Programme d'estimation du point de transition et des degrés des polynômes d'approximation d'une fonction de régression

S.R. Esterby

**ÉTUDE Nº 147, SÉRIE SCIENTIFIQUE** 

DIRECTION GÉNÉRALE DES EAUX INTÉRIEURES INSTITUT NATIONAL DE RECHERCHES SUR LES EAUX CENTRE CANADIEN DES EAUX INTÉRIEURES **BURLINGTON (ONTARIO) 1985** 

(Available in English on request)

© Ministre des Approvisionnements et Services Canada 1986 N° de cat. En36-502/147F ISBN 0-662-93696-5

# Table des matières

		Page
RÉSUMÉ	<u>.</u>	٧
ABSTRA	CT	v
INTROD	UCTION	1
ALGÒRI	THME	2
Exem Exem Exem	ES	5 5 9 11 12
ANALYS	SE	13
REMERO	CIEMENTS	13
RÉFÉRE	NCES	13
ANNEXE	I	15
	de l'utilisateur pour chaque passage du programme	2
2. Différe	ences selon l'hypothèse sur les variances	4 5
Illusi	trations	
	Organigramme général du déroulement du programme	3
Figure 3.	pour l'analyse de l'échantillon KB, exemple 1	6
Figure 4.	KB, exemple 1	8
	et la fonction de vraisemblance relative, exemple 2	10 11 12

### Résumé

On décrit un programme informatique en Fortran qui utilise une méthode de calcul pour un modèle de régression polynomiale avec un point de transition. Les hypothèses posées pour le modèle sont les mêmes que pour un seul polynôme, mais on suppose que le polynôme, et peut-être la variance, changent en un point de la série de données. L'ordre des données simples doit aussi être connu. Le degré du polynôme dans chaque segment et le point de transition sont calculés par une méthode itérative qui comporte l'emploi de la fonction de vraisemblance marginale pour le paramètre du point de transition et un test faisant appel au changement des sommes des carrés des écarts résiduels avec une diminution du degré. Les hypothèses du modèle peuvent être examinées à l'aide des valeurs des écarts résiduels et de la fonction de vraisemblance. On donne des exemples qui démontrent l'utilité du programme et indiquent les résultats obtenus avec l'imprimante et le traceur.

## **Abstract**

A Fortran computer program, which implements an estimation procedure for a polynomial regression model with a change point, is described. The assumptions of the model are analogous to those for a single polynomial but include the assumption that the polynomial, and possibly the variance, have changed at a point in the data set. The order of the data points must also be known. The degree of the polynomial in each segment and the point of change are estimated by an iterative procedure which involves the use of the marginal likelihood function for the change-point parameter and a test using the change in residual sums of squares with reduction of degree. The assumptions of the model can be examined using the residuals and the likelihood function. Examples are given to demonstrate the utility of the program and to show the output from both the printer and the plotter.

# Programme d'estimation du point de transition et des degrés des polynômes d'approximation d'une fonction de régression

S.R. Esterby

#### INTRODUCTION

On emploie souvent les méthodes de l'analyse de régression pour obtenir des renseignements sur les relations qui lient diverses variables. Dans la plupart des applications, une relation de régression unique est ajustée à un certain ensemble de données: il arrive toutefois que cette relation change de forme au cours de la période de collecte des données. La complexité de l'analyse dépend de la connaissance qu'on a de la forme de cette relation ainsi que du point de transition. Lorsque le point de la transition d'une relation de régression à l'autre est connu, l'analyse consiste à approximer chacun des segments par la méthode des moindres carrés, en imposant éventuellement aux fonctions de régression la contrainte supplémentaire qu'elles doivent avoir la même valeur au point de transition. Lorsque ce dernier est inconnu, l'analyse se complique car il faudra alors estimer le point de transition et aussi les paramètres de régression. Nous nous intéresserons, dans le présent rapport, à cette dernière situation et décrirons un algorithme fondé sur une méthode proposée par Esterby et El-Shaarawi (1981a).

Cet algorithme a été établi pour le cas d'une seule variable dépendante. On suppose que chaque segment de l'ensemble de données considéré peut être approximé par une fonction de régression unique (en appliquant la méthode des moindres carrés) et qu'en outre la relation de régression s'est modifiée à un certain point qu'il reste à déterminer. Nous supposerons que la fonction de régression est un polynôme et que l'ensemble des observations est un ensemble ordonné (par exemple par un ordre naturel dans le temps ou dans l'espace). Soit (x<sub>i</sub>, y<sub>i</sub>) un ensemble de paires d'observations indépendantes, où l'indice i désigne l'ordre de l'observation. Le modèle consiste alors en une paire de polynômes de degrés respectifs p et q, la transition du polynôme de degré p au polynôme de degré q s'effectuant après l'observation pour laquelle on a i=n1. On a donc

$$y_i = \sum_{j=0}^{p} \theta_{1j} x_i^j + e_{1i}$$
  $i = 1, 2, ..., n_1$ 

et

$$y_i = \sum_{j=0}^{q} \theta_{2j} x_i^j + e_{2i}$$
  $i = n_1+1, n_1+2...n$ 

où  $e_{1i}$  (i=1,2...n<sub>1</sub>) et  $e_{2i}$  (i=n<sub>1</sub>+1, n<sub>1</sub>+2...n) sont deux suites de variables aléatoires indépendantes que l'on suppose respectivement N(0,  $\sigma_1^2$ ) et N(0,  $\sigma_2^2$ ). Les valeurs des paramètres  $\underline{\theta}_1 = (\theta_{10}, \theta_{11} \dots \theta_{1p}), \underline{\theta}_2 = (\theta_{20}, \theta_{21} \dots \theta_{2q}),$  $\sigma_1^2$ ,  $\sigma_2^2$  et  $n_1$  sont inconnues. Dans les deux sections qui suivent, ce modèle est désigné par la notation  $M_{II}(p,q|n_1)$ . Puisque p et q sont donnés, on cherche à estimer la valeur du point de transition n<sub>1</sub> à partir de la fonction de vraisemblance marginale de n<sub>1</sub>. Dans le cas où p et q sont également inconnus, on obtient une estimation de ces deux paramètres par itération. La méthode consiste à exploiter conjointement la fonction de vraisemblance marginale de n<sub>1</sub> et un test analogue au test de réduction du degré d'un polynôme unique. On suppose que le degré de chaque polynôme est borné supérieurement. Soit (p\*, q\*) la paire des bornes supérieures des degrés des deux polynômes; à chaque étape du processus d'itération, le degré total (p + q) est réduit d'une unité, jusqu'à ce qu'aucune réduction ne soit plus possible ou que (p,q) = (0,0). Le cas où  $\sigma_1^2 = \sigma_2^2$  est également envisagé; le modèle correspondant est désigné par la notation Me(p,qln1), où les indices e et u servent à distinguer les modèles correspondant respectivement à l'inégalité des variances. Il existe en effet des différences marquées entre les fonctions de vraisemblance marginale et les méthodes d'estimation de p et q selon que les variances sont égales ou non. On trouvera une description détaillée de ces méthodes dans Esterby et El-Shaarawi (1981a) ainsi qu'une description plus détaillée des fonctions de vraisemblance dans Esterby et El-Shaarawi (1981b).

Les méthodes de l'analyse des écarts résiduels, qui sont applicables au cas où l'ensemble de données considéré peut être approximé par une fonction unique, peuvent être transposées ici : il suffit en effet d'analyser l'ensemble des

écarts résiduels correspondant à chaque segment. En outre, l'étude de la fonction de vraisemblance permet de confirmer ou d'infirmer l'hypothèse de l'existence d'un point de transition unique.

Le but principal de la méthode décrite ici consiste à estimer le point où la relation de régression change de forme dans l'hypothèse de l'existence d'un point de transition. Aucune contrainte n'est imposée aux valeurs des fonctions de régression en ce point. Selon l'objectif visé, il est possible que des méthodes différentes doivent être employées. Si l'on désire par exemple que la courbe qui approxime l'ensemble de données considéré soit une courbe lisse, on peut dans de nombreux cas approximer les données au moyen d'une spline cubique ou d'un polynôme de degré supérieur.

On trouvera une description de l'algorithme à la section suivante ainsi que des exemples d'utilisation du programme dans la troisième section. La séquence en Fortran n'a pas été listée parce que le programme est très long. On peut se procurer une copie de ce dernier ainsi que des jeux d'essai en écrivant à l'auteur. Le programme, rédigé en Fortran V, comporte un certain nombre de sousprogrammes IMSL et est exécuté sur un CDC Cyber 171 tournant sous NOS. Les programmes de traçage ont été rédigés pour un traceur CALCOMP 1036. Nous décrirons ici les principales caractéristiques du programme sans pour cela reprendre tout le développement de la méthode. Nous illustrerons son utilisation par des exemples, de façon que le lecteur qui a besoin d'un programme de cette nature puisse se faire une idée de l'utilité de celui que nous avons mis au point pour l'application particulière qu'il a en vue.

#### **ALGORITHME**

De nombreuses fonctions ont été intégrées au programme pour permettre à l'utilisateur d'analyser simplement et efficacement une ou plusieurs paires de variables en un seul passage, pourvu que toutes ces variables soient stockées dans le même fichier. On trouvera dans le tableau 1 la liste des options que l'utilisateur doit préciser. Le programme est en mesure d'estimer la valeur du point de transition n<sub>1</sub> pourvu qu'on lui fournisse les degrés des polynômes ou d'estimer les degrés des polynômes ainsi que la valeur du point de transition pourvu qu'on lui indique la borne supérieure du degré du polynôme dans chaque segment. Pour chaque paire de variables, l'utilisateur peut préciser soit la valeur du couple de degrés (p,q) soit cinq valeurs pour les couples de bornes supérieures (p\*,q\*). Pour la paire (p,q) ou pour chaque paire (p\*,q\*), l'utilisateur doit préciser si les variances sont supposées égales ou non et si des graphiques doivent être tracés ou non; il doit également préciser les deux niveaux de signification critique correspondant à chaque test d'hypothèse. Les graphiques tracés varient selon que le domaine de la variable indépendante est ordonné de facon monotone (croissante ou décroissante) ou noñ.

On trouvera représentées à la figure 1 les principales caractéristiques du programme; cette figure montre également l'effet des options choisies par l'utilisateur sur les parties du programme qui sont utilisées. Pour comprendre les détails des calculs qui sont effectués à chaque étape de traitement de l'organigramme, il faut lire l'article de Esterby et El-Shaarawi (1981a). Nous allons toutefois faire ici un certain nombre d'observations qui aideront à clarifier la

Tableau 1. Choix de l'utilisateur pour chaque passage du programme

Niveau d'application	Choix				
	Le nombre de paires de variables				
Pour chaque paire de variables	La forme selon laquelle les valeurs des variables x et y sont lues				
	Le titre				
	Si oui ou non l'ensemble des valeurs de la variable indépendante est ordonné				
	Soit 1) la valeur de (p, q) și elle est connue ou				
	<ol> <li>le nombre de paires de bornes supérieures npq des degrés des polynômes et leurs valeurs (p*, q*)</li> </ol>				
	Si tous les échantillons doivent être utilisés ou si certains (et lesquels) doivent				
	être omis				
	Les unités adoptées pour les axes de représentation graphique				
Pour la paire (p,q) ou pour chaque paire (p*,q*)	Si les variances sont égales, inégales ou s'il faut effectuer un test pour en déterminer l'égalité ou l'inégalité				
	Si des graphiques doivent être produits				
	Le niveau de signification critique pour le test d'égalité des variances				
	Le niveau de signification critique pour le test de réduction du degré				

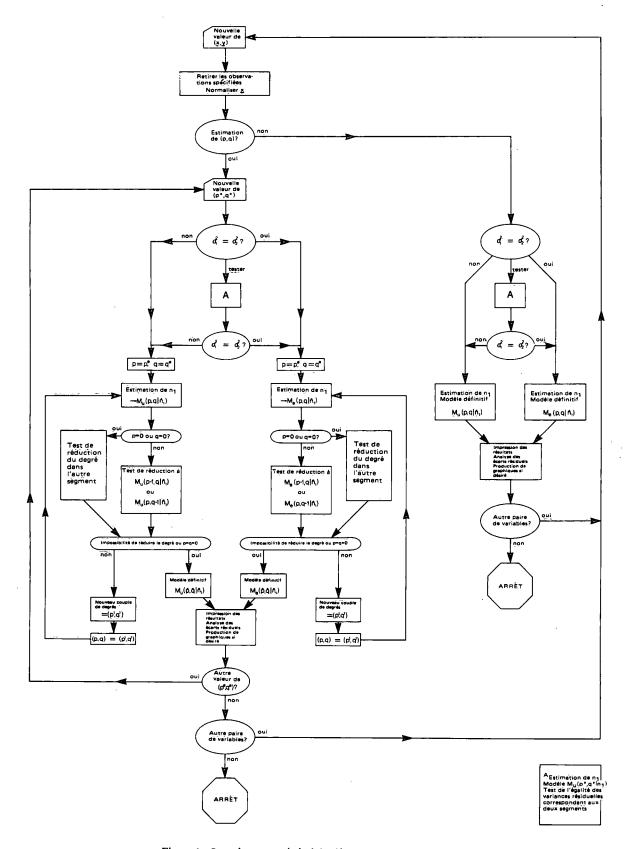


Figure 1. Organigramme général du déroulement du programme.

structure logique du programme. La méthode d'itération qui permet d'estimer p et q commence par fixer à ces deux paramètres des valeurs particulières puis, à partir de ces valeurs, attribue une valeur n1 à n1. En supposant ensuite que n<sub>1</sub> prend cette valeur n<sub>1</sub>, on effectue un test pour déterminer si le degré total (p+q) ne pourrait pas être réduit d'une unité. Dans l'affirmative, on obtient à partir de (p+q-1) une nouvelle estimation de n<sub>1</sub>. Les deux suites de la figure 1, qui correspondent respectivement au cas où les variances sont égales et inégales, sont expliquées avec plus de détails dans le tableau 2, où les différences attribuables à l'égalité ou à l'inégalité des variances sont précisées. La valeur estimée no de no est la valeur de ce paramètre qui maximise la fonction de vraisemblance marginale; on détermine cette valeur en calculant les valeurs de la fonction de vraisemblance marginale pour les valeurs de n1 qui correspondent aux valeurs de l'indice qui apparaissent dans le tableau 2.

Plusieurs des étapes de la figure 1 n'apparaissent pas dans la méthode décrite dans l'article de 1981a par Esterby et El-Shaarawi. Pour déterminer si les variances sont égales ou non (voir A à la figure 1), on commence par faire l'hypothèse que les variances en question sont inégales et on obtient ensuite une estimation du point de transition à partir de la fonction de vraisemblance marginale. On applique ensuite un test F standard pour vérifier l'hypothèse de l'égalité des variances résiduelles correspondant aux segments 1 et 2. Lors de l'estimation des degrés p et q, lorsque l'un de ces degrés s'annule dans le cadre du processus d'itération, il suffit d'appliquer la méthode de réduction à l'autre segment. On se trouve alors dans le cas d'un test de réduction de degré correspondant à l'ajustement d'un

polynôme unique à un ensemble de données. Dans ce cas, on peut utiliser directement la valeur  $\alpha$  du niveau de signification critique, alors que dans le cas où le degré peut être réduit d'une unité de deux façons, on doit employer  $\sqrt{\alpha}$ .

La mise en oeuvre de l'algorithme que nous avons mis au point comporte également plusieurs étapes additionnelles. Pour permettre l'inversion de certaines matrices, on commence par normaliser la variable indépendante. Il faut ensuite effectuer des calculs additionnels pour exprimer les paramètres de régression en fonction de la variable indépendante d'origine, connaissant la relation qui lie ces paramètres et la variable normalisée. On trouvera dans l'annexe les formules utilisées à cette fin. Dans le cas où les variances sont égales et où ni p ni q n'est nul, le dénominateur de la fonction test de réduction du degré comporte la somme des sommes des carrés des écarts résiduels correspondant aux deux segments. Les bornes inférieure et supérieure de n<sub>1</sub> peuvent donc être choisies respectivement égales à p+1 et n-(q+1), de sorte que le dégré de liberté de la somme des carrés des écarts résiduels correspondant à un segment est nul. Toutefois, lorsque cela se produit en même temps que le degré du polynôme qui correspond à l'autre segment est nul, on obtient pour la fonction test un dénominateur à 0 degré de liberté, ce qui est impossible. La solution que nous avons adoptée pour résoudre ce problème consiste à modifier la limite correspondante pour n<sub>1</sub> de façon que le nombre d'observations soit toujours supérieur d'une unité au nombre de paramètres à estimer et à conserver cette relation valide pour la totalité du processus d'itération. Dans le cas où les variances sont égales et où le nombre de paramètres est égal au nombre d'observations dans un segment, on ne calcule pas la somme des carrés des écarts

Tableau 2. Différences selon l'hypothèse sur les variances

Caractéristiques	Variances égales	Variances inégales		
Valeurs admissibles de n <sub>1</sub>	p+1n-(q+1)	p+2n-(q+2)		
Fonction de vraisemblance marginale de n <sub>1</sub> , L(n <sub>1</sub> )	$\{_{\nu_1\hat{\sigma}_1}{}^2{}_{+\nu_2\hat{\sigma}_2}{}^2{}_2\}^{-1/2(\nu_1+\nu_2-1)}$	$\frac{\Gamma(1/2\nu_1)  \Gamma(1/2\nu_2)}{\nu_1^{1/2\nu_1}  \nu_2^{1/2\nu_2}} \ \hat{\sigma}_1^{-(\nu_1-1)}  \hat{\sigma}_2^{-(\nu_2-1)}$		
Test de réduction du degré, p et q non nuls	On regroupe les sommes des carrés des écarts résiduels correspondant aux deux segments pour estimer la variance globale.	On calcule séparément les sommes des carrés des écarts résiduels correspondant à chaque segment pour estimer les variances.		
Écarts résiduels normalisés	On divise les écarts résiduels correspondant aux deux segments par $\hat{\sigma}$ dont la valeur a été déterminée à partir du regroupement des sommes des carrés des écarts résiduels	On normalise tout d'abord les écarts résiduels en les divisant par $\hat{\sigma}_1$ ou $\hat{\sigma}_2$ puis on les regroupe.		

Nota: 
$$\nu_1 = n_1 \cdot (p+1)$$
,  $\nu_2 = (n \cdot n_1) \cdot (q+1)$ ,  $\hat{\sigma}_1^2 = \sum_{i=1}^{n_1} (y_i - \sum_{j=0}^{p} \hat{\theta}_{1j} x_i^j)^2 / \nu_1$ ,  $\hat{\sigma}_2^2 = \sum_{i=n_1+1}^{n} (y_i - \sum_{j=0}^{q} \hat{\theta}_{2j} x_i^j)^2 / \nu_2$  et  $\Gamma(z)$  désigne la fonction gamma de z.

résiduels: on lui attribue plutôt la valeur théorique 0. La somme des carrés des écarts résiduels prend également cette valeur lorsque la variable dépendante est constante dans un segment. Dans le cas où les variances ne sont pas égales, l'attribution de la valeur zéro à la somme des carrés des écarts résiduels ne constitue pas une solution satisfaisante, puisque le logarithme de la somme des carrés des écarts résiduels correspondant à chaque segment intervient dans la fonction de vraisemblance marginale. Dans le cas où la variable dépendante prend une valeur constante dans un segment, on attribue arbitrairement à la somme des carrés des écarts résiduels correspondant à ce segment la valeur 0.000 000 1. Le cas où la variable dépendante a une valeur constante est donc traité de la même facon, que les variances soient égales ou non. L'utilisateur peut contourner ce choix arbitraire en indiquant tout simplement au programme d'omettre la valeur finale de la constante.

Les résultats ci-après, notamment, sont imprimés : un état récapitulatif des étapes de l'estimation de p et q, et, dans le cas du modèle définitif, la fonction de vraisemblance marginale relative de  $n_1$ , les paramètres de régression,  $\frac{\hat{\theta}}{1}$  et  $\frac{\hat{\theta}}{2}$ , la valeur estimée de  $\underline{y}$ ,  $\underline{\hat{y}}$ , ainsi que les écarts résiduels  $\underline{\hat{e}} = \underline{y} - \underline{\hat{y}}$ . La somme des carrés des écarts résiduels, le carré moyen des écarts, la somme quadratique totale corrigée ainsi que  $R^2$  (cette valeur est définie en annexe) sont également imprimés. L'analyse des écarts résiduels (Draper et Smith, 1981) s'effectue essentiellement par une méthode graphique : on trace un graphique Q-Q pour

vérifier l'hypothèse de normalité ainsi que des graphiques de représentation des écarts résiduels normalisés en fonction de  $\underline{x}$  et  $\hat{\underline{y}}$ . On applique également un test de suites aux écarts résiduels.

#### **EXEMPLES**

Des exemples visant à évaluer l'efficacité de la méthode décrite ici ont déjà été présentés dans des publications antérieures (Esterby et El-Shaarawi, 1980, 1981a,b). Les exemples que nous présentons ici visent surtout à illustrer le rôle du programme sur ordinateur ainsi que ses possibilités. L'exemple 1 est traité en détail et l'état ainsi que les graphiques produits par l'ordinateur sont présentés ci-dessous. Dans les trois autres exemples, nous n'avons présenté que les graphiques qui illustraient certaines caractéristiques particulières.

#### Exemple 1

Dans certaines conditions, la distribution verticale du pollen de type Ambrosia dans une carotte sédimentaire peut servir de moyen de datation. Les hypothèses de l'existence d'un point de transition unique et de la représentabilité de la relation de régression qui lie la concentration et la profondeur par une fonction polynômiale sont plausibles. L'ensemble d'échantillons utilisé ici est décrit par Robbins, Edgington et Kemp (1978); c'est un ensemble

Tableau 3. Choix pour les exemples présentés

	Exemple							
Option	1a	1b	2	3	4a	<b>4</b> b	4c	
Domaine de x ordonné	oui	oui	noń	oui	oui	oui	oui	
(p,q) ou		(1,0)	(1,1)	(0,0)			(0,1)	
(p*,q*)	(1,1)				(3,3)	(1,1)		
Échantillons utilisés	tous	tous	tous	tous	tous	1-44	17-49	
Variances	tester	tester	tester	égales	tester	tester	tester	
α, variances	0.05	0.05	0.05	0.05	0.05	0.05	0.05	
α, degré	0.01	0.01	0.01	0.01	0.01	0.01	0.01	
Graphiques	oui	oüi	oui	oui	non	oui	oui	
x borne inférieure	0.0	0.0	0.0	0.0		0.0	0.0	
borne supérieure	68.4	68.4	18.0	100.0		800.0	800.0	
accroissement	7.6	7.6	2.25	10.0		100.0	100.0	
y borne inférieure	0.0	0.0	0.0	-3.3		0.0	0,0	
borne supérieure	47.5	47.5	20.04	3.3		20.4	20.4	
accroissement	9.5	9.5	3.34	1.1		3.4	3.4	

Nota: Les lettres a,b,c désignent différents passages du programme correspondant à un même exemple; ainsi, 1a désigne le premier passage du programme pour l'exemple 1.

#### S.E.G.R.E.G

```
FORMAT FOR READING X AND Y: (TR6,F6.2,TL11,F5.1)
CORE KB AMBROSIA
X IS IN EITHER INCREASING OR DECREASING ORDER
P AND O WILL BE ESTIMATED
 23 PAIRS OF (X,Y) VALUES READ!
DATA IS NOT MANIPULATED
MEAN OF X VALUES = 25.69565
STANDARD DEVIATION OF X VALUES = 20.36359
  O ADJÄGENT EQUÁL VALÚES AT BEGINNING OF Y VECTOR
  O ADJACENT EQUAL VALUES AT END OF Y VECTOR
PLOT SPECIFICATIONS
                                                                                                                             9.500
                                                                                           VI THET =
                                                                                                      47.500
                                             YSTEP=
                                                         7.630
                                                                   YORIGIN=
                                                                                0.000
PAIR 1
MAXIMUM P = 1
                 MAXTMUN Q = 1
EQUALITY OF VARIANCES TO BE TESTED
PLOTS WILL BE PRODUCED
SIGNIFICANCE LEVEL FOR EQUALITY OF VARIANCES TEST = .05000
SIGNIFICANCE LEVEL FOR DETERMINATION OF DEGREES OF POLYNOMIALS = .01000
TEST FOR EQUALITY OF VARIANCES
                               DUAL RESIDUAL
MS
                          RESIDUAL
SS
   1
                      5.3674925157 .4128848397
              13
                       .1920299883
                                       . 0320048480
   2
F RATIO = 12.90067
SIGNIFICANCE LEVEL =
                           .00246486
VARIANCES ARE UNEQUAL
ESTIMATION OF P AND Q. UNEQUAL VARIANCES
                                         MODEL 2
N1HAT P Q
                                                                                      .00035
.09779 .09779 .10000 ADEQUATE
                                                                           1,13
                                                  9.496601 23.00065
THE THEORETICAL SIGNIGICANCE LEVEL IS THE SQUARE ROOT OF SPECIFIED SIGNIFICANCE LEVEL IF THERE ARE THO POSSIBLE MAYS TO REDUCE THE DEGREE: OTHERWISE IT IS THE SPECIFIED SIGNIFICANCE LEVEL:
TABLE OF RESIDUAL SUMS OF SQUARES AND MEAN SQUARES CALCULATED DURING THE ESTIMATION OF P AND Q
N1HAT P Q RSS1 RSS2 RMS1 RMS2
                   1 5.3674925157 1920290883 4128849397 0320048480 0 14.8640933333 3148875000 1.0617209524 3449839286
```

Figure 2. État produit par l'ordinateur lors du premier passage du programme pour l'analyse de l'échantillon KB, exemple 1.

```
P = 1
               Q = 1
              RESIDUAL SUMS OF SQUARES
SEGMENT 1 SEGMENT 2
                                                       RESIDUAL MEAN SQUARES
SEGMENT 1 SEGMENT 2
                                                                                            LIKELIHOOD FUNCTION
                                 .1920293883
NI HAT
            RESS1
5.3674925157
                                                       RMS1
-41288 40 397
    Y=F(THETA) Y=F(BETA)
                     36.94
                            OBSERVED
                                                                    ESTIMATED
                                                                                                 ORDERED STANDARDIZED
                                                                                                                           STANDARD NOR HAL
OBSERVATION
                                                                               RESIDUAL
SUM OF RESIDUALS, SEGMENT 1
                                     .0000000000
SUM OF RESIDUALS, SEGMENT 2
                                     .0000000000
TEST FOR RUNS IN RESIDUALS
NUMBER OF POSITIVE SIGNS
                                 10
NUMBER OF NEGATIVE SIGNS
                                 13
MEAN
                                12.3
STANDARD DEVIATION
                                2.30
N(8,1) VARIATE
                                1.82
PROBABILITY LEVEL
                                 .97
FINAL MODEL
     ALL RESIDUAL SUMS OF SQUARES AND RESIDUAL MEAN SQUARES HERE DIVICED BY
                                                                                        100.
     POINT OF CHANGE PARAMETER # 15
     DEGREE OF POLYNOMIAL, SEGMENT 1 = 1
     DEGREE OF POLYNOMIAL, SEGMENT 2 = 1
     REGRESSION PARAMETERS
      A. TRANSFORMED X
          SEGMENT 1
                              13.18205
                                                 -19.20019
                                                                       0.00000
                                                                                          0.00000
                                                                                                              0.00000
                                                                                                                                 0.00000
          SEGMENT 2
                               5.64860
                                                  -2.01965
                                                                       0.09009
                                                                                          0.00000
                                                                                                              0.00000
      B. UNTRANSFORMED X
          SEGMENT 1
                              37.40967
                                                   -.94287
                                                                       0.00000
                                                                                          0.00000
                                                                                                              0.00000
                                                                                                                                 0.00000
          SEGMENT 2
                               8.19708
                                                   -.09918
                                                                       0.00003
                                                                                                              0.00000
                                                                                                                                 0.00000
     100 R SQUARED = 86.02406
     TOTAL SUM OF SQUARES =
                               3977.92435
```

Figure 2. Suite.

de 23 couples dont le premier élément est une concentration de pollen de type Ambrosia (grains g-1 de sédiment sec) et le deuxième élément la profondeur (cm) à laquelle on constate cette concentration; la carotte sédimentaire en question provient du lac Ontario. On s'est servi du milieu de l'intervalle d'échantillon pour indiquer la profondeur de cet échantillon. L'ensemble des échantillons est ordonné verticalement de haut en bas, l'échantillon prélevé en surface se voyant attribuer le rang 1. Les diverses options choisies sont indiquées dans le tableau 3 et l'état ainsi que les graphiques produits par l'ordinateur apparaissent dans les figures 2 et 3.

L'état produit par l'ordinateur comporte suffisamment de rubriques pour être facile à comprendre (figure 2). La conclusion qui découle du test portant sur l'égalité des variances s'applique à tous les résultats successifs. La fonction de vraisemblance relative, la somme des carrés des écarts résiduels, le carré moyen des écarts, la courbe de lissage ainsi que les écarts résiduels sont imprimés uniquement pour le modèle définitif. La dernière partie de l'état présente un récapitulatif décrivant le modèle définitif. Le facteur qui y apparaît doit être appliqué à toutes les sommes de carrés d'écarts résiduels ainsi qu'à tous les carrés moyens qui figurent dans les résultats précédents.

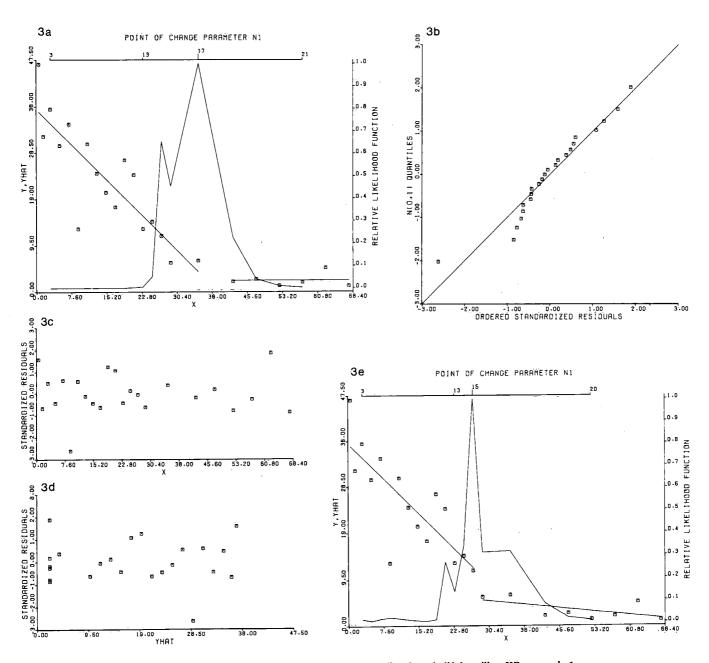


Figure 3. Graphiques produits par l'ordinateur pour l'analyse de l'échantillon KB, exemple 1.

Le modèle définitif est le modèle M<sub>u</sub>(1,1|15); il s'agit donc de deux droites de pentes différentes, une pour les échantillons 1 à 15 et une autre pour les échantillons 16 à 23; les variances qui correspondent à ces deux segments de droite sont différentes. Les variations de la concentration du pollen étant en général décrites de bas en haut, le modèle ci-dessus correspond à une augmentation du degré de concentration au taux de 0.10 grain g<sup>-1</sup> ·cm<sup>-1</sup> du dernier échantillon du bas jusqu'à l'échantillon 16, puis au taux de 0.94 grain·g<sup>-1</sup>·cm<sup>-1</sup> de l'échantillon 15 à l'échantillon prélevé à la surface. Des écarts plus grands par rapport à la moyenne ont également été constatés pour les 15 échantillons de la partie supérieure; on a en effet  $\hat{\sigma}_1^2$ =41.29 et  $\hat{\sigma}_2^2$ =3.20. Dans la figure 2, les taux et les variances sont exprimés en fonction de la valeur non transformée de x et de la somme des carrés des écarts résiduels correspondant à chaque segment. Le point de transition, l'échantillon 15, est le point où la fonction de vraisemblance relative atteint sa valeur maximum de 1. On peut se convaincre que l'échantillon 15 constitue bien le point de transition en étudiant la fonction de vraisemblance relative R(n<sub>1</sub>). Dans ce cas particulier, la valeur de la fonction de vraisemblance relative est ≥0.10 pour tous les échantillons qui appartiennent à l'intervalle fermé [12,17]. En examinant, par exemple, l'intervalle de vraisemblance relative qui correspond à la valeur 0.10, on obtient une certaine mesure du degré d'incertitude associé à l'estimation considérée.

Dans le cas où les valeurs de la variable indépendante sont ordonnées de façon croissante ou décroissante, quatre graphiques sont produits par l'ordinateur. La figure 3a représente les données, les courbes de lissage ainsi que la fonction de vraisemblance relative. Les données et les courbes de lissage sont repérées par rapport à l'axe vertical gauche et l'axe horizontal inférieur alors que la fonction de vraisemblance relative pour le paramètre du point de transition, n<sub>1</sub>, est repérée par rapport à l'axe vertical droit et l'axe horizontal supérieur. On observera que seules les valeurs entières de n<sub>1</sub> donnent lieu à une valeur non nulle de la fonction de vraisemblance relative; toutefois, pour faciliter la lecture du graphique, nous avons joint les valeurs de cette dernière. Les graphiques de représentation des écarts résiduels dans les figures 3b, c et d ne font apparaître aucune divergence significative par rapport aux hypothèses retenues par le modèle. L'écart résiduel normalisé dont la valeur absolue est la plus grande est celui qui correspond à l'échantillon 6 (valeur -2.5) et, sous des hypothèses de normalité, on s'attend à ce qu'environ 98 % des écarts résiduels appartiennent à l'intervalle [-2.5, 2.5]. En outre, le niveau de signification observé (0.97) lors du test des suites (figure 2) ne vient pas contredire l'hypothèse de l'indépendance admise par le modèle.

La connaissance générale de l'intervalle de vraisemblance auguel appartient le point de transition nous serait fort utile pour dater notre échantillon par la méthode de la distribution verticale du pollen de type Ambrosia. Or, selon Robbins et coll., on ne possède pas de définition précise du point de transition appliquée au pollen de type Ambrosia. Il est intéressant de noter que si cela désigne le point à partir duquel la concentration du pollen se met à augmenter à un taux accru, cela implique que q=0 pour le modèle décrit ici. Lorsqu'on attribue au couple (p,q) la valeur (1,0), on obtient le modèle Mu (1,017) (figure 3e) pour lequel l'intervalle de vraisemblance relative à 0.10 est l'intervalle [15, 18]. L'analyse des écarts résiduels n'a donné lieu à aucune différence marquée et on a 100R2=86.0 et 86.8 pour les modèles  $M_u(1,1|15)$  et  $M_u(1,0|17)$  respectivement. Dans Robbins et coll., le point de transition, dans le cas du pollen de type Ambrosia, a été fixé à l'échantillon 15 après analyse des données. Ce résultat correspond bien au premier modèle (deux droites) auquel nous sommes arrivés grâce à la méthode décrite ici. La comparaison des résultats obtenus par notre méthode et par la méthode d'inspection visuelle respectivement devrait aider l'analyste à mettre au point une méthode plus efficace pour la détermination du point où la fonction de répartition du pollen de type Ambrosia change de forme.

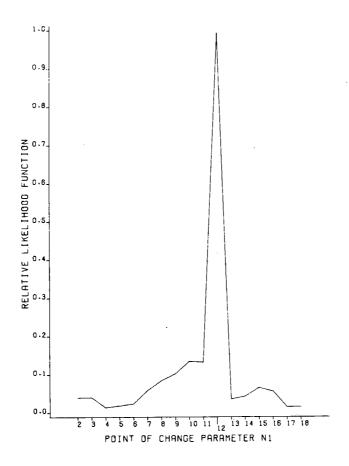
En général, la décision d'attribuer une valeur particulière au couple (p\*,q\*) ou de supposer le couple (p,q) connu dépend de la situation particulière dans laquelle on se trouve. Dans l'exemple du pollen de type Ambrosia, en attribuant une valeur supérieure à (p\*,q\*), on aurait obtenu un modèle mieux ajusté aux faibles variations des données mais qui n'aurait pas tenu compte de l'existence d'un point de transition à partir duquel le degré de concentration de ce type de pollen augmente dans les dépôts de sédiments du lac. Toutefois, à partir de l'expérience que l'auteur a acquise dans l'analyse d'autres échantillons, il est clair qu'il est impossible qu'une hypothèse unique convienne à tous les cas.

#### Exemple 2

Dans l'exemple précédent, la variable indépendante ainsi que les observations étaient rangées selon le même ordre, de sorte qu'on a supposé, pour l'établissement des modèles de régression, que les domaines de définition de la variable indépendante ne se chevauchaient pas. Bien que ce cas soit le plus courant, il arrive que les domaines de définition de la variable indépendante se chevauchent. Cela peut par exemple se produire dans l'étude de la stabilité dans le temps de la relation entre deux variables lorsque la valeur de la variable indépendante n'est pas contrôlée et que la personne qui effectue les mesures est remplacée à un certain moment par une autre personne.

Les données générées artificiellement par Quandt (1958) ont été analysées en supposant les options précisées dans le tableau 3. La figure 4 représente les deux graphiques élaborés respectivement pour décrire les données et les

courbes de lissage d'une part et la fonction de vraisemblance relative d'autre part. Des symboles distincts ont été employés pour représenter les données et les courbes de lissage correspondant aux deux segments et la fonction de



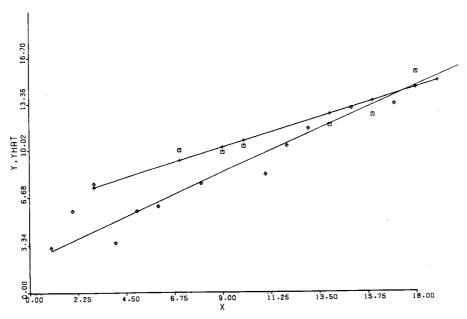


Figure 4. Représentation graphique des données avec les courbes de lissage et la fonction de vraisemblance relative, exemple 2.

vraisemblance relative est représentée de la même manière que dans le cas où un seul graphique est élaboré. Le modèle qui correspond à ce graphique est le modèle  $M_e(1,1|12)$ , avec  $\hat{y}=2.22\pm0.69x$  et  $\hat{y}=5.91\pm0.48x$  étant les équations des droites de lissage qui correspondent respectivement aux échantillons 1 à 12 et 13 à 20.

#### Exemple 3

La méthode que nous avons appliquée jusqu'à maintenant ne prévoit pas de test visant à démontrer l'existence d'un point de transition. Elle suppose à priori

l'existence de ce point. Pour voir ce qui se passe lorsqu'on applique cette méthode à un ensemble de données ne comportant pas de point de transition, nous avons appliqué le programme à un échantillon de 100 variables aléatoires générées à partir d'une distribution normale N(0,1). Les options retenues sont celles qui sont précisées au tableau 3 et les graphiques sont représentés à la figure 5. Puisque la fonction de vraisemblance relative prend une valeur supérieure ou égale à 0.30 pour toutes les valeurs possibles du point de transition, on en conclut qu'il ne peut exister aucun point de cette nature.

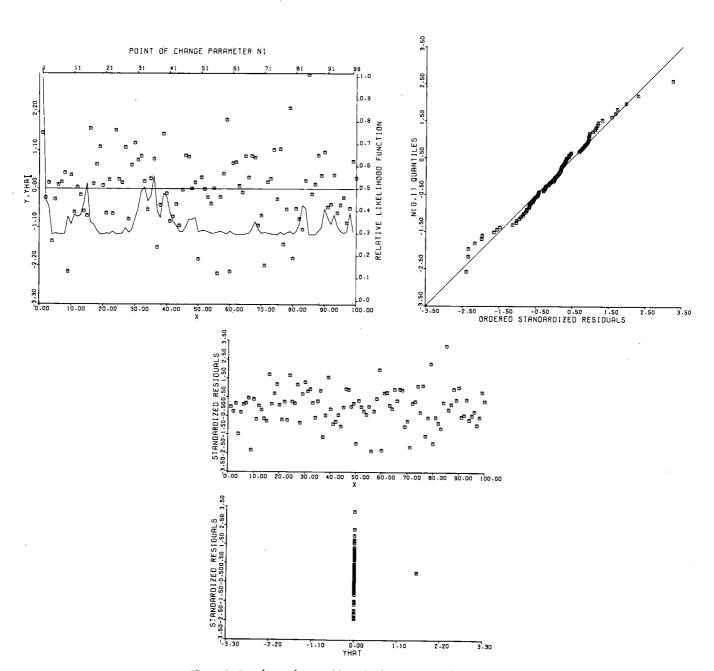


Figure 5. Représentation graphique des écarts normalisés, exemple 3.

#### Exemple 4

Bien que la méthode décrite ici ne s'applique en principe qu'au cas de l'existence d'un point de transition unique, on peut l'adapter de diverses manières au cas où il existe plusieurs points de transition. Une de ces adaptations est illustrée ici par l'exemple de la distribution verticale du pollen de type *Pinus* dans une carotte sédimentaire prélevée à Bog D Pond, Hubbard County, Minnesota. Les pourcentages ont été calculés à partir des valeurs indiquées par McAndrews (1966). Le pourcentage du pollen de type *Pinus banksiana/resinosa* décroît de haut en bas de la

carotte alors que le pourcentage du pollen de type *Pinus strobus* se comporte de façon inverse. La somme de ces deux pourcentages devrait donc présenter deux points de transition. Le modèle définitif qui représente la somme de ces deux pourcentages a été établi après trois passages du programme. On a procédé de la façon ci-après :

- 1) À partir de la totalité des 49 échantillons, on a établi le modèle M<sub>e</sub>(2,1|44).
- 2) L'ensemble de données qui correspond aux échantillons 1 à 44 est décrit par le modèle M<sub>u</sub> (0,0|16).

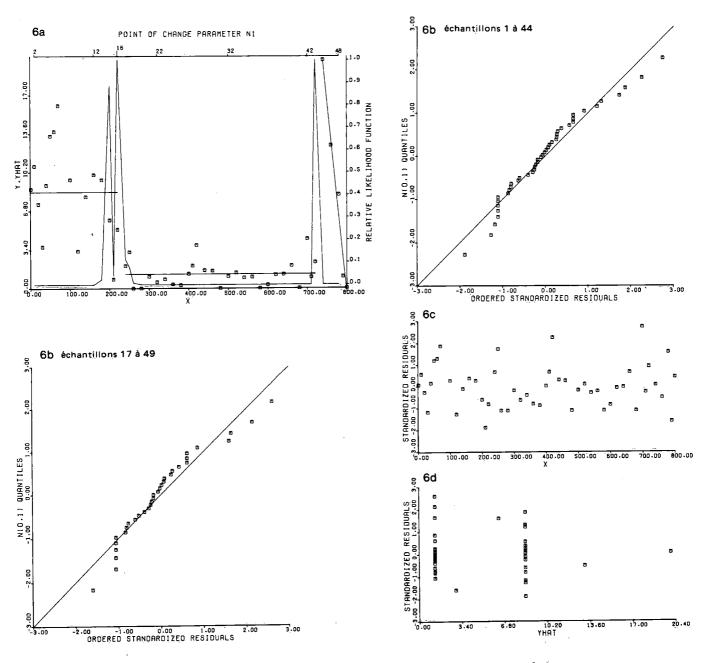


Figure 6. Graphiques tracés à partir des passages 2 et 3 du programme, exemple 4.

3) Pour vérifier que le point de transition inférieur correspond bien à l'échantillon 44, on a analysé une suite (à savoir les échantillons 17 à 49) comportant un seul point de transition. On a obtenu ainsi le modèle M. (0.1/44).

Les options retenues sont celles du tableau 3. Les graphiques produits lors des deux derniers passages du programme étant tous à la même échelle, on a pu découper et recoller les régions correspondant respectivement aux échantillons 1 à 16 et 17 à 49 (figures 6a,c,d). Le graphique Q-Q est représenté à la figure 6b. En choisissant les valeurs appropriées obtenues lors des deux derniers passages, le modèle définitif est défini par les relations ci-après

$$\hat{y}_i = 8.46$$
  $i = 1, 2 ... 16$   
= 1.25  $i = 17, 18 ... 44$   
= 273 - 0.34 x<sub>i</sub>  $i = 45, 46 ... 49$ 

avec  $\hat{\sigma}_1^2$ =16.52,  $\hat{\sigma}_2^2$ =1.27 et  $\hat{\sigma}_3^2$ =2.62. Le calcul de 100R² pour trois segments se fait de la même façon que pour deux; cela donne la valeur 100R² = 75.5. Bien que les écarts résiduels comportent encore un certain nombre de pointes, les variations importantes ont été repérées. En outre, le niveau de signification mesuré lors du test des suites appliqué au modèle à trois segments est égal à 0.20. Les valeurs nécessaires au calcul de 100R² et du niveau de signification ont été obtenues lors de l'un des trois passages du programme.

Il n'est pas toujours possible d'obtenir, dans le cas où il existe plusieurs points de transition, une estimation satisfaisante de tous ces points. Cette situation peut notamment se rencontrer lorsque les variations sont du même ordre de grandeur. Toutefois, lorsqu'on parvient, par l'étude du diagramme de dispersion, à décomposer l'ensemble de données en segments contenant chacun un seul point de transition, on peut souvent obtenir une estimation satisfaisante grâce à une série de passages du programme.

#### **ANALYSE**

Les bornes supérieures des degrés des polynômes doivent être choisies avec soin pour éviter que la courbe de lissage tienne compte de faibles variations dans les données. Comme on l'a déjà mentionné, la courbe de lissage doit être étudiée en vue de déterminer les écarts éventuels par rapport aux hypothèses. Lorsqu'on attribue à la borne supérieure du degré du polynôme une valeur trop élevée dans l'un des deux segments, l'estimation du point de

transition devient fortement dépendante de l'existence de points isolés dont la valeur diffère nettement de celle des points voisins. Cette caractéristique de la méthode est évidente lorsqu'on examine les graphiques.

Il est essentiel que la méthode permette de tenir compte de l'inégalité des variances dans les deux segments puisqu'il arrive souvent qu'une perturbation subie par un système entraîne une augmentation de la variance. La différence entre les variances dans les deux segments de la courbe de lissage n'apparaît pas dans les graphiques qui représentent les écarts résiduels puisque ces derniers sont normalisés. Les fluctuations de la variance qui ne sont pas dues simplement au fait que cette dernière prend une valeur constante différente dans chaque segment peuvent être détectées à partir de la représentation graphique des écarts résiduels et du test des suites.

Les exemples que nous avons présentés ne montrent pas comment le programme s'applique à plus d'une paire de variables. Le cas de l'exemple 4, qui a été traité en trois passages sans production de graphiques, aurait pu être traité en un seul passage avec production de graphiques. Chaque fois, ce sont les deux mêmes colonnes de données qui sont définies. La définition de paires de colonnes différentes donnerait bien à une situation plus habituelle. Ainsi, pour l'exemple 4, les colonnes de données auraient pu être constituées des trois paires: 1) profondeur et concentration du pollen de type *Pinus-banksiana/resinosa*, 2) profondeur et concentration du pollen de type *Pinus strobus* et 3) profondeur et somme des concentrations des pollens des types *Pinus-banksiana/résinosa* et *Pinus strobus*.

#### REMERCIEMENTS

L'algorithme a été codé et le programme testé par J. Hodson. Les sous-programmes de traçage de graphiques ont été rédigés à l'origine par D. Pateman; une version antérieure d'un programme semblable a été rédigée par M. Fellowes.

#### RÉFÉRENCES

Draper, N.R. et H. Smith. 1981. Applied Regression Analysis. 2e éd., John Wiley and Sons, Inc., New York.

Esterby, S.R. et A.H. El-Shaarawi. 1980. Examples of inferences about the point of change in a regression model. Non public. Institut national de recherches sur les eaux, Burlington (Ontario).

Esterby, S.R. et A.H. El-Shaarawi. 1981a. Inference about the point of change in a regression model. J. Royal Stat. Soc. Series C. 30(3): 277-285.

Esterby, S.R. et A.H. El-Shaarawi. 1981b. Likelihood inference about the point of change in a regression regime. J. Hydrol. 53: 17-30.

- McAndrews, J.H. 1966. Postglacial history of prairie, savanna and forest in Northwestern Minnesota. Mem. Torrey Bot. Club, Vol. 22, No 2, 72 p.
- Quandt, R.E. 1958. The estimation of parameters of a linear regression system obeying two separate regimes, J. Am. Stat.

Assoc. 53: 873-880.

Robbins, J.A., D.N. Edgington et A.L.W. Kemp. 1978. Compa ative <sup>210</sup> Pb, <sup>131</sup> Cs, and pollen geochronologies of sedimen from Lakes Ontario and Erie. Quat. Res. 10(2): 256-278.

#### **ANNEXE**

La normalisation de la variable indépendante par soustraction de la moyenne et division par l'écart type, où la moyenne et l'écart type sont calculés à partir de toutes les valeurs de n, s'est révélée suffisante. Le processus de normalisation adopté, analogue à l'utilisation de la matrice de corrélation pour le calcul d'une fonction de régression unique, consiste à normaliser chaque segment indépendamment. Toutefois, cette méthode augmente notablement la complexité des calculs puisque les échantillons de chaque segment varient en fonction de n<sub>1</sub>. On adopte donc, pour un segment particulier, la courbe de lissage

$$y_i = \beta_0 + \beta_1 \left(\frac{x_i - \bar{x}}{s}\right) + \dots + \beta_d \left(\frac{x_i - \bar{x}}{s}\right)^d$$

et pour exprimer les valeurs de  $\theta_0$ ,  $\theta_1$ ...  $\theta_d$  qui correspondent au modèle en fonction des variables d'origine, on calcul  $\underline{\theta} = (\theta_0, \theta_1, \dots, \theta_d)$ ' à partir de

$$\frac{\theta}{\beta} = A\underline{\beta}^*$$
où
$$\underline{\beta}^* = \begin{bmatrix} \beta_0 \\ \beta_1/s \\ \vdots \\ \beta_d/sd \end{bmatrix}$$

A est une matrice  $(d+1) \times (d+1)$  dont les éléments doivent être calculés. Ces éléments sont égaux aux termes du développement en série du binôme  $(1-\overline{x})^{j-1}$  pour j=1,2...d+1. Les j premiers éléments de la je colonne de A sont donnés par les termes de ce développement en série, en ordre inverse, et les derniers (d-j+1) éléments sont nuls. Ainsi, pour d=2, on a

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 1 & -x & x^2 \\ 0 & 1 & -2x \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1/s \\ \beta_2/s^2 \end{bmatrix}$$

Le coefficient de corrélation multiple, R<sup>2</sup>, est égal par définition à la somme des carrés attribuée à la régression divisée par la somme totale des carrés des écarts (TSS). Dans le modèle que nous étudions ici, le numérateur de R<sup>2</sup> n'est pas simplement égal à la somme des carrés attribuées à la régression calculées séparément pour chaque segment puisque cette somme ne tient pas compte de toutes les composantes du modèle. Pour le voir, il suffit de décomposer la somme totale des carrés des écarts comme suit :

$$TSS = \sum_{i=1}^{n} (y_i - \overline{y})^2 = \sum_{i=1}^{n_1} (y_i - \overline{y})^2 + \sum_{i=n_1+1}^{n} (y_i - \overline{y})^2$$
$$= TSS_1 + TSS_2$$

Mais

$$TSS_1 = \sum_{i=1}^{n_1} (y_i \cdot \overline{y}_1)^2 + \sum_{i=1}^{n_1} (\overline{y}_1 \cdot \overline{y})^2$$

et le premier terme est égal à

où l'indice 1 désigne une valeur calculée à partir des échantillons du segment 1. Les deux termes ci-dessus correspondent respectivement à la somme des carrés des écarts et à la somme des carrés attribuée à la régression pour le segment 1. En appliquant la même décomposition au segment 2, puis en recombinant et en groupant les termes correspondants des deux segments, on obtient

$$TSS = \{ResSS_1 + ResSS_2\} + \{RegSS_1 + RegSS^2\} + M$$

où  $M = n_1(\vec{y}_1 - \vec{y})^2 + (n-n_1)(\vec{y}_2 - \vec{y})^2$  et ResSS et RegSS désignent respectivement la somme des carrés des écarts et la somme des carrés attribuée à la régression. On a donc

$$100R^2 = 100 \{TSS - (ResSS_1 + ResSS_2)\}/TSS$$

Environment Canada Library, Burlington