

## MODÈLE MULTI-INDICÉ (G2, G11, H7, K1, K15)

(23 / 03 / 2020, © Monfort, Dicostat2005, 2005-2020)

(i) De façon générale, un **modèle multi-indicé** est une **représentation statistique** dans laquelle les **variables** considérées sont observées selon plusieurs **critères d'observation**, eg :

- (a) « **unité statistique** x **période d'observation** » : eg modèle de « panel » ;
- (b) « **unité statistique** x **espace** » : eg modèle « spatial » ;
- (c) « **espace** x **temps** » : eg processus spatio-temporel.

Ainsi, dans un **modèle d'analyse de la variance**, divers facteurs peuvent influencer sur une grandeur numérique donnée : les données sont alors disposées selon une « dimension » **facteur** et une dimension **observation**.

(ii) Un modèle multi-indicé se distingue d'un modèle comportant des variables « multiples », eg :

(a) un **modèle multivarié** (cf **loi multivariée**), dans lequel la liste des **variables d'intérêt** est constituée de variables de différents types : **variables qualitatives**, **variables numériques** ;

(b) un **modèle multidimensionnel** (cf **analyse multidimensionnelle**, **loi multidimensionnelle**), moins général mais plus courant que le précédent, dans lequel les variables d'intérêt sont des **vecteurs aléatoires** (très souvent réels), ie prennent leurs valeurs dans des ensembles à plusieurs dimensions (au sens algébrique) (cf **espace vectoriel**).

(iii) Chaque « observation » d'un **modèle multi-indicé** est un **tableau statistique** multidimensionnel, de la forme générale à H « dimensions » :

$$(0) \quad T = (t_l)_{l \in \mathcal{I}}, \quad \text{avec } t_l \in \mathbf{R}, \forall l = (i_1, \dots, i_H) \in \mathcal{I},$$

où  $\mathcal{I} = \prod_{h=1}^H \mathcal{I}_h$  et  $\mathcal{I}_h = \{1, \dots, n_h\}$  (multi-**indice**).

On note  $T_H(\mathbf{R})$  l'ensemble des tableaux statistiques réels H-dimensionnels.

(iv) A titre d'exemple, soit  $(\Omega, \mathcal{F}, \mathcal{P})$  un **modèle statistique** fondamental et :

$$(1) \quad \eta, \xi_1, \dots, \xi_K : \Omega \mapsto \mathbf{R}$$

une suite (ou « liste ») constituée de K **vars**  $\xi_k$ .

Un **modèle de régression multiple** linéaire tq (dans l'**espace des variables**) :

$$(2) \quad \eta = \zeta + \varepsilon = \sum_{k=1}^K b_k \xi_k + \varepsilon, \quad \text{avec } E \varepsilon = 0,$$

est alors « observé » (dans l'**espace des observations**) sous la forme :

$$(3) \quad y_I = z_I + u_I = \sum_{k=1}^K b_k x_{k,I} + u_I,$$

avec  $E u_I = 0, \forall I \in \mathcal{I}$ , dans laquelle on peut noter  $Y = (y_I)_{I \in \mathcal{I}} : \Omega \mapsto T_H(\mathbf{R})$ ,  $X_k = (x_{k,I})_{I \in \mathcal{I}} : \Omega \mapsto T_H(\mathbf{R}), \forall k \in N_K^*$ , et  $u = (u_I)_{I \in \mathcal{I}} : \Omega \mapsto T_H(\mathbf{R})$  sont  $1 + K + 1$  tableaux aléatoires. Un tel modèle est ainsi multi-indicé par  $I = (i_1, \dots, i_H) \in \mathcal{I}$ .

(v) De même, un **modèle de régression multiple** (non linéaire) tq :

$$(4) \quad \eta = f(\xi_1, \dots, \xi_K) + \varepsilon, \quad \text{avec } E \varepsilon = 0,$$

est « observé » sous la forme :

$$(5) \quad y_I = f(x_{1,I}, \dots, x_{K,I}, b) + u_I, \quad \text{avec } E u_I = 0, \quad \forall I \in \mathcal{I},$$

dans laquelle  $b \in \mathbf{R}^Q$ .

Sous l'hypothèse de non **corrélation** et d'**homoscédasticité** suivante :

$$(6) \quad C(u_I, u_J) = \delta_{IJ} \cdot \sigma_u^2, \quad \forall (I, J) \in \mathcal{I}^2,$$

on peut estimer le **paramètre**  $b$  par la **méthode des moindres carrés ordinaires**, qui équivaut à la **méthode du maximum de vraisemblance** dans le cas gaussien, où  $u_I \sim \mathcal{N}_1(0, \sigma_u^2)$  (**loi normale multidimensionnelle** centrée),  $\forall I \in \mathcal{I}$ .

Pour cela, on « empile » les observations de chaque variable  $\eta, \xi_1, \dots, \xi_K$ , eg selon l'**ordre lexicographique**. Ainsi, dans le cas de  $\eta$  :

$$(7) \quad y = (y_{11\dots 11} \quad \text{///} \quad y_{11\dots n(H)} \quad \text{///} \quad y_{11\dots 21} \quad \text{///} \quad y_{11\dots n(H-1)n(H)} \quad \text{///} \quad y_{1n(2)\dots n(H-1)n(H)} \quad \text{///} \quad y_{n(1)n(2)\dots n(H-1)n(H)}),$$

où /// dénote un changement de ligne, ce qui se représente dans le schéma ci-après :

$$y = \begin{pmatrix} y_{11} \dots 11 \\ \dots \\ y_{11} \dots n_H \\ y_{11} \dots 21 \\ \dots \\ y_{11} \dots n_{H-1} n_H \\ \dots \\ y_{1n_2} \dots n_{H-1} n_H \\ \dots \\ y_{n_1 n_2} \dots n_{H-1} n_H \end{pmatrix}$$

Ce **vecteur aléatoire** est à valeurs dans  $\mathbf{R}^N$ , avec  $N = \prod_{h=1}^H n_h$ , ce qui permet de ramener (3) à la forme standard :

$$(8) \quad y = X b + u, \quad \text{avec } E u = 0, \quad V u = \sigma_u^2 \cdot I_N,$$

et (5) à la forme classique :

$$(9) \quad y = F(b) + u, \quad \text{avec } E u = 0, \quad V u = \sigma_u^2 \cdot I_N,$$

où  $X \in M_{NK}(\mathbf{R})$  et  $F : \mathbf{R}^Q \mapsto \mathbf{R}^N$  est une fonction vectorielle (qui dépend, en général, de  $X$ ).

(vi) La « **partie certaine** »  $z_l$  de certains modèles multi-indices de forme additive :

$$(10) \quad y_l = z_l + u_l, \quad \text{avec } E u_l = 0, \quad \forall l \in \mathcal{I},$$

est souvent décomposée de façon adaptée.

Ceci est le cas des modèles (8) et (9), ainsi que :

(a) du **modèle avec constantes spécifiques**, dans lequel :

$$(11) \quad z_l = b_{1,i(1)} + \dots + b_{H,i(H)} = \sum_{h=1}^H b_{h,i(h)},$$

fréquemment utilisé en **analyse de la variance** ;

(b) du modèle plus général :

$$(12) \quad z_l = \sum_{k=1}^K b_k x_{k,l} + \sum_{h=1}^H b_{h,i(h)},$$

fréquemment utilisé en **analyse de la covariance**, etc.

(vii) Dans le cas d'un **modèle bi-indices** ( $H = 2$ ), un **modèle d'analyse de la variance** multidimensionnel peut s'écrire sous la forme :

$$(13) \quad Y_I (|I|, |T|) = X_I (|I|, |K|) b (|K|, |L|) + X_T (|L|, |T|) + U (|I|, |T|),$$

dans laquelle  $Y = (y_{it})_{(i,t)}$  est tq l'observation  $y_{it}$  de la variable endogène  $\eta$  se décompose en fonction des observations  $x_{i,i(k)}$  et  $x_{T,lt}$ , où  $X_I = (x_{i,i(k)})_{(i,k)}$  désigne la matrice des **effets inter-individuels**,  $X_T = (x_{T,lt})_{(l,t)}$  la matrice des **effets intra-individuels**,  $i \in I$  un individu (**unité statistique**) et  $t \in T$  le **temps** (en notant  $|I| = \text{Card } I$ ,  $|T| = \text{Card } T$ , etc).

(vii) Certains modèles multi-indices (10) comportent une partie aléatoire  $u_i$  décomposée de façon adaptée. Il est en effet souvent plus réaliste de supposer que le vecteur aléatoire empilé  $u$  (ou  $y$ ) possède une **matrice de dispersion** non scalaire (ie  $V u \neq \sigma_u^2 \cdot I_N$ ).

Ceci est notamment le cas du **modèle à erreurs composées**, pour lequel :

$$(14) \quad u_i = u_{i(1)} + \dots + u_{i(H)} + v_i, \quad \forall i \in \mathcal{I},$$

avec :

$$E u_{i(h)} = 0, \quad \forall i_h \in \mathcal{I}_h, \quad \forall h \in N_H^*,$$

$$(15) \quad V u_{i(h)} = \sigma_h^2, \quad \forall i_h \in \mathcal{I}_h, \quad \forall h \in N_H^*,$$

$$C (u_{g,i(g)}, u_{h,i(h)}) = 0, \quad \forall (i_g, i_h) \in \mathcal{I}_g \times \mathcal{I}_h, \quad \forall (g, h) \in N_H^* \times N_H^* (g \neq h).$$

(viii) Diverses méthodes (eg **méthode du mv** ou **méthode des mcg**) permettent de traiter ces modèles.

Les **paramètres d'intérêt**  $b$  des différentes formes précédentes (ou de leurs variantes) ne sont pas toujours estimables (donc les modèles ne sont pas toujours identifiables) (cf **estimabilité**, **identifiabilité**) et nécessitent l'introduction de **contraintes** de même nature que celles des modèles d'analyse de la variance ou de la covariance (cf aussi **singularité**).

(ix) Le « temps » est souvent l'un des « critères » selon lesquels on observe les variables  $\eta$ ,  $\xi_1, \dots, \xi_K$ . Les hypothèses (tq (18)) sont alors modifiées pour tenir compte des caractéristiques propres au temps (**relation d'ordre** total, **autocorrélation** temporelle, etc).

De même, les modèles précédents sont souvent « observés » sur des **séries temporelles** de **coupes instantanées** : eg unités statistiques « suivies » au cours du temps (données de panels).

(x) Il existe enfin des modèles « mixtes », dans lesquels, à la fois, les parties certaine  $z_i$  et aléatoire  $u_i$  sont décomposées comme ci-dessus.