

B. FICHET

**Sur les résultats de Blum en approximation et optimisation  
stochastiques multidimensionnelles**

*Statistique et analyse des données*, tome 3, n° 3 (1978), p. 57-68

[http://www.numdam.org/item?id=SAD\\_1978\\_\\_3\\_3\\_57\\_0](http://www.numdam.org/item?id=SAD_1978__3_3_57_0)

© Association pour la statistique et ses utilisations, 1978, tous droits réservés.

L'accès aux archives de la revue « Statistique et analyse des données » 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/>

SUR LES RESULTATS DE BLUM EN APPROXIMATION ET OPTIMISATION  
STOCHASTIQUES MULTIDIMENSIONNELLES

B.FICHET

Laboratoire de Physique, Faculté de Médecine.  
Université de Marseille II.

L'extension au cas multidimensionnel des processus d'approximation stochastique de Robbins-Monro et d'optimisation stochastique de Kiefer-Wolfowitz, fut l'oeuvre de Blum. On lui doit trois importants théorèmes, deux pour l'approximation et un pour l'optimisation. Depuis, de nombreux auteurs ont établi des théorèmes de convergence tant pour ces processus que pour des processus généralisés. Nous montrons par cette note que les résultats de Blum restent d'une grande actualité. Pour différentes régressions les convergences des processus de Robbins-Monro ou de Kiefer-Wolfowitz peuvent être déduites tant des théorèmes de Blum que de théorèmes postérieurs. Ainsi, à l'aide des théorèmes de Blum, nous démontrons la convergence pour l'approximation d'une régression linéaire donnée par une matrice définie négative, sous l'hypothèse de la variance uniformément bornée ; et même sans cette hypothèse à l'aide d'un exemple. Et nous établissons de même la convergence pour l'optimisation d'une régression de type parabolique.

1 - INTRODUCTION

Rappelons les processus d'approximation stochastique de Robbins-Monro et d'optimisation stochastique de Kiefer-Wolfowitz dans  $\mathbb{R}^k$ .

Pour l'approximation, on considère une famille  $\{Y_x\}$  de vecteurs aléatoires à valeurs dans  $\mathbb{R}^k$ , indicée par le paramètre  $x \in \mathbb{R}^k$ . On suppose que :  $\forall x \in \mathbb{R}^k$ ,  $E[Y_x] = M(x)$  existe, où  $M$  (inconnue) est Borel-mesurable. On se propose alors, non pas d'approcher la fonction inconnue  $M$ , mais d'estimer la solution  $\theta$  de l'équation  $M(x) = \alpha$  ( $\alpha \in \mathbb{R}^k$ , donné) ; dans ce but, on construit le processus de Robbins-Monro défini par :

(1.1)  $X_{n+1} - X_n = a_n(Y_n - \alpha)$  ( $n \in \mathbb{N}^*$ ) où :

-  $X_1 \in \mathbb{R}^k$  est donné.

-  $Y_n$  est un vecteur aléatoire ayant même loi de probabilité (conditionnelle) que  $Y_{x_n}$  quand  $x_n$  est la réalisation de  $X_n$ .

-  $\{a_n\}$  est une suite à termes positifs vérifiant :  $\sum_{n=1}^{\infty} a_n = +\infty$  ;  $\sum_{n=1}^{\infty} a_n^2 < +\infty$

Dans le processus de Markov ainsi construit, notons que  $M(X_n)$  est alors une version de  $E[Y_n | X_1, \dots, X_n] = E[Y_n | X_n]$ , pour tout  $n$ .

Pour l'optimisation, on considère une famille  $\{Y_x\}$  de variables aléatoires, indicée par le paramètre  $x \in \mathbb{R}^k$ . On suppose que :  $\forall x \in \mathbb{R}^k$ ,  $E[Y_x] = M(x)$  existe, où  $M$  (inconnue) est Borel-mesurable.  $\{u_1, \dots, u_k\}$  étant une base orthonormée de  $\mathbb{R}^k$ , pour tout réel  $c > 0$ ,

soient  $Y_{x-cu_1}^*$ ,  $Y_{x+cu_1}^*$ , ...,  $Y_{x-cu_k}^*$ ,  $Y_{x+cu_k}^*$ ,  $2k$  variables aléatoires indépendantes de même loi que, respectivement,  $Y_{x-cu_1}$ ,  $Y_{x+cu_1}$ , ...,  $Y_{x-cu_k}$ ,  $Y_{x+cu_k}$ , et définissons  $Y_{x,c}$

par :  $Y_{x,c} = [(Y_{x+cu_1}^* - Y_{x-cu_1}^*), \dots, (Y_{x+cu_k}^* - Y_{x-cu_k}^*)]$

Alors, pour estimer  $\theta \in \mathbb{R}^k$  qui maximise  $M(x)$ , on construit le processus de Kiefer-Wolfowitz défini par :

(1.2)  $X_{n+1} - X_n = (a_n/c_n)Y_n$  ( $n \in \mathbb{N}^*$ ) où :

-  $X_1 \in \mathbb{R}^k$  est donné

-  $Y_n$  est un vecteur aléatoire ayant même loi de probabilité (conditionnelle) que  $Y_{x_n, c_n}$  quand  $x_n$  est la réalisation de  $X_n$ .

-  $\{a_n\}$  et  $\{c_n\}$  sont deux suites à termes positifs vérifiant :

$$\lim_{n \rightarrow \infty} c_n = 0 ; \sum_{n=1}^{\infty} a_n = +\infty ; \sum_{n=1}^{\infty} a_n c_n < +\infty ; \sum_{n=1}^{\infty} (a_n/c_n)^2 < +\infty .$$

Similairement, Blum définit un processus par :

(1.3)  $X_{n+1} - X_n = (a_n/c_n)Z_n$  ( $n \in \mathbb{N}^*$ ) où :

$Z_n$  est un vecteur aléatoire ayant même loi de probabilité (conditionnelle) que

$Z_{x_n, c_n}$  et où :

$Z_{x,c} = [(Y_{x+cu_1}^* - Y_x^*), \dots, (Y_{x+cu_k}^* - Y_x^*)]$  est défini comme  $Y_{x,c}$

Les premiers théorèmes de convergence, pour les processus multidimensionnels cités, sont ceux de Blum [2]; nous les présentons au paragraphe suivant. Depuis, ces résultats ont été repris et améliorés - voir par exemple Macchi [13] chapitre I, et Daubeze [5] - ; des théorèmes de convergence, tant pour ces processus que pour des processus généralisés, établis sous des hypothèses faibles, peuvent être trouvés dans Dvoretzky [9], Derman et Sacks [7], Fabian [10], Venter [16], et dans les thèses de Hiriart-Urruty [12] chapitres I, II, III, et Bertran [3] chapitre I ; pour une bibliographie sur le sujet, on peut consulter Schmetterer [15], Dupac et Ivanov [8] et Wasan [17] .

Nous nous proposons de montrer par cette note, que l'approximation et l'optimisation de certaines régressions, qui se déduisent des travaux plus récents, peuvent tout aussi bien être obtenues à l'aide des théorèmes originels de Blum.

## 2 - LES THEOREMES DE BLUM

Pour l'approximation nous utilisons les notations suivantes :

- $x^+ = \frac{1}{2} [ |x| + x ]$  pour toute variable réelle  $x$ .
- $D^2$  est l'espace des fonctions de  $\mathbb{R}^k$  dans  $\mathbb{R}$ , admettant des dérivées partielles premières et secondes continues.
- $\forall f \in D^2$ ,  $\nabla f$  et  $A_f$  sont respectivement le gradient et la matrice des dérivées secondes de  $f$ .

On a alors :  $\forall f \in D^2$ ,  $f(x+h) = f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle h, A_f(x) h \rangle$   
avec :  $0 \leq t \leq 1$ ,  $t \equiv t(x, h)$ .

Nous notons encore :

- $U_f(x) = \langle \nabla f(x), M(x) - \alpha \rangle$
- $\forall a \in \mathbb{R}$ ,  $V_f(a, x) = E \{ \langle Y_x - \alpha, A_f [ x + t_1 a(Y_x - \alpha) ] (Y_x - \alpha) \rangle$   
avec  $t_1 \equiv t [ x, a(Y_x - \alpha) ]$ .

Alors, lorsque  $M(\theta) = \alpha$ , pour  $\theta$  et  $\alpha$  quelconques, les théorèmes d'approximation de Blum s'énoncent :

### Théorème 1

S'il existe  $f \in D^2$ , satisfaisant à :

- (2.1)  $\forall x \in \mathbb{R}^k$ ,  $f(x) \geq 0$
- (2.2)  $\forall \epsilon > 0$ ,  $\inf_{\epsilon \leq \|x-\theta\|} |f(x) - f(\theta)| > 0$
- (2.3)  $\forall \epsilon > 0$ ,  $\sup_{\epsilon \leq \|x-\theta\|} U_f(x) < 0$

$$(2.4) \quad \exists V < +\infty : \forall x \in \mathbb{R}^k, \forall a \in \mathbb{R}, V_f(a, x) \leq V$$

Alors la suite  $\{X_n\}$  définie par (1.1) converge presque sûrement vers  $\theta$ .

### Théorème 2

S'il existe  $f \in D^2$  satisfaisant à (2.1), (2.2) et à :

$$(2.5) \quad \forall \delta > 0 \quad \exists \lambda(\delta) > 0 : \forall a \in \mathbb{R}, \sup_{\delta \leq \|x-\theta\|} [U_f(x) + \lambda(\delta) V_f^+(a, x)] < 0$$

$$(2.6) \quad \exists \varepsilon > 0, V < +\infty : \forall a \in \mathbb{R}, \sup_{0 \leq \|x-\theta\| < \varepsilon} V_f(a, x) \leq V$$

Alors la suite  $\{X_n\}$  définie par (1.1) converge presque sûrement vers  $\theta$ .

Notons que la démonstration tient encore si (2.5) n'est satisfaite que pour  $\delta = \varepsilon$  définie par (2.6), pourvu que (2.3) soit aussi satisfaite. Signalons enfin que notre énoncé du théorème 2 diffère légèrement de celui de Blum.

Lorsque  $\theta$  est quelconque, le théorème d'optimisation de Blum s'énonce :

### Théorème 3

Si :

$$(2.7) \quad M \text{ admet des dérivées partielles premières et secondes continues}$$

$$(2.8) \quad \forall x \in \mathbb{R}^k, \sigma_x^2 = \text{var}(Y_x) \leq \sigma^2 < +\infty$$

$$(2.9) \quad \forall \varepsilon > 0, \exists \rho(\varepsilon) > 0 : \|x-\theta\| \geq \varepsilon \Rightarrow M(x) - M(\theta) \leq -\rho(\varepsilon) \text{ et } \|\nabla M(x)\| \geq \rho(\varepsilon)$$

$$(2.10) \quad \forall x \in \mathbb{R}^k, \forall i, j = 1, \dots, k \quad \left| \frac{\partial^2 M(x)}{\partial x_i \partial x_j} \right| \leq C$$

Alors la suite  $\{X_n\}$  définie par (1.3) converge presque sûrement vers  $\theta$ .

Notons que la convergence est également assurée pour le processus défini par (1.2) ; à quelques détails près, la démonstration est semblable à celle de Blum.

## 3 - APPLICATION A L'APPROXIMATION DE CERTAINES REGRESSIONS LINEAIRES.

Nous envisageons le cas d'une régression linéaire caractérisée par une matrice définie négative. Ce cas est ébauché par Blum, dans son exemple 2 - déduit du théorème 2 - et que nous présentons pour  $\theta$  et  $\alpha$  quelconques sous la forme du :

### Corollaire 1

Si :

$$(3.1) \quad M(x) = Bx \text{ ou } B \text{ est une matrice définie négative.}$$

$$(3.2) \quad \exists \varepsilon > 0, \exists C > 0 : \|x-\theta\| \leq \varepsilon \Rightarrow E \{ \|Y_x - \alpha\|^2 \} \leq C$$

$$(3.3) \quad \exists \rho(\theta) > 0 : \|x-\theta\| > \varepsilon \Rightarrow \langle x-\theta, B(x-\theta) \rangle > +\rho(\theta), E \{ \|Y_x - \alpha\|^2 \} \leq 0$$

Alors la suite  $\{X_n\}$  définie par (1.1) converge presque sûrement vers  $\theta$ .

Les conditions données par Blum dépendent encore du second membre  $\alpha$  de l'équation et de la solution  $\theta$ . Nous pouvons maintenant démontrer le :

Corollaire 2

Si (3.1) est vérifiée, et, notant  $Y_x^i$  la  $i$ ème composante de  $Y_x$ , si :

$$(3.4) \quad \exists \sigma > 0 : \forall x \in \mathbb{R}^k, \forall i = 1, \dots, k \quad (\sigma_x^i)^2 = \text{var}(Y_x^i) \leq \sigma^2$$

Alors la suite  $\{X_n\}$  définie par (1.1) converge presque sûrement vers  $\theta$ .

Démonstration :

Appliquons le corollaire 1. On a :

$$\begin{aligned} E\{\|Y_x - \alpha\|^2\} &= E\{\|Y_x\|^2\} - 2 \langle Bx, \alpha \rangle + \|\alpha\|^2 \\ &\leq k\sigma^2 + \|Bx\|^2 + 2 \|Bx\| \cdot \|\alpha\| + \|\alpha\|^2 \\ &\leq k\sigma^2 + \|B\|^2 \cdot \|x\|^2 + 2 \|B\| \cdot \|\alpha\| \cdot \|x\| + \|\alpha\|^2 \end{aligned}$$

où  $\|\cdot\|$  est la norme classique d'une application linéaire.

Dès lors,  $\forall \epsilon > 0$ , (3.2) est vérifiée, puisque :  $\|x - \theta\| \leq \epsilon \implies \|x\| \leq \epsilon + \|\theta\|$ .

Seule donc, reste à vérifier (3.3) pour un certain  $\epsilon > 0$ .

$B$  étant définie négative :  $\exists b > 0 : \forall x \in \mathbb{R}^k \quad \langle x, Bx \rangle \leq -b \|x\|^2$

d'où :  $\langle x - \theta, B(x - \theta) \rangle \leq -b \|x - \theta\|^2 \leq -b (\|x\| - \|\theta\|)^2$

Alors  $\rho(\theta)$  étant à déterminer :

$$\begin{aligned} &\langle x - \theta, B(x - \theta) \rangle + \rho(\theta) E\{\|Y_x - \alpha\|^2\} \\ &\leq -b (\|x\| - \|\theta\|)^2 + \rho(\theta) [k\sigma^2 + \|B\|^2 \cdot \|x\|^2 + 2 \|B\| \cdot \|\alpha\| \cdot \|x\| + \|\alpha\|^2] \\ &\leq [-b + \rho(\theta) \|B\|^2] \|x\|^2 + 2 [b \|\theta\| + \rho(\theta) \|B\| \cdot \|\alpha\|] \|x\| \\ &+ [-b \|\theta\|^2 + \rho(\theta) k\sigma^2 + \rho(\theta) \|\alpha\|^2] \equiv K(x) \quad (\text{posé}) \end{aligned}$$

Si nous choisissons  $0 < \rho(\theta) < b / \|B\|^2$

$K(x) \rightarrow -\infty$  quand  $\|x\| \rightarrow \infty$

Donc il existe  $\epsilon' > 0$  tel que :

$$\|x\| \geq \epsilon' \implies \langle x - \theta, B(x - \theta) \rangle + \rho(\theta) E\{\|Y_x - \alpha\|^2\} \leq K(x) < 0$$

Dès lors, la démonstration est complète avec  $\epsilon = \epsilon' + \|\theta\|$  et  $\rho(\theta)$  précédemment déterminé.  $\square$

Ainsi, pour une régression linéaire caractérisée par une matrice définie négative, la convergence est assurée quel que soit  $\alpha$  (ou  $\theta$ ), sous la seule hypothèse de la variance uniformément bornée. Wolfowitz [18] a montré que celle-ci ne pouvait être omise dans l'étude originelle de Robbins et Monro. Friedman [11] fut le premier à

démontrer la convergence sans cette hypothèse, mais en transformant le processus. L'exemple suivant montre que les résultats de Blum, bien antérieurs à ceux de Friedman, ne l'exigeaient déjà pas.

### Exemple

$\forall p \in \mathbb{R}$ , soit  $f_p$  la fonction de  $\mathbb{R}$  dans  $\mathbb{R}$  définie par :

$$\forall t \in \mathbb{R}, f_p(t) = \begin{cases} (1-m)^2 m e^{-m(1-m)t} & \text{si } t \geq 0 \\ m^2 (1-m) e^{m(1-m)t} & \text{si } t < 0 \end{cases}$$

$$\text{avec } m = \begin{cases} [p+2 - \sqrt{4+p^2}]/2p & \text{si } p \neq 0 \\ 1/2 & \text{si } p = 0 \end{cases}$$

Il est aisé de voir que lorsque  $p$  croît de  $-\infty$  à  $+\infty$ ,  $m$  décroît de 1 à 0, et que  $f_p$  est une densité de probabilité.

Si  $X$  est une variable aléatoire dont la loi de probabilité admet une densité  $f_p$ , de simples calculs montrent que :

$$E[X] = p ; \quad E[X^2] = \frac{2}{m^2(1-m)^2} ; \quad \text{Var}(X) = \frac{-4m^2 + 4m + 1}{m^2(1-m)^2}$$

Soit  $B$  une matrice définie négative d'ordre  $(k, k)$ . Notons  $B^i x$  la  $i$ ème composante de  $Bx$  ( $B^i$  est une forme linéaire).

Supposons que :

$\forall x \in \mathbb{R}^k, \forall i = 1, \dots, k$ , la loi de  $Y_x^i$  admet une densité  $f_{p_i}$ , avec  $p_i \equiv B^i x$ .

Alors les relations précédentes montrent que :

$$M(x) = E[Y_x] = Bx ; \quad \sup_{x \in \mathbb{R}^k} \text{Var}(Y_x^i) = +\infty$$

En appliquant le corollaire 1, montrons la convergence du processus (1.1) pour tout  $\theta$ .

La variation de  $m$  en fonction de  $p$ , montre que :

$$\forall A > 0, \exists r(A) : |p| \leq A \implies r(A) \leq m \leq 1 - r(A)$$

$$\text{avec } 0 < r(A) < 1/2$$

Notant  $\forall i = 1, \dots, k, M_2^i = E[(Y_x^i)^2]$ ;

$$E\{\|Y_x - \alpha\|^2\} \leq \sum_{i=1}^k M_2^i + 2 \|\alpha\| \cdot \|B\| \cdot \|x\| + \|\alpha\|^2$$

$$\text{Or : } \|x\| \leq \varepsilon \implies |p_i| \leq \|B^i\| \cdot \varepsilon \implies M_2^i \leq 2 / [r(\|B^i\| \varepsilon)]^4.$$

Dès lors,  $\forall \varepsilon > 0$  (3.2) est vérifiée.

Seule reste à vérifier (3.3) pour un certain  $\varepsilon > 0$ .

$\rho(\theta)$  étant à déterminer :

$$\begin{aligned} & \langle x-\theta, B(x-\theta) \rangle + \rho(\theta) E\{\|Y_x - \alpha\|^2\} \\ & \leq -b(\|x\| - \|\theta\|)^2 + \rho(\theta) \left[ \sum_{i=1}^k M_2^i + 2 \|\alpha\| \cdot \|B\| \cdot \|x\| + \|\alpha\|^2 \right] \\ & \leq -b \|x\|^2 + \rho(\theta) \sum_{i=1}^k M_2^i + 2 [b \|\theta\| + \rho(\theta) \|\alpha\| \cdot \|B\|] \|x\| + [-b \|\theta\|^2 + \rho(\theta) \|\alpha\|^2] \end{aligned}$$

Cherchons une borne supérieure pour  $M_2^i$ .

Un calcul simple montre que :

$$M_2^i \sim 2p_i^2 \text{ quand } p_i \rightarrow +\infty ; M_2^i \sim 2p_i^2 \text{ quand } p_i \rightarrow -\infty .$$

Soit  $\delta > 0$  donné ; alors, il existe  $A > 0$  :

$$|p_i| > A \implies M_2^i \leq 2(1+\delta) \|B^i\|^2 \|x\|^2$$

$$|p_i| \leq A \implies M_2^i \leq 2/[r(A)]^4$$

Soit  $\varepsilon' > 0$  tel que :  $2/[r(A)]^4 \leq 2(1+\delta) K^2 \varepsilon'^2$  où  $K = \max_i \|B^i\|$ .

Alors  $\|x\| \geq \varepsilon' \implies M_2^i \leq 2(1+\delta) K^2 \|x\|^2 \quad i=1, \dots, k$ .

Dès lors, la démonstration suit celle du corollaire 2. avec  $\rho(\theta) < b/2 k(1+\delta) K^2$ .  $\square$

#### Remarque

Comme exemple d'application du théorème 1, Blum envisage le cas d'une régression linéaire dans une boule, et constante à l'extérieur, soit :

$$M(x) = \begin{cases} Bx & \text{si } \|x\| \leq \rho \\ [\rho/\|x\|] Bx & \text{si } \|x\| > \rho \end{cases}$$

où  $B$  est une matrice  $(k,k)$  définie négative.

Avec  $f(x) = \|x\|^2$ , Blum démontre, sous l'hypothèse (3.4) de la variance uniformément bornée, la convergence du processus (1.1).

Mais cet exemple reste d'une portée limitée si on ne peut établir la convergence pour tout  $\theta$  vérifiant  $\|\theta\| < \rho$ . Si, par une naturelle généralisation, on prend  $f(x) = \|x-\theta\|^2$ , il est clair que (2.1) et (2.2) sont vérifiées, et que, comme dans [2], (2.4) l'est aussi. Montrons par un contre-exemple, que (2.3) peut ne pas être satisfaite si  $k > 1$ .

Contre-exemple :

Définissons B comme une matrice diagonale, d'éléments diagonaux  $\lambda_1, \dots, \lambda_k$  négatifs.

Si nous trouvons x tel que  $\langle x-\theta, M(x)-\alpha \rangle > 0$ , (2.3) ne sera pas satisfaite.

Notons  $x = (x^1, \dots, x^k)$ ,  $\theta = (\theta^1, \dots, \theta^k)$ ,

et choisissons x et  $\theta$  tels que :

$$\|x\| > \rho, \quad 0 < \rho x^1 / \|x\| < \theta^1 < x^1$$

$$(\text{ex : } k=2, \rho=x^1=1, x^2=\sqrt{3}, \theta^1=3/4, \theta^2=0)$$

Ceci entraîne :

$$\lambda_1 (x^1 - \theta^1) (\rho x^1 / \|x\| - \theta^1) > 0$$

$$(\text{ex : } -\lambda_1 / 16 > 0)$$

et puisque  $\|\theta\| < \rho$ , il existe un indice i supérieur à 1, tel que :

$$|\theta^i| < \rho |x^i| / \|x\| < |x^i|.$$

Fixons  $\lambda_j, j=2, \dots, k$  (ex :  $\lambda_2 = -1$ )

Puisque B est diagonale et que  $\|x\| > \rho$ ,

$$\frac{1}{2} U_f(x) = \lambda_1 (x^1 - \theta^1) (\rho x^1 / \|x\| - \theta^1) + \sum_{j=2}^k \lambda_j (x^j - \theta^j) (\rho x^j / \|x\| - \theta^j).$$

Si le second terme est positif ou nul (impossible si  $k=2$ ),  $U_f(x) > 0$  ;

et, s'il est négatif, il suffit de choisir :

$$\lambda_1 < - \frac{\sum_{j=2}^k \lambda_j (x^j - \theta^j) (\rho x^j / \|x\| - \theta^j)}{(x^1 - \theta^1) (\rho x^1 / \|x\| - \theta^1)} < 0 \quad (\text{ex : } \lambda_1 < -24)$$

pour avoir encore  $U_f(x) > 0$ .

□

Si  $k=1$ ,  $|\theta| < \rho$ , il est aisé de prouver la convergence ; mais c'est un résultat antérieur à l'article de Blum - voir [1] - .

## 4 - APPLICATION A L'OPTIMISATION D'UNE REGRESSION DE TYPE PARABOLOIDE

Derman [6] établit le premier la convergence du processus de Kiefer-Wolfowitz vers le point où une regression parabolique  $M(x)$  atteint son maximum. Outre que le résultat n'est qu'unidimensionnel des hypothèses restrictives sont exigées et la convergence n'a lieu qu'en probabilité. Par application du théorème 3 de Blum, antérieur aux travaux de Derman, on obtient un résultat infiniment plus fort ; c'est notre :

Corollaire 3

Si (2.8) est satisfaite et si :

$$(4.1) \quad M(x) \equiv K - K' \|A(x-\theta)\|^2$$

où  $K \in \mathbb{R}$ ,  $K' > 0$ ,  $A$  est une matrice d'ordre  $(k, k)$  régulière.

Alors les suites  $\{X_n\}$  définies par (1.2) et (1.3) convergent presque sûrement vers  $\theta$ .

Pour la démonstration, nous avons besoin du :

Lemme.

Si  $A$  est une matrice régulière d'ordre  $(k, k)$  :

$$\forall \varepsilon > 0 \quad \exists \quad q(A, \varepsilon) > 0 : \|x\| \geq \varepsilon \implies \|Ax\| \geq q(A, \varepsilon).$$

Démonstration :

S'il n'en était pas ainsi, il existerait  $\varepsilon > 0$  tel que pour tout  $q > 0$  il existerait  $x_q$  vérifiant :  $\|x_q\| \geq \varepsilon$  et  $\|Ax_q\| < q$ .

L'application  $x \mapsto \|Ax\|$  est continue. Sur le compact défini par  $\|x\| = \varepsilon$  elle atteint sa borne inférieure. Alors, puisque  $A$  est régulière, il existe  $R$  tel que :

$$\|x\| = \varepsilon \implies \|Ax\| \geq R > 0.$$

Choisissons  $q < R$ .

$$\text{Alors : } q > \|Ax_q\| = \left[ \|x_q\| / \varepsilon \right] \|A \left[ \varepsilon x_q / \|x_q\| \right]\| \geq \left[ \|x_q\| / \varepsilon \right] R \geq R$$

D'où la contradiction.

Démonstration du corollaire :

Appliquant le théorème 3., seules (2.9) et (2.10) sont à vérifier.

(2.10) est triviale.

Utilisant le lemme :

$\forall \varepsilon > 0, \|x-\theta\| \geq \varepsilon \implies M(x) - M(\theta) \leq -K' q^2(A, \varepsilon)$  ; la première partie de (2.9) est vérifiée.

Un calcul simple montre que :

$$\nabla M(x) = -2 K' A' A(x-\theta) \text{ où } A' \text{ est la transposée de } A.$$

Alors, appliquant le lemme :

$$\forall \varepsilon > 0, \|x-\theta\| \geq \varepsilon \implies \|Ax-A\theta\| \geq q(A, \varepsilon) \implies \|\nabla M(x)\| \geq 2K' q[A', q(A, \varepsilon)].$$

La seconde partie de (2.9) est prouvée et la démonstration est complète.  $\square$

## 5 - UN EXEMPLE NUMERIQUE

Nous proposons un exemple numérique, destiné à illustrer la convergence du processus de Kiefer-Wolfowitz. Les résultats ont été obtenus par simulation. Pour des exemples relevant de cas concrets, on pourra consulter par exemple [14].

Nous nous plaçons dans  $\mathbb{R}^{10}$ , de point courant  $x = (x^1, \dots, x^{10})$ . La régression est de la

forme  $M(x) = K - \sum_{i=1}^{10} (a_i x^i + b_i)^2$  avec :

$K \in \mathbb{R}$ ,  $\forall i=1, \dots, 10$   $a_i \in \mathbb{R}_+^*$ ,  $b_i \in \mathbb{R}$ .

Cette régression admet donc un maximum unique pour  $\theta = (\theta^1, \dots, \theta^{10})$  tel que :

$$\theta^i = -b_i/a_i, \quad i=1, \dots, 10.$$

On suppose que  $\forall x \in \mathbb{R}^{10}$ ,  $Y_x$  suit une loi uniforme sur  $]M(x)-a, M(x)+a[$  ( $a > 0$ ) ;

pour  $x$  donné, une observation de  $Y_x$  est obtenue par génération d'un "nombre au hasard" sur  $]0, 1[$  et transformation linéaire affine de ce nombre.

On vérifie aisément que les hypothèses du théorème 3 de Blum, sont satisfaites.

Les suites  $\{a_n\}$  et  $\{c_n\}$  intervenant dans les processus, sont les suites classiques, de type  $1/n$  et  $1/n^{1/3}$  ; toutefois, nous les "initialisons" à des valeurs correspondant à  $n$  égal à un nombre entier positif donné ; de sorte que nous prenons :

$$a_n = \frac{1}{v+n} ; \quad c_n = \frac{1}{(v+n)^{1/3}} \quad (v \in \mathbb{N} \text{ donné}).$$

Pour la régression, les valeurs numériques choisies sont :

$$K=5, \quad a=10$$

$a_i, i=1, \dots, 10$	2	1	5	8	3	7	4	9	6	5
$b_i, i=1, \dots, 10$	35	0	-4	-17	11	-83	-25	50	10	2
$\theta^i, i=1, \dots, 10$	-17,5	0,00	0,80	2,12	-3,66	11,85	6,25	-5,55	-1,66	-0,40

Nous donnons deux exemples, pour deux valeurs de  $v$  particulières. Par chacun d'eux, le vecteur  $X_1$  de début est l'origine ; les processus définis en (1.2) et (1.3) sont respectivement notés P2 et P3.

On obtient alors :

Cas  $v = 0$

n	Processus	$x_n^i, i=1, \dots, 10$									
		1	P3	0.	0.	0.	0.	0.	0.	0.	0.
9	P3	$3 \times 10^2$	-9	-90	$-1 \times 10^7$	$-3 \times 10^3$	$-3 \times 10^7$	$2 \times 10^4$	$1 \times 10^8$	$-1 \times 10^5$	$-1 \times 10^4$
3000	P3	$3 \times 10^2$	-9	-90	$-1 \times 10^7$	$-3 \times 10^3$	$-3 \times 10^7$	$2 \times 10^4$	$1 \times 10^8$	$-1 \times 10^5$	$-1 \times 10^4$

Cas v = 1000

n	Processus	$x_n^i, i=1, \dots, 10$									
1	P2	0.	0.	0.	0.	0.	0.	0.	0.	0.	0.
	P3	0.	0.	0.	0.	0.	0.	0.	0.	0.	0.
9	P2	-2,33	0,07	0,59	1,83	-1,15	9,69	2,60	-5,13	-0,97	-0,21
	P3	-1,54	-0,40	-0,22	1,36	-0,70	6,34	1,46	-4,27	-0,82	-0,43
250	P2	-16,96	0,40	0,49	2,11	-4,09	11,70	6,17	-5,63	-1,52	-0,52
	P3	-14,12	-0,41	0,81	2,20	-3,58	11,99	6,17	-5,46	-1,55	-0,37
1000	P2	-17,79	-0,30	0,69	2,10	-3,51	11,73	6,01	-5,67	-1,81	-0,21
	P3	-16,98	-0,09	0,46	2,04	-4,08	11,76	5,84	-5,51	-1,68	-0,56
4000	P2	-17,17	-1,27	0,69	2,08	-3,62	11,94	6,11	-5,71	-1,41	-0,43
	P3	-17,44	0,50	0,60	2,02	-3,96	11,97	6,03	-5,69	-1,61	-0,51

Dans le deuxième cas la convergence apparaît clairement, alors qu'elle est pratiquement inexistante dans le premier cas. Ceci montre l'importance d'un choix convenable des suites  $\{a_n\}$  et  $\{c_n\}$ . Il est certain qu'un choix entraînant un "éloignement prononcé" dès les premières itérations, peut rendre la convergence trop lente - voire, l'interdire en pratique, si les erreurs d'arrondis compensent l'évolution - ; notre premier cas en est une illustration. Pour un étude de la rapidité de convergence liée au choix de la suite  $\{a_n\}$  dans le cas de l'approximation, on pourra consulter, par exemple [4].

Si dans notre deuxième cas la convergence des processus est manifeste, nous devons toutefois convenir que cette convergence est assez lente. Notons enfin que les deux processus donnent des résultats à peu près semblables ; et n'oublions pas que le processus P3 ne nécessite que  $(k+1)$  observations à chaque itération, alors que le processus P2 en nécessite  $2k$ .

## BIBLIOGRAPHIE

- 1 - J.R. BLUM - *Approximation methods which converge with probability one.*  
Ann. Math. Stat. Vol.25 (1954) pp.382-386.
- 2 - J.R. BLUM - *Multidimensional stochastic approximation methods.*  
Ann. Math. Stat. Vol.25 (1954) pp.737-744.
- 3 - J.P. BERTRAN - Thèse d'Etat (1975) Nancy.
- 4 - CHUNG - *On a stochastic approximation method.*  
Ann. Math. Stat. Vol.25 (1954) pp.463-483.
- 5 - P. DAUBEZE - Thèse de 3ème Cycle (1974) Toulouse.
- 6 - C. DERMAN - *An application of Chung's lemma to the Kiefer-Wolfowitz stochastic approximation procedure.*  
Ann. Math. Stat. (1956) pp.532-536.

- 7 - C.DERMAN and J.SACKS - *On Dvovresky's stochastic approximation theorem.*  
Ann. Math. Stat. Vol.30 (1959) pp.601-606.
- 8 - V.DUPAC et V.V.IVANOV - *Aplikace Matematiky.*  
Svazek 22 (1977) pp.134-146.
- 9 - A.DVORESKY - *On stochastic approximation*  
Proceedings of the Third Berkeley Symposium on Mathematical Statistics  
and Probability. Vol.1, pp.39-56, University of California Press (1956).
- 10 - V.FABIAN - *Stochastic approximation methods.*  
Czechoslovak mathematical journal (1960) pp.123-159.
- 11 - S.FRIEDMAN - *On stochastic approximations.*  
Ann. Math. Stat. (1963) pp.343-346.
- 12 - J.B.HIRIART-URRUTY - *Contributions à la programmation mathématique : cas  
déterministe et stochastique.*  
Thèse d'Etat (1977) Clermont II.
- 13 - C.MACCHI - Thèse d'Etat (1972) Paris VI.
- 14 - J.M.MONNEZ - *Approximation stochastique : Le processus de Robbins-Monro. Mise à  
jour des résultats et quelques compléments.*  
Thèse de 3ème cycle (1975) Nancy I.
- 15 - L.SCHMETTERER - *L'approximation stochastique.*  
Cours de D.E.A. 1ère édition (1968), 2ème édition (1972) Clermont.
- 16 - J.H.VENTER - *On Dvovresky stochastic approximation theorems.*  
Ann. Math. Stat. Vol.37 (1966) pp.1534-1544.
- 17 - M.T.WASAN - *Stochastic approximation*  
Cambridge University Press (1969).
- 18 - J.WOLFOWITZ - *On the stochastic approximation method of Robbins and Monro.*  
Ann. Math. Stat. Vol.23 (1952) pp.457-461.