

PROCEEDINGS *of the* THIRD  
BERKELEY SYMPOSIUM ON  
MATHEMATICAL STATISTICS  
AND PROBABILITY

*Held at the Statistical Laboratory  
University of California  
December, 1954  
July and August, 1955*

VOLUME I

CONTRIBUTIONS TO THE THEORY OF STATISTICS

EDITED BY JERZY NEYMAN

For further effective support of the Symposium thanks must be given the National Science Foundation, the United States Air Force Research and Development Command, the United States Army Office of Ordnance Research, and the United States Navy Office of Naval Research.

UNIVERSITY OF CALIFORNIA PRESS  
BERKELEY AND LOS ANGELES  
1956

# ON STOCHASTIC APPROXIMATION

ARYEH DVORETZKY

HEBREW UNIVERSITY OF JERUSALEM AND COLUMBIA UNIVERSITY

## 1. Introduction

Stochastic approximation is concerned with schemes converging to some sought value when, due to the stochastic nature of the problem, the observations involve errors. The interesting schemes are those which are self-correcting, that is, in which a mistake always tends to be wiped out in the limit, and in which the convergence to the desired value is of some specified nature, for example, it is mean-square convergence. The typical example of such a scheme is the original one of Robbins-Monro [7] for approximating, under suitable conditions, the point where a regression function assumes a given value. Robbins and Monro have proved mean-square convergence to the root; Wolfowitz [8] showed that under weaker assumptions there is still convergence in probability to the root; and Blum [1] demonstrated that, under still weaker assumptions, there is not only convergence in probability but even convergence with probability 1. Kiefer and Wolfowitz [6] have devised a method for approximating the point where the maximum of a regression function occurs. They proved that under suitable conditions there is convergence in probability and Blum [1] has weakened somewhat the conditions and strengthened the conclusion to convergence with probability 1.

The two schemes mentioned above are rather specific. We shall deal with a vastly more general situation. The underlying idea is to think of the random element as noise superimposed on a convergent deterministic scheme. The Robbins-Monro and Kiefer-Wolfowitz procedures, under conditions weaker than any previously considered, are included as very special cases and, despite this generality, the conclusion is stronger since our results assert that the convergence is both in mean-square and with probability 1.

The main results are stated in section 2 and their proof follows in sections 3 and 4. Various generalizations are given in section 5, while section 6 furnishes an instructive counterexample. The Robbins-Monro and Kiefer-Wolfowitz procedures are treated in section 7. Because of the generality of our results the proofs in sections 3 and 4 have to overcome a number of technical difficulties and are somewhat involved. A special case of considerable scope where the technical difficulties disappear is discussed in section 8. This section is essentially self-contained and includes an extremely simple complete proof of the mean-square convergence result in the special case, which illustrates the underlying idea of our method. In section 8 we also find the best (unique minimax in a non-asymptotic sense) way of choosing the  $a_n$  in a special case of the Robbins-Monro scheme [they are of the form  $c/(n + c')$ ]. The concluding section 9 contains some remarks on extensions to nonreal random variables and other topics. Since the primary object of this paper is to give the general approach, no attempt has been made to study any specific procedures except the well-known Robbins-Monro and Kiefer-Wolfowitz schemes which serve as illustrations.

Research sponsored in part by the Office of Scientific Research of the Air Force under contract AF 18 (600)-442, Project R-345-20-7.

## 2. Statement of the main results

Let  $(\Omega = \{\omega\}, \mathcal{A}, \mu)$  be a probability space.  $X = X(\omega)$ ,  $Y = Y(\omega)$  and  $Z = Z(\omega)$ , as well as the same letters with primes or subscripts or both, will denote (real) random variables, and the corresponding lower-case letters will denote values assumed by the random variables.  $T_n, T'_n$  and  $T''_n$ ,  $n = 1, 2, \dots$ , will denote measurable transformations from  $n$ -dimensional real space into the reals. Instead of writing  $T_n(r_1, \dots, r_n)$  we shall often write  $T_n(r_n)$  exhibiting only the last argument.  $E\{ \}$  and  $P\{ \}$  will denote the expected value of the random variable and the probability of the event within the braces, respectively.

It is difficult to strike the proper balance between generality of result and simplicity of statement. We shall first state only a moderately general version of our results and follow it by an extension. Further generalizations will be given in section 5.

**THEOREM.** *Let  $a_n, \beta_n$  and  $\gamma_n$ ,  $n = 1, 2, \dots$ , be nonnegative real numbers satisfying*

$$(2.1) \quad \lim_{n \rightarrow \infty} a_n = 0,$$

$$(2.2) \quad \sum_{n=1}^{\infty} \beta_n < \infty$$

and

$$(2.3) \quad \sum_{n=1}^{\infty} \gamma_n = \infty.$$

Let  $\theta$  be a real number and  $T_n$ ,  $n = 1, 2, \dots$ , be measurable transformations satisfying

$$(2.4) \quad |T_n(r_1, \dots, r_n) - \theta| \leq \max [a_n, (1 + \beta_n) |r_n - \theta| - \gamma_n]$$

for all real  $r_1, \dots, r_n$ . Let  $X_n$  and  $Y_n$ ,  $n = 1, 2, \dots$ , be random variables and define

$$(2.5) \quad X_{n+1}(\omega) = T_n[X_1(\omega), \dots, X_n(\omega)] + Y_n(\omega)$$

for  $n \geq 1$ .

Then the conditions  $E\{X_1^2\} < \infty$ ,

$$(2.6) \quad \sum_{n=1}^{\infty} E\{Y_n^2\} < \infty$$

and

$$(2.7) \quad E\{Y_n | x_1, \dots, x_n\} = 0$$

with probability 1 for all  $n$ , imply

$$(2.8) \quad \lim_{n \rightarrow \infty} E\{(X_n - \theta)^2\} = 0$$

and

$$(2.9) \quad P\{\lim_{n \rightarrow \infty} X_n = \theta\} = 1.$$

The main difficulty is in proving (2.8); once this is done (2.9) follows by a simple device. In the theorem  $a_n, \beta_n$  and the restoring effect  $\gamma_n$  are assumed independent of the observations  $x_1, \dots, x_n$ . This need not be so and the following statement dispenses with this assumption.

**EXTENSION.** *The theorem remains valid if  $a_n, \beta_n$  and  $\gamma_n$  in (2.4) are replaced by non-*

negative functions  $a_n(r_1, \dots, r_n)$ ,  $\beta_n(r_1, \dots, r_n)$  and  $\gamma_n(r_1, \dots, r_n)$ , respectively, provided they satisfy the conditions: The functions  $a_n(r_1, \dots, r_n)$  are uniformly bounded and

$$(2.10) \quad \lim_{n \rightarrow \infty} a_n(r_1, \dots, r_n) = 0$$

uniformly for all sequences  $r_1, \dots, r_n, \dots$ ; the functions  $\beta_n(r_1, \dots, r_n)$  are measurable and

$$(2.11) \quad \sum_{n=1}^{\infty} \beta_n(r_1, \dots, r_n)$$

is uniformly bounded and uniformly convergent for all sequences  $r_1, \dots, r_n, \dots$ ; and the functions  $\gamma_n(r_1, \dots, r_n)$  satisfy

$$(2.12) \quad \sum_{n=1}^{\infty} \gamma_n(r_1, \dots, r_n) = \infty$$

uniformly for all sequences  $r_1, \dots, r_n, \dots$ , for which

$$(2.13) \quad \sup_{n=1, 2, \dots} |r_n| < L$$

$L$  being an arbitrary finite number.

We shall refer to the theorem and its extension together as the extended theorem. Condition (2.13) was introduced for the functions  $\gamma_n$  because of its use in applications (see section 7). Further generalizations will be given following the proof.

### 3. Proof of the theorem

Throughout the proofs, in this and the following section, we assume  $\theta = 0$ . This involves no loss of generality. Let  $m$  be any positive integer and  $A$  any set in  $\mathcal{Q}$  the definition of which can be made in terms of  $X_1, \dots, X_m$ . We shall first show that if  $a$  is any number satisfying  $a \geq a_m$  then

$$(3.1) \quad \int_A [ (|X_{m+1}| - a)^+ ]^2 d\mu \\ \leq \int_A \{ a\beta_m (1 + a\beta_m) + Y_m^2 + (1 + \beta_m)^2 (1 + a\beta_m) [ (|X_m| - a)^+ ]^2 \} d\mu$$

where, as customary,  $r^+ = (|r| + r)/2$  denotes the positive part of  $r$ .

If  $Z = Z(\omega)$  is any random variable satisfying  $|Z| \leq a$  for all  $\omega$  then clearly

$$(3.2) \quad \int_A [ (|X_{m+1}| - a)^+ ]^2 d\mu \leq \int_A (X_{m+1} - Z)^2 d\mu \\ = \int_A [T_m(X_m) - Z + Y_m]^2 d\mu.$$

If  $Z$  is defined in terms of  $X_1, \dots, X_m$  it then follows from (2.7) that

$$(3.3) \quad \int_A [T_m(X_m) - Z + Y_m]^2 d\mu = \int_A [T_m(X_m) - Z]^2 d\mu + \int_A Y_m^2 d\mu.$$

Taking  $Z = T_m(X_m)$  if  $|T_m(X_m)| \leq a$ ,  $Z = a$  if  $T_m(X_m) > a$  and  $Z = -a$  if  $T_m(X_m) < -a$  we have

$$(3.4) \quad \int_A [T_m(X_m) - Z]^2 d\mu = \int_A [ (|T_m(X_m)| - a)^+ ]^2 d\mu$$

Since  $a \geq a_m$  we note that, by (2.4), we have  $|T_m(r_m)| - a \leq (1 + \beta_m)a - a = a\beta_m$  whenever  $|r_m| \leq a$ , while otherwise  $|T_m(r_m)| - a \leq (1 + \beta_m)|r_m| - a = (1 + \beta_m)(|r_m| - a) + a\beta_m$ . Thus we have in all cases

$$(3.5) \quad 0 \leq (|T_m(r_m)| - a)^+ \leq (1 + \beta_m)(|r_m| - a)^+ + a\beta_m.$$

Using the inequality  $(u + v)^2 \leq (1 + v)u^2 + v(1 + v)$  which is valid for any  $v \geq 0$  we obtain

$$(3.6) \quad [ (|T_m(r_m)| - a)^+ ]^2 \leq (1 + a\beta_m)(1 + \beta_m)^2 [ (|r_m| - a)^+ ]^2 + a\beta_m(1 + a\beta_m).$$

Combining (3.2), (3.3), (3.4) and (3.6) we obtain (3.1).

Let now  $n \geq m$  and assume that  $a \geq \max_{m \leq j \leq n} a_j$ . Then iterating (3.1) gives immediately

$$(3.7) \quad \int_{\Delta} [ (|X_n| - a)^+ ]^2 d\mu \leq u_{m,n} \int_{\Delta} \left\{ v_{m,n} + \sum_{j=m}^{n-1} Y_j^2 + [ (|X_m| - a)^+ ]^2 \right\} d\mu$$

where

$$(3.8) \quad u_{m,n} = \prod_{j=m}^{n-1} (1 + \beta_j)^2 (1 + a\beta_j), \quad v_{m,n} = a \sum_{j=m}^{n-1} \beta_j (1 + a\beta_j)$$

[when  $n = m$  both sides of (3.7) are identical provided void sums and products are interpreted as 0 and 1, respectively].

Let  $\epsilon > 0$  be given. Choose  $a = a(\epsilon)$  so that

$$(3.9) \quad 0 < a \leq \sqrt{\frac{\epsilon}{8}}.$$

Then choose an integer  $k = k(a, \epsilon) > 1$  so that

$$(3.10) \quad \max_{j \geq k-1} a_j \leq a,$$

and

$$(3.11) \quad u_k \left( v_k + \sum_{j=k-1}^{\infty} E\{Y_j^2\} \right) \leq \frac{\epsilon}{8}$$

where

$$(3.12) \quad u_k = \prod_{j=k}^{\infty} (1 + \beta_j)^2 (1 + a\beta_j), \quad v_k = a \sum_{j=k}^{\infty} \beta_j (1 + a\beta_j);$$

such an integer  $k$  exists in virtue of (2.1) and the fact that, by (2.2) and (2.6), all infinite series and products involved are convergent.

For every  $j$  and  $\omega$  put

$$(3.13) \quad s_j = s_j(\omega) = \text{sgn } T_j(X_j)$$

where  $\text{sgn } r$  denotes, as usual, 1 when  $r > 0$ , -1 when  $r < 0$  and 0 when  $r = 0$ . For  $j > 1$  let  $B'_j$  and  $B''_j$  denote the events described by

$$(3.14) \quad B'_j = \{ \omega : \text{sgn } X_j \neq s_{j-1} \}$$

and

$$(3.15) \quad B''_j = \{ \omega : |T_{j-1}(X_{j-1})| \leq a \}.$$

Put  $B_j = B'_j \cup B''_j$  and, for  $m \geq k$ ,

$$(3.16) \quad A_m = B_m - \bigcup_{j=k}^{m-1} B_j.$$

Finally, let

$$(3.17) \quad \Gamma_n = \bigcup_{m=k}^n A_m, \quad \Delta_n = \Omega - \Gamma_n$$

for every  $n \geq k$ .

From (2.5) and (3.14) it follows that  $|X_m| \leq |Y_{m-1}|$  throughout  $B'_m$ , while from (2.5) and (3.15) we have  $|X_m| \leq a + |Y_{m-1}|$  in  $B''_m$ . Thus  $|X_m| - a \leq |Y_{m-1}|$  throughout  $B_m$  and, in particular, in  $A_m$ . Hence it results from (3.7), (3.8) and (3.12) that

$$(3.18) \quad \int_{A_m} [ (|X_n| - a)^+ ]^2 d\mu \leq u_k \int_{A_m} \left( v_k + \sum_{j=k-1}^{n-1} Y_j^2 \right) d\mu$$

whenever  $n \geq m \geq k$ . Since the sets  $A_m$  are disjoint, it follows from (3.10) on summing the inequalities (3.18) that

$$(3.19) \quad \int_{\Gamma_n} [ (|X_n| - a)^+ ]^2 d\mu \leq \frac{\epsilon}{8}.$$

As  $|X_n| \leq (|X_n| - a)^+ + a$ , it follows at once from (3.19) and (3.9) that

$$(3.20) \quad \int_{\Gamma_n} X_n^2 d\mu \leq 2 \left( \frac{\epsilon}{8} + \int_{\Gamma_n} a^2 d\mu \right) \leq \frac{\epsilon}{2}$$

for every  $n \geq k$ .

Now let us turn to  $\Delta_n$ . By (3.14) we have outside  $B'_j$

$$(3.21) \quad \begin{aligned} |X_j| &= X_j \operatorname{sgn} X_j = s_{j-1} T_{j-1}(X_{j-1}) + s_{j-1} Y_{j-1} \\ &= |T_{j-1}(X_{j-1})| + s_{j-1} Y_{j-1}, \end{aligned}$$

while outside  $B''_j$  we have for  $j \geq k$  by (3.15) and (3.10)

$$(3.22) \quad |T_{j-1}(X_{j-1})| \leq (1 + \beta_{j-1}) |X_{j-1}| - \gamma_{j-1}.$$

Hence outside  $B_j$  we have

$$(3.23) \quad |X_j| \leq (1 + \beta_{j-1}) |X_{j-1}| - \gamma_{j-1} + s_{j-1} Y_{j-1}$$

whenever  $j \geq k$ . Since  $\Delta_n$  is contained in the complement of  $B_j$  for every  $k < j \leq n$ , we can in  $\Delta_n$  iterate the inequalities (3.23) and obtain for  $\omega$  in  $\Delta_n$

$$(3.24) \quad |X_n| \leq w_{k,n} |X_k| - \sum_{m=k}^n w_{m,n} \gamma_m + \sum_{m=k}^n s_m w_{m,n} Y_m$$

where

$$(3.25) \quad w_{m,n} = \prod_{j=m}^{n-1} (1 + \beta_j).$$

Putting

$$(3.26) \quad Z_n = w_{k,n} |X_k| - \sum_{m=k}^n w_{m,n} \gamma_m, \quad Z'_n = \sum_{m=k}^n s_m w_{m,n} Y_m$$

we obtain from (3.24)

$$(3.27) \quad \int_{\Delta_n} X_n^2 d\mu \leq \int_{\Delta_n} [(Z_n + Z_n')^+]^2 d\mu \\ \leq \int_{\Delta_n} (Z_n^+ + |Z_n'|)^2 d\mu.$$

Hence

$$(3.28) \quad \int_{\Delta_n} X_n^2 d\mu \leq 2 \int_{\Delta_n} (Z_n^+)^2 d\mu + 2 \int_{\Delta_n} |Z_n'|^2 d\mu.$$

But by (3.26) and (2.7), since  $s_m$  is defined by  $T_m(X_m)$ ,

$$(3.29) \quad \int_{\Delta_n} |Z_n'|^2 d\mu \leq \int_{\Omega} |Z_n'|^2 d\mu = \sum_{m=k}^n w_{m,n}^2 E\{Y_m^2\}$$

and hence

$$(3.30) \quad \int_{\Delta_n} |Z_n'|^2 d\mu \leq \prod_{j=k}^{\infty} (1 + \beta_j)^2 \sum_{m=k}^{\infty} E\{Y_m^2\} \leq \frac{\epsilon}{8},$$

by (3.11) and (3.12).

Finally, we remark that if  $Z$  is any random variable with  $E\{Z^2\} < \infty$  then  $E\{[(Z - r)^+]^2\}$  tends to zero as  $r \rightarrow +\infty$ . By (3.7) with  $m = 1$  and  $n = k$  the random variable  $X_k$ , and hence also  $Z = |X_k| \prod_{j=k}^{\infty} (1 + \beta_j)$  have finite second moments. But,

by (3.26),  $Z_n^+ \leq Z - \sum_{m=k}^n \gamma_m$  and it follows from (2.3) and the remark made at the beginning of this paragraph that

$$(3.31) \quad \int_{\Delta_n} (Z_n^+)^2 d\mu \leq E\{(Z_n^+)^2\} < \frac{\epsilon}{8}$$

for all  $n > N = N(\epsilon, k)$ .

Combining (3.20), (3.27), (3.30) and (3.31) we have  $E\{X_n^2\} < \epsilon$  for  $n > N$ . Since  $\epsilon > 0$  is arbitrary this completes the proof of (2.8).

The proof of (2.9) will now be easily achieved. Applying (3.7) with  $A = \Omega$  we can obtain for all  $n > m$  an inequality of the form

$$(3.32) \quad E\{X_n^2\} < H\left(\max_{j \geq m} a_j, \sum_{j=m}^{\infty} \beta_j, E\{X_m^2\} + \sum_{j=m}^{\infty} E\{Y_j^2\}\right),$$

where  $H$  is an explicit function of the three exhibited variables, monotone increasing in each and tending to zero as all three of them tend to zero. The important thing for us is that  $H$  does not depend on  $X_j, T_j, Y_j$  except in the exhibited manner. In particular, the  $\gamma_n$  do not enter into (3.32), and it remains valid even if all of them are zero.

Given  $\delta > 0$  and  $\epsilon > 0$  there exists  $\eta = \eta(\delta^2\epsilon)$  such that if all three arguments of  $H$  are smaller than  $\eta$ , then  $H < \delta^2\epsilon$ . By (2.1), (2.2) and (2.8) there exists  $m = m(\eta)$  satisfying

$$(3.33) \quad H\left(\max_{j \geq m} a_j, \sum_{j=m}^{\infty} \beta_j, E\{X_m^2\} + \sum_{j=m}^{\infty} E\{Y_j^2\}\right) < \delta^2\epsilon.$$

Let this  $m$  be fixed and define  $X'_j, T'_j, Y'_j$  as follows:  $X'_j = X_j$  for  $j \leq m, T'_j = T_j$  and  $Y'_j = Y_j$  for  $j < m$ , while for other values of  $j$

$$(3.34) \quad T'_j(r_1, \dots, r_j) = \begin{cases} T_j(r_1, \dots, r_j) & \text{if } |r_j| < \delta \\ r_j & \text{if } |r_j| \geq \delta \end{cases}$$

and  $Y'_j, X'_{j+1}$  are defined recursively by

$$(3.35) \quad Y'_j = \begin{cases} Y_j & \text{if } |X'_j| < \delta \\ 0 & \text{if } |X'_j| \geq \delta \end{cases}$$

and

$$(3.36) \quad X'_{j+1} = T'_j(X'_j) + Y'_j.$$

Clearly the  $T'_n$  satisfy (2.4) with  $\gamma_n = 0$  and we also have  $E\{Y'_n|x'_1, \dots, x'_n\} = 0$  and  $E\{Y'^2_n\} \leq E\{Y^2_n\}$ . Since  $X'_m = X_m$  it follows from (3.32) and (3.33) that

$$(3.37) \quad E\{X'^2_n\} < \delta^2 \epsilon$$

for  $n \geq m$ . According to the definition of  $X'_n$  the relation  $|X'_j| \geq \delta$  for some  $j \geq m$  implies  $|X'_n| \geq \delta$  for all  $n \geq j$ . Hence we have for all  $n \geq m$

$$(3.38) \quad P\{\max_{m \leq j \leq n} |X'_j| \geq \delta\} \leq P\{|X'_n| \geq \delta\}.$$

Combining (3.37) and (3.38) we have

$$(3.39) \quad P\{\sup_{j \geq m} |X'_j| > \delta\} < \epsilon;$$

$\delta$  and  $\epsilon$  being arbitrary, (2.9) follows and the proof is completed.

#### 4. Proof of the extension

We first remark that (3.1) holds provided  $a \geq \sup a_m(r_1, \dots, r_m)$  for all  $r_1, \dots, r_m$  and  $\beta_m$  is considered as a function of  $X_1, \dots, X_m$ . Hence (3.8) also holds provided  $a \geq \sup a_j(r_1, \dots, r_j)$  for all  $m \leq j < n$  and all  $r_1, \dots, r_j$ , while  $u_{m,n}$  and  $v_{m,n}$  are the suprema of the expressions on the right side of the equalities (3.8) for all sequences  $r_1, \dots, r_n, \dots$ . Also, given any  $\epsilon$  and  $a$ , there exists according to the assumptions of the extension an integer  $k$  satisfying  $a \geq \sup a_j(r_1, \dots, r_j)$  for all  $j \geq k - 1$  and all  $r_1, \dots, r_j$  and (3.11) where  $u_k$  and  $v_k$  are defined as the suprema of the expressions on the right sides of the equalities (3.12). Therefore (3.20) holds.

Always considering  $\beta_j$  and  $\gamma_j$  as functions of  $X_1, \dots, X_j$  and replacing the infinite product in (3.30) by its sup for all sequences  $r_1, \dots, r_n, \dots$  we see that everything up to and including (3.30) carries through. Had we assumed (2.12) uniformly for all sequences  $r_1, \dots, r_n, \dots$ , (3.31) would have also followed as before; since only a weaker assumption was made, a slightly more sophisticated argument is needed for its proof.

We note that (3.32) remains valid provided the first two arguments are replaced by their suprema. Let  $M$  be a positive number and define  $X''_j, T''_j$  and  $Y''_j$  as follows:  $X''_1 = X_1, T''_j$  is given for all  $j \geq 1$  by (3.34) with  $\delta$  replaced by  $M$  and  $Y''_j, X''_{j+1}$  for  $j \geq 1$  given recursively by (3.35) with  $X'_j$  and  $\delta$  replaced by  $X''_j$  and  $M$ , and (3.36) with primes replaced by double primes. Then, exactly as in (3.37), we have

$$(4.1) \quad E\{X''^2_n\} < H_1$$



where

$$(4.2) \quad H_1 = H \left( \sup a_j, \sup \sum_{j=1}^{\infty} \beta_j, E\{X_1^2\} + \sum_{j=1}^{\infty} E\{Y_j^2\} \right)$$

the suprema being taken over all  $j$  and all sequences  $r_1, \dots, r_n, \dots$ . Hence, as in (3.38) and (3.39),

$$(4.3) \quad P\left\{ \sup_{j \geq 1} |X_j| > M \right\} \leq \sup_{n \geq 1} P\left\{ |X_n''| \geq M \right\} < \frac{H_1}{M^2}.$$

Thus the sequence  $X_n$  is bounded with probability 1.

Let us now return to  $Z_n$  as defined in (3.26). Putting

$$(4.4) \quad w_k = \sup \prod_{j=k}^{\infty} [1 + \beta_j(r_1, \dots, r_j)]$$

for all sequences  $r_1, \dots, r_j, \dots$  we have

$$(4.5) \quad Z_n \leq w_k |X_k| - \sum_{m=k}^n \gamma_m(X_1, \dots, X_m)$$

( $w_k$ , unlike  $w_{k,n}$  and  $w_{m,n}$ , is a positive constant). Since  $X_k$  has finite second moment there exists  $\zeta = \zeta(\epsilon, k) > 0$  such that

$$(4.6) \quad \int_{\Omega_1} X_k^2 d\mu < \frac{\epsilon}{16 w_k^2}$$

whenever  $\Omega_1$  is a measurable set satisfying  $P\{\Omega_1\} < \zeta$ . Let now

$$(4.7) \quad M = \sqrt{\frac{\zeta}{H_1}}$$

and denote by  $\Omega_1$  the set of  $\omega$  for which  $\sup_{n \geq 1} |X_n| > M$  and by  $\Omega_2$  the complementary set. Then by (4.3), (4.5), (4.6) and (4.7)

$$(4.8) \quad \int_{\Omega_1} (Z_n^+)^2 d\mu \leq w_k^2 \int_{\Omega_1} X_k^2 d\mu < \frac{\epsilon}{16}.$$

On  $\Omega_2$ , however,  $\sum \gamma_m(X_1, \dots, X_m)$  diverges uniformly by assumption and hence, by the argument leading to (3.31), we have

$$(4.9) \quad \int_{\Omega_2} (Z_n^+)^2 d\mu < \frac{\epsilon}{16},$$

for  $n > N' = N'(\epsilon, \zeta, k)$ . The inequalities (4.8) and (4.9) imply (3.31) and hence (2.8).

The proof of (2.9) then follows word for word as in the preceding section.

## 5. Generalizations

The requirements (2.4) are satisfied by most nonstochastic approximation schemes. Thus (2.4) (with  $\theta = 0$ ) is weaker than

$$(5.1) \quad |T_n(r_1, \dots, r_n)| \leq \max[a_n, (1 + \beta_n - \gamma_n)|r_n|]$$

with  $a_n, \beta_n, \gamma_n$  satisfying (2.1), (2.2) and (2.3). Indeed, if (2.3) is satisfied then there exists a sequence  $\rho_n, n = 1, 2, \dots$ , of positive numbers tending to zero and having the

property that  $\sum \gamma_n \rho_n$  is divergent. The second term under the maximum in (5.1) is always less than or equal to  $\max [(1 + \beta_n)\rho_n, (1 + \beta_n)|r_n| - \gamma_n \rho_n]$ . Thus (5.1) implies (2.4) with  $a_n$  replaced by  $\max [a_n, (1 + \beta_n)\rho_n]$  and  $\gamma_n$  replaced by  $\gamma_n \rho_n$ . Since these replacements do not affect conditions (2.1) and (2.3) it follows that (5.1) is subsumed under (2.4). Of course if we are interested in the rate of convergence it may be better to use (5.1) directly than to reduce it to (2.4). Some remarks in this direction will be found in section 8.

Similarly

$$(5.2) \quad |T_n(r_1, \dots, r_n)| \leq \max [a_n, (1 + \beta_n)|r_n| - \gamma_n + \delta_n]$$

with  $a_n, \beta_n, \gamma_n$  as before and  $\delta_n \geq 0$  satisfying

$$(5.3) \quad \sum_{n=1}^{\infty} \delta_n < \infty$$

is only deceptively more general than (2.4). To see this we remark that in view of (4.3) there exists a sequence of positive numbers  $\lambda_n, n = 1, 2, \dots$ , tending to infinity slowly enough so that  $\sum \lambda_n \delta_n$  is convergent. The second term under the max in (5.2) being always less than or equal to  $\max [(1 + \beta_n)/\lambda_n + \delta_n, (1 + \beta_n + \lambda_n \delta_n)|r_n| - \gamma_n]$  it follows that (5.2) implies (2.4) with  $a_n$  replaced by  $a'_n = \max [a_n, (1 + \beta_n)/\lambda_n + \delta_n]$  and  $\beta_n$  replaced by  $\beta'_n = \beta_n + \lambda_n \delta_n$ . Since these replacements do not affect conditions (2.1) and (2.2) our assertion is proved. Similar remarks to those made above concerning (2.4) apply also to the extension with  $\delta_n(r_1, \dots, r_n)$  satisfying the same requirements as  $\beta(r_1, \dots, r_n)$ .

The possibility of deducing our results under the assumption (5.2) has, however, an important consequence allowing the weakening of condition (2.7). This weakening may be useful in some applications, especially when dealing with certain rounding-off errors.

GENERALIZATION 1. *The Extended Theorem remains valid if (2.7) is replaced by*

$$(5.4) \quad \sum_{n=1}^{\infty} \sup_{x_1, \dots, x_n} |E\{Y_n | x_1, \dots, x_n\}| < \infty,$$

or even by the condition that

$$(5.5) \quad \sum_{n=1}^{\infty} E\{Y_n | x_1, \dots, x_n\}$$

be uniformly bounded and uniformly convergent for all sequences  $x_1, \dots, x_n, \dots$ .

Indeed, putting  $Y'_n = Y_n - E\{Y_n | x_1, \dots, x_n\}$  and  $T'_n(x_1, \dots, x_n) = T_n(x_1, \dots, x_n) + E\{Y_n | x_1, \dots, x_n\}$  we have  $X_{n+1} = T'_n(X_n) + Y'_n$ . If (5.4) holds, then the transformations  $T'_n$  satisfy (5.2) with  $\delta_n = \sup |E\{Y_n | x_1, \dots, x_n\}|$  which, by (5.4), satisfy (5.3), while the  $Y'_n$  satisfy the conditions imposed on  $Y_n$  in (2.6) and (2.7) since  $E\{Y'^2_n\} \leq 2E\{Y^2_n\} + 2\delta_n^2$ . A similar argument applies when (5.5) holds.

Another sometimes useful extension is the following.

GENERALIZATION 2. *Conclusion (2.9) of the Extended Theorem remains valid even without any restrictions on  $X_1$  and if (2.5) is replaced by*

$$(5.6) \quad X_{n+1} = T_n(X_n) + Y_n^*$$

with the random variables  $Y_n^*$ ,  $n = 1, 2, \dots$ , satisfying

$$(5.7) \quad P\{Y_n^* \neq Y_n \text{ for infinitely many } n\} = 0,$$

thus, in particular, when

$$(5.8) \quad \sum_{n=1}^{\infty} P\{Y_n^* \neq Y_n\} < \infty.$$

Indeed, if  $\Omega'$  denotes the set where  $|X_m| < M$  and  $\Omega''$  the set where  $Y_n^* = Y_n$  for all  $n \geq m$  then it follows from the Extended Theorem that  $P\{X_n \rightarrow \theta | \Omega' \cap \Omega''\} = 1$ . Since  $P\{\Omega' \cap \Omega''\}$  can be made arbitrarily close to 1 the result follows. This simple generalization may often be used in order to reduce the study to the case when the random variables are bounded.

Our proof was arranged in such a manner that it yields also the following.

GENERALIZATION 3. If (2.1) is replaced by

$$(5.9) \quad \overline{\lim}_{n \rightarrow \infty} a_n = a,$$

or, more generally, (2.10) by

$$(5.10) \quad \overline{\lim}_{n \rightarrow \infty} a_n(r_1, \dots, r_n) \leq a$$

uniformly for all sequences  $r_1, \dots, r_n, \dots$  then the Extended Theorem remains valid provided (2.8) and (2.9) are replaced by

$$(5.11) \quad \overline{\lim}_{n \rightarrow \infty} E\{(X_n - \theta)^2\} \leq a^2$$

and

$$(5.12) \quad P\{\overline{\lim}_{n \rightarrow \infty} |X_n| \leq a\} = 1.$$

Another type of generalization which is useful can be exemplified by the following.

GENERALIZATION 4. The Extended Theorem remains valid if the assumptions concerning  $a_n(r_1, \dots, r_n)$  are replaced by the following:  $a_1(X_1)$  is bounded with probability 1,  $a_n(X_1, \dots, X_n) \geq a_{n+1}(X_1, \dots, X_n, X_{n+1})$  with probability 1 and

$$(5.13) \quad P\{\lim_{n \rightarrow \infty} a_n(X_1, \dots, X_n) = 0\} = 1.$$

Indeed, denoting by  $a$  an upper bound in probability of  $a_1(X_1)$ , by  $\Omega_m$  the set where  $a_m(X_1, \dots, X_m) < \epsilon$  and by  $\Omega'_m$  its complement we have

$$(5.14) \quad E\{X_n^2\} = P\{\Omega_m\} E\{X_n^2 | \Omega_m\} + P\{\Omega'_m\} E\{X_n^2 | \Omega'\}$$

and (5.11), with  $\theta = 0$  for brevity, gives

$$(5.15) \quad \overline{\lim}_{n \rightarrow \infty} E\{X_n^2\} \leq \epsilon^2 + a^2 P\{\Omega'_m\}.$$

Since  $P\{\Omega'_m\} \rightarrow 0$  as  $m \rightarrow \infty$  by (5.13), we have (2.8); the proof of (2.9) is exactly the same.

The last generalization we wish to present extends the class of transformations  $T_n$ . Instead of considering transformations  $T_n$  determined by  $x_1, x_2, \dots, x_n$  we may consider random ones depending on the sample point  $\omega$ , that is, measurable mappings of  $R \times \Omega$  into  $R$ ,  $R$  being the real line. In this case  $x_1, \dots, x_n$  do not determine the value  $t_n$  as-

sumed by  $T_n(X_n)$ . However, except for this fact which necessitates a restatement of (2.7), nothing is changed in all our arguments. Hence we have

GENERALIZATION 5. *The Extended Theorem remains valid also if  $T_n, n = 1, 2, \dots$ , are random transformations provided (2.4) holds for all  $\omega$  and (2.7) is replaced by*

$$(5.16) \quad E\{Y_n \mid x_1, \dots, x_n, t_1, \dots, t_n\} = 0$$

with probability 1.

All the above generalizations may be used in conjunction. Many similar ones can easily be given.

### 6. A counterexample

The very generality of our results might lead us to suspect that even weak restrictions of the type of (2.4) or its generalizations on the  $T_n$  are entirely superfluous. In other words, one might be tempted to conjecture that whenever we have a sequence of transformations  $T_n(r_n) = T_n(r_1, \dots, r_n)$  of  $n$ -space into the reals having the property that for every  $m$  and  $r_1, \dots, r_m$  the sequence

$$(6.1) \quad r_{m+1} = T_m(r_m), \dots, r_{m+n+1} = T_{m+n}(r_{m+n}), \dots$$

converges to  $\theta$  then  $E\{X_1^2\} < \infty$ , (2.6) and (2.7) already imply (2.8) or (2.9). The following simple example shows that this is not the case.

Let  $q_n$  and  $v_n, n = 1, 2, \dots$ , be two sequences of positive numbers with  $q_n < 1$  and such that both series  $\sum q_n$  and  $\sum v_n^2/q_n$  are convergent; for instance,  $q_n = v_n = 1/(n^2 + 1)$ . Put  $s_n = v_1 + \dots + v_n$  and let  $T_n$  depend only on its last argument and be defined by  $T_n(r_n) = s_{n-1}$  for  $r_n = s_{n-1}$  and  $T_n(r_n) = 0$  otherwise. No matter what  $r_1, \dots, r_m$  are, all members of (5.1) from the second on are zero. Let now  $X_1, Y_1, \dots, Y_n, \dots$  be mutually independent with  $X_1 = 0$  and  $Y_n$  assuming the two values  $v_n$  and  $-(1 - q_n)v_n/q_n$  with probabilities  $1 - q_n$  and  $q_n$ , respectively. Clearly  $E\{Y_n\} = 0$  and  $E\{Y_n^2\} < 3v_n^2/q_n$ , and thus (2.6) and (2.7) are satisfied. On the other hand, the probability that  $X_{n+1} = s_n$  for every  $n$  equals the probability that  $Y_n = v_n$  for every  $n$  and, being equal to  $\prod (1 - q_n)$ , is positive. Hence not only (2.8) and (2.9) fail to hold but  $X_n$  does not even converge in probability to zero. [In this example the  $T_n$  are discontinuous; this is easily remedied. All we have to do is to define  $T_n(r_n) = s_{n-1}$  for  $r_n = s_{n-1}$ ,  $T_n(r_n) = 0$  for  $r_n \leq s_{n-2}$  or  $r_n \geq s_n$  and by linear interpolation for the remaining values of  $r_n$ .]

### 7. The Robbins-Monro and Kiefer-Wolfowitz procedures

In this section we deal with a very special case of the general theory. It will be shown that specializing the general results will, without further ado, improve the best results previously obtained for the specific procedures.

Let  $Z_u$  be a one-parameter family of random variables, the parameter space being the real line, and assume that

$$(7.1) \quad f(u) = E\{Z_u\}$$

exists for every  $u$ . The Robbins-Monro and Kiefer-Wolfowitz procedures are concerned with finding, under suitable assumptions, the location of the root  $f(u) = 0$  and of the maximum of the regression function  $f(u)$ .

The Robbins-Monro procedure is based on a sequence of positive numbers  $a_n$ ,  $n = 1, 2, \dots$ , satisfying

$$(7.2) \quad \sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty.$$

Then, starting with an arbitrary  $x_1$ , it defines recursively a sequence  $x_n$ ,  $n = 1, 2, \dots$ , by

$$(7.3) \quad x_{n+1} = x_n - a_n Z_n$$

where  $Z_n$  is an observation on the random variable  $Z_{x_n}$ .

The Kiefer-Wolfowitz procedure is based on two sequences of positive numbers  $b_n$  and  $c_n$  satisfying

$$(7.4) \quad \sum_{n=1}^{\infty} b_n = \infty, \quad \lim_{n \rightarrow \infty} c_n = 0, \quad \sum_{n=1}^{\infty} \left(\frac{b_n}{c_n}\right)^2 < \infty.$$

Then, starting with an arbitrary  $x_1$ , it defines recursively a sequence  $x_n$ ,  $n = 1, 2, \dots$ , by

$$(7.5) \quad x_{n+1} = x_n + \frac{a_n}{c_n} (z'_n - z''_n)$$

where  $z'_n$  and  $z''_n$  are observations on  $Z_{x_n+c_n}$  and  $Z_{x_n-c_n}$ , respectively.

RESULT 1. *If the  $Z_u$  have uniformly bounded variances and if the regression function  $f(u)$  is measurable and satisfies*

$$(7.6) \quad |f(u)| < A|u| + B < \infty$$

for all  $u$  and suitable  $A$  and  $B$ , and

$$(7.7) \quad \inf_{1/k < u - \theta < k} f(u) > 0, \quad \sup_{1/k < \theta - u < k} f(u) < 0$$

for all integers  $k$ ; then the Robbins-Monro sequence (7.3) converges to  $\theta$  both in mean square and with probability 1.

Indeed, an underlying probability space can be constructed in which  $x_n$  is an observation on the random variable  $X_n$ ; then  $X_1 = x_1$  and

$$(7.8) \quad X_{n+1} = X_n - a_n f(X_n) + Y_n$$

with

$$(7.9) \quad Y_n = a_n [Z_{X_n} - f(X_n)].$$

From (7.1) we have (2.7) while the assumption

$$(7.10) \quad E\{Z_u^2\} \leq \sigma^2 < \infty$$

for all  $u$  gives  $E\{Y_n^2\} \leq \sigma^2 a_n^2$  and thus, by (7.2), (2.6) holds. Assume, for simplicity of writing,  $\theta = 0$  and let  $\rho_n$ ,  $n = 1, 2, \dots$ , be a sequence of positive numbers tending to zero and for which  $\sum \rho_n a_n = \infty$ ; and let  $\eta_n$  be also a null sequence of positive numbers having the property that  $\inf_{\eta_n < |u - \theta| < 1} |f(u)| > \rho_n$ . By (7.6) we have  $|u| - a_n |f(u)| > -Ba_n$

for all  $n > n_0$ , while given any  $L > 0$  we have  $|u| - a_n |f(u)| < |u| - a_n \rho_n$  for all  $\eta_n < |u| < L$  and  $n > n_L$ . Thus the transformations  $T_n(r_1, \dots, r_n) = r_n - a_n f(r_n)$  occurring in (7.8) satisfy, for large  $n$ , condition (2.4) with  $a_n = \max(\eta_n, Ba_n)$  and  $\gamma_n = a_n \rho_n$ . Since  $\sum a_n \rho_n$  is divergent, the result follows from the Extended Theorem. (The argument

could be somewhat simplified by using generalization 3; the introduction of the sequences  $\rho_n$  and  $\eta_n$  could then have been avoided.)

*Remark 1.* Condition (7.6) is necessary in order to dampen the restoring effect of  $-a_n z_n$ . This is illustrated by the following simple example:  $a_n = 1/n$ ,  $Z_u = f(u) = u|u|$  with probability 1. Taking  $x_1 = 3$  we have  $x_2 = 3 - 3^2 = -6$ ,  $x_3 = -6 + 6^2/2 = 12$ ,  $\dots$  and it is easily verified that  $|x_n| \rightarrow \infty$ . Condition (7.6) can be shown to be the only one of its type that will eliminate this phenomenon for all sequences  $a_n$ ; for any specific sequence this condition can, of course, be somewhat relaxed. However, it should be emphasized that in practice, condition (7.6) causes no trouble. Indeed, in all practical situations one knows in advance that the root  $\theta$  lies in some finite interval  $(C_1, C_2)$ . Then, provided  $f(u)$  is bounded in  $(C_1, C_2)$  one can then replace  $Z_u$  by, say,  $+1$  for  $u > C_2$  and by  $-1$  for  $u < C_1$ . [Such a replacement also substitutes a stronger version of (7.7); it is no longer necessary to consider the possibility that  $|f(u)|$  may become very small as  $|u| \rightarrow \infty$  and result 1 would in this case follow directly from the theorem.]

We now proceed to deal with the Kiefer-Wolfowitz scheme. We denote by  $\bar{D}f(u)$  and  $Df(u)$  the upper and lower derivatives of  $f(u)$ :

$$(7.11) \quad \begin{aligned} \bar{D}f(u) &= \overline{\lim}_{0 \neq h \rightarrow 0} \frac{f(u+h) - f(u)}{h}, \\ Df(u) &= \underline{\lim}_{0 \neq h \rightarrow 0} \frac{f(u+h) - f(u)}{h}. \end{aligned}$$

**RESULT 2.** *If the  $Z_u$  have uniformly bounded variances and if the regression function  $f(u)$  satisfies*

$$(7.12) \quad |f(u+1) - f(u)| < A|u| + B < \infty$$

for all  $u$  and suitable  $A$  and  $B$ , and

$$(7.13) \quad \sup_{1/k < u - \theta < k} \bar{D}f(u) < 0, \quad \inf_{1/k < \theta - u < k} Df(u) > 0$$

for all integers  $k$ , then the Kiefer-Wolfowitz sequence (7.5) converges to  $\theta$  both in mean square and with probability 1.

Indeed, putting  $X_1 = x_1$  and

$$(7.14) \quad X_{n+1} = X_n + \frac{b_n}{c_n} [f(X_n + c_n) - f(X_n - c_n)] + Y_n$$

with

$$(7.15) \quad Y_n = \frac{b_n}{c_n} [Z_{X_n + c_n} - f(X_n + c_n) - Z_{X_n - c_n} + f(X_n - c_n)],$$

we see, by (7.1), that  $Y_n$  satisfies (2.7). Also, since (7.10) holds, we have  $E\{Y_n^2\} < 2\sigma^2 b_n^2 / c_n^2$  and hence, by (7.4), condition (2.8) is also satisfied. Thus, again assuming  $\theta = 0$ , all that remains to be shown is that the transformations  $T_n(r_1, \dots, r_n) = r_n + b_n [f(r_n + c_n) - f(r_n - c_n)] / c_n$  satisfy (2.4). Since  $c_n$  tends to zero we have from (7.12) the inequality  $|f(u + c_n) - f(u - c_n)| < A|u| + A + B$  for  $n > n_0$ . Noticing that  $u$  and  $f(u + c_n) - f(u - c_n)$  have different signs for  $|u| > c_n$  and remembering that  $b_n / c_n \rightarrow 0$  by (7.4), we see that given  $\rho > 0$  we have  $|T_n(r_n)| < \max[\rho + b_n(A\rho + A + B) / c_n, r_n]$  for all  $n > n_0$ . Also, given any  $L > \rho$  we have  $|f(u + c_n) - f(u - c_n)| > 2\gamma c_n$

where  $\gamma = \min \left[ - \sup_{\rho/2 < u < L+1} Df(u), \inf_{\rho/2 < -u < L+1} Df(u) \right]$  for  $\rho < |u| < L$  and  $n > n_L$ .

Hence  $T_n$  satisfies, for large  $n$ , (2.4) with  $a_n = 2\rho$ ,  $\beta_n = 0$ , and  $\gamma_n \geq 0$  satisfying furthermore  $\gamma_n(r_1, \dots, r_n) > 2\gamma b_n$  for  $|r_n| < L$ . Since  $\rho$  is arbitrary our result follows from generalization 3. (The use of generalization 3 could have been avoided as in the proof of result 1; for the sake of variety we illustrated both methods.)

*Remark 2.* Like (7.6) in the previous result, condition (7.12) here has no practical importance and is necessary for the same reasons.

The conclusion that  $x_n$  converges to  $\theta$  with probability 1 was proved by Blum [1] in the case of the Robbins-Monro procedure under exactly the assumptions made by us. He also proved the same conclusion for the Kiefer-Wolfowitz procedure under the following stronger assumptions that  $f(u)$  satisfies (7.12) with  $A = 0$  and the condition obtained from (7.13) on replacing  $1/k < u - \theta < k$  and  $1/k < \theta - u < k$  by  $1/k < u - \theta < \infty$  and  $1/k < \theta - u < \infty$  under the sup and inf signs, respectively (Blum, following Kiefer and Wolfowitz, formulates his assumptions somewhat differently but they are easily seen to be equivalent to those stated here). Blum's results contain those due to Robbins-Monro [7], Wolfowitz [8], Kiefer-Wolfowitz [6] and Kallianpur [5]. Besides the stronger conclusion in both cases, our result 2 allows such regression functions as  $f(u) = -u^2$  and  $f(u) = \exp(-u^2)$  which do not satisfy Blum's conditions.

### 8. A special case

The method of proof of sections 2 and 3 can be adapted to give explicit bounds for  $E\{(X_n - \theta)^2\}$ , etc. Here we shall do it only for a special case, furnishing an extremely simple proof of the theorem in this case.

ASSUMPTION. *The transformations  $T_n$  of (2.5) satisfy*

$$(8.1) \quad |T_n(r_1, \dots, r_n) - \theta| \leq F_n |r_n - \theta|,$$

$F_n, n = 1, 2, \dots$ , being a sequence of positive numbers satisfying

$$(8.2) \quad \prod_{n=1}^{\infty} F_n = 0.$$

Putting  $V_n^2 = E\{(X_n - \theta)^2\}$  and  $\sigma_n^2 = E\{Y_n^2\}$  we have at once from (2.5) and (8.1)

$$(8.3) \quad V_{n+1}^2 \leq F_n^2 V_n^2 + \sigma_n^2.$$

On iteration we have

$$(8.4) \quad V_{n+1}^2 \leq \sigma_n^2 + \sigma_{n-1}^2 F_n^2 + \dots + \sigma_m^2 F_{m+1}^2 F_{m+2}^2 \dots F_n^2 \\ + \dots + \sigma_1^2 F_2^2 F_3^2 \dots F_n^2 + V_1^2 F_2^2 \dots F_n^2.$$

This is the estimate of  $E\{X_{n+1} - \theta\}^2$ . To prove (2.8) we merely have to remark that by (8.4)

$$(8.5) \quad V_{n+1}^2 \leq \sum_{j=m}^n \sigma_j^2 \cdot \max_{m \leq k \leq n} \prod_{j=k+1}^n F_j^2 + \left( V_1^2 + \sum_{j=1}^{m-1} \sigma_j^2 \right) \cdot \max_{1 \leq k < m} \prod_{j=k+1}^n F_j^2.$$

Because of (8.2) all partial products  $\prod_{j=r}^n F_j^2$  are uniformly bounded by a finite number  $A$ ,

say. Given any  $\epsilon > 0$  choose  $m$  so large that  $A \sum_{j=m}^{\infty} \sigma_j^2 < \epsilon/2$ , then (8.5) gives

$$(8.6) \quad V_{n+1}^2 \leq \frac{\epsilon}{2} + \left( V_1^2 + \sum_{j=1}^{\infty} \sigma_j^2 \right) \cdot \max_{1 \leq k \leq m} \prod_{j=k+1}^n F_j^2.$$

With  $m$  being fixed, the max term in (8.5) tends to zero as  $n \rightarrow \infty$  by (8.2) and hence  $V_{n+1}^2 \leq \epsilon$  for all sufficiently large  $n$ . This proves (2.8), and (2.9) can be deduced thence by a simplified version of the argument at the end of section 3. [If it is assumed that all  $F_n$  are  $\leq 1$  then the writing can be somewhat abbreviated since the first max in (8.5) is simply 1 and the second occurs at  $k = m - 1$ .]

We shall illustrate the use of (8.3) by proving the following minimax result on the Robbins-Monro procedure.

RESULT 3. *If the  $Z_u$  satisfy (7.10) and if the regression function  $f(u)$  is measurable and satisfies*

$$(8.7) \quad 0 < A \leq \frac{f(u) - \theta}{u - \theta} \leq B < \infty$$

and if it is known that

$$(8.8) \quad |x_1 - \theta| \leq C \leq \sqrt{\frac{2\sigma^2}{A(B-A)}};$$

then if we use in the Robbins-Monro procedure the sequence

$$(8.9) \quad a_n = \frac{AC^2}{\sigma^2 + nA^2C^2}$$

we shall have

$$(8.10) \quad E\{(X_n - \theta)^2\} \leq \frac{\sigma^2 C^2}{\sigma^2 + (n-1)A^2C^2}.$$

For any other sequence  $a_n$  there are  $Z_u$  and  $x_1$  satisfying all the above conditions for which (8.10) does not hold.

By (7.8) the transformation  $T_n(r_1, \dots, r_n)$  in the Robbins-Monro scheme is  $r_n - a_n f(r_n)$ . Hence (taking  $\theta = 0$  throughout the proof), we have from (8.7)

$$(8.11) \quad |T_n(r_1, \dots, r_n)| \leq |r_n| \cdot \sup_{u \neq 0} \left| 1 - a_n \frac{f(u)}{u} \right| \\ \leq |r_n| \cdot \max(1 - Aa_n, Ba_n - 1).$$

Thus if  $a_n \rightarrow 0$ , (8.1) holds for large  $n$  with  $F_n = 1 - Aa_n$ . Therefore if the  $a_n$  satisfy (7.2) the assumption is verified and hence  $E\{X_n^2\}$  tends to zero. As the sequence (8.9) clearly satisfies (7.2) the conclusion holds in this case. [So far no use was made of (8.8). Also all the above follows directly from result 1.]

From (8.11) it follows that if

$$(8.12) \quad a_n \leq \frac{2}{A+B}$$

then  $T_n$  satisfies (8.1) with  $F_n = 1 - Aa_n$ ; hence we have in this case according to (8.3)

$$(8.13) \quad V_{n+1}^2 \leq (1 - Aa_n)^2 V_n^2 + a_n^2 \sigma^2$$



by (7.9) and (7.10). The minimum of the right side of (8.13) is achieved at

$$(8.14) \quad a_n = \frac{A V_n^2}{\sigma^2 + A^2 V_n^2},$$

and for this choice of  $a_n$  we have

$$(8.15) \quad V_{n+1}^2 \leq \frac{\sigma^2 V_n^2}{\sigma^2 + A^2 V_n^2}.$$

Also, by (8.8),  $V_1^2 \leq C^2$ , this and the recursion formula (8.15) give (8.10); and substituting the right side of (8.10) for  $V_n^2$  in (8.14) we obtain (8.9). Moreover, if the  $a_n$  thus computed satisfy (8.12), and if  $x_1 = C$ , and  $f(u) = Au$  for all  $u$ , and the equality sign always holds in (7.10), we have an equality sign also in (8.15). Thus our last assertion will be proved if we show that the  $a_n$  given by (8.9) satisfy (8.12), but this is evident since they form a monotone sequence and, by (8.8),  $a_1 = AC^2/(\sigma^2 + A^2C^2) \leq 2/(A + B)$ .

It is also easy to dispose of the case when  $C$  does not satisfy (8.8). We merely have to start with  $a_1 = 2/(A + B)$  and keep using this value until, for the first time, we have  $V_m$  not larger than the right side of (8.8). After that we define  $a_{m+n-1}$  by the right side of (8.9) with  $V_m$  replacing  $C$ .

## 9. Concluding remarks

In all the preceding we have dealt with real random variables. Our methods carry over, however, to more general situations. Since then it may be impossible to multiply or square, one has either to operate with a pair of adjoint spaces or with the norm. We shall show how the latter can be done in the case treated in the beginning of the last section. Suppose  $X_1$  and  $Y_n$  assume values in a normed linear space  $\mathcal{N}$  with  $\|r\|$  denoting the norm of  $r$ . Let  $\theta$  be an element of  $\mathcal{N}$  and  $T_n(r_1, \dots, r_n)$  be measurable transformations from the  $n$ th Cartesian power of  $\mathcal{N}$  into  $\mathcal{N}$  and assume that  $\|T_n(r_1, \dots, r_n) - \theta\| \leq F_n \cdot \|r_n - \theta\|$  with  $F_n > 0$  satisfying (8.2). Then, if we put  $X_{n+1} = T_n(X_n) + Y_n$ , the assumptions  $E\{\|X_1\|^2\} + \sum E\{\|Y_n\|^2\} < \infty$  and

$$(9.1) \quad E\{\|\varphi(x_1, \dots, x_n) + Y_n\|^2\} \leq E\{\|\varphi(x_1, \dots, x_n)\|^2\} + E\{\|Y_n\|^2\}$$

for every measurable function  $\varphi(x_1, \dots, x_n)$  imply  $E\{\|X_n\|^2\} \rightarrow 0$  and  $\|X_n - \theta\| \rightarrow 0$  with probability 1. Except for substituting norm instead of absolute value not one word in the proof is changed. We replaced the condition (2.7) by (9.1) since, in the case of real variables, (2.7) was used solely to have  $E\{\varphi_n Y_n\} = 0$  and thus obtain (9.1) with the sign of equality. As, in general, we cannot multiply we assumed (9.1) to start with. This condition is related to orthogonality, and in many important cases may be deduced from relations similar to (2.7). What has been said here about the special case treated in section 7 can be suitably extended to cover the general case. So far as we know, the only treatment of nonreal random variables in this connection is Blum's study [2] of the Robbins-Monro and Kiefer-Wolfowitz schemes for random variables assuming values in finite dimensional Euclidean space.

Chung [3] studied for special cases of the Robbins-Monro procedure the asymptotic distribution of  $X_n - \theta$ . Under the general assumptions of our theorem nothing can, of course, be asserted about asymptotic distributions. Such assertions would certainly necessitate assumption of lower bounds on the effectiveness of the transformations  $T_n$

and very much else. Our methods, however, do give bounds for the second moment of  $X_n - \theta$  and, assuming higher moments for  $X_1$  and  $Y_n$ , can be used to obtain bounds for the corresponding moments of  $X_n - \theta$ . In this connection the inequalities of Doob (see chapter 8, section 3 in [4]) can be useful. (Doob's theorems may also serve to give a modified proof of our main results.)

The general theory embraces naturally many other schemes besides those of Robbins-Monro and Kiefer-Wolfowitz. It may also be modified to yield methods of obtaining confidence intervals and the like.



*Note added in proof.* Since submitting the paper the author became aware of the following two studies.

(a) D. L. BURKHOLDER, "On a certain class of stochastic approximation processes," Mimeograph Series No. 129, Institute of Statistics, University of North Carolina (1955).

(b) C. DERMAN, "An application of Chung's lemma to the Kiefer-Wolfowitz stochastic approximation procedure," to appear in *Annals of Math. Stat.*

The most relevant results of these papers are a proof of the probability 1 part of Result 2 of section 7 in (a), and studies of the asymptotic distribution in special cases of the Kiefer-Wolfowitz procedure in both (a) and (b).

#### REFERENCES

- [1] J. R. BLUM, "Approximation methods which converge with probability one," *Annals of Math. Stat.*, Vol. 25 (1954), pp. 382-386.
- [2] ———, "Multidimensional stochastic approximation methods," *Annals of Math. Stat.*, Vol. 25 (1954), pp. 737-744.
- [3] K. L. CHUNG, "On a stochastic approximation method," *Annals of Math. Stat.*, Vol. 25 (1954), pp. 463-483.
- [4] J. L. DOOB, *Stochastic Processes*, New York, John Wiley and Sons, 1953.
- [5] G. KALLIANPUR, "A note on the Robbins-Monro stochastic approximation method," *Annals of Math. Stat.*, Vol. 25 (1954), pp. 386-388.
- [6] J. KIEFER and J. WOLFOWITZ, "Stochastic estimation of the maximum of a regression function," *Annals of Math. Stat.*, Vol. 23 (1952), pp. 462-466.
- [7] H. ROBBINS and S. MONRO, "A stochastic approximation method," *Annals of Math. Stat.*, Vol. 22 (1951), pp. 400-407.
- [8] J. WOLFOWITZ, "On the stochastic approximation method of Robbins and Monro," *Annals of Math. Stat.*, Vol. 23 (1952), pp. 457-461.