>3</sup> INSERM U170, 16 Avenue Paul Vaillant Couturier, 94807 Villejuif Cedex, France

#### SUMMARY

Modified tests of association based on the correlation coefficient or the covariance between two spatially autocorrelated processes are presented. These tests can be used both for lattice and nonlattice data. They are based on the evaluation of an effective sample size that takes into account the spatial structure.

For positively autocorrelated processes, the effective sample size is reduced. A method for evaluating this reduction via an approximation of the variance of the correlation coefficient is developed. The performance of the tests is assessed by Monte Carlo simulations. The method is illustrated by examples from geographical epidemiology.

#### 1. Introduction

The calculation of correlation coefficients for variables that are observed at a variety of different spatial locations has suggested intriguing hypotheses about the relationship between environmental factors and disease pathologies (Armstrong and Doll, 1975; Hoover and Fraumeni, 1975). It is well known that the null distribution of r, the product moment correlation coefficient, is influenced by spatial and temporal autocorrelation (Student, 1914; Bartlett, 1935; Richardson and Hémon, 1981).

In a preliminary report, Clifford and Richardson (1985) have suggested a method of approximating the critical values of r. Their method depends on an estimate of the variance of r. In this paper we investigate a general method of obtaining such an estimate and show that the test which results is related to a test based on the standardised covariance between the processes. Theoretical properties of the distribution of r are reviewed in Section 2. In Section 3 we report on an empirical simulation study of the performance of these procedures and in Section 4 we apply our methods to test the association between cigarette consumption, industrial risk factors, and deaths from lung cancer for French départements.

Some authors have devised measures of association that involve transformation of the data before standard techniques are applied. Student (1914) advocated a form of what would now be called prewhitening, and more recent work, such as that of Haugh (1976) with time series and Davies and Jowett (1958) can be seen as developments of this local approach, in that filters are applied to the data in order to reduce temporal or spatial

Requests for reprints should be addressed to the second author.

Key words: Correlation coefficient: Geographical epidemiology; Monte Carlo simulations: Significance tests: Spatial processes.

autocorrelation. However, the application of this approach to irregularly spaced data presents a number of problems that have not been thoroughly explored and information is lost by filtering techniques.

The problem of testing the association between autocorrelated variables can also be tackled by regression techniques (Ord, 1975) when a relationship between a dependent variable and a set of independent variables is postulated. Cook and Pocock (1983), in their study of the association between water hardness and heart disease, considered a regression model with stationary autocorrelated errors, the autocorrelation decreasing exponentially with distance. Mardia and Marshall (1984) have demonstrated the asymptotic properties of this approach. In general, these techniques involve a substantial amount of computing time and depend on an explicit parametric model for the autocovariance.

A nonparametric index of association has been proposed by Tjøstheim (1978) and generalised by Hubert and Golledge (1982). However, they assess the significance of the index with reference to a randomisation distribution that involves the spatial redistribution of one of the variables, the other staying fixed. This full set of permutations clearly does not preserve the autocorrelation of the permuted variable and hence the significance of the index is not correctly assessed.

#### 2. Modified Tests of Association

#### 2.1 Basic Properties of r

We are interested in data sets which consist of a set A of N locations numbered from 1 to N and a set of pairs of observations  $(X_{\alpha}, Y_{\alpha}), \alpha \in A$ , where each pair is indexed by its location. The correlation coefficient is then given by

$$r = \frac{S_{XY}}{S_X S_Y},\tag{2.1}$$

where  $s_{XY} = N^{-1} \sum_A (X_\alpha - \overline{X})(Y_\alpha - \overline{Y})$  is the sample covariance,  $s_X^2 = N^{-1} \sum_A (X_\alpha - \overline{X})^2$ ,  $s_Y^2 = N^{-1} \sum_A (Y_\alpha - \overline{Y})^2$ , and where  $\overline{X} = N^{-1} \sum_A X_\alpha$  and  $\overline{Y} = N^{-1} \sum_A Y_\alpha$ . If, conditional on X, the elements of Y are normal i.i.d. variables or conditional on Y, the elements of X are normal i.i.d., then r has the standard null distribution with p.d.f.

$$f_N(r) = \frac{(1 - r^2)^{(N-4)/2}}{B[\frac{1}{2}, \frac{1}{2}(N-2)]}, \quad r \le 1,$$
(2.2)

where B is the beta function.

Critical values of r are usually obtained from t-tables since  $(N-2)^{1/2}r/(1-r^2)^{1/2}$  has a t-distribution with N-2 degrees of freedom under these assumptions. This is also the t-statistic that is calculated in testing the significance of the linear regression either of Y on X or of X on Y.

#### 2.2 The Standardised Covariance

For independent samples the correlation coefficient can be thought of as a standardised covariance. If the elements of Y are i.i.d., then conditional on X = x the sample covariance  $s_{XY}$  has mean zero and variance  $N^{-2} \sum_{i} (x_{\alpha} - \bar{x})^2 \sigma_Y^2$  where  $\sigma_Y^2$  is the variance of the elements of Y. To standardise  $s_{XY}$ , the unknown quantity  $\sigma_Y^2$  is replaced by an unbiased estimate  $\sum_{i=1}^{n} (Y_{\alpha} - \bar{Y})^2 / (N - 1)$  and  $s_{XY}$  is divided by the resulting estimate of its standard deviation. This leads to the expression  $(N - 1)^{1/2} r$ , whose significance would be approximately assessed with reference to tables of the normal distribution, relying on a central limit theorem to justify the approximation. Thus, a test based on the standardised covariance is equivalent to a test based on r.

In the general case, we suppose now that X and Y are independent but that both X and Y are multivariate normal vectors with constant means and variance-covariance matrices  $\Sigma_X$  and  $\Sigma_Y$ , respectively. The conditional variance  $s_{XY}$  is given by

$$N^{-2} \sum_{\alpha,\beta} (X_{\alpha} - \bar{X})(X_{\beta} - \bar{X}) \operatorname{cov}(Y_{\alpha}, Y_{\beta}), \qquad (2.3)$$

which is equal to  $N^{-2} \sum_{\alpha,\beta} (X_{\alpha} - \overline{X})(X_{\beta} - \overline{X}) \operatorname{cov}(Y_{\alpha} - \overline{Y}, Y_{\beta} - \overline{Y})$  since  $\sum_{A} (X_{\alpha} - \overline{X}) = 0$ . Replacing  $\operatorname{cov}(Y_{\alpha} - \overline{Y}, Y_{\beta} - \overline{Y})$  by the unbiased estimate  $(Y_{\alpha} - \overline{Y})(Y_{\beta} - \overline{Y})$  gives the expression  $s_{XY}^2$ , as a trivial estimate of the conditional variance of  $s_{XY}$ . It is clear that no progress can be made until some plausible restrictive structure is imposed on  $\Sigma_{Y}$ .

#### 2.3 Structure for $\Sigma_X$ and $\Sigma_Y$

We will assume that the set of locations A is a subset of some larger set  $\Omega$ . For stationary processes  $\Omega$  can, in principle, be arbitrarily large. For real spatial data,  $\Omega$  is finite, perhaps equal to A itself. We assume that the set of all ordered pairs of elements of  $\Omega$  can be divided into strata  $S_0, S_1, S_2, \ldots$ , so that the covariances within strata are constant, i.e.,  $\operatorname{cov}(X_{\alpha}, X_{\beta}) = C_{X}(k)$  and  $\operatorname{cov}(Y_{\alpha}, Y_{\beta}) = C_{Y}(k), (\alpha, \beta) \in S_k, k = 0, 1, \ldots$ . Of course, if  $(\alpha, \beta) \in S_k$ , then  $(\beta, \alpha) \in S_k$  for consistency. For stationary processes the stratification is indexed by directional lags. For isotropic processes the number of strata is reduced in the lattice case and when the data are irregularly spaced the strata can be indexed by a discrete distance function. The general formulation is flexible enough to permit spatially inhomogeneous variances and other aspects of nonstationarity. With this structure (2.3) becomes

$$N^{-2} \sum_{k} N_{k} \left[ \frac{\sum_{A_{k}} (X_{\alpha} - \bar{X})(X_{\beta} - \bar{X})}{N_{k}} \right] C_{Y}(k), \qquad (2.4)$$

where  $N_k$  is the cardinality of  $A_k$ ,  $A_k = (A \times A) \cap S_k$ , and the summation over  $A_k$  is for all ordered pairs  $(\alpha, \beta) \in A_k$ .

For stationary processes, Clifford and Richardson (1985) have suggested using

$$\hat{C}_{\mathbf{Y}}(k) = \sum_{A_k} (Y_{\alpha} - \overline{Y})(Y_{\beta} - \overline{Y})/N_k$$
(2.5)

as an estimate of  $C_{Y}(k)$  for values of k corresponding to small spatial lags and shrinking  $\hat{C}_{Y}(k)$  to zero otherwise. This was partially for computational convenience. Here we propose to consider the inclusion of all lags. The resulting estimate of the conditional variance of  $s_{XY}$  is therefore

$$N^{-2} \sum_{k} N_{k} \hat{C}_{X}(k) \hat{C}_{Y}(k).$$
 (2.6)

It does not rely on an arbitrary notion of what constitutes a small lag, it is invariant to shifts in the mean, and it is unbiased for periodic processes. It has the additional advantage of being symmetric in X and Y, so that it is also the estimate of the conditional variance of  $s_{XY}$  given Y. Note that when  $S_0 = A \times A$ , (2.6) reduces to  $N^{-1}s_X^2s_Y^2$ .

#### 2.4 The Standardised Covariance and the Modified t-Test

Using the estimate (2.6), the standardised covariance, W, becomes

$$W = Ns_{YY} \left[ \sum_{k} N_k \hat{C}_X(k) \hat{C}_Y(k) \right]^{-1/2}$$
(2.7)
If we consider the correlation coefficient, in Appendix 1 it is shown that to the first order,

$$\sigma_r^2 \approx \frac{\operatorname{var}(s_{YY})}{\operatorname{E}(s_Y^2)\operatorname{E}(s_Y^2)}.$$
(2.8)

We therefore take as our estimate of  $\sigma_r^2$ 

$$\hat{\sigma}_r^2 = \frac{\sum N_k \hat{C}_{\rm X}(k) \hat{C}_{\rm Y}(k)}{N^2 s_{\rm Y}^2 s_{\rm Y}^2}.$$
(2.9)

We consider a modified *t*-test for the correlation coefficient by approximating the critical values of *r* by percentage points of the p.d.f.  $f_{\hat{M}}(r)$ , where  $\hat{M} = 1 + \hat{\sigma}_r^{-2}$  and *f* is defined by (2.2). The quantity  $(\hat{M})$  *M* can be thought of as an (estimated) "effective sample size" that takes into account the spatial autocorrelation in the variables X and Y. Note that when  $S_0 = A \times A$ ,  $\hat{M} = N + 1$ . For positively autocorrelated processes, the estimated effective sample size is typically less than N. If one of the processes has negative autocorrelation it is, in principle, possible that the effective sample size will be larger than N.

Comparing (2.7) and (2.9), we see that

$$W = (\hat{M} - 1)^{1/2} r, \qquad (2.10)$$

where  $\hat{M} = 1 + \hat{\sigma}_r^{-2}$ .

In Section 3 we consider the performance of the test obtained by assuming that W has a N(0, 1) distribution under the null hypothesis and compare this test with that obtained by assuming that r has p.d.f.  $f_{\hat{M}}(r)$ , that is, using a *t*-statistic with  $\hat{M} - 2$  in place of N - 2.

#### 3. Monte Carlo Simulations

The performance of the modified *t*-test or the test based on the standardised covariance has been assessed by Monte Carlo simulation in both a lattice and a nonlattice case.

#### 3.1 The Lattice Case

Method We simulated stationary first-order isotropic simultaneous autoregressive processes defined by

$$X_{s,t} = a(X_{s-1,t} + X_{s+1,t} + X_{s,t-1} + X_{s,t+1}) + \varepsilon_{s,t}, \qquad (3.1)$$

where  $\{\varepsilon_{s,t}\}$  is a sequence of i.i.d. N(0, 1),  $0 < |a| < \frac{1}{4}$ . This class of processes has been widely discussed in the modelling of spatial patterns (Whittle, 1954; Besag, 1974; Cliff and Ord, 1975) and is relatively easy to simulate (see Appendix 2).

The simulation was performed on a DPS7 C.I.I. Honeywell Bull. A spectral decomposition similar to that given by Besag in the discussion of Bartlett's (1978) paper for autonormal processes was used to generate processes on a  $26 \times 26$  lattice with zero on the boundary and then restricted to middle  $12 \times 12$ ,  $16 \times 16$ , and  $20 \times 20$  squares. This restriction was sufficient to render negligible the influence of the boundary condition.

For each simulation, a 26 × 26 field of i.i.d. N(0, 1) variables was first generated using a polar algorithm. Values of *a* equal to 0, .0945, .165, .2099, .2364 were chosen to give values of the nearest neighbour autocorrelation,  $\rho_X(1)$ , equal to 0, .2, .4, .6, and .8, respectively.

For each value of a, 500 pairs  $(X_{\alpha}, Y_{\alpha}), \alpha \in A$ , were generated. Taking advantage of the lattice structure, each process X can also be rotated by 90° or reflected with respect to the diagonals, leading to eight copies (four rotations and/or two reflections) of the original process. These copies have the same set of estimated autocovariances,  $\{\hat{C}_X(k)\}$ , which therefore is calculated only once. Thus, 4,000 trials, dependent in groups of eight, were obtained for each value of  $\rho_X(1)$  and  $\rho_Y(1)$ .

The proposed statistics are denoted by W and  $t_{M-2}$ , where  $\bar{M}$  is the estimated effective sample size.  $\hat{M} = 1 + \hat{\sigma}_r^{-2}$ . If expression (2.6) for the estimated variance of  $s_{MT}$  gave an inadmissible negative estimate, it was replaced by the product  $N^{-1}s_M^2s_T^2$ , which is the estimate in the case of no autocorrelation. The standard *t*-statistic for *r* based on *N* observations is denoted by  $t_{N-2}$ . For the three statistics, W,  $t_{M-2}$ , and  $t_{N-2}$ , a 5% nominal rejection level was chosen. For testing  $t_{M-2}$ , the integer part of  $\hat{M}$  was taken. The empirical variance of the rejection levels was estimated by first averaging the rejection indicator function over the eight dependent pairs  $(X_{\alpha}, Y_{\alpha})$  (where  $X_{\alpha}$  is obtained by rotations/reflections) and then calculating the empirical variance of this average over the 500 independent simulations.

**Results** Figure 1 illustrates the poor performance of the standard *t*-test in the presence of positive autocorrelation. The observed Type I error rates are represented for the three lattice sizes together with their 95% confidence interval. Note that the observed Type I error rates are significantly larger than 5%, even for low values of  $\rho_X(1)$  and  $\rho_Y(1)$ .



Figure 1. Standard test of the correlation coefficient: 95% confidence intervals for the proportion of Type I errors for a 5% nominal test in the case of two mutually independent simultaneous autoregressive processes on a lattice. For each value of the nearest neighbour autocorrelation for X or Y.  $\rho_X(1)$  or  $\rho_X(1)$ , the confidence intervals are plotted for three lattice sizes:  $12 \times 12$ ,  $16 \times 16$ ,  $20 \times 20$ .

When both X and Y are highly autocorrelated, observed Type I error rates vary between 25% and 55%, thus clearly showing that the testing procedure needs to be modified.

For the two proposed tests, the observed Type I error rates ranged from 4.2% to 5.85% for the *W* test (Figure 2) and from 4.1% to 5.9% for the  $t_{M-2}$  test, and are thus close to the nominal 5% level. Note that when one of the processes has no autocorrelation, *W* and  $t_{M-2}$  perform as well as the standard *t*-test. The confidence interval did not include the 5% nominal level in only two cases  $[20 \times 20, \rho_X(1) = .2, \rho_Y(1) = .2$  and  $16 \times 16, \rho_X(1) = .6, \rho_Y(1) = .8]$ . Inadmissible estimates were rare: they did not occur more than twice per lattice size in all the trials except in the case  $\rho_X(1) = 0, \rho_Y(1) = .8$ , where there were three negative estimates for the  $12 \times 12$  lattice.



Figure 2. Test based on the standardised covariance W: 95% confidence intervals for the proportion of Type I errors for a 5% nominal test in the case of two mutually independent simultaneous autoregressive processes on a lattice.

A comparison of the empirical variance of r (averaged over the 4,000 trials) and of the average  $v_r$  of the estimated variance  $\hat{\sigma}_r^2$  given by (2.9), was also made. For small to moderate autocorrelations, there is practically no difference between the empirical variance of r and  $v_r$ . As the autocorrelation increases,  $v_r$  is consistently too low. That (2.6) is negatively biased as an estimate of the variance of  $s_{XY}$  can be easily seen in the case  $\Sigma_X = \Sigma_Y$ .

Quantile plots of the distribution of the W statistic show that it has short tails compared with the normal distribution in the extreme case of  $\rho_X(1) = \rho_Y(1) = .8$ . The departure from normality is confirmed by a Kolmogorov-Smirnov test, which is significant at the 5% level but not at the 1% level. The tendency for short tails also occurs in other cases of  $\rho_X(1)$  and  $\rho_Y(1)$ , but is more marked for higher autocorrelation. Nevertheless, the departure from normality of W and the order of the negative bias in  $v_r$  do not seem to be quantitatively altering the levels of the W and  $t_{M-2}$  tests even in the strongly autocorrelated cases.

#### 3.2 Nonlattice Case

*Choice of the model* The structure of the network and the type of spatial dependence were both chosen in order to be similar to that of examples in geographical epidemiology which will be discussed in the next section.

The coordinates of the points of the network were identified with the geographical locations of the administrative centers ("préfectures") of French départements. The variables considered are the mortality rate for lung cancer (LC) for men, the cigarette sales (CS) per inhabitant, and the percentage of metal workers (MW) and the percentage of textile workers (TW) with respect to the male working population in each département.

The spatial structure of these variables was investigated by means of a variogram. In this analysis, N = 82 locations were retained after grouping the départements around Paris into one area. The distances between the centres of départements were partitioned into 15 classes of 50-km intervals each. This gives 15 strata  $S_1, \ldots, S_{15}$ : the stratum  $S_0 = \{(\alpha, \alpha) | \alpha \in A\}$ . The number of strata chosen should take into account a balance between the sampling fluctuations of the estimated autocovariances when the number of strata is large and the introduction of bias when the number of strata is small. In making this judgement it is helpful to compute the observed semivariogram for several cases. For

the present data set we have found that the results were not sensitive to the choice of the number of strata.

The observed variogram of LC, i.e., the plot of

$$N_k^{-1} \sum_{(\alpha,\beta) \in S_k} (X_{\alpha} - X_{\beta})^2, \qquad k = 1, \dots, 15,$$

against the average distance,  $d_k$ , for départements in  $S_k$ , is shown in Figure 3. The variograms of CC and MW were similar and exhibited also an upward trend of fairly linear shape with increasing distance. Hence, a disc model for the covariance matrix (Ripley, 1981, p. 55) seemed appropriate and we chose it to simulate a spatially dependent process on this irregular network.



**Figure 3.** Variogram of the lung cancer mortality (LC). Fifteen classes of distance are considered: the number of ordered pairs in each class is (82: 400: 582: 674: 764: 822: 812: 726: 630: 476: 304: 172: 94; 58; 40). The abscissa corresponds to the average distance in kilometres within each class.

*Method* The average distance between départements in the first stratum  $S_1$  is approximately 40 km. The parameters of the disc models for X and Y were chosen to be such that the autocorrelation at distance 40 km is equal to .2, ..., .9. We also denote these by  $\rho_X(1)$  or  $\rho_X(1)$ .

For each chosen value of  $\rho_X(1)$ , the matrix  $\Sigma_X$  was triangularised ( $\Sigma_X = \mathbf{LL}^T$ ) and then a realisation of X was obtained by first generating a vector of N i.i.d. N(0, 1) and then multiplying this vector by L. In each case 500 pairs of mutually independent processes (X, Y) were simulated and the statistics W and  $t_{M-2}$  calculated for each pair (X, Y) as in Sections 2.2 and 2.3, where the autovariances  $\hat{C}_X(k)$  are defined as

$$\hat{C}_{\mathbf{X}}(k) = N_k^{-1} \sum_{(\alpha,\beta) \in S_k} (X_\alpha - \bar{X})(X_\beta - \bar{X}).$$

The same procedure as in the lattice case was adopted for an inadmissible negative estimate for the variance of  $s_{XY}$ . A 5% nominal level was chosen to assess the performance of W,  $t_{M-2}$ , and  $t_{N-2}$ .

#### Biometrics, March 1989

**Results** Figure 4 demonstrates the increased proportion of Type I errors of the standard *t*-test in the case of positive autocorrelation of both processes. Figure 5 shows the performance of the *W* test, which is indistinguishable from the performance of the  $t_{M-2}$  test. All observed Type I error rates, even in the most highly autocorrelated case, were close to 5%. No inadmissible negative estimate,  $\hat{\sigma}_r^2$ , arose. Kolmogorov–Smirnov tests performed on the distribution of *W* were significant in two cases (.2 × .6 and .9 × .9).



Figure 4. Standard test of the correlation coefficient: 95% confidence intervals for the proportion of Type I errors for a 5% nominal test in the case of two mutually independent processes generated by a disc model on a network of 82 points. (The parameters  $\rho_N(1)$  and  $\rho_N(1)$  of the disc models for X and Y respectively are equal to the autocorrelation at 40 km.)



Figure 5. Test based on the standardised covariance W:95% confidence intervals for the proportion of Type I errors for a 5% nominal test in the case of two mutually independent processes generated by a disc model on a network of 82 points. (The parameters  $\rho_X(1)$  and  $\rho_Y(1)$  of the disc models for X and Y respectively are equal to the autocorrelation at 40 km.)

#### 4. Examples

Our examples concern the relationship between lung cancer, smoking, and industrial factors. We calculate W and  $t_{M-2}$ .

#### 4.1 The Data

For 82 départements we considered male lung cancer mortality rate (LC) standardised over the age 35-74 and over a 2-year period, 1968-1969; cigarette sales per inhabitant (CS) in 1953 (a 15-year time lag was chosen to account for the delay between exposure and the onset of the pathology); and demographic data on the percentage of employed males in the metal industry (MW) and the textile industry (TW) recorded by census in 1962.

#### 4.2 Results

Using the standard test based on r, there is a highly significant positive association both between LC and CS, and between LC and MW. The association between LC and TW, on the other hand, is less strong but still significant at the 1% level (Table 1). The W and  $t_{M-2}$  statistics for these three examples are shown in Table 1. One can see a substantial reduction of the degrees of freedom when the autocorrelation is taken into account.

Table 1Comparison of the significance levels for tests of the association between lung cancer mortality ratesand several risk factors given by standard test, W, and  $t_{M-2}$  tests

	r	$l_{\lambda-2}^{a}$	W	Ŵ	$t_{\tilde{M}-2}$
Cigarette sales per inhabitant (1953) (CC)	.76	$10.48 P \approx 10^{-21}$	2.94 P = .0032	15	4.22 P = .001
% male workers in metal industry (1962) (MW)	.63	$7.16 P \approx 10^{-11}$	2.48 P = .0136	16	3.00 P = .01
% male workers in textile industry (1962) (TW)	.28	2.57 $P = .01$	1.51 P = .13	30	1.52 P = .15

<sup>a</sup> Standard test (t-transformation with 80 d.f.).

For CS and MW the CS effective sample size is about 20% of the original sample size. Consequently, the significance levels are reduced but even after this "adjustment" these two factors are statistically significantly associated with lung cancer. For TW the significance disappears after adjustment.

These results can be considered in good agreement with current knowledge concerning life style and occupational risk factors for lung cancer (Schottenfeld and Fraumeni, 1982).

#### 5. Discussion

Our study of the correlation coefficient has been motivated by its widespread use in epidemiology, where relatively small data sets of up to 100 irregularly spaced points are encountered. In this context, positively autocorrelated X and Y are most commonly observed.

We have confirmed both theoretically and empirically that the uncritical use of the correlation coefficient for testing association between positively autocorrelated processes leads to an inflated proportion of Type I errors, and have shown theoretically that the magnitude of this inflation is consistent with a reduction in the effective sample size. We have then investigated how the correlation coefficient behaves when it is adjusted to take account of this effect. We have done this in two related ways, the W and  $t_{\dot{M}-2}$  tests, and

have shown that when adjustment is made the Type I error rate is much closer to the nominal level. These methods do not require the identification of a particular parametric model for the type of spatial autocorrelation and they cope equally well with regularly and irregularly spaced points. Isotropy need not be assumed. The strata can be defined by orientation as well as distance in the spirit of Granger (1969) because equations (2.7) and (2.9) can be easily adapted to any partition of the set of pairs of locations. Only simple calculations of autocorrelations, which can be done on small computers, are involved.

The detection of association between spatial data sets is not a simple problem. In this paper we have investigated one particular approach. As well as having reservations about the assumption of stationarity, in a detailed statistical analysis it should be recognised that association can exist simultaneously at a number of different geographical scales. The correlation coefficient is a single omnibus statistic that averages the scale-dependent association. Thus, for example, it is possible that negative association at small scales is swamped by positive association at large scales. We have not attempted to explore the ways in which these scale-dependent associations can be separated out. This is an important area of research.

#### **ACKNOWLEDGEMENTS**

The authors wish to thank Annie Mollié and Nicole Le Moual for computing assistance. This work was supported by Euratom Contract BI6-0126-F.

#### Résumé

Des tests d'association modifiés. portant sur le coefficient de corrélation ou la covariance entre deux processus spatiaux autocorrelés, sont proposés. Ces tests peuvent être utilisés à la fois pour des données observées sur un lattice ou sur un domaine irrégulier. Il sont fondés sur l'évaluation d'un nombre d'observation corrigé qui prend en compte la structure spatiale de chaque processus.

Le nombre d'observations corrigé est inférieur à la taille de l'échantillon quand l'autocorrélation de chaque processus est positive. Une méthode pour évaluer cette réduction par l'intermédiaire d'une approximation de la variance du coefficient de corrélation est développée. La performance des tests est évaluée par des simulations de Monte-Carlo. La méthode est illustrée par des exemples d'épidémiologie géographique.

#### References

- Armstrong, B. and Doll, R. (1975). Environmental factors and cancer incidence and mortality in different countries, with special reference to dietary practices. *International Journal of Cancer* **15**, 617-631.
- Bartlett, M. S. (1935). Some aspects of the time-correlation problem in regard to tests of significance. Journal of the Royal Statistical Society 98, 536–543.
- Bartlett, M. S. (1978). Nearest neighbour models in the analysis of field experiments. Journal of the Royal Statistical Society, Series B 40, 147–174.
- Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society, Series B* 36, 192–236.
- Cliff, A. D. and Ord, J. K. (1975). Model building and the analysis of spatial pattern in human geography. Journal of the Royal Statistical Society, Series B 37, 297-348.
- Clifford, P. and Richardson, S. (1985). Testing the association between two spatial processes. *Statistics and Decisions*, Supp. issue 2, 155–160.
- Cook, D. G. and Pocock, S. J. (1983). Multiple regression in geographic mortality studies with allowance for spatially correlated errors. *Biometrics* 39, 361-371.
- Davies, H. M. and Jowett, G. H. (1958). The fitting of Markoff serial variation curves. Journal of the Royal Statistical Society. Series B 20, 120-142.
- Granger, C. W. J. (1969). Spatial data and time series analysis. In London Papers in Regional Science 1. Studies in Regional Science, A. J. Scott (ed.), 1–24. London: Pion.
- Haugh, L. D. (1976). Checking the independence of two covariance stationary time series: A univariate residual cross-correlation approach. *Journal of the American Statistical Association* **71**, 378-385.

Hoover, R. and Fraumeni, J. F. (1975). Cancer mortality in U.S. counties with chemical industries. Environmental Research 9, 196–207.

Hubert, L. J. and Golledge, R. G. (1982). Measuring association between spatially defined variables: Tjostheim's index and some extensions. *Geographical Analysis* 14, 273-278.

Mardia, K. V. and Marshall, R. J. (1984). Maximum likelihood estimation of models for residual covariance in spatial regression. *Biometrika* 71, 135-146.

Ord. K. (1975). Estimation methods for models of spatial interaction. Journal of the American Statistical Association 70, 120-126.

Richardson, S. and Hémon, D. (1981). On the variance of the sample correlation between two independent lattice processes. *Journal of Applied Probability* 18, 943-948.

Ripley, B. D. (1981). Spatial Statistics. New York: Wiley.

Schottenfeld, D. and Fraumeni, J. F. (1982). Cancer Epidemiology and Prevention. Philadelphia: W. B. Saunders.

Student (W. S. Gosset) (1914). The elimination of spurious correlation due to position in time or space. *Biometrika* 10, 179-181.

Tjøstheim, D. (1978). A measure of association for spatial variables. *Biometrika* 65, 109–114.

Whittle, P. (1954). On stationary processes in the plane. Biometrika 41, 434-449.

Received June 1987; revised June 1988.

#### APPENDIX 1

#### The Variance of r

Let  $\Sigma_{\xi}$  and  $\Sigma_{\eta}$  denote the variance–covariance matrices of the vectors with elements  $\xi_{\eta} = X_{\eta} - \overline{X}$  and  $\eta_{\eta} = Y_{\eta} - \overline{Y}$ .

If we write

$$r = \xi^{\mathrm{T}} \eta / (\xi^{\mathrm{T}} \xi \eta^{\mathrm{T}} \eta)^{1/2}$$

then

$$\operatorname{var}(r) = \operatorname{E} \operatorname{tr} \left\{ \frac{\xi \xi^{\mathsf{T}}}{\xi^{\mathsf{T}} \xi} \frac{\eta \eta^{\mathsf{T}}}{\eta^{\mathsf{T}} \eta} \right\} = \operatorname{tr} \left( \operatorname{E} \frac{\xi \xi^{\mathsf{T}}}{\xi^{\mathsf{T}} \xi} \operatorname{E} \frac{\eta \eta^{\mathsf{T}}}{\eta^{\mathsf{T}} \eta} \right)$$
(A.1)

by the independence of  $\xi$  and  $\eta$ .

Expanding  $\eta^{T} \eta$  about its expectation,  $k_{1} = tr(\Sigma_{\eta})$ , we have

$$\frac{\eta \eta^{\mathsf{T}}}{\eta^{\mathsf{T}} \eta} = \frac{\eta \eta^{\mathsf{T}}}{k_1} \left[ 1 - \frac{\eta^{\mathsf{T}} \eta - k_1}{k_1} + \frac{(\eta^{\mathsf{T}} \eta - k_1)^2}{k_1^2} - \cdots \right].$$
(A.2)

Note that in the case  $\Sigma_N = I$  we have  $k_1 = N - 1$  so that we would be expanding in inverse powers of N - 1. To evaluate the expectation of (A.2) we calculate

$$\mathrm{E}(\eta\eta^{\mathrm{T}})^{2} = 2\Sigma_{\eta}^{2} + k_{\mathrm{I}}\Sigma_{\eta}$$

and

$$\mathsf{E}(\boldsymbol{\eta}\boldsymbol{\eta}^{\mathsf{T}})^{3} = 8\boldsymbol{\Sigma}_{\boldsymbol{\eta}}^{3} + 2k_{2}\boldsymbol{\Sigma}_{\boldsymbol{\eta}} + 4k_{1}\boldsymbol{\Sigma}_{\boldsymbol{\eta}}^{2} + k_{1}^{2}\boldsymbol{\Sigma}_{\boldsymbol{\eta}},$$

where  $k_2 = \operatorname{tr}(\Sigma_{\pi}^2)$ .

If we make the plausible assumption that the maximum eigenvalue of  $\Sigma_{\eta}$  is bounded above by  $\lambda_{\rm B}$  as  $N \to \infty$ , then using the inequality tr( $\Sigma_{\eta}^2$ )  $\leq \lambda_{\rm B}$  tr( $\Sigma_{\eta}$ ), we have the second-order asymptotic expression

$$\mathsf{E}\left(\frac{\eta\eta^{\mathsf{T}}}{\eta^{\mathsf{T}}\eta}\right) = \frac{\Sigma_{\eta}}{k_1} - 2\frac{\Sigma_{\eta}^2}{k_1^2} + 2\frac{\Sigma_{\eta}k_2}{k_1^3}.$$
 (A.3)

Finally, substituting (A.3) and a similar expression for  $E(\xi\xi^T/\xi^T\xi)$  into (A.1), we have to the first order

$$\sigma_r^2 = \frac{\operatorname{tr}(\boldsymbol{\Sigma}_{\eta}\boldsymbol{\Sigma}_{\xi})}{k_1 j_1}, \qquad (A.4)$$

and to the second order

$$\sigma_r^2 = \frac{\operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{\eta}}\boldsymbol{\Sigma}_{\boldsymbol{\xi}})}{k_1 j_1} - \frac{2\operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{\eta}}^2 \boldsymbol{\Sigma}_{\boldsymbol{\xi}})}{k_1^2 j_1} + \frac{2\operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{\eta}}\boldsymbol{\Sigma}_{\boldsymbol{\xi}})}{k_1 j_1} \left(\frac{k_2}{k_1^2} + \frac{j_2}{j_1^2}\right) - \frac{2\operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{\eta}}\boldsymbol{\Sigma}_{\boldsymbol{\xi}}^2)}{j_1^2 k_1}.$$
 (A.5)

where  $j_n = \operatorname{tr}(\Sigma_{\varepsilon}^n)$ .

In the case  $\Sigma_{\pi} = \Sigma_{\xi}$  the second-order term in (A.5) becomes

$$-4\left\{\frac{\operatorname{tr}(\Sigma_{\pi}^{3})}{k_{1}^{3}} - \frac{\operatorname{tr}(\Sigma_{\pi}^{2})\operatorname{tr}(\Sigma_{\xi}^{2})}{k_{1}^{4}}\right\} = -4\left\{\sum_{i=1}^{N-1} p_{i}^{3} - \left(\sum_{i=1}^{N-1} p_{i}^{2}\right)^{2}\right\}$$
$$= -4\left\{\sum_{i=1}^{N-1} \left(p_{i} - \sum_{i=1}^{N-1} p_{i}^{2}\right)^{2} p_{i}\right\} \le 0,$$

where  $p_i = \mu_i / \sum_{j=1}^{N-1} \mu_j$ ,  $\{\mu_i\}_{i=1}^{N-1}$  being the eigenvalues of  $\Sigma_{\pi}$  and  $\Sigma_{\xi}$  (matrices of rank  $\leq N-1$ ), i.e., the first-order approximation is no larger than the second-order approximation.

Note that in the case where (i)  $\Sigma_{\xi}$  and  $\Sigma_{\eta}$  commute, (ii)  $\Sigma_{\eta}$  is proportional to an idempotent matrix, and (iii) the eigenvalues of  $\Sigma_{\xi}$  are zero whenever the associated eigenvalues of  $\Sigma_{\eta}$  are zero, then the first-order term in (A.4) is exactly equal to  $\sigma_r^2$ . This follows from the remark that in this case r has p.d.f.  $f_M(r)$  with effective sample size  $M = 1 + \operatorname{rank}(\Sigma_{\eta})$ , where  $f_M(r)$  is the standard distribution of r defined in (2.2).

#### Appendix 2

#### Simulation of Simultaneous Autoregressive Process

For a process given by (3.1) with zero boundary, i.e.,  $X_{s,0} = X_{0,t} = X_{n+1,t} = X_{s,n+1} = 0$  (s,  $t = 0, 1, \dots, n + 1$ ), the spectral decomposition of  $\Sigma_{\infty}$  can be deduced from the expression given by Besag [discussion of a paper by Bartlett (1978)]. The process can be simulated by computing  $\mathbf{P}^{\mathsf{T}}\operatorname{diag}(\delta_t^{1/2})\mathbf{Z}$ , where Z is a vector of independent N(0, 1) variables and where  $\mathbf{P}^{\mathsf{T}}\operatorname{diag}(\delta_t)\mathbf{P}$  is the spectral decomposition of  $\Sigma_{\infty}$ . If we denote the matrix  $\{X_{s,t}\}_{s,t=1}^{n}$  by X and rearrange the vector of Z into an  $n \times n$  matrix whose elements are  $Z_{t_t}$  ( $i, j = 1, 2, \dots, n$ ), then, exploiting the properties of P, we find that X is simulated by

$$\mathbf{X} = \mathbf{Q} \{ Z_{ii} \lambda_{ij} \} \mathbf{Q}.$$

where

$$\lambda_{ij} = 1 - 2a \{ \cos[\pi i/(n+1)] + \cos[\pi j/(n+1)] \}$$

and

$$Q_{i} = [2/(n+1)]^{1/2} \sin[\pi i j/(n+1)], n \text{ even.}$$

The computational cost of a single realisation is therefore of the order of  $n^3$  operations or  $N^{3/2}$  since  $N = n^2$ . In contrast, a single triangulation  $LL^T$  of  $\Sigma_X$  requires an order of  $N^3$  operations and each realisation involves an order of  $N^2$  operations.

# TESTING ASSOCIATION BETWEEN SPATIAL PROCESSES

### TESTING ASSOCIATION BETWEEN SPATIAL PROCESSES

## Sylvia RICHARDSON\*, Peter CLIFFORD+

## 1. Introduction

This paper is concerned with a topic in spatial statistics, that of testing for association between two spatial processes. This question arises frequently in various fields and examples abound in geography and regional sciences (Cliff and Ord, 1981), geology (Malin and Hyde, 1982), sociology (Doreian, 1981)... In the epidemiology of chronic diseases, etiological clues are sometimes sought by studying joint geographic variations of environmental risk factors and disease rates (Doll 1980). The examples discussed in the later part of this paper will be taken from this field.

Throughout the data will consist of a set A of N locations and pairs of variables  $(X_{\alpha}, Y_{\alpha})$ ,  $\alpha \in A$ , indexed by their location. These variables will be spatially autocorrelated.

In the first part a method for approximating the critical values of the product moment correlation coefficient  $r_{XY}$  will be summarised. This method has been described in detail in Clifford, Richardson, Hémon (1989).

In the following section we give some complementary results on the performance of the tests for small domains and on their power. We also discuss the influence of the choice of the partition of the covariance structure. Finally we give some examples, comparing our results with those of Monte Carlo tests and giving pivotal confidence intervals for the regression coefficient.

## 2. Modified tests of associations

We have devised modified tests of association based either on  $s_{XY}$ : the empirical covariance between pairs of observations  $(X_{\alpha}, Y_{\alpha})$ ,  $\alpha \in A$ , or based on  $r_{XY}$ : the corresponding empirical correlation coefficient. We shall use the notation :

 $\overline{X} = N^{-1}(\Sigma X_{\alpha}), \quad s_{XY} = N^{-1} \Sigma (X_{\alpha} - \overline{X})(Y_{\alpha} - \overline{Y}), \quad s_X^2 = N^{-1} \Sigma (X_{\alpha} - \overline{X})^2$ 

and similarly for Y and  $s_Y^2$ .

<sup>\*</sup> INSERM U.170, 16 Avenue Paul Vaillant-Couturier 94807 VILLEJUIF Cedex, France

<sup>+</sup> Mathematical Institute, University of Oxford,

<sup>24-29</sup> St Giles, Oxford OXI 3LB

### Method

Suppose that X and Y are independent but that both X and Y are multivariate normal vectors with constant means and variance-covariance matrices  $\Sigma_X$  and  $\Sigma_Y$  respectively. A stratified structure for  $\Sigma_X$  and  $\Sigma_Y$  is imposed. Pairs in A x A are divided into strata S0,S1, S2,... such that the covariances within strata remain constant,

ie 
$$\operatorname{cov} (X_{\alpha}, X_{\beta}) = C_X(k)$$
 if  $(\alpha, \beta) \in S_k$ 

An estimate of the conditional variance of  $s_{XY}$  is then derived :

$$N^{-2} \sum_{k} N_k \hat{C}_X(k) \hat{C}_Y(k) , \qquad (1)$$

 $N_k$  is the number of pairs in strata  $S_k$  and  $\hat{C}_X(k)$  (respectively  $\hat{C}_Y(k)$ ) is the estimated autocovariance :

$$\hat{C}_{X}(k) = \sum_{\substack{X \\ S_{k}}} (X_{\alpha} - \overline{X}) (X_{\beta} - \overline{X}) / N_{k}$$

Thus the estimate takes into account the autocorrelation of both X and Y.

Further it can be shown that, to the first order, the variance of  $r_{XY}$ ,  $\sigma^2 r$  is :

$$\sigma_{r}^{2} = \frac{Var(s_{XY})}{E(s_{X}^{2}) E(s_{Y}^{2})}$$
, (2)

which leads to the following estimate :

$$\hat{\sigma}_{r}^{2} = \frac{\sum N_{k} \hat{C}_{X}(k) \hat{C}_{Y}(k)}{N^{2} s_{X}^{2} s_{Y}^{2}}$$
(3)

Note that an equivalent expression for  $Var(s_{XY})$  is N<sup>-2</sup> tr  $(\Sigma_{\xi} \Sigma_{\eta})$  where  $\Sigma_{\xi}$  and  $\Sigma_{\eta}$  are the variance-covariance matrices of the centered vectors X -  $\overline{X}$  and Y -  $\overline{Y}$ .

In the classical non autocorrelated case, when either  $\Sigma_X$  or  $\Sigma_Y = I$ , it can be shown that the approximation given by (2) is exact and that  $r_{XY}$  follows a t-distribution with N-2 d.f ( $t_{N-2}$ ). Further :  $N = 1 + (\sigma^2_T)^{-1}$ 

In general, an estimated effective sample size,  $\hat{M}$ , is defined by the relationship

$$\hat{M} = 1 + (\hat{\sigma}^2 r)^{-1}$$

where  $\hat{\sigma}_{T}^{2}$  is given by (3). A modified t-test :  $t_{M-2}^{A}$  is proposed which rejects the null hypothesis of no association when :

$$|(\hat{M}-2)^{1/2} r(1 - r^2)^{-1/2}| > t^{\alpha} \hat{M}-2$$

where  $t^{\alpha}_{M-2}$  is critical value of the **t**-statistic with  $\hat{M}-2$  d.f.

Equivalently a standardised covariance can be used :

W = N s<sub>XY</sub> (
$$\Sigma$$
 N<sub>k</sub>  $\hat{C}_X(k)$   $\hat{C}_Y(k)$ )<sup>-1/2</sup>

and tested as a standard normal relying upon central limit theorems for spatially dependent variables.

## Results on the performance of W and tM-2

The performance of these tests under the null hypothesis of stochastic independence between X and Y was first assessed by Monte Carlo simulations for two models :

a) X and Y were generated on 3 lattice sizes (12x12, 16x16, 20x20) as nearest neighbour isotropic autoregressive gaussian processes, with 1st order autocorrelation  $\rho_X(1)$  and  $\rho_Y(1)$  ranging from 0.2 to 0.8;

b) X and Y were generated on the grid of the administrative centres of the French départements as gaussian variables with a disc model for their autocovariance (Ripley 1981) and arbitrarily defined 1st order autocorrelation ranging from 0.2 to 0.9.

In both cases and for several levels of autocorrelation in X or Y, 500 trials were executed with a nominal rejection level of 5%.

For the two statistics,  $t_{M-2}$  and W, the type I errors found did not vary in any systematic way with the level of autocorrelation and fluctuated around the nominal 5% level (the confidence interval excluded 5% in only 2 cases out of 45 simulated in (a)). This contrasted with the performance of the standard t-test procedure based on N-2 d.o.f. where the type I error increased systematically with the autocorrelation to reach values around 50 % in the highly autocorrelated cases (Clifford, Richardson, Hémon, 1989).

## 3. Comparison of W and $t_{M-2}$ on small lattices

Based on our first results the performance of the statistics  $t_{M-2}$  and W seemed indistinguishable. The number of sample points was reasonably large and a difference between their respective performances could be better highlighted by studying smaller samples.

Results from further simulations on smaller lattices of sizes 6x6, 8x8, 10x10 are shown in Figure 1. The type I errors of the W statistic exhibit a systematic downward trend with increasing autocorrelation which is not apparent for the modified  $t_{M-2}$  statistic. Consequently, for small domains, the use of the modified  $t_{M-2}$  statistic is preferable to that of W. We also note that the convergence of W to normality is slower as the autocorrelation increases. Figure 2 shows a Q.Q. plot of W, i.e. quantiles of a standard normal distribution against the sample quantiles of W for 500 independent trials, in the case  $\rho_X(1) = \rho_Y(1) = 0.8$  and a 12x12 lattice. We note that the distribution of the W statistic has short tails compared to the normal distribution. The departure from normality is confirmed by a Kolmogorov-Smirnov test which is significant at the 5% level but not at the 1% level.

## 4. Power of the modified tests

The power of the modified tests was assessed under a simple alternative hypothesis of a linear regression between Y and X :  $H_1$  : Y = aX + W, X ~ N( $\mu_X$ ,  $\Sigma_X$ ), W ~ N( $\mu_W$ , $\Sigma_W$ ) and X and W independent. It is difficult to calculate theoretically the power of the modified tests because their distribution under  $H_1$  is not precisely known. Their power can be assessed by simulations.

Two independent spatially autocorrelated processes X and W were generated on the grid of the administrative centers of French départements as gaussian variables with a disc model for their autocovariance. Without loss of generality  $\sigma^2_X = \sigma^2_W$  was chosen and hence the correlation  $\rho_{XY}$  between X and Y was only dependent on the parameter a. Five hundred trials were carried out for several levels of autocorrelation in X and W and for the values  $\rho_{XY} = 0.2$ and 0.4. The grid contained N = 82 points. Results for higher values of  $\rho_{XY}$  are not reported because the power of the tests was very close to 1. In these simulations, the power of W and  $t_{M-2}$ , was evaluated with a 5 % nominal level.

On the other hand it might be interesting to calculate the power  $\pi_T(s_{XY})$  of a test on the covariance  $s_{XY}$ , similar to W but where the estimate : N<sup>-2</sup>  $\sum N_k \stackrel{\frown}{C}_X(k) \stackrel{\frown}{C}_Y(k)$ , of the variance of  $s_{XY}$  is replaced by its theoretical value under H1 :

N<sup>-2</sup> tr 
$$(\Sigma_{\xi} \Sigma_{\eta}) = N^{-2} [a^2 tr(\Sigma_{\xi}^2) + tr(\Sigma_{\xi} \Sigma_{\theta})].$$

where  $\Sigma_{\theta}$  denotes the variance-covariance matrix of the centered vector W -  $\overline{W}$  and  $\Sigma_{\xi}$  and  $\Sigma_{\eta}$  are defined as before.

In order to carry out the calculation of  $\pi_T(s_{XY})$ , it is necessary to suppose that the distribution of  $\sqrt{N} s_{XY}$  under H1 is approximately normal. This approximation can be justified by central limit theorems if appropriate hypotheses are placed on the rate of decrease of the autocovariances of X and W as a function of the lag.

The expectation and the variance of  $s_{XY}$  under H1 are given by :

$$\begin{split} &E_{H1} \left( Ns_{XY} \right) = a \ tr \ (\Sigma_{\xi}) \\ &V_{H1} \left( Ns_{XY} \right) = 2a^2 \ tr \ (\Sigma_{\xi}^2) + tr \ (\Sigma_{\xi} \Sigma_{\theta}). \end{split}$$

The traces of the matrices  $\sum_{\xi} \sum_{\theta}$  or their product can be expressed in terms of  $\sum_{X}$  and  $\sum_{W}$  and thus evaluated for specific models for  $\sum_{X}$  and  $\sum_{W}$ .

The power  $\pi_T(s_{XY})$  of the test of the covariance using the statistic Ns<sub>XY</sub> tr( $\Sigma_{\xi} \Sigma_{\eta}$ )<sup>-1/2</sup> for a bilateral test of nominal level  $\alpha$  is thus equal to 1 - [ $\Phi(M_1) - \Phi(M_2)$ ] with :

$$M_{1} = \frac{C_{\alpha} \left[a^{2} \operatorname{tr} (\Sigma_{\xi}^{2}) + \operatorname{tr} (\Sigma_{\xi}\Sigma_{\theta})\right]^{1/2} - a \operatorname{tr}(\Sigma_{\xi})}{\left[2a^{2} \operatorname{tr} (\Sigma_{\xi}^{2}) + \operatorname{tr} (\Sigma_{\xi}\Sigma_{\theta})\right]^{1/2}}$$
$$M_{2} = \frac{-C_{\alpha} \left[a^{2} \operatorname{tr} (\Sigma_{\xi}^{2}) + \operatorname{tr} (\Sigma_{\xi}\Sigma_{\theta})\right]^{1/2} - a \operatorname{tr}(\Sigma_{\xi})}{\left[2a^{2} \operatorname{tr} (\Sigma_{\xi}^{2}) + \operatorname{tr} (\Sigma_{\xi}\Sigma_{\theta})\right]^{1/2}}$$

 $\Phi$  being the N(0,1) distribution function and  $C_{\alpha}$  being such that  $P\{|N(0,1)| > C_{\alpha}\} = \alpha$ .

It is also interesting to be able to compare the power observed by simulations to a reference value. Along these lines, we thus also calculated the power of the classical test of  $r_{XY}$ ;  $\pi_{N*}(r)$ , in a case which would be compatible with the observed empirical variance of  $r_{XY}$ ,  $v_e$ , estimated by the Monte Carlo simulations. Recall that in the case of non autocorrelated variables X and Y and large samples, the variance of  $r_{XY}$  is approximately equal to  $(1 - \rho^2_{XY})^2 / N-1$  for a sample of N observations. For autocorrelated X and Y, we thus computed an approximately equivalent sample size, N\*:

$$N^* = 1 + (1 - \rho^2 X Y)^2 / v_e$$

This number N\* was used to compute the reference value,  $\pi_{N*}(r)$ , power of the classical test of  $r_{XY}$  based on N\* observations.

A summary of the results is given in Table1a and 1b. In those tables, the observed power of W is only given since it is almost identical to that of  $t_{M-2}$ . Overall all the powers, whether observed or calculated, are close. The only differences are seen a few cases of high autocorrelation for either X and W.

In summary we can say that the "theoretical power"  $\pi_T(s_{XY})$  gives a good approximation of the observed power in most cases and does not require Monte-Carlo simulations. Furthermore, the modified W and  $t_{M-2}$  tests have comparable power to that of a classical test based on an "equivalent" number of observations N\*.

## 5. Choice of the partition for the covariance structure of X and Y.

Computations needed in order to apply the modified tests to a data set are straight forward but rely upon a choice of strata  $\{S_0, S_1, S_2, ...\}$  of AxA, on each of which the covariance of the two processes is assumed to be constant.

In Table 2 the performance of the  $t_{M-2}$  statistic is investigated for different choices of strata in the case of the irregular grid network and gaussian disc models. Six partitions were defined, ranging from 5 to 21 strata and corresponding to different discretisations of the distance between pairs of locations (assuming isotropy).

Overall the type I error was close to 5 % for most of the partitions. As the autocorrelation increased, the type I error for the 5-strata partition was inflated whereas there was a stability in the observed error rate when the number of strata increased. The results for the W statistic were similar.

This confirmed that a balance has to be reached between choosing too few classes which can bias expression (1) or too many resulting in less precise estimates of the autocovariances. Clearly in the case simulated where the autocorrelation decreases smoothly with distance, the performance of the  $t_{M-2}$  statistic is robust to various choices of partitions.

## 6. Confidence interval for regression coefficient

Once we have a test for independence it is possible to construct a confidence interval for the regression coefficient b, in the model

$$Y = a1 + bX + Z$$

where Z is a process independent of X and 1 is a vector with unit elements. The confidence interval for b is the set of values of b which we would not actually reject i.e the set of b such that Y - bX has no significant correlation with X. Defining :

 $f_{\alpha} = X_{\alpha} - \overline{X}$  and  $g_{\alpha} = Y_{\alpha} - \overline{Y}$  the standardised covariance between Y - bX and X is :

$$W_{b} = \frac{(g^{T} f - b f^{T} f)}{\sqrt{\sum N_{k} \hat{C}_{X}(k) \hat{C}_{Y-bX}(k)}}$$

where  $\hat{C}_{Y-bX}(k) = \hat{C}_Y(k) + b^2 \hat{C}_X(k) - 2b N_k^{-1} \sum_{S_k} f_{\alpha} g_{\beta}$ .

This standardised covariance can be used as a pivot to obtain a confidence interval for b. We do not reject the null hypothesis when

$$|W_b| < C_\alpha$$
 where :  $P\{|N(0,1)| > C_\alpha\} = \alpha$ .

Therefore the confidence interval is :

$$\{b: (g^{T}f - bf^{T}f)^{2} \le C_{\alpha}^{2} \sum N_{k} \hat{C}_{X}(k) \hat{C}_{Y-bX}(k)\}$$
or b - T<sub>2</sub> ± (T<sub>2</sub><sup>2</sup> + T<sub>1</sub> - T<sub>1</sub>T<sub>3</sub> - 2bT<sub>2</sub> + b<sup>2</sup>T<sub>3</sub>)<sup>1/2</sup> (1 - T<sub>3</sub>)<sup>-1</sup>
(4)

where 
$$b = s_{XY} / s_X^2$$
;  $T_1 = d_X \sum_k N_k \hat{C}_Y(k) \hat{C}_X(k)$ ;  $T_2 = d_X \sum_k N_k \hat{C}_X(k) \hat{C}_{XY}(k)$   
 $T_3 = d_X \sum_k N_k \hat{C}_X^2(k)$  and where  $\hat{C}_{XY}(k) = \sum_k f_{\alpha} g_{\alpha} / N_k$  and  $d_X = C_{\alpha}^2 N^{-2} s_X^{-4} \cdot k$ 

## 7. Examples

Our examples concern the relationship between lung cancer, smoking and industrial factors. We calculate W,  $t_{M-2}$  and a confidence interval for the regression coefficient for each example. To provide an additional check on the performance of our tests, we also carried out a Monte Carlo test based on the disc model for cases in which such a model is plausible.

## The data

For 82 départements we considered male lung cancer mortality rate (LC) over a 2 year period, 1968-1969, standardised over the age 35-74, cigarette sales per inhabitant (CS) in 1953 (a fifteen year time lag was chosen to account for the delay between exposure and the onset of the pathology) and demographic data on the percentage of employed males in the metal industry (MW) and the textile industry (TW) recorded by census in 1962.

The coordinates of the points of the network were identified with the geographical locations of the administrative centers ("préfectures") of French départements. The spatial structure of these variables was investigated by means of a variogram. In this analysis, N = 82 locations were retained after grouping the départements around Paris into one area. The distances between the centres of départements were partitionned into 15 classes of 50 kilometres intervals each. This gives 15 strata S<sub>1</sub>, ..., S<sub>15</sub>.

The observed variograms, i.e the plot of  $N_k^{-1} \Sigma$   $(X_{\alpha} - X_{\beta})^2$ ,

$$(\alpha,\beta) \in S_k$$

k = 1,... 15, against the average distance,  $d_k$ , for départements in S<sub>k</sub>, for the four variables considered are shown in Figure 3. Note that the last two classes contain few pairs and hence have a large variability. Three of the variograms (LC, CS, MW) exhibit clearly an upward trend with increasing distance. Up until the 10th class, the shape of this trend is fairly linear with increasing distance, thus compatible with the disc model for the covariance matrix discussed by Ripley (1981). The variogram for TW gives little indication of any spatial autocorrelation.

## **Results**

Using the standard test base on  $r_{XY}$  there is a highly significant positive association both between LC and CS and between LC and MW. The association between LC and TW on the other hand is less strong but still significant at the 1 % level (Table 3). The W and t $\hat{M}$ -2 statistics for these 3 examples are shown in Table 3. One can see a substantial reduction of the degrees of freedom when the autocorrelation is taken into account. We note that this occurs also for the case LC-TW which is surprising if one recalls the shape of the variogram of TW. A possible explanation for this is that the geographically small départements, which are over-represented in the first few strata, are atypical for this particular variable (TW). For CS and MW the effective sample size is about 20 % of the original sample size. Consequently the significance levels are reduced but even after this "adjustement" these two factors are statistically significantly associated with lung cancer. For TW the significance disappears after adjustement.

The first two examples were also investigated by a Monte Carlo test. A disc model for the covariance was fitted by maximum likelihood. The parameters were found by direct search and corresponded to autocorrelation  $\rho$  (1) of 0.91, 0.85 and 0.91 respectively for LC, CS and MW.

For each example, 1000 pairs of mutually independent variables, with covariance given by the estimated disc model, were generated and the observed correlation coefficient was ranked among the 1000 generated coefficients. The significance levels obtained are given in Table 3. The agreement between them and the significance levels of the W or  $t_{M-2}$  tests is better for CS than for MW; this is possibly due to a better fit of the disc model for the CS variable. Confidence intervals, given by (4), for the regression coefficients were also calculated. We note that they are not symmetric.

## 8. Discussion

In this paper we have studied the properties of modified tests of the empirical correlation coefficient between two spatial processes. We have shown that by a simple adjustment, correct level of significance can be reached and that the power under a simple linear alternative is compatible with that of a standed test in an equivalent situation. These tests can be applied both to regularly and irregularly spaced points and can be considered as a first step in an analysis of association when detailed spatial modelling is not suitable.

The performance did not vary much when different strata of equal covariance were chosen. On small lattices, the modified  $t_{M-2}$  statistic is better than the standardised covariance.

Application of these tests to data may also give pivotal confidence interval for the regression coefficient. Furthermore, if one is prepared to model the observed covariance structure, Monte Carlo tests of association can be performed. In the examples investigated, the results from the two types of testing procedures were close.

It would be interesting to develop distribution-free tests of association based on permutations and to compare their performance to that of the proposed modified tests in the case of non gaussian spatial distributions.

## REFERENCES

- Cliff, A.D. and J.K. Ord (1981), Spatial Processes. Models and Applications. London, Pion.

- Clifford, P. Richardson, S. and Hémon, D. Assessing the significance of the correlation between two spatial processes, (1989) to appear in Biometrics.

- Doll, R. (1980). The epidemiology of cancer. Cancer, 45, 2475-2485.

- Doreian, P. (1981), Estimating linear models with spatially distributed data, in Sociological Methodology, edited by Leinhardt and Samuel. San Francisco, Jossey-Bass Publishers, 359-388.

- Malin, S.R.C. and Hide, R. (1982). Bumps on the core-mantle boundary : geomagnetic and gravitational evidence revisted. Phil. Trans. R. Soc. Lon A, 306, 281-289.

- Ripley, B.D. (1981), Spatial Statistics. New York, Wiley.









QQ plot (sample quantiles (Y) against quantiles of a standard normal distributions (X)) of 500 trials for the W statistic with two mutually independent simultaneous autogressive processes (12x12 lattice,  $\rho_X(1) = \rho_Y(1) = 0.8$ )

VARIOGRAMS OF THE VARIABLES CONSIDERED IN &7



Fifteen classes of distance are considered. The number of pairs in each class is : 82, 400, 582, 674, 764, 822, 812, 726, 630, 476, 304, 172, 94, 58, 40.

## Table 1a

Power of the modified tests: results concerning the testing of the correlation between X and Y = aX + W where both X and Y are spatially autocorrelated, X and W independent and of equal variance and a is chosen so that the correlation  $\rho_{XY}$  between X and Y takes the value 0.2.

		ρ <sub>xy</sub> = 0.2					
ρ <sub>W</sub>	$\rho_X$	0.0	0.2	0.4	0.6	0.8	
	N*	79	82	78	79	84	
0.0	π <sub>N*</sub> (r) power of W π <sub>T</sub> (s <sub>XY</sub> )	0.42 0.42 0.44	0.44 0.44 0.44	0.42 0.43 0.43	0.42 0.39 0.42	0.44 0.35 0.37	
		76	74	67	71	63	
0.2		0.41 0.42 0.44	0.40 0.42 0.41	0.36 0.39 0.39	0.38 0.34 0.36	0.34 0.30 0.32	
		80	78	66	57	54	
0.4		0.42 0.44 0.44	0.42 0.36 0.39	0.36 0.37 0.36	0.32 0.34 0.32	0.30 0.21 0.27	
		78	75	61	46	35	
0.6		0.41 0.48 0.45	0.40 0.36 0.39	0.33 0.36 0.34	0.25 0.28 0.27	0.20 0.20 0.20	
		78	66	54	39	19	
0.8		0.41 0.52 0.48	0.36 0.48 0.40	0.30 0.42 0.33	0.22 0.28 0.23	0.12 0.12 0.13	

500 simulations were carried out. The observed power of W is compared with  $\pi_{N*}(r)$  and  $\pi_{T}(s_{XY})$  (cf § 4). The observed power of  $t_{M-2}$  is almost identical to that of W. Standard deviations and confidence intervals for the observed proportions can be calculated according to binomial sampling.

## Table 1b

Power of the modified tests: results concerning the testing of the correlation between X and Y = aX + W where both X and Y are spatially autocorrelated, X and W independent and of equal variance and a is chosen so that the correlation  $\rho_{XY}$  between X and Y takes the value 0.4.

		$p_{xy} = 0.4$					
ΡW	$\rho_X$	0.0	0.2	0.4	0.6	0.8	
	N*	83	87	81	70	63	
0.0	π <sub>N*</sub> (r) power of W π <sub>T</sub> (s <sub>XY</sub> )	0.97 0.97 0.94	0.97 0.96 0.93	0.96 0.96 0.92	0.93 0.92 0.89	0.90 0.87 0.74	
		82	73	76	68	55	
0.2		0.96 0.96 0.94	0.94 0.96 0.91	0.95 0.94 0.90	0.90 0.92 0.85	0.86 0.83 0.70	
		84	67	61	63	49	
0.4		0.97 0.97 0.94	0.92 0.92 0.91	0.89 0.91 0.87	0.90 0.88 0.81	0.81 0.79 0.65	
		83	69	66	44	30	
0.6		0.97 0.96 0.94	0.93 0.95 0.90	0.92 0.93 0.85	0.77 0.81 0.74	0.58 0.56 0.54	
		68	52	51	33	20	
0.8		0.92 0.97 0.95	0.83 0.93 0.91	0.82 0.89 0.84	0.63 0.76 0.68	0.39 0.42 0.38	

500 simulations were carried out. The observed power of W is compared with  $\pi_{N^*}(r)$  and  $\pi_T(s_{XY})$  (cf § 4). The observed power of  $t_{M-2}$  is almost identical to that of W. Standard deviations and confidence intervals for the observed proportions can be calculated according to binomial sampling.

# PERCENTAGE OF TYPE I ERRORS OF THE $t_{M-2}^{A}$ STATISTIC FOR DIFFERENT PARTITIONS OF THE COVARIANCE STRUCTURE

Number of strata in each partition	5	9	13	15	17	21
$\rho_X^{(1)} = \rho_Y^{(1)} = 0$	0.056	0.052	0.05	0.056	0.054	0.054
$\rho_{\rm X}^{(1)} = \rho_{\rm Y}^{(1)} = 0.2$	0.04	0.04	0.036	0.032	0.034	0.03
$\rho_X^{(1)} = \rho_Y^{(1)} = 0.4$	0.062	0.046	0.05	0.05	0.052	0.05
$\rho_{\rm X}^{(1)} = \rho_{\rm Y}^{(1)} = 0.6$	0.066	0.058	0.054	0.056	0.052	0.052
$\rho_{\rm X}^{(1)} = \rho_{\rm Y}^{(1)} = 0.8$	0.068	0.058	0.048	0.04	0.046	0.046

	r	t <sub>N-2</sub> *	W	Ŵ	t <sub>M-2</sub>	Monte + Carlo	95%CI for ¥
Cigarette sales per inhabitant (1953) (CS)	0.76	10.48 P≃10 <sup>-21</sup>	2.94 p=0.0032	15	4.22 p=0.001	2/1000	0.78 [0.54;0.88]
<pre>% male workers in metal industry (1962) (MW)</pre>	0.63	7.16 p≃10 <sup>-11</sup>	2.48 p=0.0136	16	3.00 p=0.01	45/1000	0.29 [0.11;0.36]
<pre>% male workers in textile industry (1962) (TW)</pre>	0.28	2.57 p≅0.01	1.51 p=0.13	30	1.52 p=0.15	-	0.18 [-0.07;0.37]

Table 3: Comparison of the significance levels for tests of the association between lung cancer mortality rates and several risk factors given by standard test, W and  $t_{M-2}^{2}$  tests and Monte Carlo simulations,  $\hat{\gamma}$  is the estimated regression coefficient.

- \* standard test (t transformation with 80 d.f.)
- + Monte Carlo significance levels over 1000 simulations.
# A METHOD FOR TESTING THE SIGNIFICANCE OF GEOGRAPHICAL CORRELATIONS WITH APPLICATION TO INDUSTRIAL LUNG CANCER IN FRANCE

# A method for testing the significance of geographical correlations with application to industrial lung cancer in France

## Sylvia RICHARDSON

# SUMMARY

This paper discusses some of the problems encountered when using tests of significance in geographical epidemiology studies where the variables analysed will typically exhibit some spatial autocorrelation.

A test of partial correlations between spatially autocorrelated variables is presented. This test is based on an evaluation of an effective sample size which takes into account the spatial structure. Its performance is assessed via Monte Carlo simulations. The method proposed is applied to studying the relationship between male lung cancer rate and specific industries.

## Key words :

significance tests, partial correlation, spatial processes, geographical epidemiology, lung cancer.

#### 1. Introduction

Ecological correlation studies aim to analyse the association between a set of variables defined on groups. Among these, geographical correlation studies are particularly concerned with variables measured as averages over geographical units and the study of their joint variations.

In epidemiology, studies of the geographical distribution of incidence or mortality rates for particular diseases have been widely used for obtaining some clues about the etiology of those diseases<sup>1,2,3</sup>. In particular, the correlation between the spatial variation of exposure indicators, like dietary practice or industrial employment, and health indicators have been investigated<sup>4,5</sup>. Ecological analysis has also been used to check epidemiological hypotheses suggested by other approaches such as animal experiments, case-control or cohort studies<sup>6,7</sup>.

Examples of geographical correlation studies are also found in other domains like regional science, for instance when relating agricultural output and road accessibility<sup>8</sup>, or sociology with reference for instance to voting behaviour and socio-economic factors<sup>9</sup>.

When considering results from such studies, it is necessary to be aware of the methodological problems which arise from their design. These problems fall broadly into two categories : problems of interpretation due to the nature of the data considered, and statistical problems connected with the spatial structure of the data. Indeed, the data units can rarely be considered as independent replicates and often present some spatial autocorrelation.

In epidemiology, the limitations which stem from the nature of the data have been reviewed by many authors<sup>10-13</sup>. They are related to numerous factors : the potential non uniformity in diagnosis and registration of causes of death, the difficulty in adequately defining the population at risk, the difficulty in finding data on health and exposure on the same geographical units, the use of average

exposure when exposure is heterogeneous within the population and the difficulty in considering appropriate time lags between exposure and disease.

A discussion of the biases encountered when trying to estimate quantitatively relative risks from ecological studies is given in Richardson, Stücker and Hémon<sup>14</sup> and Greenland and Morgenstern<sup>15</sup>.

Furthermore, geographical data often exhibit a regular component in these spatial variations analogous to a trend or a gradient. Such gradients, if they are present in both the exposure and the health indicator rates, make the interpretation of the joint variations difficult since many potential confounder variables may also show the same regular spatial gradient. This is clearly illustrated in the study of cardiovascular mortality and water hardness of Pocock et al<sup>16</sup> where a northwest to southeast gradient is observed for both variables which is also related to climatic and industrial variables in Great Britain. At an individual level, this situation would be equivalent to having a confounding factor whose distribution nearly parallels that of the risk factor investigated. This has suggested a further step in the analysis : correlating the residual variations in the exposure and the health variables after removing a trend component or gradient which is obtained by regression of each of the variables on the coordinates of the geographical units. This is equivalent to computing the partial correlation between the exposure and the health variables after conditioning on the coordinates of the geographical units. A significant correlation between both the non-adjusted and adjusted observations would strengthen a causal interpretation.

Similarly one might want to adjust the relationship investigated on some geographically distributed confounding factors. Nevertheless, the residual variations after the removing of confounding factors and/or a spatial gradient might still be spatially autocorrelated. In a multifactorial chronic disease for instance, this residual autocorrelation might arise from the spatial structure of other unidentified or unmeasurable risk factors. In the case of autocorrelation among the residuals the classical methods for testing the significance of the association are not applicable.

In this paper, we will describe a method for assessing the significance of a partial correlation coefficent in the presence of spatial autocorrelation. This method is an extension of a modified test of correlation developped in Clifford, Richardson and Hémon<sup>17</sup> to the case of several variables. The proposed method is outlined in section 3, and some results on its performance are given which were obtained by Monte Carlo simulations. This section is preceded by a review of some available methods in geographical correlation studies where, in particular, the modified test of correlation between two spatially autocorrelated variables is briefly described as well as the fitting of multivariate regression models with a variance-covariance spatially parametrised error matrix. In a final section, the test proposed is applied to an example concerning the relationship between lung cancer mortality for men in France and particular industries. This relationship is investigated at the level of the French départements and adjustment is made on cigarette sales. The results clearly illustrate the interest of performing an adjustment for spatial autocorrelation. The effect of adjusting a linear spatial gradient is also discussed on the examples and through some simulations.

#### 2. <u>*Review*</u>

The data usually considered in geographical correlation studies is a set A of N locations and variables indexed by their locations. These variables will typically exhibit some degree of spatial autocorrelation. A first step in the analysis is often the calculation and testing of correlation coefficients between pairs of variables  $(X_{\alpha}, Y_{\alpha}), \alpha \in A$ .

The consequences of neglecting the autocorrelation in the X and Y variables were first pointed out by Bartlett<sup>18</sup> in the case of stationary time series. When the autocorrelations are positive, the significance level of classical tests

are overestimated. This phenomenon is more accentuated for spatial series. For instance, when  $X_{\alpha}$  and  $Y_{\alpha}$  are both spatially autocorrelated with nearest neighbour autocorrelation around 0.8, the observed significance level of a 5 % nominal test of the correlation coefficent can be increased up to 50 % <sup>17</sup>.

The situation is similar in regression analysis. The effect of neglecting the autocorrelation in the dependent variables or in the error term has been pointed out by several authors, Johnston<sup>19</sup> for time series, Cliff and Ord<sup>8</sup> for spatial series. Inflated values of the t and F statistics are again found when there is positive autocorrelation even though the estimation of the regression coefficient is not biased.

A modified test of the correlation coefficient  $r_{XY}$  has been proposed by Clifford, Richardson and Hémon<sup>17</sup> for testing the association between a pair of variables (X<sub> $\alpha$ </sub>,Y<sub> $\alpha$ </sub>). It consists in modifying the degrees of freedom of the classical test on  $r_{XY}$ . An effective sample size is calculated based on an estimation of the variance of  $r_{XY}$  which takes the internal autocorrelations of X and Y into account.

Suppose that X and Y are independent but that both X and Y are multivariate normal vectors with constant means and NxN variance-covariance matrices  $\Sigma_X$  and  $\Sigma_Y$  respectively. It can be shown that, to the first order, the variance of  $r_{XY}$ ,  $\sigma^2_r$  is :

$$\sigma^{2}r = \frac{1}{E(s_{X}^{2})} \frac{1}{E(s_{Y}^{2})}$$
(1)

where  $s_{XY}$  denotes the empirical covariance between pairs of observations ( $X_{\alpha}, Y_{\alpha}$ ),  $\alpha \epsilon A$  and  $s_X^2$  (respectively  $s_Y^2$ ) the empirical variances of X (respectively Y).

To be able to estimate var (s<sub>XY</sub>), one needs to impose a stratified structure for  $\Sigma_X$  and  $\Sigma_Y$ . More precisely pairs in A x A are divided into strata S0,S1, S2,... such that the covariances within strata remain constant, i.e. cov (X<sub> $\alpha$ </sub>, X<sub> $\beta$ </sub>) = C<sub>X</sub>(k), if ( $\alpha$ , $\beta$ )  $\varepsilon$  S<sub>k</sub>.

An estimate of the conditional variance of sXY is then derived :

$$N^{-2} \Sigma N_k \hat{C}_X(k) \hat{C}_Y(k) , \qquad (2)$$

where N<sub>k</sub> is the number of pairs in strata S<sub>k</sub> and  $\hat{C}_X(k)$  (respectively  $\hat{C}_Y(k)$ ) is the estimated autocovariance :

Thus the estimate takes into account the autocorrelation of both X and Y.

Equations (1) and (2) lead to the following estimate of the variance of r :

In the classical non autocorrelated case, when either  $\Sigma_X$  or  $\Sigma_Y = I$ , it can be shown that the approximation given by (1) is exact and that  $(N-2)^{1/2} r_{XY} / (1-r^2_{XY})^{1/2}$  follows a t-distribution with N-2 degrees of freedom (d.f.)  $(t_{N-2})$ . Further : N=1+( $\sigma^2_r$ )<sup>-1</sup>

In general an estimated effective sample size,  $\overset{\frown}{M}$  , is defined by the relationship

$$\hat{M} = 1 + (\hat{\sigma}^2 r)^{-1}$$

where  $\delta^2_r$  is given by (3). A modified t-test :  $t_{M-2}$  is proposed which rejects the null hypothesis of no association when :

$$|(\hat{M}-2)^{1/2} r(1 - r^2)^{-1/2}| > t^{\alpha} \hat{M}_{-2}$$
 (4)

where  $t^{\alpha}_{M-2}$  is critical value of the t-statistic with  $\hat{M}$ -2 degrees of freedom.

Monte Carlo simulations carried out both for lattice and non lattice models show that the  $t_{M-2}$  statistic has indeed a correct significance level<sup>17</sup>.

In the case of several variables, classical regression analysis can be extended. The spatial structure is taken into account either by considering a mixed simultaneous regressive-autoregressive model :

$$Y = L(Y) + X\beta + \varepsilon$$
 (5)

where  $\{\varepsilon_{\alpha}\}$  are i.i.d. N(0,  $\sigma^2$ ), and L(Y)<sub> $\alpha$ </sub> represents some linear function of the values of Y in neighbouring locations  $\gamma, \gamma \neq \alpha$ ; or by considering a regression model with autocorrelated error term :

$$Y = X\beta + U \tag{6}$$

where the spatial autocorrelation is reflected in the structure of  $\Sigma_U$ , the variancecovariance matrix of the error U, which is parametrised in terms of a finite number of parameters  $\theta$ .

The interpretation of  $\beta$  is essentially different for the regressions (5) and (6). In model (5),  $\beta$  represents the effect of  $X_{\alpha}$  on  $Y_{\alpha}$  when the influence of neighbouring values of  $Y_{\alpha}$  has already been substracted. When relating the incidence of a non contagious disease to environmental factors, model (5) is not plausible since disease rates in neighbouring areas do not have a direct causal influence and model (6) should be chosen. In model (6), it is assumed that the autocorrelation of the variable Y comes from an underlying structural variable U.

The estimation of models (5) and (6) with an autocorrelation structure modelled via a contiguity matrix of weights W, was first discussed by  $Ord^{20}$ . Mardia and Marshall<sup>21</sup> have studied the asymptotic properties of the maximum likelihood estimators for the model (6). Assuming that Y is a gaussian process, they give conditions which ensure the consistency and the asymptotic normality of  $(\hat{\beta}, \hat{\theta})$  the maximum likelihood estimators of  $(\beta, \theta)$ . Conditionally upon  $\hat{\theta}$ , the variance-covariance matrix for  $\hat{\beta}$  can be calculated and hence tests of the regression coefficients can be performed. Because of the conditionality, the standard errors of  $\hat{\beta}$  might be underestimated, which leads some authors to prefer a Bayesian approach<sup>22</sup>.

In using model (6), one is left with choosing how to parametrise  $\Sigma_U$  and to perform a numerical maximisation of the likelihood function which can lead to

involved computing and is fraught with difficulties. A study of Warnes and Ripley<sup>23</sup> showed that the Fisher scoring technique usually ensures convergence only to the nearest local maximum. Moreover Ripley<sup>24</sup> reports some simulation results where the global maximum found is well away from the true value. Models used for  $\Sigma_U$  have been basically of two kinds, either involving directly the contiguity matrix<sup>20</sup>, or supposing that the elements of  $\Sigma_U$  depend upon the distance between pairs of locations with a functional form suggested by the variogram of the ordinary least square regression residuals. This is the approach chosen by Cook and Pocock<sup>25</sup> in their study of cardiovascular mortality and environmental factors (in particular water hardness). They fitted an exponentially decreasing function of distance for  $\Sigma_U$  and still found a statistically significant link with water hardness.

In summary, generalised regression where the variance-covariance matrix of the errors follows a specific spatial model involves some degree of arbitrariness and is often computationally difficult in practice. We have thus found it interesting to develop a simple test of partial correlation which does not depend on parametric modelling.

### 3. <u>A modified test for partial correlation</u>

Our aim is to show how the modified  $t_{M-2}$  test described in §2 can be simply extended to test partial correlations. For the sake of clarity the method is going to be described and assessed for testing the association between two variables ( $Y_{\alpha}, Z_{\alpha}$ ) adjusted on a third one,  $X_{\alpha}$ ,  $\alpha \in A$ . Its generalisation to any number of adjustment variables is straightforward.

We suppose that the 3N vector  $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix}$  follows a multivariate normal distribution with mean  $\begin{pmatrix} \mu_X \\ \mu_Y \\ \mu_Z \end{pmatrix}$  and variance-covariance matrix given by :

$$\begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} & \Sigma_{XZ} \\ \Sigma_{YX} & \Sigma_{YY} & \Sigma_{YZ} \\ \Sigma_{ZX} & \Sigma_{ZY} & \Sigma_{ZZ} \end{pmatrix}.$$

Then the joint distribution of (Y,Z) conditional on X is also a multivariate normal with marginals given by N( $\mu_{Y} - \sum_{XY} \sum_{XX} - 1\mu_{X}$ ,  $\sum_{YY} - \sum_{XY} \sum_{XX} - 1\sum_{YX}$ ) and N( $\mu_{Z} - \sum_{XZ} \sum_{XX} - 1\mu_{X}$ ,  $\sum_{ZZ} - \sum_{XZ} \sum_{XX} - 1\sum_{ZX}$ ) respectively. We can therefore test the hypothesis :  $\rho_{YZ,X} = 0$ , where  $\rho_{YZ,X}$  is the partial correlation between Y and Z conditional on X, by testing that the correlation between the residuals of the regression of Y on X and of Z in X is zero. Hence, the method outlined in §2 equations (1) to (4) can be extended to test the partial correlation by considering the relevant variance-covariances in the conditional distributions of Y and Z given above and carrying out the same adjustment of the degrees of freedom (d.f.). This adjustment will now take into account spatial autocorrelation in the conditional distributions. The same approach can be used if conditionally on X, (Y,Z) is multivariate normal even if X is not a random variable as it is the case when X is a vector of spatial coordinates.

In practice, this implies using the modified  $t_{M-2}$  statistic on the residuals of the linear regression of Y on X and of Z on X respectively. These residuals need therefore to be estimated. We are proposing to do this by ordinary least squares (O.L.S.), thus ignoring the autocorrelation, since the O.L.S regression estimates are unbiased.

In summary the following steps are followed :

1 - regress Y on X by O.L.S giving estimated residuals  $\hat{U}$ 

2 - regress Z on X by O.L.S giving estimated residuals  $\hat{V}$ 

3 - test the correlation coefficient between  $\hat{U}$  and  $\hat{V}$  using the modified test statistic  $t_{\hat{M}-2}$  given in formula (4) with

$$\hat{M} = (\text{var } r_{\hat{U}}, \hat{v})^{-1} + 1$$
 (7).

If the residuals U and V are non autocorrelated, the variance of the partial correlation coefficient is  $(N-2)^{-1}$ , hence  $\hat{M} = N-1$  as usual.

In order to check whether the estimation of the residuals by O.L.S could influence the performance of the test, a Monte-Carlo simulation was carried out to check the type I errors and to give some indication on the power of the proposed method.

## Simulation model

In order to stay close to the epidemiological examples which will be analysed in the last section, the simulation was carried out on an irregular grid of points defined by the administrative centers of the French *département* (N = 82). Spatial dependence was introduced directly on the variance-covariance matrix of the multivariate normal distributions by considering a disc model. In this model, the covariance between the values of the variable in 2 areas (represented by their grid points) is defined as proportional to the intersection area of 2 discs of common radius R, centered on the grid points (see Ripley<sup>26</sup> p. 55). By varying the length of R, models with different degrees of autocorrelation can be generated. The shape of the covariance function of this model shows a fairly linear decrease with distance. This shape is similar to the observed variograms ; i.e. the plot of N<sub>k</sub>-1  $\Sigma$  (X<sub> $\alpha$ </sub> - X<sub> $\gamma$ </sub>)<sup>2</sup> against the average distance between ( $\alpha$ , $\gamma$ )  $\epsilon$  S<sub>k</sub>

locations in  $S_k$ , of a number of variables such as lung cancer mortality rate or proportion of workers in specific industries (Figure 1). These will be considered in the next section. The parameter R of the disc model is chosen so that the autocorrelation for a distance of 40 km between points is equal to 0.2, 0.5 or 0.8. We denote this autocorrelation by  $\rho(1)$  indexed by the name of the variable.

To simulate a process X on A, for each chosen value of  $\rho_X(1)$ , a NxN matrix  $\Sigma_X$  is generated following the disc model.  $\Sigma_X$  was then triangularised,  $\Sigma_X = LL^T$ ,

and a realisation of X distributed as  $N(0,\Sigma_X)$  was obtained by first generating a vector of N i.i.d. N (0,1) and then premultiplying this vector by L.

The distances between the centers of *départements* were partitioned into 15 classes of 50 km intervals each. This gives 15 strata S<sub>1</sub>, ... S<sub>15</sub>; the stratum  $S_0 = \{(\alpha, \alpha), \alpha \in A\}.$ 

## -Figure 1-

We now describe the model for partial correlation which was chosen. First three mutually independent autocorrelated variables, X, V and W are generated. Then we define : Z = aX + V, Y = cX + dZ + W. (8) For simplification, the autocorrelation structure of X,V and W are chosen to be identical. Therefore  $\rho_{YZ,X}$  depends only on the value of d :  $\rho_{YZ,X} = d/(d^2 + 1)^{1/2}$ . The test statistic  $t_{M-2}$  is computed as described in §2 steps 1 to 3 and formula (7) using { $\hat{C}_U(k)$ } and { $\hat{C}_V(k)$ }, k = 1,... 15, the estimated autocovariances in each of the 15 strata of the O.L.S residuals :  $\hat{U}$  and  $\hat{V}$ , from the regression of Y on X and Z on X respectively.

Five thousand simulations were carried out for  $\rho_{YZ,X} = 0.0$ , 0.2 and 0.4 and  $\rho_X(1) = \rho_V(1) = \rho_W(1)$  equal to 0.0, 0.2, 0.5 and 0.8. The values of a and c were fixed. Several sets of a and c were tried and were found not to influence the results.

#### RESULTS

Type I errors of the modified statistic were first checked by setting d=0 in (8) and choosing 5 % and 1 % nominal levels. The results are summarised in Table 1. At the 5 % nominal level the performance of the modified  $t_{M-2}$  statistic is clearly satisfactory as the value 5 % belongs to all the confidence intervals and there is no systematic variation with increasing autocorrelation. At the 1 % level the modified  $t_{M-2}$  statistic becomes slightly conservative for the highest autocorrelation but is otherwise satisfactory. As an illustration, the empirical

variance of  $r_{YZ,X}$  is also shown and one can see that it increases significantly with  $\rho(1)$ . As expected, the type I errors of the non adjusted standard test of partial autocorrelation are greatly inflated for high values of  $\rho(1)$ . The results are similar to those obtained in an earlier work for the test of the single correlation coefficient<sup>17</sup>.

In choosing in practice the number of strata for the partition of AxA, a compromise has to be reached between choosing to few strata resulting in a bias in formula (2) or too many leading to imprecise estimates of the autocovariances in each strata. Further simulations were carried out to check the robustness of the performance of the  $t_{M-2}$  statistic with respect to the choice of partitions. Using the same random number generator seed as in Table 1 and 5000 trials, we evaluated the type I error for isotropic partitions ranging from 5 classes of around 150 km each to 21 classes of around 35 km each, in the most highly autocorrelated case ( $\rho_X(1) = \rho_Y(1) = 0.8$ ). The observed type I errors stayed overall close to their nominal level and decreased progressively from 6.4 % [5.7 %-7.1 %] for 5 strata to 4.8 % [4.2%-5.4%] for 21 strata for a 5 % nominal level and from 1.1 % [0.86 % - 1.45 %] to 0.56 % [0.35 % - 0.76 %] for a 1 % nominal level. This indicates that at the 5 % level the test is slightly over significant when too few strata are used and that a finer partition with number of strata between 11 and 17 is preferable. In a case (like the one simulated) where the autocorrelation decreases smoothly with distance, one can thus see that the performance of the  $t_{M-2}$  statistic is robust to various choices of partitions.

#### <u>Table 1</u>

The power of the modified test was then investigated for value of  $\rho_{YZ,X}$  equal to 0.2 and 0.4 and a 5 % nominal significance level. Higher values of  $\rho_{YZ,X}$  led to a power close to 100 % and were less informative.

In order to have a reference value for the power the modified  $t_{M-2}$  statistic, we also calculated the power of the standard test of  $\rho_{YZ,X}$  in a case which would be compatible with the observed empirical variance of  $r_{YZ,X}$ ,  $v_e$ , estimated by

Monte-Carlo simulations. Recall that in the case of non autocorrelated variables X, Y and Z, the variance of  $r_{YZ,X}$  is approximatively equal to  $(1 - \rho^2_{YZ,X})^2 / N-2$  for a sample of N observations when N is large. For autocorrelated X, Y and Z we thus computed an equivalent sample size, N<sup>\*</sup>:

$$N^{\star} = 2 + (1 - \rho^2 \gamma_{Z,X})^2 / v_e \qquad (9)$$

This number N\* was used to compute the reference value  $\pi$ , power of the standard test of  $\rho_{YZ,X}$  based on N\* observations.

A summary of the results is given in Table 2. On the first line is indicated  $v_e$ , the empirical variance of  $r_{YZ,X}$ . One can see that  $v_e$  is smaller for  $\rho_{YZ,X} = 0.4$  than for  $\rho_{YZ,X} = 0.2$  for the same autocorrelations. This is similar to the case of non correlated variables. Further  $v_e$  increases noticeably with increasing  $\rho(1)$ . N\* and  $\pi$  are given respectively on the 2nd and 3rd lines. Note that as expected N\* is fairly constant for the same autocorrelation. The power of the test  $t_{M-2}$  is given on the 4th line. Clearly the test  $t_{M-2}$  perform satisfactorily as its power is comparable to  $\pi$ .

#### -Table 2-

## §4 Examples

We consider in this section examples concerning the relationship between male lung cancer and some industrial factors after adjustment on smoking. Among the different cancer sites, male lung cancer has been frequently associated with industrial exposure<sup>27,28</sup> and the estimation of the proportion of lung cancer due to occupational risk factors has been the subject of recent debates<sup>29</sup>. To illustrate and discuss further the method proposed we present some results on 4 industrial branches : metal, general engineering, mining and textiles.

## The data

Male lung cancer mortality rate has been standardised over the age 35-74 and over a 2 year period, 1968-1969. The data were provided by the French National Institute for Health and Medical Research (INSERM) at the scale of the French *départements*. Cigarettes sales data were compiled by the French Nationalised Tobacco Company (SEITA). To take into account the time lag between smoking and the onset of a lung pathology, cigarettes sales per inhabitant were recorded in 1953. Demographic data on the percentage of employed males in the metal industry, in general engineering, in the textile industry and in mining were taken from the 1962 census (INSEE). After the grouping of the *départements* around Paris into one area and the exclusion of 4 others owing to the poor quality of the data, N = 82 locations were retained.

## Results for the modified t<sub>M-2</sub> test

Separate results for the 4 industrial branches are presented in Table 3. Simple correlations, partial correlations after adjustment on smoking or after both adjustment on smoking and a linear gradient were tested by the standard methods and using the modified  $t_{M-2}$  test.

Generally one can see that by including more explanatory variables the autocorrelation in the residuals decreases and so the effective sample size  $\hat{M}$  increases. As expected, the difference between the standard and the  $t_{\hat{M}-2}$  tests, which can be substantial, is more marked where there is more autocorrelation.

The results obtained are quite different for the 4 industries. For metal and general engineering industry, there is a statistically significant association with male lung mortality rate after adjustment on cigarette sales both with and without the removal of a linear gradient. For mining, the significant association after adjustment on smoking is not observed after the removal of a linear gradient. For the textile industry, no significant association is found before adjustment with the modified  $t_{M-2}$  statistic whilst this would not be the case if only standard test were

carried out. After adjustment on both cigarette sales and a linear gradient, the association is also clearly non significant. Note that in all the examples except mining the adjustment solely on cigarette sales leads to slightly larger  $t_{M-2}$  statistic.

#### <u>-Table 3-</u>

In Table 3, the  $t_{M-2}$  statistic is calculated using the 15 isotropic strata defined in §3. To further study the influence of the number of strata we have calculated as an example values of the  $t_{M-2}$  statistic for metal industry adjusted on cigarette sales (2nd line in Table 3) for 9 isotropic partitions ranging from 5 classes to 21 classes. The values of  $t_{M-2}$  varied little, from 3.61 for 5 classes to 3.34 for 21 classes ; similarly the effective sample size  $\hat{M}$  varied only from 37 (5 classes) to 32 (21 classes). When these calculations were carried out for the other examples there is again little change in the results. This is in agreement with the simulations described earlier.

In a further analysis, it was thought interesting to also include other industrial factors in the adjustments. We chose the metal industry and general engineering for this purpose since they are the only two factors clearly significantly related to lung cancer. The results are presented inTable 4. Taking metal industry, with an adjustment on general engineering, or taking general engineering, with an adjustment on metal industry, did not overall modify the results found in Table 3 and the relationships are confirmed. It is interesting to note that for mining, adjustment on both industries had a similar effect to that of adjusting for a gradient with the association with lung cancer mortality becoming non significant and the residual autocorrelation disappearing. For the textile industry, the adjustment for both industries gives a similar result to that for cigarette sales with no clear evidence of any link with lung cancer mortality.

## <u>-Table 4-</u>

The inclusion of a linear gradient, whilst strengthening the interpretation of the association of lung cancer mortality with metal or general enginnering industry has on the contrary, weakened the observed association with mining or textile. We now discuss in some detail the problems connected with the existence of gradients.

### Gradient influence

The adjustment of gradients always involves some degree of arbitrariness in the definition of the appropriate trend. Furthermore, the observed variations in the residuals may be of a smaller order of magnitude than the original variations thus leading to less precision in the estimates. Hence, one needs to be cautious when interpretating a non significant result after de-trending.

We have investigated the effect of omitting to adjust a gradient on the performance of the standard and modified  $t_{M-2}$  test of the correlation coefficient. Simple models were chosen for two variables  $X_{\alpha}$  and  $Y_{\alpha}$ ,  $\alpha \in A$ . The centred and standardised coordinates of a point  $\alpha$  in A will be denoted by i $\alpha$  and j $\alpha$ . Define :

 $X_{\alpha} = U_{\alpha} + ai_{\alpha} + bj_{\alpha}$  and  $Y_{\alpha} = V_{\alpha} + a'i_{\alpha} + b'j_{\alpha}$  (10) with  $U_{\alpha}$  and  $V_{\alpha}$  mutually independent variables,  $U_{\alpha} \sim N(0, \Sigma_U)$ ,  $V_{\alpha} \sim N(0, \Sigma_V)$ ,  $\Sigma_U$ and  $\Sigma_V$  generated by a disc model with equal parameter  $\rho_{(1)}$ .  $X_{\alpha}$  and  $Y_{\alpha}$  defined by (10) will not be stationary and the autocorrelation between 2 points will depend on their coordinates.

The effect of the linear gradient on the expectation of the covariance between X and Y leads to a term proportional to aa' + bb' involving  $\sum i_{\alpha}^2$  and  $\sum j_{\alpha}^2$ , corrected by a term involving  $\sum i_{\alpha}j_{\alpha}$ . The correction term is small since our centered and standardised coodinates are nearly orthogonal. Hence the effect is mainly related to aa' + bb' which is maximum when a' = a , b' = b and minimum when a' = -b ; b' = a. The effect of the gradient on the correlation coefficient is more complicated since a,a',b,b' are also involved in the denominator and a first order approximation for the expectation of r<sub>XY</sub> can be calculated. For the sake of simplicity in the following discussion we relate the results to aa' + bb'. Values of the first order approximation are indicated in brackets in Table 5. Five hundred simulations were carried out for different choices of aa' + bb' (keeping a + b = a' + b'), 3 levels of autocorrelation and a nominal 5 % significance level. The results are presented in Table 5.

#### -Table 5-

For the standard test and  $\rho_{(1)} = 0$ , type I errors are clearly increasing with aa' + bb'. The situation is similar when some autocorrelation is present  $(\rho_{(1)} = 0.4)$  with overall higher type I errors. When the autocorrelation is high  $(\rho_{(1)} = 0.8)$  the gradient influence on the overall type I error decreases in regard to the autocorrelation effect.

The modified  $t_{M-2}$  test is based on the estimated autocovariances  $\{\hat{C}_X(k)\}\)$ and  $\{\hat{C}_Y(k)\}\)$  even in the case  $\rho_{(1)} = 0$ . Due to the linear trend, these autocovariance functions are quite regular and reach moderately large negative values for high k. From formula (2), it is easy to see that the adjustment in the degrees of freedom for  $t_{M-2}$  is strongest when  $\{\hat{C}_X(k)\}\)$  and  $\{\hat{C}_Y(k)\}\)$  are quite similar, whether positive or negative. This happens clearly in the extreme case a' = a, b'= b of parallel gradients ; but also to a slightly lesser extent in the case aa' + bb' = 0 since values of the variable  $\{Y_{\alpha}\}\)$  would be equal to those of  $\{X_{\alpha}\}\)$  if the set of coordinates were rotationally invariant by  $\pi/2$  which is nearly the case in France. This explains that for aa' + bb' = 0, the d.f. are overadjusted and the statistic  $t_{M-2}$  is conservative. For the other values of aa' + bb', the correction of the d.f. reduces the type I error in comparison to that of the standard test but nevertheless it does not always stay close to 5%. When  $\rho_{(1)} = 0.4$  or 0.8, the intrinsic autocorrelation prevails more and more on the gradient so that irrespective of the gradient, all the type I errors are around 5% for  $\rho_{(1)} = 0.8$ .

In conclusion it is worth calculating  $t_{M-2}$  both before and after de-trending. Even in the presence of a gradient,  $t_{M-2}$  may realise a satisfactory adjustment if the intrinsic autocorrelation in the variables is strong.

## Concluding remarks

In this paper, we have outlined a simple test for partial correlations between spatially autocorrelated variables which is particularly suitable for ecological correlation studies. The performance of this test was shown to be satisfactory by Monte Carlo simulations. The application of this method to data sets is straightforward and only requires simple computing. The examples analysed concern the relationship between male lung cancer mortality rate and some industrial factors. They were chosen in order to illustrate that geographical epidemiology studies, when analysed with concern for the spatial structure, can identify risk factors which are related to those found in individual epidemiological studies and that this can be done quite easily with the modified test proposed. It would be interesting to compare our results with those of a generalised regression with spatially modelled error structure.

At the level of individual epidemiological studies, working in the metal industry, ship building and motor vehicule construction has been recognised to present a carcinogenic risk for lung cancer due to possible exposure to arsenic, chromium, benzo(a) pyrene, nickel or asbestos. Concerning the mining industry, the association is specifically reported for arsenic, iron-ore, asbestos or uranium mining but not for coal mining. Finally, there is no reported relationship between lung cancer and the textile industry<sup>29</sup>. Hence, the associations observed in our results using the modified  $t_{M-2}$  statistic with the metal industry and general engineering are particularly interesting as well as the non-association with the textile industry. Using standard methods would lead somewhat to misleading results concerning the textile industry. For the mining industry, the association is bordeline as it can be accounted for by a simple spatial gradient structure or by other industrial factors. As the epidemiological evidence does not involve all types of mining, a more detailed indicator of exposure would need to be analysed.

It is also interesting to note that part of the autocorrelation in the dependent variable is explained by the inclusion of confounding factors in the model. In this light, the gradient can be thought of as a proxy confounding variable. Finally one has to recall that as discussed in Greenland<sup>15</sup> confounding at the geographical level by a covariate can occur under broader conditions than for individual epidemiological studies and that geographic adjustment might be insufficient to control for it.

# **Acknowledgements**

The author wishes to thank Chantal Guihenneuc and Virginie Lasserre of the Laboratoire de Statistiques Médicales, Université de Paris V for computing assistance and Evelyne Przybilski for secretarial assistance. This work was supported by a EURATOM contract n° BI6-126F.

- (1) Doll, R. The epidemiology of cancer. Cancer, 45,2475-85 (1980).
- (2) Doll, R. ed. The geography of disease. British Medical Bulletin. Published for the British Council by Churchill Livingstone (1984).
- (3) Hutt, M.S. and Burkitt, D.P. The geography of non-infectious disease. Oxford University Press Oxford (1986)
- (4) Armstrong, B., Doll, R. Environmental factors and cancer incidence and mortality in different countries, with special reference to dietary practices. Int J Cancer, 15, 617-631 (1975).
- (5) Blot, W.J. and Fraumeni, J.F. Geographic patterns of oral cancer in the United States : etiological implications. J Chron Dis, 30, 745-757 (1977)
- (6) Graham, S. Methodological problems in ecologic studies of the asbestos cancer relationship. Envir Res, 25, 35-49 (1981)
- (7) Cantor, K.P., Hoover, R., Mason, T.J, Mc Cabe, L.J. Associations of cancer mortality with halomethanes in drinking water. J Natl Cancer Inst, 61, 979-985 (1978).
- (8) Cliff, A.D. and Ord, J.K. Spatial Processes. Models and Applications. London, Pion (1981).
- (9) Doreian, P. "Estimating linear models with spatially distributed data", in Sociological Methodology, edited by Leinhardt and Samuel. San Francisco, Jossey-Bass Publishers, 359-388 (1981).
- (10) Breslow, N.E. and Enstrom, J.E. Geographic correlations between cancer cancer mortality rates and alcohol-tobacco consumption in the United States. J Natl Canc Inst, 53, 631-639 (1974).
- (11) Stravraky, K.M. The role of ecologic analysis in studies of the etiology of disease : a discussion with reference to large bowel cancer. J Chron Dis, 29, 435-444 (1976).
- (12) Davies, J.M. and Chilvers, C. The study of mortality variations in small administrative areas of England and Wales, with special reference to cancer. J Epidem C, 34, 87-92 (1980).
- (13) Morgenstern, H. Uses of ecological analysis in epidemiologic research. Am J Publ Health, 72, 1336-44 (1982).
- (14) Richardson, S., Stücker, I., Hémon, D. Comparison of relative risks obtained in ecological and individual studies : some methodological considerations. Int J Epidemiol, 16, 111-120 (1987).
- (15) Greenland, S., Morgenstern, H. Ecological bias, confounding, and effect modification. Int J Epidemiol, 18, 269-274 (1989).

- (16) Pocock, S.J., Cook, D.G. and Shaaper, A.G. Analysing geographic variation in cardiovascular mortality : methods and results. J.R.Statist. Soc A, 145, 313-341 (1982).
- (17) Clifford, P., Richardson, S. and Hémon, D. "Assessing the significance of the correlation between two spatial processes" Biometrics, 45(1), 123-134 (1989).
- (18) Bartlett, M.S. Some aspects of the time-correlation problem in regard to tests of significance. J. R. Statist. Soc., 98, 536-543 (1935).
- (19) Johnston, J. Econometric Methods, 2nd edition. New York, McGrawHill (1972).
- (20) Ord, K. Estimation methods for models of spatial interaction. J. Am. Statist. Ass, 70, 120-126 (1975).
- (21) Mardia, K.V. and Marshall, R.J. Maximum likelihood estimation of models for residual covariance in spatial regression. Biometrika, 71, 135-146 (1984)
- (22) Hepple, L.W. Bayesian analysis of the linear model with spatial dependence, in Exploratory and Explanatory Statistical Ananlysis of Spatial Data pp 179-199 (Bartels, C.P.A and Ketellapper, R.H., Eds). Martinus Mijhoff : Boston (1979).
- (23) Warnes, J.J. and Ripley, B.D. Problems with likelihood estimation of covariance functions of spatial Gaussian processes. Biometrika, 74(3), 640-642 (1987).
- (24) Ripley, B.D. Statistical inference for spatial processes. Cambridge University Press, Cambridge (1988).
- (25) Cook, D.G. and Pocock, S.J. Multiple regression in geographic mortality studies with allowance for spatially correlated errors. Biometrics, 39, 361-371 (1983).
- (26) Ripley, B.D. Spatial statistics. New York, Wiley (1981).
- (27) Pastorino, U., Berrino, F., Gervasio, A., Pesenti, V., Riboli, E., Crosignomi, P. Proportion of lung cancers due to occupational exposure. Int J Cancer, 33, 231-237 (1984).
- (28) Benhamou, S., Benhamou, E., Flamant, R. Occupational risk factors of lung cancer in a French case-control study. Br J Ind Med 45, 231-233, (1988).
- (29) Simonato, L., Vineis, P. and Fletcher, A. Estimates of the proportion of lung cancer attributable to occupational exposure. Carcinogenesis, 9, 1159-1165 (1988).

# <u>Table 1</u>

Type I errors (per cent) of the  $t_{M-2}$  statistic for testing at 5 % and 1 % nominal rejection levels the partial correlation coefficient  $r_{YZ,X}$  between Z = aX + U and Y = cX + W after adjusting on X.

5000 simulations are carried out for several levels of autocorrelation in X,U and W ( $\rho_X(1)=\rho_U(1)=\rho_W(1)$ ), with a = 2c = 0.204. As reference, the empirical variance of  $r_{YZ,X}$  and the type I errors (per cent) for the classical test of  $r_{YZ,X}$  based on N = 82 observations are also indicated.

autocorrelation $\rho(1)$ for variables X,U and W	0.0	0.2	0.5	0.8	
empirical variance of ryz.x	0.0128	0.0137	0.0176	0.0522	
% type I errors for t <sub>M-2</sub> [95 % CI] 5 % norninal level	5.5 % [4.87% - 6.13%]	5.44 % [4.81% - 6.07%]	5.54 % [4.91% - 6.17%]	5.16 % [4.55% - 5.77%]	
% type I errors for the t test * of ryz x based on N observations [95 % Cl] 5 % nominal level	5.56 % [4.93% - 6.19%]	6.5 % [5.81% - 7.18%]	10.48 % [9.63% -11.33%]	36.42 % [35.09% - 37.75%]	
% type I errors for t <sub>M-2</sub> [95 % CI] 1 % nominal level	1.14 % [0.85% - 1.43%]	1.2 % [0.9% - 1.5%]	0.92 % [0.65% - 1.18%]	0.72 % [0.48%- 0.95%]	
% type I error for the t test *of ryz x based on N observations [95 % CI] 1 % nominal level	1.22 % [0.92% -1.52%]	1.76 % [1.40% - 2.10%]	3.06 % [2.58% - 3.53%]	22.96 % [21.79% - 24.13%]	

\* test of  $\sqrt{N-3}$   $_{YZ,X}$  / (1 -  $r^2{}_{YZ,X})^{1/2}$  as a  $t_{N-3}$  distribution

Power of the modified  $t_{M-2}$  statistic for testing at a 5 % nominal rejection level the partial correlation coefficient  $r_{YZ,X}$  between Z = aX + U and Y = cX + dZ + W after adjusting on X.

500 simulations are carried out for several levels of autocorrelation in X,U and W ( $\rho_X(1) = \rho_U(1) = \rho_W(1)$ ), with a = 2c = 0.204. N\* defined in (9) is an equivalent sample size base on v<sub>e</sub>, the observed variance of r<sub>YZ,X</sub> and  $\pi$  is the power of a standard test of r<sub>YZ,X</sub> based on N\* observations.

ργΖ.Χ	ρ(1)	0.0	0.2	0.5	0.8
0.2	ve	0.0111	0.0121	0.0156	0.0465
	N*	85	78	61	22
	π	44.8%	40.4 %	33.3 %	10.8 %
	power of $t_{\dot{M}-2}$	43.4 %	41.7 %	34.6 %	13.6 %
0.4	ve	0.0087	0.0094	0.0120	0.0372
	N*	83	77	61	21
	π	96.4 %	95 %	89.1 %	41.7 %
	power of $t_{M-2}$	96 %	95.4 %	91 %	46.8 %

Comparison of the significance levels for tests of simple and partial correlations between male lung cancer mortality rates and industrial risk factors (adjusted on cigarette sales or on cigarette sales and a linear gradient) given by standard tests and  $t_{M-2}$  tests.

	r	standard test (N = 82)	Modified $t_{M-2}$ test	Â
		(values of the t <sub>N-2</sub> statistic and significance level)	(values of the t <sub>M-2</sub> statistic and significance level)	
Metal industry workers				
no adjustment	0.63	7.16 p < 10 <sup>-9</sup>	3.00 p = 0.01	16
adjustment on cigarette sales	0.52	5.46 p < 10 <sup>-6</sup>	3.46 p = 0.002	34
adjustment on cigarettes sales and a linear gradient	0.35	3.34 p = 0.002	2.72 p = 0.01	55
General engineering workers				
no adjustment	0.43	4.23 p < 10 <sup>-4</sup>	2.72 p = 0.01	35
adjustment on cigarette sales	0.37	3.58 p < 10 <sup>-3</sup>	2.86 p = 0.007	53
adjustment on cigarettes sales and a linear gradient	0.26	2.41 p = 0.02	2.32 p = 0.02	76
Mine workers				
no adjustment	0.33	3.16 p = 0.003	2.37 p = 0.02	47
adjustment on cigarette sales	0.24	2.26 p = 0.03	2.42 p = 0.02	94
adjustment on cigarettes sales and a linear gradient	0.14	1.27 NS	1.26 NS	81
Textile industry workers				
no adjustment	0.28	2.57 p = 0.01	1.52 p = 0.14	30
adjustment on cigarette sales	0.26	2.40 p = 0.02	1.91 p = 0.07	53
adjustment on cigarettes sales and a linear gradient	0.11	1.02 NS	0.96 NS	73

Significance levels of partial correlations between male lung cancer mortality rates and industrial factors after respective adjustments on metal industry and/or general engineering given by the standard tests and  $t_{M-2}$  tests.

	r	standard test (N = 82) (values of the t <sub>N-2</sub> statistic and significance level)	Modified $t_{M-2}$ test (values of the $t_{M-2}$ statistic and significance level)	M
Metal industry workers				
adjustment on general engineering	0.58	6.33 p < 10 <sup>-7</sup>	3.32 p = 0.004	24
adjustment on cigarette sales, general engineering and a linear gradient	0.34	3.28 p =0.002	2.62 p = 0.02	53
General engineering workers				
adjustment on metal industry	0.33	3.09 p =0.003	2.81 p = 0.007	68
adjustment on cigarette sales, metal industry and a linear gradient	0.25	2.33 p =0.03	2.20 p = 0.04	73
Mine workers				
adjustment on metal industry and general engineering	0.17	1.5 NS	1.54 NS	86
adjustment on cigarette sales, metal industry, general engineering and a linear gradient	0.11	0.96 NS	0.98 NS	86
Textile industry				
adjustment on metal industry and general engineering	0.23	2.07 p=0.05	1.88 p=0.07	68
adjustment on cigarette sales, metal industry, general engineering and a linear gradient	0.18	1.63 NS	1.51 NS	70

Proportion of type I errors for the standard test and the  $t_{\dot{M}-2}$  test of the correlation coefficient between the variables  $X_{\alpha} = U_{\alpha} + ai_{\alpha} + bj_{\alpha}$  and  $Y_{\alpha} = V_{\alpha} + a'i_{\alpha} + b'j_{\alpha}$  where  $U_{\alpha}$  and  $V_{\alpha}$  are mutually independent spatially autocorrelated disc processes with equal parameter  $\rho(1)$ .

500 simulations are carried out in each case. The nominal rejection levels of the standard and  $t\dot{M}$ -2 test chosen is 5 %.

% type	errors	a' = 0.4 b' = -0.1 aa'+bb' = 0 (0.00)*	a' = 0.35 b' = 0.05 aa'+bb' = 0.075 (0.043)*	a' = 0.3 b' = 0.2 aa'+bb' = 0.11 (0.091)*	a' = 0.2 b' = 0.3 aa'+bb' =0.14 (0.119)*	a' = 0.1 b' = 0.4 aa'+bb'=0.17 (0.143)*
ρ(1) = 0	standard t	4.6 %	6.8 %	13.2 %	16 %	26.4 %
	tM-2	1.8 %	2.8 %	6.2 %	7.8 %	14 %
ρ(1) = 0.4	standard t	10.4 %	12.2 %	16.8 %	21.8 %	28.6 %
	tM-2	3.8 %	4.8 %	5.2 %	8.6 %	11 %
ρ(1) = 0.8	standard t	35.4 %	46.2 %	43.6 %	50 %	46%
	tŵ-2	2.6 %	4 %	5.2 %	5.2. %	5 %

#### Linear gradient slopes for Y

[the linear gradient slopes for X are fixed with a=0.1 and b= 0.4]

\* first order approximation for the expected value of  $r_{XY}$ 

Figure 1





Fifteen classes of distance for the plot of  $N_k^{-1} \sum (X_{\alpha} - X_{\gamma})^2$  against the average distance between locations in  $S_k$  are considered.

The number  $N_k$  of pairs in each class is : 82, 400, 582, 674, 764, 822, 812, 726, 630, 476, 304, 17 $\pmb{g}$ , 94, 58, 40.
# COMPARISON OF RELATIVE RISKS OBTAINED IN ECOLOGICAL AND INDIVIDUAL STUDIES : SOME METHODOLOGICAL CONSIDERATIONS

# Comparison of Relative Risks Obtained in Ecological and Individual Studies: Some Methodological Considerations

# SYLVIA RICHARDSON\*†, ISABELLE STÜCKER\* AND DENIS HÉMON\*

Richardson S (INSERM, U170, 16 av. Paul Vaillant-Couturier, 94807 Villejuif Cedex, France), Stücker I and Hémon D. Comparison of relative risks obtained in ecological and individual studies: some methodological considerations. *International Journal of Epidemiology* 1987, **16**: 111–120.

This paper is concerned with the problem of estimating relative risks from ecological correlation studies. In the first part, some of the biases encountered when analysing aggregated data are discussed and in particular attention is focused on the shape of the dose-response relationship obtained from aggregated and non-aggregated data, on the need for extrapolation and on the scale of aggregation. In the second part some empirical observations are made on these points by means of four examples concerning the relative risk between smoking and different pathologies. The estimates of relative risks derived from French geographical data and from case control or cohort studies are compared. The performance of ecological studies is discussed with respect to the strength of the risk factor considered, the geographical distribution of counfounding factors and the adjustment of different models.

Studies of the geographical distribution of incidence or mortality rates for particular diseases have been widely used by epidemiologists in formulating hypotheses about the aetiology of those diseases.<sup>1-5</sup> These studies, where the variables concerned are averages measured on groups of people, are often referred to as 'ecological correlation studies'. How to interpret and to assess the statistical significance of an observed, ecological correlation has been discussed by many authors.<sup>6-13</sup> Our aim is to compare the quantification of the doseresponse relationship obtained by geographical studies to that assessed at the individual level, an approach first used by Béral.<sup>14</sup>

# ESTIMATION OF RELATIVE RISK FROM ECOLOGICAL DATA: THEORETICAL CONSIDERATIONS

Comparison of the Shape of the Dose-Response Relationship between Aggregated and Non-Aggregated Data Consider the relation:

$$f(\mathbf{x}) = \mathbf{P}(\mathbf{D} \mid \mathbf{X} = \mathbf{x}) \tag{1}$$

which models the probability of contracting the disease

D if an individual is exposed to the risk factor X taking the value x. Suppose that inside a group G the distribution of X has mean  $\mu_G$  and that it is only possible to observe the relationship between  $\mu_G$  and the mean level of the disease in the group  $Y_G$ :

 $Y_G = g(\mu_G) = E_X [P(D | X)].$  (2) It is assumed here that  $E_X [P(D | X)]$  can be indexed solely by  $\mu_G$ . Let us compare in some examples estimations of the relative risk (**RR**) derived from (1) and (2). (a) Firstly let us consider the case where f is linear,

$$f(x) = \alpha + \beta x$$
 (3)  
then clearly g is also linear with the same coefficients

$$Y_G = \alpha + \beta \mu_G. \tag{4}$$

As suggested by Béral.<sup>14</sup> (4) can be used to estimate the relative risk RR( $x_0$ ) for an average exposure  $\mu_G = x_0$ : RR( $x_0$ ) = ( $\beta/\alpha$ )  $x_0 + 1$ .

$$\mathbf{RR}(\mathbf{x}_0) = (\mathbf{p}/\alpha) \mathbf{x}_0 + 1$$

 $RR(x_0)$  can therefore be equally estimated from equation (3) or (4), that is from individual or ecological studies. So in the linear case with one risk factor there is no mathematical bias resulting from using aggregated data.

(b) Secondly if f is convex, then using Jensen's inequality, it follows that

$$Y_G = g(\mu_G) = E_X(f(X)) \ge f(\mu_G).$$

In particular in the classical exponential case,  

$$f(x) = c e^{ix}, \gamma > 0$$
 (6).

$$Y_G = g(\mu_G) = E_X(ce^{\gamma x}) = c \exp(\gamma \mu_G) \times b(\gamma, X)$$
 (7)

<sup>\*</sup> INSERM U170 Unité de Recherches Epidémiologiques et Statistiques sur l'Environnement et la Santé, 16 av. Paul Vaillant-Couturier, 94807 Villejuif Cedex, France.

<sup>\*</sup> Laboratoire de Statistiques Médicales, Université de Paris V 45 rue des Saints Pères, 75006 Paris, France.

where the bias coefficient  $b(\gamma, X)$  can be expressed in terms of the cumulants of the distribution of X and will in general depend on  $\gamma$ ,  $\mu_G$  and the parameters of f.<sup>15</sup> Consequently (7) will lead to a biased estimate of RR(x<sub>0</sub>) unless  $b(\gamma, X)$  cancels from the numerator and the denominator of RR(x<sub>0</sub>). This is the case for instance when X follows a normal distribution ( $\mu_G$ ,  $\sigma_{\mu}^2$ ) under the condition that  $\sigma_{\mu}^2$  does not vary over the range of  $\mu$ , a condition which would rarely hold.

In general the bias in the estimate of the relative risk will be more pronounced the larger the values of  $\mu$ , the wider its range and the more pronounced the convexity of f.

#### Average response at zero exposure

In most cases, the range of the distribution of the ecological risk factor will not include zero. An extrapolation towards a non-exposed group will be necessary in order to compute an ecological relative risk. This extrapolation is subject to errors and even if two models are indistinguishable in terms of fit on the range of the ecological risk factor, they might lead to quite different risk estimates at the non-exposed level. The problem of extrapolation can also occur in individual studies where only a few cases have low exposure. In both individual and ecological studies information about relative risks can still be derived by considering non-zero reference exposure within the range of the risk factor.

## Scale of aggregation

The geographical units on which the ecological analysis is based are divisible. Hence ecological studies can often be performed at several levels which may not give identical results. Checking the stability of the results in relation to the geographical scale of the analysis is important for their interpretation.<sup>16-22</sup>

Assuming that the variables can be considered at two levels, either ungrouped or averaged over groups (symbolized by G), the total variation of a variable (var (X)) can be decomposed into the average variation A = E[var(X | G)] within each group and the variation B = var [E(X | G)] of the mean between groups:

var(X) = A+B = E[var(X | G)] + var[E(X | G)].Similarly for the covariance:

$$cov(X,Y) = A'+B' = E[cov(X, Y | G)] + cov$$
  
[E(X | G), E(Y | G)].

When the slope of the individual association of X and Y does not vary from one area to another, several situations can occur:

(a) A'/A = B'/B (Figure 1a). In this case the slope of the between group variation of X and Y is equal to the within groups one. Hence (A'+B')/(A+B) = A'/A and so the evaluation of the slope is the same using total or intra-group variation.

(b)  $A'/A \neq B'/B$  (Figures 1b and c). In this case there is an 'ecological cross level bias' caused by the presence of a confounding factor which alters the intra-group variation of X and Y either by amplifying it (Figure 1b) or by reducing it (Figure 1c). This confounding factor can influence the analysis either at the ungrouped or at the grouped level.

The same reasoning applies to situations where small areas can be grouped into larger regions. When the slope of the relationship changes from one level to another, one may consider this as the result of a confounding factor which varies with the level of the analysis.

# ESTIMATING RELATIVE RISKS FROM ECOLOGICAL DATA: SOME EXAMPLES In the following examples French ecological relative

risks between cigarette consumption and lung,



FIGURE 1 Aggregation and ecological bias.

bladder, oesophagus cancer and coronary heart disease (CHD) mortality are compared with those estimated in individual studies. Confidence intervals for the relative risk are calculated by Fieller's method.<sup>23</sup>

## Mortality Rates

Throughout we consider the mortality in France for two time periods 1968-1969 and 1974-1976 corresponding to national census. The mortality rates have been standardized by the direct method using ten-year age groups from 35-44 to 65-74 and the 1968 population. From the 86 administrative areas (départements)

for which both cigarette consumption and mortality rates were available for those dates, four were excluded from the analyses because the proportion of deaths due to unspecified causes was too large.<sup>24</sup>

# Cigarette Consumption

Damiani and Massé have estimated tobacco consumption per age and per sex for the 22 French regions<sup>25,26</sup> under the assumption that individual and ecological relative risks are equal. Their estimates cannot therfore be used in the present study which aims at discussing that assumption.



FIGURE 2 Relative risk for lung cancer as a function of current number of cigarettes smoked per day estimated in eight cohort studies and in ecological studies. Cohort studies as quoted in:28

- 1 Best et al., 1961
- 2 Cederlof et al., 1975
- 3 Doll et al., 1972
- 4 Hammond and Horn, 1958
- 5 Dorn et al., 1974
- 6 Hammond et al., 1976
- 7 Weir et al., 1970
- 8 Hirayama, 1977

It is necessary to consider an average time lag of 15 years between cigarette sales and mortality rates. The number of cigarettes sold per inhabitant for each French 'département' was recorded in 1953 and 1957 by the French nationalized tobacco company.<sup>27</sup> To correlate mortality rates in the age range 35–74 to cigarette sales which occurred 15 years before, we assessed retrospectively the proportion of the total cigarette sales (in 1953 or 1957) which is attributable solely to men aged 20–59. We were able to estimate this proportion using the results of a large survey carried out in 1960 on smoking habits in France.<sup>28</sup> These estimates will be used throughout the rest of the study and referred to as cigarette consumption for men.

# Lung Cancer and Cigarette Consumption

In Figure 2, are represented relative risks for lung cancer mortality (ICD 162) among men as a function of

current consumption of cigarettes per day estimated in eight major cohort studies.<sup>29</sup> The relative risks corresponding to 5, 10, 15 or 20 cigarettes per day estimated geographically in France by a linear model between the standardized mortality rate for lung cancer in men in 1968 and cigarette consumption for men in 1953 are also shown in this figure.

The values corresponding to daily consumptions of 15 or 20 cigarettes are far outside the range of the average consumption in French 'départements', which is 2 to 10 cigarettes per day (Figure 3). The comparison between cohort estimates and ecological estimates should only be judged on the 2–10 cigarettes per day range where there is a good agreement with a linear model fitted to the French data.

In this figure are also represented the same relative risks estimated by a linear model from international data and from US data. Correlation between national



FIGURE 3 Lung cancer mortality and average cigarette consumption per day in 82 French départements: fitting of linear (--) or exponential (--) model.

lung cancer mortality and cigarette sales has been discussed by many authors.<sup>1,24,30</sup> Data on lung cancer mortality in 1975 or 1976 for 19 countries (agestandardized over the range 35-74 using the European population as reference) were taken from WHO statistics<sup>31</sup> and the corresponding national cigarette sales per adult (in 1957) from Lee.<sup>32</sup> The US data represented come from Fraumeni.<sup>33</sup> The relative risks estimated for the international and US data sets are much smaller than those given by the cohort studies or the French ecological data. This could be the consequence of the different ways of collecting data. Furthermore, US and international data on cigarettes sales and mortality rates were not restricted in this study to the most relevant age and sex groups thus resulting possibly in a dilution effect. In Table 1 estimates for different time lags and different models are given.

For the purpose of illustration, a consumption of 10 cigarettes per day was chosen, but the results are similar within the range 2–10 cigarettes per day. All correlation coefficients corresponding to this table are highly significant, and vary between 0.75 and 0.79. All ecological estimates given in Table 1 are within the range of the individual relative risks cited in the US Surgeon General's Report.<sup>29</sup> Further, estimates corresponding to 1953 and 1957 cigarette consumption are very close together. On the other hand, estimates based on 1968 death rates lie closer to the median value of relative risk obtained from individual studies than those for 1975 (Figure 2).

Similarly the estimates given by the linear model correspond better to cohort studies in the 2 to 10 cigarettes range than those given by the exponential model (Figure 2 and Table 1). This is understandable when looking at the fairly linear shape of the individual doseresponse relationship in this range. Nevertheless there is a very good fit of both models on the French data (Figure 3). The exponential estimates of the relative risk are lower than the linear ones because the baseline mortality rates for non-smokers given by the two models are different.

Health and exposure variables often exhibit a regular component in their spatial variations analogous to a time trend which makes the interpretation of the joint variations difficult since many potential confounding variables might also show the same regular spatial gradient.<sup>34</sup> This suggests a further step in the analysis, the multiple regression of the health variable on the exposure variable and on trend components expressed in terms of the coordinates of the geographical units. A causal interpretation would be strengthened if a statistically significant correlation between exposure and health indicators is observed both with and without adjustment for a spatial trend. Table 2 shows that for both the linear and the exponential model, the regression coefficients between lung cancer mortality and cigarette consumption after the adjustment of a linear gradient are hardly modified.

## Bladder Cancer and Cigarette Consumption

In Figure 4, the relative risk for daily cigarette consumption and bladder cancer mortality (ICD 188) estimated from the case-control studies of Wynder.<sup>35</sup> Cartwright.<sup>36</sup> Howe.<sup>37</sup> and Vineis<sup>38</sup> are represented. In contrast to Figure 2, only two studies clearly demonstrate a dose-response relationship. As for lung cancer the geographical estimates of the relative risk for the French data analysed and for the US data from Lee<sup>32</sup> are shown in Figure 4, but in this case the relative risk estimates based on French ecological data are higher than those of the individual studies.

In Table 3 we note that the relative risks do not vary between 1968 and 1975 and that the linear model gives only slightly larger estimates than the exponential one. Correlation coefficients corresponding to this table were all statistically significant and vary between 0.52 and 0.65.

In the case of bladder cancer we therefore observe some discrepancy between ecological estimates and

TABLE 1	Lung cancer	r mortality	and cigarette	consumption:	estimation o	f the relative	e risk	corresponding	to 10	) cigarettes	per	day,	from .	French
					ecologica	l data								

	Linear	model	Exponential model			
	1953 <sup>(a)</sup>	1957(b)	1953(2)	1957(5)		
1968	8.2	8.2	6.0	5.1		
Death rate	(5.8 : 11.3)(*)	(5.7 : 11.6)(*)	(4.2 : 8.6)(*)	(3.6 : 7.2)(*)		
1975	5.3	5.2	4.7	4.1		
Death rate	(4.1 : 6. <sup>-</sup> )(*)	(4.0 : 6.6)(*)	(3.5 : 6.2)(*)	(3.1 : 5.4)(*)		

(\*) 95% confidence intervals.

Evaluation of cigarette consumption (a) in 1953. (b) in 1957.

# INTERNATIONAL JOURNAL OF EPIDEMIOLOGY

Fitted model	Estimated slope	% variance explained	Significant(*) explanatory variables
Lung cancer	······································		
$y = \alpha + \beta X + \varepsilon$	$7.76 \times 10^{-7}$	57.3%	Cigarette consumption
$y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$7.84 \times 10^{-7}$	67.9%	Cigarette consumption, latitude
$\log y = \alpha + \beta X + \varepsilon$	$1.09 \times 10^{-3}$	55.3%	Cigarette consumption
$\log y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$1.10 \times 10^{-3}$	63.7%	Cigarette consumption, latitude
Bladder cancer			
$y = \alpha + \beta X + \varepsilon$	$1.07 \times 10^{-7}$	30.5%	Cigarette consumption
$y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$0.71 \times 10^{-7}$	35.4%	Cigarette consumption, longitude
$\log y = \alpha + \beta X + \varepsilon$	$0.97 \times 10^{-3}$	27.2%	Cigarette consumption
$\log y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$0.61 \times 10^{-3}$	33.1%	Cigarette consumption, longitude
CHD			
$y = \alpha + \beta X + \varepsilon$	$5.37 \times 10^{-7}$	22.1%	Cigarette consumption
$y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$5.46 \times 10^{-7}$	32.9%	Cigarette consumption, latitude
$\log y = \alpha + \beta X + \varepsilon$	$0.34 \times 10^{-3}$	21.1%	Cigarette consumption
$\log y = \alpha + \beta X + a_1 U + a_2 V + \varepsilon$	$0.35 \times 10^{-3}$	31.2%	Cigarette consumption, latitude

 TABLE 2
 Simple and partial correlations between lung cancer, bladder cancer and CHD mortalities and cigarette consumption after adjustment for a linear trend.

y = standardized mortality (men) 1968.

x = cigarette consumption per inhabitant 1953. V: latitude.

(\*) p < 0.05.

individual studies. A possible explanation may be found through failure to take into account ecological confounding factors which would be positively correlated with the distribution of cigarette consumption. Thus the slope of the regression between bladder cancer mortality and cigarette consumption would also include part of the effect of these confounding factors and lead to an overestimation of the relative risk. A



FIGURE 4 Relative risk for bladder cancer as a function of current number of cigarettes smoked per day estimated in four case-control studies and in ecological studies. 1 Vineis<sup>38</sup> - 2 Howe<sup>37</sup> - 3 Wynder<sup>35</sup> - 4 Cartwright<sup>36</sup>

U: longitude.

good candidate for such a confounding factor is the exposure to some industrial risk factors.<sup>39</sup> Indeed, we know from consumer surveys<sup>28</sup> that cigarette consumption is higher among industrial workers. When a linear gradient is adjusted on ecological data (Table 2) the slope of the dose-response relationship varies more than for lung cancer, thus pointing to the likely influence of confounding factors having an heterogeneous spatial distribution.

## Coronary Heart Disease and Cigarette Consumption

Figure 5 shows the relative risks for coronary heart disease (ICD 410-414) with respect to the current number of cigarettes smoked per day estimated in five cohort studies,<sup>40</sup> together with the ecological estimates for the French data analysed or for the US data taken from Friedman.<sup>41</sup> Estimates are given both for the standardized mortality and for the age group 45-54.

In Figure 5 one sees a good agreement between the individual and ecological studies and in Table 3 it is shown that the estimates for the 45-54 age group are higher than those for the 35-74 one. This tallies with the Framingham study<sup>42</sup> where the effect of cigarette consumption is more pronounced among middle-aged men than for those over 65. Just as for lung cancer the estimates for 1968 seem to give a better fit than those for 1975. All the correlation coefficients corresponding to Table 3 are statistically significant and similar for the two periods but are of a lower order of magnitude (from 0.35 to 0.47) than for lung or bladder cancer. The

TABLE 3 Ecological relative risks (French data) estimates for an average consumption of 10 cigarettes per day

Cause and year of death	Linear model	Exponential model
Bladder cancer 1968 (a)	5.4 (3.3 : 8.3)(*)	5.0 (2.8 : 8.9)(*)
Bladder cancer 1975 (b)	5.3 (3.6 : 7.6)(*)	4.6 (3.0 : 6.9)(*)
CHD 1968 (a)	1.8 (1.4 : 2.2)(*)	1.8 (1.4 : 2.2)(*)
CHD 1968: age (45–54) only (a)	2.4 (1.6 : 3.2)(*)	2.5 (1.7 : 3.7)(*)
СНД 1975 (b)	1.5 (1.2 : 1.8)(*)	1.5 (1.2 : 1.9)(*)
CHD 1975: age (45–54) only (b)	2.0 (1.5 : 2.6)(*)	2.0 (1.4 : 2.7)(*)

(\*) 95% confidence interval.

Evaluation of cigarette consumption (a) in 1953, (b) in 1957.

linear and exponential models give very close estimates due to a less pronounced convexity of the exponential curve. Adjustment of a linear trend does not modify the slopes of the dose-response relationships (Table 2).

## Oesophagus Cancer and Cigarette Consumption

The consumption of alcoholic beverages and of tobacco products has long been suspected of increasing the risk of oesophagus cancer (ICD 150).43-46 The ecol-



FIGURE 5 Relative risk for coronary heart disease as a function of current number of cigarettes smoked per day estimated in five cohort studies and in ecological studies. Cohort studies as quoted in:40

- 1 Veterans Administration study 2 British Physicians study
- 3 American Cancer Society study
- 4 Framingham study
- 5 Pooling Project

ogical analysis presented concerns 1968 death rates and cigarette consumption for the year 1953. A grouping of the deaths from alcoholism (ICD 291, 303) and cirrhosis of the liver (ICD 571) was used as an indicator of alcoholic beverage consumption and these causes will be referred to as alcoholic mortality.

The simple correlation between standardized oesophagus cancer mortality for men and cigarette consumption is negative and statistically significant (r = -0.28) and the same result holds for the partial correlation coefficient when alcoholic mortality is taken into account, whether an additive model (r = -0.27) or a multiplicative one (r = -0.28) is fitted. Looking at Figure 6, a group of eight contiguous 'départements' situated in Normandy or Brittany clearly stands out by their very high standardized

oesophagus mortality and their low cigarette consumption. When these eight 'départements' are excluded the correlation coefficients between oesophagus cancer mortality and cigarette consumption becomes non-significant for the additive and the multiplicative model.

Thus cigarette consumption does not appear to play a part in the variations by 'départements' of the male mortality for oesophagus cancer in France. This is not due to a parallelism in the geographical trends of alcoholic mortality and cigarette consumption since in France they display nearly independent variations at the 'département' level although they are correlated among individuals.

In a recent analysis of the geographical pathology of cancer of the oesophagus, Day<sup>47</sup> points out that there



FIGURE 6 Oesophagus cancer mortality and average daily cigarette consumption in 82 French départements: regression line (a) including all départements. (b) excluding eight contiguous départements \* situated in the west of France.

are some geographical contrasts, for instance within Normandy and Brittany, that cannot be explained solely in terms of variations of alcohol or cigarette consumption. It seems that even in industrial countries nutritional deficiencies which could interact with alcohol consumption might play an important part in the aetiology of the disease. Hence finding the precise role played by smoking is a difficult task, the more so as the dose-response relationship is weaker with respect to smoking than with respect to alcohol consumption.<sup>48</sup> In this context ecological analysis in France at the 'département' level does not contribute to our understanding of the aetiology of oesophagus cancer.

### CONCLUSIONS

In this paper the biases that one encounters when trying to estimate a relative risk from an ecological study are discussed.

For smoking and lung cancer, a good match between individual and ecological relative risk estimates is obtained. In this case smoking is an overwhelming risk factor with a fairly linear dose-response relationship in the low and medium dose range. Furthermore, time lags were included and a careful evaluation of the relevant exposure was possible. Using a time lag is necessary but its magnitude is difficult to interpret since both mortality and consumption have been averaged over different age groups. This example also illustrates that the choice of the model has consequences on the estimates because of the frequent need for extrapolation to zero.

In the case of bladder cancer and cigarette consumption the ecological relative risk found is higher than the one given by individual studies. This example brings out the importance of considering the geographical distribution of confounding factors. This is also underlined by the modification of the regression coefficients when spatial trends are adjusted.

As in the case of bladder cancer smoking is not necessarily the dominant risk factor involved in CHD, but here the geographical estimates show a better agreement with the individual studies. The difference between the two examples probably stems from the different geographical distribution of possible confounding factors. Whilst cigarette consumption and industrial exposure are very likely to be positively correlated, this is presumably not the case for the other important risk factors for CHD such as diet. In this example of a multifactorial disease the closeness of the geographical and individual estimates of the relative risks might thus be partially coincidental.

The last example on smoking and oesophagus cancer illustrates that, in the case of a complex multifactorial

disease, ecological analysis techniques are likely to fail to bring out aetiological factors when performed with routinely collected data at a regional level.

## ACKNOWLEDGMENTS

This work was supported in part by the Radiation Protection Programme of the Commission of the European Communities, contract BIO-F-515-82-F.

The authors would like to thank Dr F Hatton, who kindly provided the mortality data used in the present study, Miss N Lincot for her technical assistance and Mrs M Guerrois for typing the manuscript.

#### REFERENCES

- <sup>1</sup> Doll R. The epidemiology of cancer. Cancer 1980; 45: 2475-85.
- <sup>2</sup> Armstrong B. Doll R. Environmental factors and cancer incidence and mortality in different countries, with special reference to dietary practices. Int J Cancer 1975; 15: 617-31.
- <sup>3</sup>Blot WJ, Mason TJ, Hoover R, Fraumeni JF. Cancer by county: etiologic implications. Origins of human cancer 1977.
- <sup>4</sup> Graham S. Methodological problems in ecologic studies of the asbestos cancer relationship. *Envir Res* 1981; 25: 35-49.
- <sup>5</sup> Cantor KP, Hoover R, Mason TJ, McCabe LJ. Associations of cancer mortality with halomethanes in drinking water. J Natl Canc Inst 1978; 61: 979-85.
- <sup>6</sup> MacMahon B. Pugh TF. Epidemiology: principles and methods. Boston, Little Brown, 1970.
- <sup>\*</sup> Breslow NE, Enstrom JE. Geographic correlations between cancer mortality rates and alcohol-tobacco consumption in the United States. J Natl Canc Inst 1974: 53: 631-9.
- <sup>8</sup> Stravraky KM. The role of ecologic analysis in studies of the etiology of diseases: a discussion with reference to large bowel cancer. J Chron Dis 1976; 29: 435–44.
- <sup>o</sup> Davies JM, Chilvers C. The study of mortality variations in small administrative areas of England and Wales, with special reference to cancer. J Epidemiol Community Health 1980; 34: 87-92.
- <sup>10</sup> Richardson ST, Hémon D. On the variance of the sample correlation between two independent lattice processes. J Appl Prob 1981; 18: 943-8.
- <sup>11</sup> Cook DG, Pocock SJ. Multiple regression in geographic mortality studies with allowance for spatially correlated errors. *Biometrics* 1983; 39: 361-71.
- <sup>12</sup> Cliff AD, Ord JK. Spatial processes, models and applications. London, Pion, 1981.
- <sup>13</sup> Morgenstern H. Uses of ecologic analysis in epidemiologic research. Am J Publ Health 1982; 72: 1336-44.
- <sup>14</sup> Béral V, Chilvers C, Fraser P. On the estimation of relative risk from vital statistical data. J Epidemiol Community Health 1979; 33: 159-62.
- <sup>15</sup> Kendall M, Stuart A. The advanced theory of statistics, vol 1, London. Charles Griffin, 1977.
- <sup>16</sup> Robinson WS. Ecological correlations and the behaviour of individuals. Am Sociol Rev 1950; 15: 351-57.
- <sup>17</sup> Goodman LA. Some alternatives to ecological correlation. Am J Sociol 1959; 64: 610-25.
- <sup>18</sup> Thomas EN, Anderson DL, Additional comments on weighting values in correlation analysis of areal data. Annals of the Association of American Geographers 1965; 55: 492-505.
- <sup>19</sup> Alker HR. A typology of ecological fallacies. In: M. Dogan and S.

Rokkan (Eds) Quantitative ecological analysis in the social sciences, Cambridge, Ma, MIT Press, 1969.

- <sup>20</sup> Clark WAV, Avery KL. The effects of data aggregation in statistical analysis. *Geographical Analysis* 1976; 8: 428–33.
- <sup>21</sup> Firebaugh G. A rule for inferring individual level relationship from aggregate data. Am Sociol Rev 1978; **43:** 557-72.
- <sup>22</sup> Openshaw S. Ecological fallacies and the analysis of areal census data. Envir Pl-A 1984; 16: 17-31.
- <sup>23</sup> Finney DJ. Statistical methods in biological assay, London, Charles Griffin, 1952.
- <sup>24</sup> Doll R, Peto R. Avoidable risks of cancer in the U.S. J Natl Canc Inst 1981; 66: 1191-1308.
- <sup>25</sup> Damiani P, Massé H. Mortalité par cause et tabac: application d'un modèle de liaison et évaluation de la consommation de tabac par sexe et par âge. Journal de la Société de Statistique de Paris 1980; 121: 81-89.
- <sup>26</sup> Damiani P, Massé H. Liaison de la mortalité par cause avec l'ensemble des consommations de tabac et d'alcool. Journal de la Société de Statistique de Paris 1981; 122: 174-81.
- <sup>27</sup> Atlas of tobacco consumption in France, SEITA 1953–1957.
- <sup>28</sup> Etude du marché français du tabac, Enquête générale 1960, SEMA.
- <sup>29</sup> U.S. Department of Health and Human Services. The health consequences of smoking. *Cancer: A report of the Surgeon General*. Department of Health and Human Services, Public Health Service, Office on Smoking and Health, 1982.
- <sup>30</sup> Stocks P. Cancer mortality in relation to national consumption of cigarettes. solid fuel, tea and coffee. Br J Canc 1970; 24: 215– 25.

<sup>31</sup> World Health Statistics. Geneva, World Health Organization. 1979.

- <sup>32</sup> Lee PN. Tobacco consumption in various countries. London, Tobacco Research Council, 1975.
- <sup>33</sup> Fraumeni J. Cigarette smoking and cancers of the urinary tract: geographic variation in the United States. J Natl Canc Inst 1968; 41: 1205-11.
- <sup>34</sup> Lazar P. Geographical correlations between disease and environmental exposures. Proceedings of the European Symposium on Medical Statistics, Rome, 1980.

- <sup>35</sup> Wynder EL. Goldsmith R. The epidemiology of bladder cancer. Cancer 1977; **40**: 1246–68.
- <sup>36</sup> Cartwright RA, Adib R, Appleyard I et al. Cigarette smoking and bladder cancer: an epidemiological inquiry in West Yorkshire. J Epidemiol Community Health 1983; 37: 256-63.
- <sup>37</sup> Howe GR, Burch JD, Miller AB et al. Tobacco use, occupation, coffee, various nutrients, and bladder cancer. J Natl Canc Inst 1980; 64: 701-13.
- <sup>38</sup> Vineis P, Frea B, Uberti E, Ghisetti V, Terracini B. Bladder cancer and cigarette smoking in males: a case-control study. *Tumori* 1983; 69: 17-22.
- <sup>39</sup> Vineis P, Segnan N, Costa G, Terracini B. Evidence of a multiplicative effect between cigarette smoking and occupational exposures in the aetiology of bladder cancer. *Cancer Let* 1981; 14: 285-90.
- <sup>40</sup> Kannel WB, Framingham MD. Update on the role of cigarette smoking in coronary artery disease. Am Heart J 1981; 101: 319-28.
- <sup>41</sup> Friedman GD. Cigarette smoking and geographic variation in coronary heart disease mortality in the United States. J Chron Dis 1967; 20: 769–79.
- <sup>42</sup> Dawber TR. *The Framingham study*. Cambridge, Mass, Harvard University Press, 1980.
- <sup>43</sup> Wynder EL, Bross IJ. A study of etiological factors in cancers of the oesophagus. *Cancer* 1961; 14: 389–413.
- <sup>44</sup> Schwartz D, Flamant R, Lellouch J et al. Alcool et cancer. Résultats d'une étude rétrospective. Rev Fr Etudes Clin Biol 1962; 7: 590-604.
- <sup>45</sup> Tuyns AJ, Péquignot G, Jensen OM. Le cancer de l'oesophage en Ille et Vilaine en fonction des niveaux de consommation d'alcool et de tabac. *Bulletin du cancer* 1977; 64: 1: 45–60.
- <sup>46</sup> Tuyns A. Oesopheal cancer in non-smoking drinkers and in nondrinking smokers. Int J Canc 1983; 32: 443-4.
- <sup>47</sup> Day NE. The geographic pathology of cancer of the oesophagus. British Medical Bulletin 1984; 40: 329-34.
- <sup>44</sup> Breslow NE. Day NE. The analysis of case-control studies. Lyon: IARC Scientific Publications no 32, 1980.

(Revised version received April 1986)