

## ON CONVERGENCE OF THE KIEFER-WOLFOWITZ APPROXIMATION PROCEDURE

BY J. H. VENTER

*Potchefstroom University*

**1. Summary.** It is shown that the Kiefer-Wolfowitz procedure for estimating the maximum of a regression surface converges almost surely to the maximum if this is the only stationary point of the surface and some other conditions are satisfied. The result is in a sense stronger than those existing presently.

**2. Introduction and preliminaries.** Let  $R^p$  denote  $p$ -dimensional Euclidian space; the elements of  $R^p$  will be thought of as column vectors;  $u_i$  is the vector with  $j$ th component  $\delta_{ij}$ , the Kronecker delta;  $(\cdot, \cdot)$  and  $\|\cdot\|$  denote the usual inner product and norm. Let  $f$  be a real valued function on  $R^p$ . Whenever it exists, we write  $D(x) = \sum u_i f_i(x)$  where  $f_i(x)$  is the derivative of  $f$  in the direction specified by  $u_i$  at  $x$ ; we write  $H(x)$  for the matrix with  $i, j$ th element  $f_{ij}(x)$ , the derivative of  $f_i(x)$  in the direction specified by  $u_j$ .

Suppose  $f$  achieves its maximum at  $\theta \in R^p$  and that for each  $x \in R^p$  random variables  $Y(x) = f(x) + Z(x)$  with  $EZ(x) = 0$  can be observed. The problem is to estimate  $\theta$  on the basis of the observed  $Y(x)$ 's at a number of suitably chosen  $x$ 's. The Kiefer-Wolfowitz (KW) procedure [1], [2], [4], [5], [6], to this effect is as follows: Let  $\{a_n\}$  and  $\{c_n\}$  be positive sequences, let  $X_1$  be arbitrary and supposing that after the  $(n - 1)$ st step we have an estimate  $X_n$  of  $\theta$ , observe the random variables  $Y(X_n \pm c_n u_i)$ ,  $i = 1, \dots, p$ ; put

$$(1) \quad \Delta_n = \sum_{i=1}^p u_i (2c_n)^{-1} [Y(X_n + c_n u_i) - Y(X_n - c_n u_i)]$$

and take

$$(2) \quad X_{n+1} = X_n + a_n \Delta_n$$

as the next estimate of  $\theta$ .

Convergence of  $\{X_n\}$  to  $\theta$  in the almost sure (a.s.) sense has been discussed in the references above and elsewhere. Throughout this paper it will be assumed that  $f$  has bounded second order partial derivatives. Under this assumption the key condition on  $f$  needed by these authors (see e.g. [5], p. 401) essentially amounts to

$$(3) \quad (x - \theta, D(x)) < 0 \quad \text{for all } x \neq \theta.$$

In order to see the geometric implication of this condition, consider the derivative of  $f$  along the line through  $\theta$  in the direction specified by the unit vector  $\lambda = \sum \lambda_i u_i$ . The equation of the line is  $x - \theta = r\lambda$  and the derivative becomes  $df(\theta + r\lambda)/dr = \sum \lambda_i f_i(\theta + r\lambda)$ . Hence  $r df(\theta + r\lambda)/dr = \sum r \lambda_i f_i(\theta + r\lambda) =$

Received 29 September 1966.

$(r\lambda, D(\theta + r\lambda))$  and it follows that (3) is equivalent to the condition:

(4) for each  $r \neq 0$  and for each unit vector  $\lambda$   $r \frac{df(\theta + r\lambda)}{dr} < 0$ ,

i.e., along each line through  $\theta$ ,  $f$  is unimodal with maximum at  $\theta$ . We will refer to this condition as linear unimodality.

Since the KW procedure may be thought of as an estimated steepest ascent procedure, one would expect that convergence would occur under conditions weaker than linear unimodality. Calling  $x$  stationary if  $D(x) = 0$ , we show below that if  $\theta$  is the only stationary point of  $f$  and some other technical conditions hold, then  $X_n$  converges to  $\theta$  a.s.

**3. Main result.** By  $\theta$  being a local maximum of  $f$  we will mean that for some  $\epsilon > 0$  we have  $f(x) < f(\theta)$  for all  $x \neq \theta$  in the set  $\{x: \|x - \theta\| < \epsilon\}$ .  $K_0, K_1, K_2, \dots$ , will denote finite positive constants chosen to suit the context in which they appear.

**THEOREM.** *Let the following conditions hold:*

- (a)  $f$  has bounded second order derivatives over  $R^p$ ;
- (b)  $f$  is bounded above over  $R^p$ ;
- (c)  $\theta$  is a local maximum of  $f$ ;
- (d) for each  $x \in R^p, D(x) \neq 0$  if  $x \neq \theta$ ;
- (e)  $E \|Z(x)\|^2 < K_0$  for all  $x \in R^p$ ;
- (f)  $\sum a_n = \infty, \sum a_n c_n < \infty, \sum (a_n/c_n)^2 < \infty, c_n \rightarrow 0$  and  $a_n < K_1 c_n^2$  for all  $n$ .

Let  $\{X_n\}$  be defined by (1) and (2). Then, with probability one, either  $X_n \rightarrow \theta$  or  $\|X_n\| \rightarrow \infty$ .

**PROOF.** Without loss of generality we may take  $\theta = 0$ . By Taylor's theorem

$$(5) \quad (2c_n)^{-1}\{f(X_n + c_n u_i) - f(X_n - c_n u_i)\} \\ = f_i(X_n) + \frac{1}{4}c_n\{f_{ii}(X_n + \phi'_{in}c_n u_i) - f_{ii}(X_n - \phi''_{in}c_n u_i)\}$$

where  $0 \leq \phi'_{in}, \phi''_{in} \leq 1$ . From (1)

$$(6) \quad \Delta_n = D(X_n) + L_n + U_n$$

where

$$(7) \quad L_n = \frac{1}{2}c_n \sum_i u_i \{f_{ii}(X_n + \phi'_{in}c_n u_i) - f_{ii}(X_n - \phi''_{in}c_n u_i)\},$$

$$(8) \quad U_n = (2c_n)^{-1} \sum_i u_i \{Z(X_n + c_n u_i) - Z(X_n - c_n u_i)\}.$$

It follows that

$$(9) \quad E\{U_n | X_n\} = 0 \text{ a.s.};$$

$$(10) \quad E\|U_n\|^2 \leq c_n^{-2}K_2;$$

$$(11) \quad \|L_n\| \leq c_n K_3;$$

where we have used conditions (a) and (e). Again from Taylor's theorem and (2)

$$(12) \quad f(X_{n+1}) = f(X_n) + a_n(\Delta_n, D(X_n)) + \frac{1}{2}a_n^2(\Delta_n, H_n\Delta_n)$$

where  $H_n \equiv H(X_n + \phi_n a_n \Delta_n)$  with  $0 \leq \phi_n \leq 1$ . By Schwarz's inequality, condition (a), (6) and (11),

$$(13) \quad |(\Delta_n, H_n\Delta_n)| \leq \|\Delta_n\|^2 K_4 \leq 4\{\|D(X_n)\|^2 + c_n^2 K_3^2 + \|U_n\|^2\} K_4.$$

Similarly

$$(14) \quad |(L_n, D(X_n))| \leq \|L_n\| \|D(X_n)\| \leq c_n K_3 \{1 + \|D(X_n)\|^2\}.$$

Substituting (6) into (12) and applying the bounds derived above together with condition (f), it follows that

$$(15) \quad f(X_{n+1}) = f(X_n) + a_n A_n \|D(X_n)\|^2 + a_n(U_n, D(X_n)) + B_n$$

where

$$(16) \quad |A_n| \leq K_5, A_n \rightarrow 1 \text{ a.s. and } \sum E|B_n| < \infty.$$

From (9)  $E\{(U_n, D(X_n)) | X_n\} = 0$  a.s. Hence if  $\sum a_n^2 E\{(U_n, D(X_n))^2 | X_n\} < \infty$  for some sample sequence  $\{X_n\}$  then according to Lemma 10 of [3] we may take it that  $\sum a_n(U_n, D(X_n))$  converges. Iteration of (15) then shows that if  $\sum a_n \|D(X_n)\|^2 = \infty$ , we must have  $f(X_n) \rightarrow \infty$  which is impossible in view of condition (b). Hence we may take it that

$$(17) \quad \sum a_n \|D(X_n)\|^2 < \infty,$$

$$(18) \quad \{f(X_n)\} \text{ converges}$$

for sample sequences such that  $\sum a_n^2 E\{(U_n, D(X_n))^2 | X_n\} < \infty$ . On the other hand, if this series diverges, then according to Lemma 8 of [3] we may assume that

$$\left\{ \sum_1^n a_k(U_k, D(X_k)) \right\} / \left\{ \sum_1^n a_k^2 E\{(U_k, D(X_k))^2 | X_k\} \right\} \rightarrow 0.$$

Hence using Schwarz's inequality, (10) and condition (f) we get

$$\left\{ \sum_1^n a_k(U_k, D(X_k)) \right\} / \left\{ \sum_1^n a_k A_k \|D(X_k)\|^2 \right\} \rightarrow 0$$

and now iteration of (15) again shows that if (17) does not hold then  $f(X_n) \rightarrow \infty$ . It follows that (17) and (18) may be taken to be true generally. Further, from (2) and (6)

$$(19) \quad X_{n+1} - X_n = a_n D(X_n) + a_n L_n + a_n U_n.$$

It follows from (9), (10), condition (f) and Lemma 10 of [3] that  $\sum a_n U_n$  converges a.s. and hence  $a_n U_n \rightarrow 0$  a.s.; from (11)  $a_n L_n \rightarrow 0$  and from (17)  $a_n D(X_n) \rightarrow 0$  a.s. Hence

$$(20) \quad X_{n+1} - X_n \rightarrow 0 \text{ a.s.}$$

Now, supposing that we have a particular sample sequence  $\{X_n\}$  for which the results above hold, we have four possibilities, viz.

$$(i) \quad 0 = \liminf \|X_n\| < \limsup \|X_n\|;$$

- (ii)  $0 < \liminf \|X_n\| \leq \limsup \|X_n\| < \infty$ ;
- (iii)  $0 < \liminf \|X_n\| < \limsup \|X_n\| = \infty$ ;
- (iv) the statement of the theorem holds for  $\{X_n\}$ .

To rule out case (i), choose  $\epsilon$  such that  $0 < \epsilon < \limsup \|X_n\|$  and such that  $f(x) < f(0)$  for all  $x \neq 0$  in the set  $\{x: \|x\| < \epsilon\}$ . Then, due to (20) there must be at least one limit point of  $\{X_n\}$  in the set  $\{x: \frac{1}{2}\epsilon < \|x\| < \epsilon\}$ , say  $x^*$ . Let  $\{X_{n_k}\}$  be a subsequence converging to  $x^*$ . Then  $f(x^*) = \lim_{k \rightarrow \infty} f(X_{n_k}) = \lim_{n \rightarrow \infty} f(X_n)$ , due to (18) and continuity of  $f$ . There is also a subsequence,  $\{X_{n_m}\}$  say, converging to 0. Hence  $f(0) = \lim_{m \rightarrow \infty} f(X_{n_m}) = \lim_{n \rightarrow \infty} f(X_n)$ , from which follows the contradiction  $f(0) = f(x^*)$ .

In case (ii) condition (d) implies that  $\liminf \|D(X_n)\| > 0$ ; this contradicts (17) and condition (f).

In order to rule out (iii) put  $d = \liminf \|X_n\|$  and let  $c > d$ . Let  $\epsilon > 0$  satisfy  $0 < d - \epsilon < d + \epsilon < c$ . Put

$$(21) \quad \alpha = \inf \{ \|D(x)\| : d - \epsilon \leq \|x\| \leq c \}.$$

By continuity of  $D(x)$  and condition (d),  $\alpha > 0$ . Let  $\{X_{n_k}\}$  be a subsequence of  $\{X_n\}$  such that  $d - \epsilon \leq \|X_{n_k}\| \leq d + \epsilon$  for all  $k$ . For each  $k$  let  $m_k$  be the smallest positive integer such that  $\|X_{n_k+m_k}\| \geq c$ ; in the present circumstances  $1 \leq m_k < \infty$ . Iterating (19), taking norms and using the various bounds established above, we get

$$\|X_{n_k+m_k}\| \leq \|X_{n_k}\| + \sum_{j=0}^{m_k-1} a_{n_k+j} \|D(X_{n_k+j})\| + \delta_{n_k}$$

where  $\delta_{n_k} \rightarrow 0$  as  $k \rightarrow \infty$ . Hence

$$(22) \quad \sum_{j=0}^{m_k-1} a_{n_k+j} \|D(X_{n_k+j})\| \geq c - d - \epsilon - \delta_{n_k}.$$

Iteration of (15) together with arguments similar to those used in establishing (17) and (18) now show that for all  $k$  large enough

$$f(X_{n_k+m_k}) \geq f(X_{n_k}) + \frac{1}{2}\alpha \sum_{j=0}^{m_k-1} a_{n_k+j} \|D(X_{n_k+j})\| + \delta'_{n_k}$$

where  $\delta'_{n_k} \rightarrow 0$  as  $k \rightarrow \infty$ . Substituting (22) we find

$$(23) \quad f(X_{n_k+m_k}) \geq f(X_{n_k}) + \frac{1}{2}\alpha(c - d - \epsilon - \delta_{n_k}) + \delta'_{n_k}$$

and letting  $k \rightarrow \infty$  the left hand side of this inequality tends to  $\lim f(X_n)$  according to (18) whereas the right hand side tends to  $\lim f(X_n) + \frac{1}{2}\alpha(c - d - \epsilon) > \lim f(X_n)$ , yielding a contradiction.

The theorem therefore follows.

**4. Discussion.** 1. The "usual" choice  $a_n = an^{-1}$ ,  $c_n = cn^{-\gamma}$  with  $a, c > 0$  and  $0 < \gamma < \frac{1}{2}$  satisfies condition (f).

2. In the one dimensional case ( $p = 1$ ) the possibility  $|X_n| \rightarrow \infty$  can also be ruled out under the conditions of the theorem, as follows. There are three cases, viz. (a)  $X_n \rightarrow -\infty$ , (b)  $\liminf X_n = -\infty$ ,  $\limsup X_n = +\infty$ , (c)  $X_n \rightarrow +\infty$ . In case (b) there must be a limit point of  $\{X_n\}$  at  $\theta$  due to (20) and the argument

of case (i) above becomes applicable. Further, the assumptions imply that  $D(x) < 0$  for  $x > \theta$ . Hence in case (c), from (19), for all  $n$  large enough  $X_{n+1} - X_n < a_n L_n + a_n U_n$  and iteration together with the convergence of  $\sum a_n L_n$  and  $\sum a_n U_n$  show that  $\{X_n\}$  is bounded above, contradicting the assumed  $X_n \rightarrow +\infty$ . Similarly for case (a).

3. It seems unlikely that the possibility  $\|X_n\| \rightarrow \infty$  can be ruled out under the conditions of the theorem in the multi-dimensional case ( $p > 1$ ). We have been able to construct an example satisfying all conditions except boundedness of second order derivatives in which  $\|X_n\| \rightarrow \infty$ .

4. A simple additional condition sufficient for ruling out  $\|X_n\| \rightarrow \infty$  is

(g)  $\lim_{t \rightarrow \infty} \inf \{ \|D(x)\| : \|x - \theta\| > t \} > 0$ , for this implies  $\liminf \|D(X_n)\| > 0$  if  $\|X_n\| \rightarrow \infty$ , which would be impossible in view of (17) and condition (f).

5. Condition (g) is not satisfied by functions such as  $f(x) = \exp \{-\|x - \theta\|^2\}$  for which one would expect convergence. Such cases are covered by the following weak linear unimodality condition.

(g') For some  $T$   $(x - \theta, D(x)) \leq 0$  for all  $x$  such that  $\|x - \theta\| > T$ .

Convergence under this additional condition is proved as follows. Suppose  $\theta = 0$  and  $\|X_n\| \rightarrow \infty$ . From (2)

$$(24) \quad \|X_{n+1}\|^2 = \|X_n\|^2 + 2a_n(X_n, \Delta_n) + a_n^2 \|\Delta_n\|^2.$$

From (6), (11), Schwarz and (g')

$$(25) \quad (X_n, \Delta_n) = (X_n, D(X_n)) + (X_n, L_n) + (X_n, U_n) \leq c_n K_3 \{1 + \|X_n\|^2\} + (X_n, U_n)$$

for all  $n$  large enough. Also from (6) and (11)

$$\|\Delta_n\|^2 \leq 4\{\|D(X_n)\|^2 + c_n^2 K_3^2 + \|U_n\|^2\}$$

and substitution into (24) shows that

$$(26) \quad \|X_{n+1}\|^2 \leq \{1 + 2a_n c_n K_3 + V_n \|X_n\|^{-1}\} \|X_n\|^2 + F_n$$

where  $V_n = 2a_n(X_n, U_n)\|X_n\|^{-1}$  and  $\sum |F_n| < \infty$  a.s. Hence  $E\{V_n | X_n\} = 0$  and  $E\{V_n^2 | X_n\} \leq 4a_n^2 c_n^{-2} K_2$  according to (9) and (10). By condition (f) and Lemma 10 of [3] it follows that  $\sum V_n$  converges and hence also  $\sum V_n \|X_n\|^{-1}$  converges. Also  $EV_n^2 \leq 4a_n^2 c_n^{-2} K_2$  and hence  $\sum V_n^2 < \infty$  which implies  $\sum V_n^2 \|X_n\|^{-2} < \infty$ . Lemma 1b of [7] is applicable and iteration of (26) shows that  $\{\|X_n\|\}$  is bounded, a contradiction.

6. The argument just given shows that condition (g) can be weakened as follows

(g'')  $\lim_{t \rightarrow \infty} \inf \{ \|x - \theta\| \|D(x)\| : \|x - \theta\| > t \} > 0$ ,

for, in this case there exists  $K_5 > 0$  such that for all  $n$  large enough ( $\theta = 0$ ),

$$(X_n, D(X_n)) \leq \|X_n\| \|D(X_n)\| \leq K_5 \|X_n\|^2 \|D(X_n)\|^2$$

if  $\|X_n\| \rightarrow \infty$ , and using this in (25) an alternative to (26) with  $2K_5 a_n \|D(X_n)\|^2$  added to the coefficient of  $\|X_n\|^2$  is obtained. (17) ensures that the rest of the argument goes through.

## REFERENCES

- [1] DÛPAC, V. (1957). On the Kiefer-Wolfowitz approximation method. *Sel. Transl. Math. Statist. Prob.* **4** 43-69.
- [2] DVORETZKY, A. (1956). On stochastic approximation. *Proc. Third Berkeley Symp. Math. Statist. Prob.* **1** 39-55. Univ. of California Press.
- [3] DUBINS, L. E. and FREEDMAN, D. A. (1965). A sharper form of the Borel-Cantelli lemma and the strong law. *Ann. Math. Statist.* **36** 800-807.
- [4] KIEFER, J. and WOLFOWITZ, J. (1952). Stochastic estimation of the maximum of a regression function. *Ann. Math. Statist.* **23** 462-466.
- [5] SACKS, J. (1958). Asymptotic distributions of stochastic approximation procedures. *Ann. Math. Statist.* **29** 373-405.
- [6] SCHMETTERER, L. (1960). Stochastic approximation. *Proc. Fourth Berkeley Symp. Math. Statist. Prob.* **1** 587-609. Univ. of California Press.
- [7] VENTER, J. H. (1966). On Dvoretzky stochastic approximation theorems. *Ann. Math. Statist.* **37** 1534-1544.