

L. LEBART

104

La formule précédente demande, pour le calcul effectif de  $\hat{\beta}$ , l'inversion de matrices ayant autant de lignes qu'il y a d'observations, (matrices  $88 \times 88$  pour le cas des statistiques départementales. Ces matrices faisant l'objet d'estimations assez approchées, le résultat des calculs risque de n'avoir pas grand sens. Il semble donc plus efficace de retenir pour  $\hat{\beta}$  l'estimation centrée :

$$\hat{\beta} = (Z'Z)^{-1} Z'X$$

en tenant compte du fait que cette estimation n'est plus à dispersion minimale, et en estimant la nouvelle dispersion de  $\hat{\beta}$ .

$$V(\hat{\beta}) = (Z'Z)^{-1} Z'V(\epsilon)Z(Z'Z)^{-1} \quad (1)$$

c) Estimation de la matrice des covariances des résidus.

Si la variable endogène  $X$  et les variables exogènes  $Z_1, Z_2, Z_p$  sont susceptibles d'avoir été générées par un modèle stationnaire\*, la variable  $\epsilon = X - Z\beta$  est également stationnaire, (la liaison entre les observations  $Z_{i1}$  et  $Z_{kj}$  des variables  $Z_1$  et  $Z_k$  dépendant de  $l, k$  et de  $d(i, j)$ , sans dépendre ni de  $i$ , ni de  $j$ ).

On peut donc supposer que la matrice des covariances du vecteur résiduel est du type :  $V(\epsilon) = v_0 I + v_1 M_1 + \dots + v_e M_e$ .

Pour les échantillons importants, on estimera  $v_0, \dots, v_e$  à partir des résidus observés (comme précédemment,  $v_0$  est la variance résiduelle,  $v_i$  la covariance des deux résidus correspondant à des sommets distants de  $i$  sur le graphe,  $M_i$  la Matrice associée au graphe des distances «  $i$  »).

Le biais résultant de l'utilisation de  $\hat{\beta}$  au lieu de la vraie valeur de  $\beta$  est d'autant plus réduit que la taille «  $n$  » de l'échantillon est importante : (si  $\hat{\epsilon}$  désigne le vecteur de résidus observés :  $\hat{\epsilon} = Z(\beta - \hat{\beta}) + \epsilon$ ).

Par suite, quand  $n \rightarrow \infty$ ,  $\hat{\beta}$  tend en probabilité vers  $\beta$ , et  $\hat{\epsilon}$  tend en probabilité vers  $\epsilon$ .

\* c'est-à-dire tel que la covariance de deux observations ne dépend que de leur distance sur le graphe.

ANALYSE STATISTIQUE DE LA CONTIGUÏTE

105

Sans faire intervenir les valeurs numériques des résidus on peut estimer  $v_i(\epsilon)$  à partir de  $v_i(X)$  et des  $v_i(Z_k)$  à l'aide de la relation  $E(\hat{\epsilon}\hat{\epsilon}') = E(X - Z\hat{\beta})(X - Z\hat{\beta})'$ .

(Notons que les tests de significativité des coefficients de contiguïté  $C(\epsilon)$  ne s'appliquent plus rigoureusement, car dans l'hypothèse où les résidus théoriques ont une matrice de covariance  $\sigma^2 I$ , les résidus observés ont

une matrice de covariance :  $\sigma^2 (I - \frac{U}{n} - Z(Z'Z)^{-1}Z')$ , d'où une inextricable complication de la loi de  $C(\hat{\epsilon})$ .

On aboutit donc à une estimation du type :

$$V(\epsilon) = v_0 I + v_1 M_1 + \dots + v_e M_e$$

Par suite, en remplaçant  $V(\epsilon)$  par sa valeur dans la formule (1)

$$V(\hat{\beta}) = (Z'Z)^{-1} Z'(v_0 I + v_1 M_1 + \dots + v_e M_e)Z(Z'Z)^{-1}$$

soit :

$$(2) \quad V(\hat{\beta}) = v_0 (Z'Z)^{-1} + v_1 (Z'Z)^{-1} (Z'M_1 Z) (Z'Z)^{-1} + \dots + v_e (Z'Z)^{-1} (Z'M_e Z) (Z'Z)^{-1}$$

Notons que  $v_0 (Z'Z)^{-1}$  n'est autre que la matrice des covariances de  $\hat{\beta}$  dans l'hypothèse de non contiguïté des résidus.

La formule précédente montre que si les variables exogènes sont « non contiguës », la zone de confiance des coefficients de régression est inchangée, quelle que soit la contiguïté des résidus.

Si les variables exogènes dépendent significativement du graphe, les éléments de  $Z'M_k Z$  ont des valeurs élevées pour les premières valeurs de  $k$ .

Le terme rectificatif est alors donné par la formule (2).

## d) Exemple d'application numérique.

Nous nous proposons d'étudier la dépendance qui peut exister entre un indicateur de dépenses départementales  $D$ , et le taux de population urbaine  $U$  d'un même département.

$$D_i = \alpha U_i + \beta + e_i$$

La M. de C. du vecteur  $\begin{pmatrix} D \\ U \end{pmatrix}$  est  $V = \begin{pmatrix} 653,273 & 332,715 \\ 332,715 & 271,839 \end{pmatrix}$

On a donc :  $D_i = 1,224 U_i + \beta + \epsilon_i$  avec un coefficient de corrélation  $\rho = 0,7897$ .

On trouve :  $\sigma_\alpha^2 = 0,01053$ , sans faire d'hypothèse sur la liaison des résidus.

Ce qui donne approximativement la zone de confiance au seuil 0,05 :

$$1,02 \leq \alpha \leq 1,43$$

En fait, le coefficient de contiguïté des résidus est  $C(\hat{\epsilon}) = 0,396$  ; bien que le test de contiguïté ne s'applique plus rigoureusement puisqu'il s'agit de résidus observés, la taille de l'échantillon (88) et la valeur très faible trouvée pour  $C(\hat{\epsilon})$  (à 7 écart-types de la moyenne) nous conduisent à considérer que les résidus sont significativement contigus. En ne retenant que le terme rectificatif qui fait intervenir des liaisons au niveau 1 :

On trouve  $\sigma_\alpha^2 = 0,01053 + 0,0083 = 0,01883$  ( $\sigma_\alpha = 0,137$ ).

La nouvelle zone de confiance est alors :  $0,950 \leq \alpha \leq 1,498$ , pour le même seuil de 0,05.

D'une manière générale, on notera que, comme dans le cas des séries temporelles, la précision des coefficients est souvent surestimée lorsqu'il existe des liaisons directes.

## B - Perte d'information et « pseudo-taille » de l'échantillon.

Le modèle théorique le plus simple, susceptible de générer des réalisations dépendantes du graphe est le modèle « stationnaire », où il existe

entre deux observations un coefficient d'autocorrélation qui ne dépend que de leur distance sur le graphe.

Ainsi, le vecteur des observations d'une variable  $X$  a pour matrice des covariances :

$$V = v_0 I + v_1 M_1 + \dots + v_k M_k$$

(les différentes observations ont même moyenne  $\mu$ , et même variance  $v_0$ ).

$I$  est la matrice unité.

$M_1, \dots, M_k$  les matrices associées aux graphes  $G_1, \dots, G_k$  ( $m_{u,ij}$   $\epsilon$ )

$M_u = 1$  si  $d(i, j) = u$ ,  $m_{u,ij} = 0$  si  $d(i, j) \neq u$ .

$\bar{\mu}$  est le vecteur dont toutes les composantes sont égales à  $\mu$ .

Calculons les quantités d'informations relatives aux paramètres  $\mu$  et  $v_u$  ( $u = 1, \dots, k$ ).

La densité de probabilité du vecteur  $X$  s'écrit :

$$P(X) = \frac{1}{(2\pi)^{\frac{n}{2}} |V|} \exp \left\{ -\frac{1}{2} (X - \bar{\mu})' V^{-1} (X - \bar{\mu}) \right\}$$

$$\text{Log } P(X) = -\frac{n}{2} \text{Log } 2\pi + \frac{1}{2} \text{Log } |V^{-1}| - \frac{1}{2} \text{tr} (V^{-1} (X - \bar{\mu}) (X - \bar{\mu})')$$

Posons  $V^{-1} = C = (c_{ij})$ ,

$$(X - \bar{\mu}) (X - \bar{\mu})' = D(\mu),$$

$$\text{et } \frac{1}{2} (\text{Log } |C| - \text{tr} (CD)) = L.$$

$$\text{On a : } \frac{\partial L}{\partial v_u} = \sum_{i,j} \frac{\partial L}{\partial c_{ij}} \cdot \frac{\partial c_{ij}}{\partial v_u} = \text{tr} \left( \frac{\partial L}{\partial C} \cdot \frac{\partial C}{\partial v_u} \right)$$

où la matrice  $\frac{\partial L}{\partial C}$  a pour terme générique  $\frac{1}{2} \left( \frac{\partial L}{\partial c_{ij}} \right)$  et  $\left( \frac{\partial L}{\partial c_{ii}} \right)$  sur la diagonale :

$$\text{La matrice } \frac{\partial L}{\partial C} \text{ est telle que } \begin{cases} \frac{\partial L}{\partial c_{ij}} = \left( \frac{\text{cof}(c_{ij})}{|C|} - d_{ij} \right) \\ \frac{\partial L}{\partial c_{ii}} = \frac{1}{2} \left( \frac{\text{cof}(c_{ii})}{|C|} - d_{ii} \right) \end{cases}$$

$$\text{Donc } \frac{\partial L}{\partial C} = \frac{1}{2} (C^{-1} - D) = \frac{1}{2} (V - D)$$

$$\text{On a aussi : } \frac{\partial C}{\partial v_u} = \frac{\partial V^{-1}}{\partial v_u} = -V^{-1} \frac{\partial V}{\partial v_u} V^{-1} = -V^{-1} M_u V^{-1}$$

$$\text{donc : } \frac{\partial L}{\partial v_u} = -\frac{1}{2} \text{tr} (V - D) V^{-1} M_u V^{-1}$$

$$\text{D'autre part : } \frac{\partial L}{\partial \mu} = \text{tr} (C \cdot \frac{\partial D}{\partial \mu})$$

$$\text{Calculons : } \frac{\partial^2 L}{\partial v^2} \text{ et } \frac{\partial^2 L}{\partial \mu^2} :$$

$$\begin{aligned} \frac{\partial V^{-1} M_u V^{-1}}{\partial v_u} &= \frac{\partial V^{-1}}{\partial v_u} \cdot M_u V^{-1} + V^{-1} \frac{\partial M_u V^{-1}}{\partial v_u} \\ &= -V^{-1} M_u V^{-1} M_u V^{-1} - V^{-1} (M_u \cdot V^{-1} M_u V^{-1}) \\ &= 2 V^{-1} (M_u V^{-1})^2 \end{aligned}$$

Par suite :

$$\frac{\partial^2 L}{\partial v_u^2} = -\frac{1}{2} \text{tr} (2 DV^{-1} - I) (M_u V^{-1})^2$$

Si les observations sont indépendantes  $v = \sigma^2 I$  et  $M_u = M_0 = I$ , l'équation  $\frac{\partial L}{\partial \sigma^2} = 0$  redonne  $\sigma^2 = \frac{1}{2} \sum (x_i - \mu)^2$  et l'on a bien

$$\frac{\partial^2 L}{\partial \sigma^2} = \frac{-1}{2\sigma^4} (2n - 1) < 0$$

$$\text{On a de même } \frac{\partial^2 L}{\partial \mu^2} = -\sum_{ij} c_{ij} = -\text{tr} (V^{-1} U)$$

(avec  $u_{ij} = 1$  pour tout  $i$  et  $j$ ).

Par suite :

$$I_{v_u} = E \left( \frac{-\partial^2 L}{\partial v_u^2} \right) = \frac{1}{2} \text{tr} (M_u V^{-1}) \text{ et } I_\mu = E \left( -\frac{\partial^2 L}{\partial \mu^2} \right) = \dots = + \text{tr} (V^{-1} U)$$

$I_{v_u}$  et  $I_\mu$  étant les quantités d'information de Fischer, relatives aux paramètres  $v_u$  et  $\mu$ .

Dans le cas où  $V = \sigma^2 (I + \rho M)$ , cas où la contiguïté ne se manifeste qu'au niveau (1), on a :

$$I_m = \frac{1}{\sigma^2} \text{tr} (I + \rho M)^{-1} U = \frac{1}{\sigma^2} \text{tr} (I - \rho M + \rho^2 M^2 - \rho^3 M^3 + \dots) U$$

Or,  $\text{tr} M^\alpha U$  est le nombre «  $c_\alpha$  » de chemins de longueur  $\alpha$  du graphe.

$$\text{Donc, } I_m = \frac{1}{\sigma^2} (n - c_1 \rho + c_2 \rho^2 - c_3 \rho^3 + \dots)$$

Si  $\rho = 0$ , on retrouve la quantité connue  $\frac{n}{\sigma^2}$ . On constate ce qui est intuitif, que la quantité d'information sur la moyenne, les autres paramètres étant connus, est plus élevée si les liaisons (caractérisées par  $\rho$ ) sont « négatives ». Elle est, en général, plus faible si les liaisons sont « positives ».

On peut définir la pseudo-taille  $n_1$  de l'échantillon comme l'effectif qui suffirait à donner la même information sur la moyenne, si les observations étaient indépendantes.

$n_1$  est donc tel que :

$$\frac{n_1}{\sigma^2} = I_n = \frac{1}{\sigma^2} (n - c_1 \rho + c_2 \rho^2 \dots)$$

Soit :

$$n_1 = n - \sum c_i \rho^i = \text{tr} (I + \rho M)^{-1} U$$

Supposons ainsi que le graphe soit formé de doublets disjoints : l'échantillon est donc formé de couples d'observations corrélées. Le nombre de chemin de longueur  $k$  est constant et égal à  $n$  quelque soit  $k$ .

$$\text{On a donc ici : } n_1 = n (1 - \rho + \rho^2 - \rho^3 \dots) = \frac{n}{(1 + \rho)}$$

Pour  $\rho = 1$ , c'est-à-dire lorsque les  $n$  observations se décomposent en  $\frac{n}{2}$  couples contenant la même observation, on a bien :  $n_1 = \frac{n}{2}$ .

Supposons au contraire que le graphe soit complet (tout couple de sommet est relié par une arête). Toutes les observations sont corrélées entre elles de la même façon. Le calcul direct de l'inverse de  $V = I + \rho M$  est immédiat, et l'on trouve :

$$n_1 = \frac{n}{1 + (n-1)\rho}$$

Si  $\rho = 1$ , tout se passe comme si l'on n'avait qu'une seule observation.

La pseudo-taille de l'échantillon est  $n_1 = 1$ .

Pour des graphes plus complexes, le calcul de  $n_1 = \text{tr} \left\{ (I + \rho M)^{-1} U \right\}$  doit être fait automatiquement ; dans ce cas, on peut estimer  $\rho$  par le complément à 1 du coefficient de contigüité ( $\rho = 1 - c$ ).

Les valeurs trouvées dans le cas des statistiques départementales ( $n_1$  varie de 40 à 80 pour les variables des exemples précédents, où  $n=88$ ) doivent inciter à une grande prudence en ce qui concerne l'interprétation des tests.

## BIBLIOGRAPHIE

- [ 1 ] BERGE (C.) - *Théorie des graphes et ses applications*. (Dunod 1963).
- [ 2 ] CHARTIER (F.) - *L'estimation statistique dans le cas d'observations non indépendantes. Etude d'un cas particulier I.S.U.P.* Vol. 1 Fasc. 2 - 1952.
- [ 3 ] FERREZ (J.) - *Inégalité régionale des possibilités d'accès à l'Education*. (Publ. O.C.D.E. 1961).
- [ 4 ] GEARY (R.C.) - A general expression for the moments of certain symmetrical functions of normal sample. (*Biometrika* XXV - p. 184 - 1933).  
*The contiguity ratio and statistical mapping « the incorporated statistician »* Vol. 5, n° 3, p. 115-45, 1954).
- [ 5 ] GURLAND (J.) - Quadratic forms in normally distributed random variables - *Sankhya* - Vol. 17 (1956).
- [ 6 ] ISARD (W.) - *Methods of regional analysis. An introduction to regional science* (The MIT Press - Cambridge - Massachusett).
- [ 7 ] JORESKO (K.G.) - *Statistical estimation in factor analysis. A new technique and its foundations* (Almqvist & Wiksell - Stockholm, Uppsala - 1963).
- [ 8 ] LEDERMANN (S.) - *Cours d'analyse multivariate*, 1966, I.S.U.P. (Ronéo).
- [ 9 ] MALINVAUD (E.) - *Méthodes statistiques de l'économétrie*. Dunod 1964.
- [ 10 ] TRANQUILLI (G.B.) - *Les lois jointes des statistiques et de l'écart carré moyen des échantillons normaux* (thèse 1963).
- [ 11 ] Von NEUMANN (J.) - Distribution of the ratio of the mean square successive difference to the variance (*Annals of math. stat* Vol. 12 - 1941).

- [12] *Cahiers de l'I.S.E.A. – Série L – N° 1 à 13, en particulier n° 10 – Bibliographie de science économique régionale.*
- [13] *Bulletin de l'Institut International de Statistique - 1958 - Tome 36 (Divers articles sur l'économie régionale).*
- [14] «L'espace économique français» (1955, 1962) et diverses remises à jour. (n° spécial de la revue « Etudes et conjoncture »).