

# REVUE DE STATISTIQUE APPLIQUÉE

J. ULMO

## **Problèmes et programmes de régression**

*Revue de statistique appliquée*, tome 19, n° 1 (1971), p. 27-39

[http://www.numdam.org/item?id=RSA\\_1971\\_\\_19\\_1\\_27\\_0](http://www.numdam.org/item?id=RSA_1971__19_1_27_0)

© Société française de statistique, 1971, tous droits réservés.

L'accès aux archives de la revue « *Revue de statistique appliquée* » (<http://www.sfds.asso.fr/publicat/rsa.htm>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme  
Numérisation de documents anciens mathématiques  
<http://www.numdam.org/>

# PROBLÈMES ET PROGRAMMES DE RÉGRESSION

J. ULMO

*Après un rappel des résultats classiques de l'étude des modèles linéaires, on étudie le problème du choix des variables indépendantes à faire intervenir dans une équation de régression.*

*On montre que la suppression à tort d'un certain nombre de variables indépendantes conduit en général à des estimateurs biaisés des coefficients de régression relatifs aux variables restantes, puis on passe en revue les principales procédures de choix des variables indépendantes à faire intervenir dans une équation de régression, quand celle-ci est utilisée pour la prévision par "interpolation" (c'est-à-dire quand on se limite au domaine couvert par l'ensemble des valeurs prises par les variables indépendantes dans l'échantillon utilisé).*

## PLAN

	Pages
I - Introduction : Rappel de résultats classiques de l'étude des modèles linéaires.....	27
II - Comparaison des régressions obtenues sur K échantillons (analyse de la covariance).....	32
III - Choix des variables explicatives ou indépendantes à faire intervenir dans une équation de régression.....	33
Références bibliographies.....	39

## I - INTRODUCTION : RAPPEL DE RESULTATS CLASSIQUES DE L'ETUDE DES MODELES LINEAIRES

### 1/ Définition d'un modèle de régression linéaire

C'est un modèle du type

$$(1) \quad y = \underline{x}'\underline{\beta} + \varepsilon$$

où  $\underline{x}' = (x^1 \dots x^p)^{(1)}$  représente un ensemble de variables dites indépendantes ou explicatives qu'on traite comme si elles étaient non aléatoires, que ceci soit

réalisé ou non, et  $y$  est une variable dépendante ou à expliquer,  $\underline{\beta} = \begin{Bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{Bmatrix}$  est

-----  
(1) Pour faciliter la compréhension du texte les vecteurs (ou leurs transposés) seront généralement représentés par des lettres soulignées.

le vecteur des coefficients de régression de  $y$  sur les  $\underline{x}^\alpha$  ( $\alpha = 1 \dots p$ ) et c'est lui qu'on cherche à estimer à partir d'un  $n$  échantillon de réalisations  $(y_i, \underline{x}_i')$ ,  $i = 1, 2 \dots n$ .

$\varepsilon$  est une variable aléatoire de moyenne nulle, en sorte que  $E^{\mathbf{x}'}(y) = \mathbf{x}'\beta$ .  $E^{\mathbf{x}'}(y)$  = espérance de  $y$  conditionnée par  $\underline{x}'$ .

## 2/ Méthode d'estimation des moindres carrés

Quand on dispose de  $n$  réalisations de  $(y_i, \underline{x}_i')$  le modèle correspondant à (1) s'écrit

(2)  $y = X\beta + \varepsilon$  où  $\underline{y}$  et  $\underline{\varepsilon}$  sont des vecteurs de  $\mathbb{R}^n$  et  $X = \|\underline{x}_i^\alpha\| \begin{cases} \alpha = 1 \dots p \\ i = 1 \dots n \end{cases}$  est une matrice  $(n, p)$ .

La méthode des moindres carrés consiste à chercher un estimateur  $\underline{b}$  de  $\beta$  tel que  $\underline{y}^* = X\underline{b}$  soit le plus voisin possible de  $y$  au sens d'une norme euclidienne de  $\mathbb{R}^n$  lieu de  $\underline{y}$ .

Dans la méthode des moindres carrés ordinaire la norme considérée est la norme habituelle définie par le produit scalaire  $\|\underline{y}\|^2 = \underline{y}'\underline{y} = \sum_{i=1}^n y_i^2$ .

$\underline{b}$  est donc tel que  $\underline{y}^* = X\underline{b}$  et  $\|\underline{y} - \underline{y}^*\|^2 = \|\underline{y} - X\underline{b}\|^2 = (\underline{y} - X\underline{b})'(\underline{y} - X\underline{b})$  soit minimum

Rien n'empêche d'introduire une norme définie par une forme quadratique définie positive de matrice  $M$  quelconque, soit  $\|\underline{y}\|_M^2 = \underline{y}'M\underline{y}$ .

Les seules métriques  $M$  utilisées dans la pratique sont celles qui sont caractérisées par une matrice diagonale.  $M = D_{w_i}$  de poids  $w_i$ . On a alors la méthode des moindres carrés pondérée qui consiste à définir  $\underline{b}$  par la condition

$$(\underline{y} - X\underline{b})'M(\underline{y} - X\underline{b}) = \sum_{i=1}^n w_i (y_i - y_i^*)^2 \text{ minimum avec } y_i^* = \underline{x}_i' \underline{b}$$

En fait on montre que pour conserver les propriétés classiques de la méthode des moindres carrés on doit prendre  $M = K\Sigma^{-1}$  où  $K$  est une constante si  $\Sigma = \text{cov } \underline{y} = \text{cov } \underline{\varepsilon}$ .

On suppose généralement que  $\text{cov } y = \sigma^2 I_n$  ou  $\text{cov } y = D_{\sigma_i^2}$  matrice diagonale d'élément  $\sigma_i^2$ .

## 3 - Interprétation géométrique de la méthode

$\underline{y}^* = X\underline{b}$  est la projection "M orthogonale" de  $\underline{y} \in \mathbb{R}^n$  sur le sous-espace linéaire  $V$  engendré par les colonnes de  $X$  tandis que  $\underline{e} = \underline{y} - \underline{y}^*$  vecteur écart résiduel est M orthogonal à  $V$  c'est-à-dire est élément de  $V^\perp$  sur lequel il est la projection de  $\underline{y}$ .

On a donc

$$\boxed{\begin{matrix} \underline{y} = \underline{y}^* + \underline{e} \\ \underline{e} \in V^\perp \end{matrix}}$$

(1)  $V^\perp$  (ou  $V$  "orthogonal") désigne le sous espace supplémentaire de  $V$  dans  $\mathbb{R}^n$  qui lui est M orthogonal

et  $y'My = y^*My^* + e'Me$  car  $y^*Me = 0$

soit

(3)  $\boxed{\|y\|^2 = \|y^*\|^2 + \|e\|^2}$  théorème de Pythagore généralisé.

Dans le cas où  $M = I$ , auquel on se bornera par la suite, (3) s'écrit

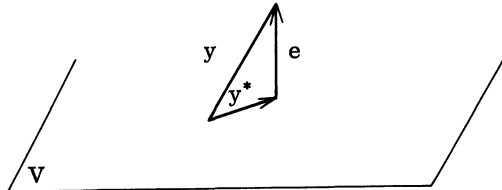
$$y'y = y^*y^* + e'e$$

soit

$$\boxed{\sum_i y_i^2 = \sum_i y_i^{*2} + \sum_i e_i^2}$$

avec

$$\sum_i e_i^2 = \sum_i (y_i - y_i^*)^2$$



#### 4 - Détermination de $\underline{b}$

Quand  $X_{(n,p)}$  est de rang  $p$ , l'application  $f$  associée à  $X$  est injective,  $V$  est de dimension  $p$  et  $\underline{b} = f^{-1}(y^*)$  est défini de façon unique.

Si  $X$  est de rang  $r$  inférieur à  $p$  il y a indétermination non pas sur  $y^* = X\underline{b}$  mais sur  $\underline{b}$  qui appartient à la variété  $f^{-1}(y^*)$  qu'on sait être de dimension  $p - r$ .

C'est ce qui se produit quand parmi les variables indépendantes considérées certaines sont qualitatives.

Soit  $A$  une variable qualitative prenant  $K$  modalités. On lui associe alors  $K$  variables binaires  $x^{\alpha_e}$ ,  $e = 1, 2, \dots, K$  prenant les valeurs 1 ou 0 suivant que  $A$  prend ou non la modalité  $A_k$ . Les  $K$  colonnes de  $X$  relatives à l'ensemble des  $x^{\alpha_e}$  forment alors un ensemble de rang  $K - 1$  puisque la somme des éléments d'une ligne pour ces  $K$  colonnes est égale à 1. On lève généralement l'indétermination en imposant aux composantes  $\beta_{\alpha_e}$ ,  $e = 1, \dots, K$  de  $\beta$

correspondant aux  $x$  de satisfaire à une contrainte, par exemple à  $\sum_{e=1}^K \alpha_e = 0$  ce qui revient à prendre pour origine des effets des modalités  $A_e$ , la moyenne de ces effets.

#### Equations normales et expression de $\underline{b}$

La dérivation par rapport à  $\underline{b}$  de  $(y - Xb)'M(y - Xb)$  et l'annulation de cette dérivée conduit à  $X'M(y - y^*) = 0$  soit  $X'MX\underline{b} = X'My$

$X'M(y - y^*) = 0$  traduit la M orthogonalité de  $y - y^* = e$  et des colonnes de X donc de  $y - y^*$  et de V.

Expression de  $\underline{b}$

Dans la suite on supposera que  $M = I$

(Sinon on pourra toujours écrire  $M = A'A$  et il suffira de poser  $\underline{z} = \underline{A}y$   $U = AX$  pour être ramené au modèle  $\underline{z} = U\underline{\beta} + \underline{\eta}$  ( $\underline{\eta} = A\underline{\varepsilon}$ ) que l'on pourra traiter avec la métrique unité. En effet  $\underline{z}'\underline{z} = y'A'Ay = y'My$ )

Si X est de rang p il en est de même de  $X'X = S$  et on obtient :

$$\underline{b} = (X'X)^{-1} X'y$$

Si X est de rang  $r < p$ , on ajoute les contraintes  $F \underline{b} = 0$  avec F de  $\begin{matrix} (t,p) \\ \end{matrix}$   $p - r$  et telle que  $\text{rang de } G = \begin{vmatrix} X \\ F \end{vmatrix} = p$ .

On a alors  $\begin{cases} X'X\underline{b} = X'y \\ F'F\underline{b} = 0 \end{cases}$

soit  $G'G\underline{b} = X'y$  où  $G'G$  est de rang p.

On en déduit

$$\underline{b} = (G'G)^{-1} X'y$$

$\underline{b}$  est un estimateur de  $\underline{\beta} = (G'G)^{-1} X'E(y)$

qui satisfait à  $\begin{cases} E(y) = X\underline{\beta} \\ F\underline{\beta} = 0 \end{cases}$

5 - Propriétés de  $\underline{b}$  et des écarts résiduels

a) Théorème de Markov

1/  $\underline{\beta}$  étant défini par  $\begin{cases} E(y) = X\underline{\beta} \\ F\underline{\beta} = 0 \end{cases}$  où F est de rang  $p - r$

et telle que  $G = \begin{vmatrix} X \\ F \end{vmatrix}$  soit de rang p.

$\underline{b}$  est un estimateur sans biais de  $\underline{\beta}$  ( $E(\underline{b}) = \underline{\beta}$ ) tel que  $F\underline{b} = 0$

2/  $\underline{b}$  est fonction linéaire de  $\underline{y}$

3/ Si  $\text{cov } \varepsilon = \sigma^2 I_n$ , parmi les estimateurs  $\hat{\underline{\beta}}$  de  $\underline{\beta}$  satisfaisant à 1. et 2.,  $\underline{b}$  est tel que ses composantes ont des variances minimales. Plus généralement on a  $\text{cov } \hat{\underline{\beta}} - \text{cov } \underline{b}$  définie positive si  $\hat{\underline{\beta}} \neq \underline{b}$ . On montre que si X est de rang p :  $\text{cov } \underline{b} = \sigma^2 (X'X)^{-1}$  et dans le cas général :  $\text{cov } \underline{b} = \sigma^2 (G'G)^{-1} X'X (G'G)^{-1}$ , avec  $G'G = X'X + F'F$

b) On montre aussi que  $s^{*2} = \frac{e^2}{n-r} = \frac{\bar{y}^2 - \overline{Xb}^2}{n-r}$  est un estimateur sans biais de  $\sigma^2$  et que  $\underline{e}$  et  $\underline{b}$  sont non corrélés ( $\underline{e}$  et  $X\underline{b}$  sont les projections orthogonales de  $\underline{y}$  sur  $V^\perp$  et  $V$ ).

c) Propriétés complémentaires quand on suppose que la distribution de  $\underline{\varepsilon}$  donc de  $\underline{y}$  est normale de matrice variance  $\sigma^2 I_n$

1/  $\underline{b}$  a une distribution normale de rang  $r$  caractérisée par sa moyenne  $\underline{\beta}$  et sa matrice variance  $\text{cov } \underline{b}$  (cf a) 3.).

2/  $\underline{e}$  a une distribution normale de rang  $n-r$ , de moyenne nulle.

3/  $\underline{e}$  et  $\underline{b}$  sont indépendants en probabilité.

4/  $\underline{e}^2 \in \sigma^2 \chi^2(n-r)$  et  $\|X(\underline{b} - \underline{\beta})\|^2 = \|y^* - E(y)\|^2 \in \sigma^2 \chi_r^2$ ,

ces deux variates étant indépendantes en probabilité.

On peut aussi dire que  $\overline{Xb}^2 = y^{*2} \in \sigma^2 \chi_r^2(\overline{X\beta}^2)$  où  $\chi_r^2(\lambda)$  désigne une variable chi deux non centrée à  $r$  degrés de liberté de paramètre de non centralité  $\lambda$ .

## 6 - Test d'une sous-hypothèse linéaire sur les coefficients de régression

Ce qui suit n'est valable que sous l'hypothèse où  $\underline{\varepsilon}$  donc  $\underline{y}$  suit une distribution normale de matrice variance  $\sigma^2 I_n$ .

On se propose de tester l'existence de certaines relations linéaires entre les composantes de  $\underline{\beta}$ .

On peut toujours mettre ces relations sous la forme générale :

$$(4) \quad H(\underline{\beta} - \underline{\theta}) = 0$$

où  $H$  est une matrice  $(q, p)$  donnée de rang  $q$  et  $\underline{\theta}$  est donné.

$q$  est l'ordre de l'hypothèse qu'on désire tester.

Comme exemples les plus fréquemment rencontrés on peut citer l'hypothèse  $\underline{\beta}_2 = 0$  où  $\underline{\beta}_2$  est une composante de  $\underline{\beta} = \begin{pmatrix} \underline{\beta}_1 \\ \underline{\beta}_2 \end{pmatrix}$  dans la décomposition

de  $\underline{\beta}$ . suivant deux sous-espaces de  $\mathbb{R}^p$  ayant respectivement  $p-q$  et  $q$  dimensions. Ce test est alors appelé "test de signification" d'un ensemble de  $q$  coefficients de régression partielle.

Plus généralement on peut tester l'hypothèse  $\underline{\beta}_2 = \underline{\beta}_{20}$  où  $\underline{\beta}_{20}$  est fixé.

On rencontre également assez souvent le test de l'hypothèse d'égalité entre certaines composantes  $\beta_j, \beta_j \in \mathbb{R}^a$ , de  $\underline{\beta}$ . On en verra un exemple très important en II.

L'hypothèse (4) signifie que  $\underline{\beta}$  n'est pas quelconque dans  $\mathbb{R}^p$  mais qu'il est situé dans une variété linéaire  $B^{p-q}$  dimension  $p-q$  de  $\mathbb{R}^p$ .

Ce qu'on teste en fait c'est que  $E^x(y) = X\underline{\beta}$  appartient à la variété  $V^{p-q}$  image par  $X$  de  $B^{p-q}$  qui est bien entendu contenue dans l'image  $V$  de  $\mathbb{R}^p$ .

On est donc amené à tester l'hypothèse :

$$H_1 : (1) \underline{y} = X\underline{\beta} + \underline{\varepsilon}$$

avec (4)  $H(\underline{\beta} - \underline{\theta}) = 0$

comme "sous-hypothèse" de  $H_0$  : (1)  $\underline{y} = X\underline{\beta} + \underline{\varepsilon}$

La méthode générale de test repose sur la comparaison des sommes de carrés résiduelles,

$\bar{e}^2 = Q_{E/H_0}$  et  $\bar{e}_1^2 = Q_{E/H_1}$  attachées à l'estimation de  $\underline{\beta}$  par la méthode des moindres carrés sous les hypothèses  $H_0$  et  $H_1$  respectivement.

Quand  $\theta = 0$ ,  $V_1$  est un sous-espace de  $R^n$  et l'interprétation géométrique de la méthode montre que si  $\underline{b}_{H_1}$  est l'estimateur des moindres carrés de  $\underline{\beta}$  sous  $H_1$ ,  $\underline{y}_{H_1} = X\underline{b}_{H_1} = \text{proj } \underline{y}$  sur  $V_1$

et  $\underline{e}_1 = \text{proj } \underline{y}$  sur  $V_1^\perp$ .

Comme  $V_1 \subset V$ ,  $V_1^\perp \supset V^\perp$ . On a donc  $e = \text{proj } e_1$  sur  $V^\perp$ . (1)

C'est à partir de cette remarque qu'on montre que la statistique

$$G_{H_1/H_0} = \frac{Q_{E/H_1} - Q_{E/H_0}}{q} : \frac{Q_{E/H_0}}{n-p}$$

est distribuée comme  $F$  à  $q$  et  $n-p$  degrés de liberté si, et seulement si l'hypothèse  $H_1$  est réalisée.

On peut montrer que le test le plus puissant basé sur  $G_{H_1/H_0}$  est un test unilatéral à droite.

Si  $\alpha$  est le seuil que l'on s'est fixé, on rejettera donc  $H_1$  si

$$G_{H_1/H_0} > F_{q, n-p}(1 - \alpha), \text{ où } F(1 - \alpha) \text{ est défini par } P[F < F(1 - \alpha)] = 1 - \alpha$$

Remarque : Test de signification d'un coefficient de régression partielle  $\beta_j$

Dans le test de  $H_1 : \beta_j = 0$ ,  $q = 1$  et  $G_{H_1/H_0}$  est sous  $H_1$ , distribuée comme  $F_{1, n-p}$ .

On montre que  $G_{H_1/H_0} = (b_j / s^* b_j)^2$  où  $b_j$  et  $s_b^*$  sont respectivement l'estimateur des moindres carrés de  $\beta_j$  et l'écart type estimé de  $b_j$  (obtenu en remplaçant  $\sigma$  inconnu dans  $\sigma_{b_j}$  par  $s^*$ ) sous l'hypothèse  $H_0$ .

L'application de la méthode générale est donc équivalente à l'application du test  $t$  de Student à  $n-p$  degrés de liberté à  $b_j / s_b^*$ .

## II - COMPARAISON DES REGRESSIONS OBTENUES SUR K ECHANTILLONS (ANALYSE DE LA COVARIANCE)

On suppose que les formules de régression comportent un terme constant c'est-à-dire que la première colonne de  $X$  est formée de 1.

(1) Quand  $\theta \neq 0$  on se ramène au cas précédent (1) sous la forme

$$\underline{y} - X\underline{\theta} = X(\underline{\beta} - \underline{\theta}) + \underline{\varepsilon}$$

soit (1')  $\underline{z} = X\underline{\gamma} + \underline{\varepsilon}$  avec  $\underline{z} = \underline{y} - X\underline{\theta}$  connu  $\in R^n$

$$\text{et } \underline{\gamma} = \underline{\beta} - \underline{\theta} \in R^p$$

(4) devient alors (4') :  $H\underline{\gamma} = 0$

et on est amené à tester (4') dans le modèle (1')

On aura alors  $\underline{b}_{H_1} = \underline{g}_{H_1} + \underline{\theta}$ , si  $\underline{g}_{H_1}$  est l'estimateur des moindres carrés de  $\underline{\gamma}$ .

On est alors amené à se poser les problèmes suivants : test du parallélisme des hypersurfaces de régression :  $H_1$ , ou test de la confusion de ces hypersurfaces :  $H_2$

1/ de prime abord, ou

2/ en les supposant parallèles (c'est ce problème qui constitue "l'analyse de la covariance").

Le modèle de départ correspondant à l'hypothèse  $H_0$  est :

$$H_0 \quad \begin{cases} \underline{y}_1 = X_1 \underline{\beta}_1 + \underline{\varepsilon}_1 & \underline{\beta}_1, \underline{\beta}_2 \dots \underline{\beta}_k \in \mathbb{R}^p \\ \underline{y}_2 = X_2 \underline{\beta}_2 + \underline{\varepsilon}_2 & \underline{y}_i, \underline{\varepsilon}_i \in \mathbb{R}^{n_i} \\ \dots \\ \underline{y}_k = X_k \underline{\beta}_k + \underline{\varepsilon}_k & i = 1, 2 \dots K \end{cases}$$

Si on fait l'hypothèse que  $\varepsilon_i \in N(0, \sigma^2 I_{n_i})$  et que les  $K$  échantillons sont indépendants, ces problèmes rentrent dans le cadre de la méthode générale de test d'une sous-hypothèse linéaire sur les coefficients de régression (cf. I - 6), mais il ne semble pas qu'ils aient fait l'objet d'un programme sauf en ce qui concerne le test de confusion des hypersurfaces supposées parallèles ( $H_2$  comme sous-hypothèse de  $H_1$ ).

En effet on peut représenter l'ensemble des observations par le modèle :

$$M : \underline{y} = X\underline{\beta} + \underline{\varepsilon} \quad \underline{\varepsilon} \in N(0, \sigma^2 I_N)$$

avec

$$y = \begin{pmatrix} \underline{y}_1 \\ \vdots \\ \underline{y}_k \end{pmatrix} \quad \underline{\varepsilon} = \begin{pmatrix} \underline{\varepsilon}_1 \\ \vdots \\ \underline{\varepsilon}_k \end{pmatrix} \quad \begin{cases} y \in \mathbb{R}^N \\ \varepsilon \in \mathbb{R}^N \\ N = \sum_i n_i \end{cases} \quad \underline{\beta} = \begin{pmatrix} \underline{\beta}_1 \\ \vdots \\ \underline{\beta}_k \end{pmatrix} \quad \beta \in \mathbb{R}^{pk}$$

$$(N, pk) \quad X = \begin{pmatrix} X_1 & \dots & 0 \\ & X_2 & \\ \vdots & & \\ 0 & & X_k \end{pmatrix}$$

et le test de confusion des hypersurfaces de régression est le test de l'hypothèse  $H : \underline{\beta}_1 = \underline{\beta}_2 = \dots = \underline{\beta}_k$  tandis que le test du parallélisme est celui de l'hypothèse  $H_2 : \beta_{1j} = \beta_{2j} = \dots = \beta_{kj} \quad \forall j \neq 1$

### III - CHOIX DES VARIABLES EXPLICATIVES OU INDEPENDANTES A FAIRE INTERVENIR DANS UNE EQUATION DE REGRESSION

#### 1 - Evaluation du biais introduit dans les estimateurs des $q$ coefficients de régression restant quand on supprime $p - q$ variables indépendantes dans un modèle de régression linéaire

Considérons le modèle de régression

$$(1) \quad E^{x'}(y) = x'_1 \beta_1 + x'_2 \beta_2 = \pi$$



où  $\underline{\beta}_1$  est de dimension  $q$  et  $\underline{\beta}_2$  de dimension  $p - q$

et le modèle simplifié (2)  $E^{x'}(y) = x'_1 \beta_1 = \pi_1$

On trouve que si  $\underline{b}_1^0$  est l'estimateur des moindres carrés de  $\underline{\beta}_1$  dans le modèle (2) on a  $E(\underline{b}_1^0) = \underline{\beta}_1 + (X'_1 X_1)^{-1} X'_1 X_2 \underline{\beta}_2$  où  $X_1$  et  $X_2$  sont respectivement les matrices  $n, q$  et  $n, p - q$  dont les lignes sont les valeurs de  $\underline{x}'_1$  et  $\underline{x}'_2$  pour le  $n$  échantillon considéré ( $X = \begin{pmatrix} |X_1| & |X_2| \end{pmatrix}$  et  $\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$ )

Le biais commis en estimant  $\underline{\beta}_1$  par  $\underline{b}_1^0$  est donc

$$\underline{B} = (X'_1 X_1)^{-1} X'_1 X_2 \underline{\beta}_2$$

Il est nul si et seulement si :  $\underline{\beta}_2 = 0$  c'est-à-dire si le modèle (2) convient

$$\text{ou } X'_1 X_2 = 0$$

c'est-à-dire si les colonnes de  $X_1$  sont orthogonales aux colonnes de  $X_2$ . (on dit si les variables composant  $X_1$  sont orthogonales aux variables composant  $X_2$ )

Ce biais dépend donc du plan de relevé des données par l'intermédiaire de  $X_1$  et  $X_2$

On a

$$X_1 \beta = \underbrace{X_1 (X'_1 X_1)^{-1} X'_1 X_2 \beta_2}_{\text{proj}(X_2 \beta_2) \text{ sur } V_1}$$

$V_1 =$  sous espace de  $R^n$  engendré par les colonnes de  $X_1$

Ce biais est donc l'image réciproque par  $f_1$  associée à  $X_1$  de la projection orthogonale sur  $V_1$  sous espace engendré par  $X_1$ , de  $X_2 \beta_2$  composante de  $E^{x'}(y)$  négligée dans le modèle (2).

On montre aussi que  $\underline{b}_1^0 = \underline{b}_1 + (X'_1 X_1)^{-1} X'_1 X_2 \underline{b}_2$

si  $b_1$  et  $b_2$  sont les estimateurs des moindres carrés de  $\beta_1$  et  $\beta_2$  dans le modèle (1).

## 2 - Conclusions de l'étude

1/ La suppression à tort d'une ou plusieurs variables indépendantes conduit à des estimateurs biaisés pour les coefficients de régression des autres variables si celles-ci ne sont pas orthogonales aux variables supprimées dans le  $n$  échantillon considéré (ou non corrélées avec elles s'il s'agit de variables centrées).

2/ L'étude de l'influence sur  $E(y)$  d'une variable ou d'un groupe de variables considéré isolément exige que cette variable ou ce groupe de variables soit orthogonal à chacune des autres dans l'échantillon considéré. S'il n'en n'est pas ainsi il se peut que l'influence de ces variables puisse être négligée par rapport aux influences des autres quand on se limite à l'ensemble des valeurs prises par les variables indépendantes dans l'échantillon, mais rien ne permet d'affirmer qu'il en sera ainsi pour d'autres valeurs, ou même d'autres combinaisons de valeurs des variables indépendantes. C'est ce qui interdit la suppression de variables indépendantes non orthogonales aux autres

pour la prévision par extrapolation, si on ne veut pas risquer d'introduire un biais dont on ne peut pas connaître la valeur.

De même, dans la prévision par extrapolation de l'influence de certaines variables indépendantes sur  $E(y)$  on ne peut envisager que la suppression éventuelle d'un groupe de variables qui seraient orthogonales aux autres dans l'échantillon si on veut prévoir sans biais l'influence des variables restant.

Aussi semblerait-il intéressant de procéder à une analyse en composantes principales sur l'ensemble des variables indépendantes (ou explicatives) pour rechercher des groupes de variables non corrélées (ou plus généralement orthogonales d'un groupe à l'autre si on ne considère pas des variables centrées), ou substituer aux variables initiales de nouvelles variables deux à deux non corrélées ou plus généralement orthogonales. C'est le problème de la "régression orthogonalisée" que nous n'aborderons pas ici.

### 3 - Les principales procédures de choix des variables indépendantes à faire intervenir dans une équation de régression utilisée pour la prévision par interpolation

On supposera pour simplifier le langage que la formule de régression recherchée comporte un terme constant, ce qui revient à considérer une première variable indépendante  $x^0$  constante, et  $p - 1$  variables indépendantes non constantes au sein desquelles on se propose de choisir les variables à conserver. Il importe de ne pas perdre de vue qu'il n'y a pas de procédure statistique unique pour procéder à ce choix et que le jugement personnel devra intervenir chaque fois.

#### a) Essai de toutes les régressions possibles

Ceci conduit à essayer  $2^{p-1}$  régressions et n'est donc praticable que pour  $p - 1 < 10$  ;

On est amené à choisir "une meilleure régression" dans chaque groupe de régressions à  $r$  variables indépendantes et ceci pour toutes les valeurs possibles de  $r$  soit  $r = 1, 2, \dots, p - 1$  et à comparer ensuite les meilleures régressions obtenues pour les diverses valeurs de  $r$  afin de décider quelle est la plus petite valeur de  $r$  que l'on retiendra.

Pour procéder à ce choix on utilise deux types de critères :

1/ Le carré du coefficient de corrélation multiple  $R^2$ , ou la variance résiduelle empirique  $s^2 = e^2/n$

$R^2$  est la réduction relative de variance introduite par le remplacement de la variance marginale empirique de  $y$  par sa variance résiduelle empirique dans la régression considérée.

$$R^2 = \frac{\sum_1 (y_1 - \bar{y})^2 - \sum_1 (y_1 - y_1^*)^2}{\sum_1 (y_1 - \bar{y})^2}$$

puisque

$$e^2 = \sum_1 (y_1 - y_1^*)^2$$

On compare donc la variance résiduelle de  $y$  attachée à une régression à sa variance en l'absence de toute régression. Comme la variance marginale ne dépend que de l'échantillon, le critère de  $R^2$  maximum équivaut au critère de variance résiduelle empirique minimum.

On notera que quand  $r$  augmente par adjonction de nouvelles variables explicatives  $e^2$  ne peut que diminuer donc  $R^2$  ne peut qu'augmenter ; donc la meilleure régression à  $r + 1$  termes sera toujours au moins aussi bonne que la meilleure régression à  $r$  termes.

Pour choisir la valeur de  $r$  à adopter on peut porter les valeurs de  $R^2$  Max ou de  $e^2$  Min/ $n$  en fonction de  $r$ . On retiendra la plus petite valeur de  $r$  pour laquelle  $R^2$  Max ou  $e^2$  Min/ $n$  se stabilise.

## 2/ Le critère $C_K$ de Mallows

$K$  est le nombre total des termes introduits dans la régression, donc  $K = r + 1$  s'il y a  $r$  variables explicatives non constantes. Par définition  $C_K = K +$  somme des carrés des erreurs systématiques de prévision commises en chacun des  $n$  points de l'échantillon, ces erreurs étant comptées en écarts types liés.

$K$  est la valeur moyenne de la somme des carrés des erreurs aléatoires de prévision aux mêmes points, comptée elle aussi en écarts types liés.

On montre que  $C = \bar{e}_K^2 / \sigma^2 - (n - 2K)$  en désignant par  $\bar{e}_K^2$  la somme des carrés résiduelle attachée à la régression à  $K$  termes considérée.

Pratiquement on ne connaît pas  $\sigma^2$  et on l'estime par  $s^{*2} = \frac{\bar{e}_p^2}{n - p}$  estimateur de  $\sigma^2$  dans le modèle complet de régression à  $p$  termes, en sorte qu'on obtient  $C_p = p$

On compare donc en fait la somme des carrés résiduelle des modèles à  $K$  termes à la somme des carrés résiduelle du modèle "complet" à  $p$  termes.

En l'absence d'erreur systématique  $C_K = K$  ; par suite les valeurs de  $C_K$  doivent être comparées à  $K$  et on retiendra le plus petit nombre  $K$  tel que  $C_K$  Min soit voisin de  $K$ .

### Remarques importantes

- Le meilleur groupe de  $r + 1$  variables ne contient pas nécessairement le meilleur groupe de  $r$  variables, en sorte que même sous l'hypothèse de normalité de la distribution de  $\underline{e}$  il n'est pas toujours possible de tester la signification de l'influence de la  $(r + 1)^{\text{ème}}$  "meilleure" variable.

- Par contre toujours sous l'hypothèse de normalité pour la distribution de  $\underline{e}$  on peut tester la signification de l'ensemble des  $p - (r + 1)$  variables non introduites c'est-à-dire tester l'hypothèse  $\beta_2 = 0$  dans le modèle (1) de (I,6) où  $x'$  est l'ensemble des  $p - (r + 1)$  variables non introduites. Quand on a calculé les coefficients  $C_K$ , la méthode générale indiquée en (I,6) conduit à appliquer un test  $F$  à  $p - K$  et  $n - p$  degrés de liberté, unilatéral

$$\text{à } 1 + \frac{C_K - K}{p - K}$$

## b) Elimination progressive (Backward Elimination)

Cette méthode consiste à :

1/ Calculer une formule de régression comportant toutes les variables.

2/ Calculer la valeur de toutes les statistiques

$$\frac{b_j}{s_{b_j}^*}, j = 1, 2, \dots, p - 1$$

qui permettent de tester la signification des régressions partielles  $\beta_j$  au moyen du test t de Student (cf I, 6 : Remarque)

3/ Tester la signification de la régression partielle la moins significative (celle qui correspond à la plus petite valeur de  $b_j/s_{b_j}^*$ ) par exemple de  $\beta_L$ .

4/ - Si cette régression partielle est significative, conserver la formule trouvée.

- Si elle n'est pas significative, supprimer la variable  $x_L$  et recommencer le processus.

Cette méthode est particulièrement recommandée à ceux qui désirent avoir la formule de régression sur l'ensemble des variables indépendantes. Elle donne par ailleurs la meilleure régression à  $p - 1$  termes.

Si on ne tient pas spécialement à garder ou du moins à essayer l'ensemble des  $p - 1$  variables indépendantes, la régression pas à pas paraît préférable.

### Remarques

1/ Le critère de choix de la variable à éliminer à chaque étape est en fait le carré du coefficient de corrélation partielle de  $y$  avec cette variable.

Il garde un sens concret même si les hypothèses de validité du test appliqué ne sont pas réalisées.

2/ Il faut se garder d'éliminer simultanément plusieurs variables qui auraient des régressions partielles non significatives ; c'est en fait souvent parce que ces variables sont fortement corrélées qu'aucune d'elles n'a d'influence significative.

3/ On peut envisager, une fois la formule de régression finale obtenue, de tester la signification de l'ensemble des variables éliminées.

## c) Régression pas à pas ou "Stepwise regression"

(cf B.M.D. 02 - R)

C'est une méthode d'adjonction progressive des variables.

Les étapes de la méthode sont les suivantes :

1/ Choisir la variable  $x_j$  la plus fortement corrélée avec  $y$ , soit  $x_1$  cette variable.

Tester la signification du coefficient de régression totale  $b_1$  à un seuil  $\alpha_1$  fixé.

S'il n'est pas significatif, conclure à l'absence de régression significative de  $y$  sur les  $x_j$ .

Si l'influence de  $x_1$  est significative passer à

2/ Introduire successivement chacune des variables  $x_j$ ,  $j \neq 1$ , en plus de  $x_1$ .

Retenir la variable qui donne un coefficient de corrélation multiple  $R^2_{y/x_1, x_j}$  le plus grand possible c'est-à-dire une somme des carrés résiduelle  $\bar{e}^2$  la plus petite possible.

Soit  $x_2$  cette variable. Tester la signification du coefficient de régression partielle  $\beta_2$  dans la régression de  $y$  sur  $x_1$  et  $x_2$  au risque  $\alpha_1$ .

Si l'influence de  $x_2$ , après  $x_1$  n'est pas significative, garder la seule variable explicative  $x_1$ .

Si elle est significative passer à

3/ Tester la signification de  $x_1$ , après  $x_2$  au risque  $\alpha_2$  ( $\alpha_2 \geq \alpha_1$ ).

a) Si  $x_1$  n'est pas significatif après  $x_2$ , garder  $x_2$  seule et reprendre à l'étape 2. Le fait que  $\alpha_2 \geq \alpha_1$  implique que  $x_1$  ne peut avoir après  $x_2$  une influence significative au risque  $\alpha_1$  et ne peut par suite pas être réintroduite après  $x_2$

b) Si  $x_1$  est significatif après  $x_2$ , garder  $x_1$  et  $x_2$  et passer à

4/ Introduire successivement chacune des variables  $x_j$ ,  $j \geq 3$  en plus de  $x_1$  et  $x_2$  et retenir la variable donnant la plus forte corrélation multiple.

Soit  $x_3$  cette variable. Si son influence après  $x_1$  et  $x_2$  n'est pas significative au risque  $\alpha_1$ , garder la formule avec  $x_1$  et  $x_2$ . Si l'influence de  $x_3$  après  $x_1$  et  $x_2$  est significative, tester au risque  $\alpha_2$  la signification de celle des variables  $x_1$  et  $x_2$  qui est la moins significative ( $b_j/s_{b_j}^*$  le plus petit).

a) Si les deux variables  $x_1$  et  $x_2$  ont des influences significatives, poursuivre le processus en essayant d'introduire une quatrième variable.

b) Si  $x_1$  par exemple n'est pas significative, la supprimer et passer à l'étape 4.

On arrête le processus quand il n'y a plus aucune variable à introduire ou à supprimer.

La méthode est aisée à mettre en oeuvre si on fait apparaître les  $s_{b_j}^*$  à chaque étape, c'est-à-dire si on calcule les éléments diagonaux de  $(X'X)^{-1}$  puisque le test de signification d'une régression partielle  $\beta_j$  est le test  $t$  appliqué à  $b_j/s_{b_j}$ . Le plus souvent on ne change pas de valeur critique pour  $t$  quand le nombre  $K$  des variables augmente. C'est sans importance quand  $n$  est grand devant  $K$ .

### Remarques

1/ Cette méthode conduit à renoncer à toute régression dans le cas où un ensemble de variables serait collectivement significatif sans qu'aucune d'elle considérée isolément ne le soit.

On peut se prémunir contre ce risque en calculant la régression sur l'ensemble des  $p - 1$  variables candidates et en testant la signification de la ré-

