

A REPRESENTATION THEOREM FOR GENERALIZED
ROBBINS-MONRO PROCESSES AND APPLICATIONS

David Ruppert¹
Department of Statistics
University of North Carolina at Chapel Hill

Running title: Robbins-Monro Processes

Summary

Many generalizations of the Robbins-Monro process have been proposed for the purpose of recursive estimation. In this paper, it is shown that for a large class of such processes, the estimator can be represented as a sum of possibly dependent random variables, and therefore the asymptotic behavior of the estimate can be studied using limit theorems for sums. Two applications are considered: robust recursive estimation for autoregressive processes and recursive nonlinear regression.

¹ Supported by National Science Foundation Grant MCS81-00748.

1. Introduction. Stochastic approximation was introduced with Robbins and Monro's (1951) study of recursive estimation of the root, θ , to $R(x) = 0$, when the function R from \mathbb{R} (the real line) to \mathbb{R} is unknown, but when for each x , an unbiased estimator of $R(x)$ can be observed. Their procedure lets $\tilde{\theta}_1$ be an initial estimate of θ and defines $\tilde{\theta}_n$ recursively by

$$\tilde{\theta}_{n+1} = \tilde{\theta}_n - a_n Y_n,$$

where a_n is a suitable positive constant and Y_n is an unbiased estimate of $R(\tilde{\theta}_n)$. Blum (1954) extended their work to the multidimensional case: that is, when R maps \mathbb{R}^p (p -dimensional Euclidean space) to \mathbb{R}^p .

A situation where the classical Robbins-Monro (RM) procedures are inadequate occurs when there is a different function (now mapping \mathbb{R}^p to \mathbb{R}) at each stage of the recursion. Letting R_n be the function at the n^{th} stage, there may be multiple solutions to

$$(1.1) \quad R_n(x) = 0,$$

but it is assumed that there exists a unique θ solving (1.1) for all n . At the n^{th} stage, we can, for any x in \mathbb{R}^p , observe an unbiased estimate of $R_n(x)$. In this context, we will define a generalized Robbins-Monro (GRM) procedure to be a process satisfying

$$(1.2) \quad \tilde{\theta}_{n+1} = \tilde{\theta}_n - a_n Y_n U_n,$$

where $\tilde{\theta}_1$ is an arbitrary initial estimate of θ , $\{a_n\}$ is a sequence of positive constants, Y_n is an unbiased estimate of $R_n(\tilde{\theta}_n)$, and U_n is a suitably chosen random vector in \mathbb{R}^p .

GRM processes are used by Albert and Gardner (1967) for recursive estimation of nonlinear regression parameters. They assume that one observes a process $\{Y_n: n = 1, 2, \dots\}$ such that $EY(n) = F_n(\theta)$, where F_n is a known map of \mathbb{R}^p to \mathbb{R} and θ is an unknown parameter vector which is to be estimated. If $R_n(x) = F_n(x) - F_n(\theta)$, then θ solves (1.1) for all n .

Another example of a GRM process appears in Campbell (1979), who considers robust recursive estimation for autoregressive processes. Suppose $\theta = (\theta_1, \dots, \theta_p)'$ is a vector parameter and the process $\{Y_n: n = 1, 2, \dots\}$ satisfies

$$Y_n = \sum_{k=1}^p Y_{n-k} \theta_k + u_n,$$

where $\{u_n\}_{n=-\infty}^{\infty}$ are iid. Assume ψ , which maps \mathbb{R} to \mathbb{R} , satisfies $E\psi(s+u_1) = 0$ if and only if $s = 0$. Then for $x = (x_1, \dots, x_p)'$, let

$$R_n(x) = \psi\left(Y_n - \sum_{k=1}^p Y_{n-k} x_k\right).$$

In this case, because the Y 's are random for each x , $R_n(x)$ is a random variable. (Albert and Gardner's setup allows $R_n(x)$ to be either random or fixed; the situation is analogous to regression where the independent variables may be random or fixed.) Except in degenerate situations,

$$P\{\text{there exists } x \neq 0 \text{ satisfying } R_n(x) = 0 \text{ for all } n\} = 0.$$

A third example of a GRM is given by Ruppert (1979, 1981a).

The classical Robbins-Monro process has been studied quite extensively,

but little beyond consistency has been established for GRM procedures. Albert and Gardner (1967) investigate asymptotic distributions for $p = 1$ but state that "owing to its considerable complexity, the question of large-sample distribution for the vector estimates is not examined." Campbell (1979) does not treat asymptotic distributions, and Ruppert (1981a) proves asymptotic normality of $n^{1/2}(\tilde{\theta}_n - \theta)$ (only in his special case) under restrictions which imply that for each fixed x , $\{R_n(x): n = 1, 2, \dots\}$ are iid, that $\{(Y_n - R_n(\tilde{\theta}_n)): n = 1, 2, \dots\}$ are iid, and that $\{R_n(x): n = 1, 2, \dots \text{ and } x \in \mathbb{R}^p\}$ is independent of $\{(Y_n - R_n(\tilde{\theta}_n)): n = 1, 2, \dots\}$.

The classical studies of stochastic approximation procedures -- for example, Chung (1954), Sacks (1958), and Fabian (1968) -- assume that the sequence $Y_n - R(\tilde{\theta}_n)$ forms a martingale difference sequence, but recent authors, including Ljung (1978) and Kushner and Clark (1978), have felt that this assumption can be unduly restrictive in applications. Ljung (1978) proves consistency, and Kushner and Clark (1978) give a detailed treatment of consistency and asymptotic distributions for stochastic approximation methods, under assumptions considerably weaker than that the errors are martingale differences. Ruppert (1978) proves almost sure representations for multidimensional Robbins-Monro and Kiefer-Wolfowitz (1952) processes with weakly dependent errors. These representations can be used to establish almost sure behavior, i.e. laws of the iterated logarithm and invariance principles, as well as asymptotic distributions. Invariance principles are useful for the study of stopping rules; see Sielken (1973).

For GRM procedures, the results of Kushner and Clark (1978) -- in particular, their theorems 2.4.1, 2.5.1, and 2.5.2 -- are sufficient to prove consistency under weak conditions on $(Y(n) - R_n(\tilde{\theta}(n)))$, but they do not appear adequate to handle asymptotic distributions.

In this paper, the techniques of Ruppert (1978) are used to show that GRM processes can be almost surely approximated by a sequence of partial sums of random vectors, and therefore, central limit theorems, invariance principles, laws of the iterated logarithm, and other results for such sums can also be established for GRM processes. The procedures of Campbell (1979) and Albert and Gardner (1967) are used as examples. In future work, these results will be applied to the procedure of Ruppert (1979, 1981a).

2. Notation and assumptions. All random variables are defined on a probability space (Ω, \mathcal{F}, P) . All relations between random variables are meant to hold with probability 1. $(\mathbb{R}^k, \mathcal{B}^k)$ is k -dimensional Euclidean space with the Borel σ -algebra. Let $X^{(i)}$ be the i^{th} coordinate of the vector X , $A^{(ij)}$ be the i, j^{th} entry of the matrix A , A' be the transpose of A , I_k be the $k \times k$ identity matrix, and $\|A\|$ be the Euclidean norm of the matrix A , i.e. $\|A\| = (\sum_{i=1}^n \sum_{j=1}^m (A^{(ij)})^2)^{1/2} = \text{trace}(A'A)^{1/2}$. If A is a square matrix, then $\lambda_*(A)$ and $\lambda^*(A)$ denote the minimum and maximum of the real parts of the eigenvalues of A , respectively.

Let $I(A)$ be the indicator of the set A . Throughout, K will denote a generic positive constant.

For a square matrix M , $\exp(M)$ is given by the usual series definition:

$$\exp(M) = \sum_{n=0}^{\infty} M^n (n!)^{-1},$$

and for $t > 0$,

$$t^M = \exp((\log t)M).$$

Let p be a fixed integer. For convenience, we will often abbreviate I_p by I .

Notice that $t^I = tI$, $t^M s^M = s^M t^M = (st)^M$, $(t^{-1})^M = t^{-M}$, $Mt^M = t^M M$, and $(d/dt)t^M = Mt^{M-I}$.

We write $X_n = o(Y_n)$ or $X_n \ll Y_n$ if the random vectors X_n and Y_n satisfy $\|X_n\| \leq X \|Y_n\|$ for a (finite) random variable X .

We will be needing the following assumptions.

A1. $M_n(X, \omega)$ ($= M_n(X)$) is a $(\mathbb{R}^p \times \Omega, \mathcal{B}^p \times \mathcal{F})$ to $(\mathbb{R}^p, \mathcal{F})$ measurable transformation, $A_n(\omega)$ ($= A_n$) is a random matrix, v is a nonnegative Borel function on \mathbb{R}^p , θ is in \mathbb{R}^p , and B_n is a nonnegative random variable such that

$$(2.1) \quad v(x) = o(x - \theta) \text{ as } x \rightarrow \theta$$

and

$$(2.2) \quad \|M_n(x) - A_n(x - \theta)\| \leq B_n v(x)$$

for all x in \mathbb{R}^p and all n .

A2. $\tilde{\theta}_{n+1} = \tilde{\theta}_n - n^{-1}(M_n(\tilde{\theta}_n) + W_n)$, where W_n is a random vector in \mathbb{R}^p .

$$A3. \quad \tilde{\theta}_n \rightarrow \theta.$$

A4. A is a $p \times p$ matrix such that $\lambda_*(A) > \frac{1}{2}$ and C_1 is a constant such that

$$(2.3) \quad \sum_{k=1}^{\infty} k^{-1} (A_k - A) \text{ converges}$$

and

$$(2.4) \quad \limsup_{m \rightarrow \infty} \sup_{m \leq n} \sum_{k=m}^n k^{-1} (||A_k - A|| - C_1) \leq 0 .$$

A5. Suppose

$$(2.5) \quad \sup_n E(||A_n||^2 + B_n^2) < \infty ,$$

$$(2.6) \quad \sum_{k=1}^{\infty} k^{-\frac{1}{2}-\epsilon} W_k \text{ converges for all } \epsilon > 0 ,$$

and for a constant C_2 ,

$$(2.7) \quad \limsup_{m \rightarrow \infty} \sup_{m \leq n} \sum_{k=m}^n k^{-1} (||B_k|| - C_2) \leq 0 .$$

A6. There exists a constant C_3 such that for all $n > m$, α real, and matrices Q_i such that $||Q_i|| \leq 1$, we have

$$(2.8) \quad E\{\max || \sum_{i=m}^k i^\alpha Q_i (A_i - A) ||^2 : m \leq k \leq n\} \leq C_3 \sum_{i=m}^n i^{2\alpha}$$

and

$$(2.9) \quad E\{\max || \sum_{i=m}^k i^{-1} W_i ||^2 : m \leq k \leq n\} \leq C_3 \sum_{i=m}^n i^{-2} .$$

A7. For some $\delta > 0$,

$$v(x) = o(||x||^{1+\delta}) \text{ as } x \rightarrow 0 .$$

REMARKS. The relationship between (1.2) and A2 is that $M_n(x) = a^{-1} R_n(x) U_n$ (a is some positive number), $W_n = a^{-1} (Y_n - R_n(\tilde{\theta}_n)) U_n$, and $a_n = a n^{-1}$. More general a_n sequences could be treated, but they would not lead to better rates for convergence of $\tilde{\theta}_n$ to θ . McLeish (1975) has investigated the convergence of the series $\sum_{n=1}^{\infty} c_n X_n$, where the c_n are constants and $\{X_n\}$ is a sequence of random variables satisfying certain mixing and moment conditions. As we will see in Section 4, his results can be used in the verification of A4 and A5 in some specific case. Note that (2.4) is a weaker assumption than that $\sum_{k=1}^{\infty} (||A_k - A|| - C_1) k^{-1}$ converges (and similarly for (2.7)). Also McLeish's theorem 1.6, which is a generalization of a martingale inequality due to Doob (1953), gives conditions which imply A6.

3. Main results.

LEMMA 3.1. *Let M be a $p \times p$ square matrix such that $a < \lambda_*(M) \leq \lambda^*(M) < b$ for real numbers a and b . Then there exists a norm $|| \cdot ||_*$ on \mathbb{R}^p such that*

$$e^{at} ||X||_* \leq ||e^{Mt} X||_* \leq e^{bt} ||X||_*$$

for all real t and all X in \mathbb{R}^p .

PROOF. The proof is given by Hirsch and Smale (1974, page 146). \square

LEMMA 3.2. *Assume A1 to A5 hold. Then $n^\epsilon (\tilde{\theta}_n - \theta) \rightarrow 0$ for all $\epsilon < \frac{1}{2}$.*

PROOF. Without loss of generality, we take $\theta = 0$. Fix $\epsilon < \frac{1}{2}$ and define $\phi_n = (n-1)^\epsilon \tilde{\theta}_n$. Then from (2.2) and A2, we obtain

$$\phi_{n+1} = [I + n^{-1}\{\epsilon I - A_n + u_n\}]\phi_n - n^{\epsilon-1} W_n ,$$

where the random matrix u_n satisfies

$$(3.1) \quad \|u_n\| = o(n^{-1} + n^{-1}\|A_n\|) + o(B_n) .$$

Define $D_n = A_n - \epsilon I - u_n$. Suppose $n \geq m$. Then iterating (3.1) back to m and using the definitions

$$D = A - \epsilon I ,$$

$$\mathcal{D}_m^n = \sum_{k=m}^n k^{-1}(D_k - D), \quad K_m^n = \sum_{k=m}^n k^{-1} ,$$

and

$$\omega_m^n = \sum_{k=m}^n k^{\epsilon-1} W_k ,$$

we see that

$$(3.2) \quad \begin{aligned} \phi_{n+1} = \phi_m - [& K_m^n D \phi_m + \mathcal{D}_m^n \phi_m + D \sum_{k=m}^n k^{-1} (\phi_k - \phi_m) \\ & + \sum_{k=m}^n k^{-1} (D_k - D) (\phi_k - \phi_m) + \omega_m^n] . \end{aligned}$$

Choose t_0 and a norm $\|\cdot\|_*$ such that

$$\|(I - Dt)x\|_* \leq (I - \lambda_*(D)t/2) \|x\|_*$$

for all t in $[0, t_0]$. (This can be done by Lemma 3.1 and a simple calculation.) Let

$$||D||_* = \sup\{||Dx||_* : ||x||_* = 1\} .$$

For $x > 0$ and $y > 0$, define

$$(3.3) \quad \rho(x,y) = (x - y(x+2)) / (||D||_*(1+x) + y + C_1 x + xy)$$

(C_1 is given by A4) and

$$(3.4) \quad r(x,y) = (\rho(x,y) - y)(\lambda_*(D)/2 - y - ||D||_* x - xy - C_1 x) - (2y + xy) .$$

Choose $\eta > 0$ and $L > 0$ such that

$$t_0/2 > \rho(L,\eta) > \eta \text{ and } r(L,\eta) > 0 .$$

Then using (2.3), (2.4), (2.6), and (2.7), choose $N > \max\{2/t_0, \eta^{-1}\}$ so large that if $n > m \geq N$, then

$$(3.5) \quad ||\mathcal{D}_m^n||_* \leq \eta(1 + K_m^n) ,$$

$$(3.6) \quad \sum_{k=m}^n k^{-1} ||D_k - D||_* \leq (\eta + C_1) K_m^n + \eta ,$$

and

$$(3.7) \quad ||w_m^n||_* < \eta^2 .$$

Define $n(1) = N$ and $n(i+1) = \inf\{k \geq n(i) : K_{n(i)}^k \geq \rho(L,\eta)\}$ for $i = 1, 2, \dots$. Note that

$$(3.8) \quad K_{n(i)}^{n(i+1)} \leq t_0$$

since $\rho(L,\eta) \leq t_0/2$ and $(n(i+1))^{-1} \leq N^{-1} \leq t_0/2$. Since $\rho(L,N) > \eta > N^{-1}$,

$n(i+1) > n(i)$ for all $i \geq 1$. Let

$$(3.9) \quad \eta_i = \sup\{\|\omega_{n(i)}^k\|^{\frac{1}{2}} : n(i) \leq k \leq n(i+1)\}.$$

By (3.7), $\eta_i < \eta$ for $i \geq 1$, and by (2.6), $\eta_i \rightarrow 0$.

We now will prove that for each i and each k in $\{n(i), \dots, n(i+1)\}$,

$$(3.10) \quad \|\phi_k - \phi_{n(i)}\| \leq L \max\{\|\phi_{n(i)}\|, \eta_i\}.$$

The proof will be by induction on k for each fixed i . Suppose that (3.10) holds for all ℓ less than or equal to some $k < n(i+1)$. Then by (3.2), (3.5), (3.6), and (3.9), and the fact that $\kappa_{n(i)}^k < \rho(L, \eta)$ for $k < n(i+1)$,

$$\begin{aligned} \|\phi_{k+1} - \phi_{n(i)}\|_* &\leq \kappa_{n(i)}^k \|D\|_* \|\phi_{n(i)}\|_* + \eta(\kappa_{n(i)}^k + 1) \|\phi_{n(i)}\|_* \\ &\quad + ((C_1 + \eta + \|D\|_*) \kappa_{n(i)}^k + \eta) \cdot \sup\{\|\phi_\ell - \phi_{n(i)}\|_* : n(i) \leq \ell \leq k\} + \eta_i^2 \\ &\leq \{\rho(L, \eta) [\|D\|_*(1+L) + \eta + C_1 L + \eta L] + \eta(L+2)\} \cdot \max\{\|\phi_{n(i)}\|_*, \eta_i\} \\ &\leq L \max\{\|\phi_{n(i)}\|_*, \eta_i\}, \end{aligned}$$

which completes the proof of (3.10). It follows from (3.2), (3.5) to (3.8), and (3.10) that

$$\begin{aligned} &I\{\|\phi_{n(i)}\|_* > \eta_i\} \|\phi_{n(i+1)}\|_* \\ &\leq \|\phi_{n(i)}\|_* \{1 + 2\eta + \eta L - \kappa_{n(i)}^{n(i+1)-1} [\lambda_*(D)/2 - \eta - \|D\|_* L - \eta L - C_1 L]\}. \end{aligned}$$

By the definition of $n(i+1)$,

$$\kappa_{n(i)}^{n(i+1)-1} \geq \rho(L, n) - (n(i+1))^{-1} \geq \rho(L, n) - \eta,$$

and therefore

$$(3.11) \quad I\{\|\phi_{n(i)}\| > \eta_i\} \|\phi_{n(i+1)}\|_* \leq \|\phi_{n(i)}\|_* (1 - r(n, L)).$$

Now choose $\gamma_i > 0$ such that $\gamma_i > \eta_i$ and $\sum \gamma_i = \infty$. Then by (3.10) and (3.11), we have

$$\|\phi_{n(i+1)}\|_* \leq \max\{\gamma_i(L+1), \|\phi_{n(i)}\|_* - r(n, L)\gamma_i\}.$$

An application of Derman and Sack's (1959) lemma 1, with $a_i = \gamma_i(L+1)$, $b_i = 0$, $\delta_i = 0$, and $c_i = r(n, L)\gamma_i$, proves that $\|\phi_{n(i)}\|_* \rightarrow 0$, and then by (3.10) we can conclude that $\|\phi_n\|_* \rightarrow 0$. \square

THEOREM 3.1. *Assume A1 to A7. Then there exists $\epsilon > 0$ such that*

$$n^{1/2}(\tilde{\theta}_{n+1}^{-\theta}) = -n^{-1/2} \sum_{k=1}^n (k/n)^{A-I} W_k + o(n^{-\epsilon}).$$

PROOF. Again we take $\theta = 0$. By (2.2), A2, and A7,

$$\tilde{\theta}_{n+1} = \tilde{\theta}_n - n^{-1}(A\tilde{\theta}_n + \kappa_n \tilde{\theta}_n + \xi_n + W_n),$$

where

$$\kappa_n = A_n - A \text{ and } \xi_n = o(B_n \|\tilde{\theta}_n\|^{1+\delta}).$$

By A6, $E\|\kappa_n\|^2 \leq C_3$ for all n , whence

$$\sum_{n=1}^{\infty} n^{-1-\varepsilon} E \|\kappa_n\| < \infty$$

for all $\varepsilon > 0$. Thus for all $\varepsilon > 0$,

$$\sum_{n=1}^{\infty} n^{-1-\varepsilon} \|\kappa_n\| < \infty .$$

(Here and elsewhere in the proof, we make use of the fact that for nonnegative random variables X_n , $\sum_1^{\infty} X_n$ converges if $\sum_1^{\infty} EX_n$ converges.) Then since by Lemma 3.2, $\tilde{\theta}_n \ll n^{-\frac{1}{2}+\varepsilon}$ for all $\varepsilon > 0$, all $\varepsilon > 0$,

$$\sum_{n=1}^{\infty} n^{-\frac{1}{2}-\varepsilon} \|\kappa_n \tilde{\theta}_n\| < \infty .$$

Also, we can prove in a similar fashion that

$$\sum_{n=1}^{\infty} n^{-\frac{1}{2}} \|\xi_n\| < \infty .$$

Therefore for each $\varepsilon > 0$,

$$\sum_{n=1}^{\infty} n^{-\frac{1}{2}-\varepsilon} (\kappa_n \tilde{\theta}_n + \xi_n + W_n) \text{ converges}$$

by (2.6). Now applying Theorem 3.1 of Ruppert (1978) with $\tau = 0$, $X_n = \tilde{\theta}_n$, $f(x) = Ax$, $\beta_n = 0$, $\beta = 0$, $e_n = (\kappa_n \tilde{\theta}_n + \xi_n + W_n)$, and $\gamma = \frac{1}{2}$, we obtain

$$n^{\frac{1}{2}} \tilde{\theta}_{n+1} = -n^{-\frac{1}{2}} \sum_{k=1}^n (k/n)^{A-I} (\kappa_k \tilde{\theta}_k + \xi_k + W_k) + O(n^{-\varepsilon})$$

for some $\varepsilon > 0$. The proof will be completed by showing that for some

$\epsilon > 0$,

$$(3.12) \quad n^{-\frac{1}{2} + \epsilon/3} \sum_{k=1}^n (k/n)^{A-I} (\kappa_k \tilde{\theta}_k + \xi_k) = o(1) .$$

We will assume that $A-I$ has only one eigenvalue, λ . This involves no loss in generality since, by a change of basis, we can put $A-I$ in real canonical form (Hirsch and Smale (1974, page 130)) so that

$$(A-I) = \text{diag}(M_1, \dots, M_q) , \quad q \leq p ,$$

where the square matrices M_1, \dots, M_q each have exactly one eigenvalue.

Then we can apply the proof q times. By Lemma 3.1, there exists Q_k and Q_k^* such that $\sup_k (||Q_k|| + ||Q_k^*||) < \infty$ and

$$k^{A-I} = k^{\lambda + \epsilon/3} Q_k$$

$$k^{I-A} = k^{-\lambda + \epsilon/3} Q_k^* .$$

Then (3.12) holds if

$$n^{-\frac{1}{2} - \lambda + 2\epsilon/3} \sum_{k=1}^n Q_k^* k^{\lambda + \epsilon/3} Q_k (\kappa_k \tilde{\theta}_k + \xi_k) = o(1) ,$$

and so by Kronecker's lemma it is sufficient to find $\epsilon > 0$ such that

$$(3.13) \quad \sum_{k=1}^{\infty} k^{-\frac{1}{2} + \epsilon} Q_k \xi_k \text{ converges}$$

and

$$(3.14) \quad \sum_{k=1}^{\infty} k^{-\frac{1}{2} + \epsilon} Q_k \kappa_k \tilde{\theta}_k \text{ converges} .$$

Now since $||\xi_k|| \ll B_n ||\tilde{\theta}_n||^{1+\delta}$ and $||\tilde{\theta}_n|| \ll n^{-\frac{1}{2}+\epsilon}$ for all $\epsilon > 0$, (3.13) holds for some $\epsilon > 0$. It remains to prove (3.14).

Now fix a , ϵ , and ϵ' such that

$$\begin{aligned} a &> 3/2, \\ 0 &< \epsilon < \frac{1}{4}, \end{aligned}$$

and

$$0 < 2a\epsilon' < \epsilon.$$

Define $\phi(n)$ to be the greatest integer less than or equal to n^a , $\eta(i) = \{j: \phi(i) \leq j \leq \phi(i+1) - 1\}$, and $\psi_k = Q_k \kappa_k k^{-\frac{1}{2}+\epsilon'}$. Then

$$\sum_{k=1}^{\infty} \psi_k \tilde{\theta}_k = \sum_{k=1}^{\infty} \sum_{j \in \eta(k)} \psi_j \tilde{\theta}_j.$$

We will prove (3.14) by showing that

$$(3.15) \quad \sum_{k=1}^{\infty} \left\| \sum_{j \in \eta(k)} \psi_j \tilde{\theta}_j \right\| < \infty.$$

Now define

$$\begin{aligned} z_{1,i} &= \sum_{j \in \eta(i)} \psi_j \left(\sum_{\ell=\phi(i)}^{j-1} \ell^{-1} W_{\ell} \right), \\ z_{2,i} &= \sum_{j \in \eta(i)} \psi_j \left(\sum_{\ell=\phi(i)}^{j-1} \ell^{-1} M_{\ell}(\tilde{\theta}_{\ell}) \right), \end{aligned}$$

and

$$z_{3,i} = \tilde{\theta}_{\phi(i)} \sum_{j=\eta(i)} \psi_j.$$

Then since $\tilde{\theta}_j = \tilde{\theta}_{\phi(i)} - \sum_{\ell=\phi(i)}^{j-1} \ell^{-1} (M_\ell(\tilde{\theta}_\ell) + W_\ell)$ for $j > \phi(i)$, we will have proved (3.15) if we prove that $\sum_{i=1}^{\infty} \|z_{m,i}\| < \infty$ for $m = 1, 2$, and

3. Now

$$\begin{aligned} \sum_{i=1}^{\infty} \|z_{1,i}\| &\leq \sum_{i=1}^{\infty} \left\| \sum_{j \in \eta(i)} \psi_j \right\| \left\{ \max_{j \in \eta(i)} \left\| \sum_{\ell=\phi(i)}^{j-1} \ell^{-1} W_\ell \right\| \right\} \\ &\leq \left(\sum_{i=1}^{\infty} i^{-\varepsilon} \left\| \sum_{j \in \eta(i)} \psi_j \right\|^2 \right)^{1/2} \left(\sum_{i=1}^{\infty} i^\varepsilon \max_{j \in \eta(i)} \left\| \sum_{\ell=\phi(i)}^{j-1} \ell^{-1} W_\ell \right\|^2 \right)^{1/2} \\ &< \infty, \end{aligned}$$

since by A6

$$(3.16) \quad E \left\| \sum_{j \in \eta(i)} \psi_j \right\|^2 \ll i^{-1+2a\varepsilon'},$$

$$2a\varepsilon' < \varepsilon,$$

$$i^\varepsilon E \max_{j \in \eta(i)} \left\| \sum_{\ell=\phi(i)}^{j-1} \ell^{-1} W_\ell \right\|^2 \ll i^{\varepsilon-a-1},$$

and $\varepsilon < a$ by choice of ε and a .

By (2.1), (2.2), and Lemma 3.2,

$$\ell^{-1} \|M_\ell(\tilde{\theta}_\ell)\| \ll (\|A_\ell\| + B_\ell) \ell^{-3/2 + \varepsilon'}.$$

Therefore,

$$(3.17) \quad \sum_{i=1}^n \|z_{2,i}\| \ll \sum_{i=1}^n \left\| \sum_{j \in \eta(i)} \psi_j \right\| \left[\sum_{j \in \eta(i)} (\|A_j\| + B_j) \right] i^{a(-3/2 + \varepsilon')}.$$

By (2.5),

$$E\left[\sum_{j \in \eta(i)} (||A_j|| + B_j)\right]^2 \ll (\#\eta(i))^2 \ll i^{2(a-1)},$$

where #A is the cardinality of the set A. Therefore, using (3.16), the expectation of the i^{th} summand on the RHS of (3.17) is bounded by

$$K i^{a(-3/2 + \epsilon')} i^{-1/2 + a\epsilon'} i^{a-1} = K i^{-a/2 - 3/2 + 2a\epsilon'}.$$

Since $\epsilon' < 1/4$, the expectation of the RHS of (3.17) has a finite limit.

Therefore, $\sum_{i=1}^{\infty} ||z_{2,i}|| < \infty$.

Since $\sup_i ||i^{(\frac{1}{2} - \epsilon')a} \tilde{\theta}_{\phi(i)}|| < \infty$, $\sum_{i=1}^{\infty} ||z_{3,i}||$ converges if the following converges:

$$(3.18) \quad \sum_{i=1}^{\infty} i^{(\epsilon' - \frac{1}{2})a} ||\sum_{j \in \eta(i)} \psi_j||.$$

Using (3.16) and $-(a+1)/2 + 2a\epsilon' < -(a+1)/2 + \epsilon < -1$, (3.18) converges. \square

4. Applications. In this section, we elaborate upon two examples of generalized Robbins-Monro processes that were mentioned briefly in the introduction.

ROBUST RECURSIVE ESTIMATION FOR AUTOREGRESSIVE PROCESSES.

Let $\{Y_n\}$ be a strictly stationary autoregressive process such that

$$Y_n = \sum_{k=1}^p Y_{n-k} \theta^{(k)} + u_n,$$

where $p < \infty$ and $\{u_n\}$ is an iid sequence with distribution function F and satisfying $Eu_n^4 < \infty$ and $Eu_n = 0$. There are, of course, many methods, e.g. least squares, of estimating $\theta = (\theta^{(1)}, \dots, \theta^{(p)})'$ based on the sample (Y_1, \dots, Y_n) . If the Y_n are being observed rapidly and we need to update our estimate after receiving each new observation, then recursively defined estimators are very attractive. Campbell (1979) has defined a class of recursive estimators which are robust (that is, relatively insensitive to the tail behavior of the distribution of u_n and to the presence of spurious observations), but she only examined the consistency of these procedures. We will look at a class of procedures including Campbell's, investigate its asymptotic distributions, and find the procedure which minimizes the asymptotic variance-covariance matrix. Let $Z_n = (Y_{n-1}, \dots, Y_{n-p})'$. Define $\tilde{\theta}_n$ recursively by

$$(4.1) \quad \tilde{\theta}_{n+1} = \tilde{\theta}_n - n^{-1} g(Z_n) \psi(Z_n' \tilde{\theta}_n - Y_n) D Z_n,$$

where g and ψ are maps from \mathbb{R}^p to \mathbb{R} and \mathbb{R} to \mathbb{R} respectively, and D is a $p \times p$ matrix. For robustness, we require that

$$(4.2) \quad \sup_{x \in \mathbb{R}} |\psi(x)| < \infty$$

and

$$\sup_{x \in \mathbb{R}^p} \|g(x)x\| < \infty.$$

These boundedness assumptions will also simplify later proofs. Campbell uses only $D = I$, but she considers a more general subsequence, $\{a_n\}$, in place of n^{-1} in expression (4.1). We will demonstrate that $D = I$ is asymptotically optimal only under rather special circumstances.

Moreover, as we will see, certain of our recursive estimators are asymptotically equivalent to their nonrecursive competitors, bounded influence regression estimates (Hampel (1978)), so there appears to be no advantage in using other a_n .

First, we will verify that A1 to A7 hold. Let

$$F_n = \sigma(u_k : k \leq n - 1) ,$$

$$M_n(x) = \left[\int_{-\infty}^{\infty} \psi(Z_n'(X-\theta) - u) dF(u) \right] g(Z_n) DZ_n$$

$$= E^n \psi(Z_n'(X-\theta) - u_n) g(Z_n) DZ_n ,$$

and

$$W_n = \psi(Z_n'(\tilde{\theta}_n - \theta) - u_n) g(Z_n) DZ_n - M_n(\tilde{\theta}_n) .$$

Then (4.1) is equivalent to A2.

We will also assume that:

B1. ψ has a bounded Radon-Nikodym derivative $\dot{\psi}$.

B2. ψ has a bounded second derivative $\ddot{\psi}$ except at a finite collection of points $\{\alpha_1, \dots, \alpha_J\}$.

B3. $E\psi(-u_1) = 0$.

B4. F has a density f satisfying

$$\int_{-\infty}^{\infty} |f(x) - f(x+\alpha)| dx \leq K|\alpha|^\beta \text{ for all } \alpha \text{ and some } \beta > 0 .$$

B5. The polynomial $x^p - \sum_{r=1}^p x^{p-r} \theta^{(r)}$ has all roots inside the unit circle.

B6. Define $c = \int_{-\infty}^{\infty} \dot{\psi}(-u) dF(u)$, $S_n = g(Z_n) Z_n Z_n'$, $A_n = c D S_n$, $S = E S_n$, and $A = E A_n$. Assume $\lambda_*(A) > \frac{1}{2}$.

All ψ functions in the Princeton robustness study (Andrews (1972)), for example, Huber's

$$\psi(x) = \min(k, |x|) \text{sign } x$$

and Andrews'

$$\begin{aligned} \psi(x) &= \sin \frac{x}{k} \quad \text{if } |x| \leq \pi k \\ &= 0 \quad \text{otherwise,} \end{aligned}$$

satisfy B1 and B2, and they satisfy B3 if F is symmetric about zero.

Since $\{W_n, F_{n+1}\}$ is a martingale difference sequence and is uniformly bounded in L_2 , (2.6) and (2.9) hold by the martingale convergence theorem and Doob's inequality (Doob (1953, Chapter VII, Theorem 3.4)).

Theorem 3.2 of Pham and Tran (1980), B4, and the assumption that $E\mu_4^4 < \infty$ imply that Y_n is strong mixing with mixing coefficients $\alpha(n)$ which are $O(\rho^n)$ as $n \rightarrow \infty$ for some ρ in $(0,1)$.

Therefore, if h is a Borel function on \mathbb{R}^q and $V_k = h(Y_k, \dots, Y_{k+q-1})$, then $\{V_n\}$ is also strong mixing, and its mixing coefficients are $o(\rho^{n/2})$.

Define

$$(4.3) \quad B_n = ||Z_n||^2.$$

Since $EY_n^4 < \infty$ and $g(x)$ is bounded, A_n (see B6) and B_n are strictly stationary and have finite second moments, so (2.5) holds. Define $C_1 = E||A_n - A||$ and $C_2 = EB_n$. Then $\{A_n\}$, $\{||A_n - A|| - C_1\}$, and $\{B_n - C_2\}$ are mixingales of size $-q$ for all $q > 0$. (See McLeish (1975) for a definition and discussion of mixingales.) Therefore, by Corollary 1.8 of McLeish, (2.3), (2.4), and (2.7) hold. Therefore A5 holds, and since $\lambda_*(A) > \frac{1}{2}$ is assumed, A4 holds. By Theorem 1.6 of McLeish, A6 holds.

By B1, B2, and B4,

$$||M_n(x) - A_n(x-\theta)|| = o(B_n ||x - \theta||^2)$$

as $x \rightarrow \theta$, so that A1 and A7 hold with $v(x) = ||x - \theta||^2$ and B_n given by (4.3). One can show that $\tilde{\theta}_n \rightarrow \theta$, i.e. that A3 holds, using Theorem 2.4.1 of Kushner and Clark (1978) with $X_n = \tilde{\theta}_n - \theta$, $a_n = n^{-1}$, $\xi_n = (Y_{n-1}, \dots, Y_{n-p}, u_n)$, $h_0 = 0$, $B_n \equiv 0$, and $h(x, y) = Dg(y_1)y_1 \psi(y_1'x - y_2)$, where $y' = (y_1'y_2')$, $y_1 \in \mathbb{R}^p$, $y_2 \in \mathbb{R}$, and $x \in \mathbb{R}^p$. Kushner and Clark's assumptions A2.4.2 and A2.4.3' can be verified using the mixing properties of Y_n which we have established. Use $\theta(x) = K||x||$, $g_1(x, y) = 1$, and $g_2(\xi) = ||\xi||$ in A2.4.3'. The assumption that $\tilde{\theta}_n - \theta$ is bounded is established using the fact that

$$E^n \|\tilde{\theta}_{n+1} - \theta\|^2 \leq \|\tilde{\theta}_n - \theta\|^2 + n^{-2} K$$

and Theorem 1 of Robbins and Siegmund (1971).

We have now verified A1 to A7, so by Theorem 3.1,

$$n^{1/2}(\tilde{\theta}_n - \theta) = -n^{-1/2} \sum_{k=1}^n (k/n)^{\text{CDS-I}} DW_k + o(n^{-\epsilon})$$

for some $\epsilon > 0$. With this approximation, we could investigate the asymptotic behavior of $\tilde{\theta}_n$ rather thoroughly. Here, however, we will restrict attention to asymptotic normality. Define

$$\ddagger_n = [E\psi(-u_1)] [g(Z_n)^2 Z_n Z_n']$$

and

$$\ddagger = E\ddagger_n .$$

Since $\tilde{\theta} \rightarrow \theta$ and W_n is bounded uniformly in n ,

$$(4.4) \quad E^n W_n W_n' - \ddagger_n \rightarrow 0$$

by the bounded convergence theorem. Now $(\ddagger_n - \ddagger)$ is a mixingale of size $-q$ for all $q > 0$, so by Corollary 1.8 of McLeish (1975),

$$\sum_{k=1}^{\infty} k^{-(1/2+\epsilon)} Q_k (\ddagger_n - \ddagger) Q_k' \text{ converges}$$

for any $\epsilon > 0$ and any sequence of uniformly bounded matrices Q_k . Thus by (4.4) and the same argument used to conclude (3.12) from (3.13) and

(3.14), we see that

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n (k/n)^{cDS-I} (E^k W_k W_k') ((k/n)^{cDS-I}), \\
 &= \lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n (k/n)^{cDS-I} \frac{1}{n} ((k/n)^{cDS-I}), \\
 &= \int_0^1 (t^{cDS-I} D \frac{1}{n} D' (t^{cDS-I}))' dt \\
 &= V(c,D,S) \quad , \quad \text{say} .
 \end{aligned}$$

(The integral does not diverge at 0 since $\lambda_*(cDS - I) > -\frac{1}{2}$.) By the martingale central limit theorem of, for example, Brown (1971, Theorem 2),

$$n^{\frac{1}{2}}(\tilde{\theta}_n - \theta) \xrightarrow{D} N(0, V(c,D,S)) .$$

(A functional central limit theorem can also be established by the same theorem of Brown.) The Lindeberg condition is easily established because W_n is uniformly bounded. If ψ and g are fixed so that c and S are fixed, then $D = (cS)^{-1}$ is optimal, that is,

$$V(c,D,S) - V(c, (cD)^{-1}, S)$$

is p.s.d. This can be seen by noting that

$$\begin{aligned}
 \int_0^1 t^{cDS-I} D \frac{1}{n} D' (cS)^{-1} dt &= (cDS)^{-1} \int_{t=0}^{t=1} t^{cDS} D \frac{1}{n} D' (cS)^{-1} \\
 &= (cS)^{-1} \frac{1}{n} (cS)^{-1}
 \end{aligned}$$

so that

$$\int_0^1 (t^{cDS-I} D - (cS)^{-1}) \frac{1}{n} (t^{cDS-I} D - (cS)^{-1})' dt = V(c,D,S) - V(c, (cS)^{-1}, S) ,$$

and the integrand of the LHS is p.s.d.

Note that

$$(4.5) \quad V(c, (cS)^{-1}, S) = c^{-2} S^{-1} \dagger S^{-1} \\ = \frac{E\psi^2(-u_1)}{(E\dot{\psi}(-u_1))^2} (Eg(Z_1)Z_1Z_1')^{-1} (Eg^2(Z_1)Z_1Z_1') (Eg(Z_1)Z_1Z_1')^{-1} .$$

Also, if $D = (cS)^{-1}$, then

$$n^{1/2}(\tilde{\theta}_n - \theta) = -n^{-1/2} \sum_{k=1}^n (cS)^{-1} W_k + o(n^{-\epsilon}) \\ = n^{1/2} \sum_{k=1}^n (cS)^{-1} \psi(u_n) g(Z_n) DZ_n + o_p(1) ,$$

so by Carroll and Ruppert (1981, Theorem 1), $\tilde{\theta}_n$ is asymptotically equivalent to the bounded influence regression estimate, $\tilde{\theta}$, which solves

$$0 = \sum_{i=1}^n \psi(Y_i - Z_i' \tilde{\theta}) g(Z_i) Z_i = 0 .$$

Of course, c and S will generally not be known *a priori*. One can, however, estimate them by, e.g.,

$$\tilde{c}_n = n^{-1} \sum_{i=1}^n \dot{\psi}(Y_i - Z_i' \tilde{\theta}_i)$$

and

$$\tilde{S}_n = n^{-1} \sum_{i=1}^n g(Z_i) Z_i Z_i'$$

and then use $(\tilde{C}_n S_n)^{-1}$ for D in (4.1). These estimates can be defined recursively. The asymptotic behavior of the resulting stochastic approximation process is still an open problem, but once consistency is established, then general results of Section 3 should be applicable.

If one is not concerned with robustness, then for any value of ψ , $g(x) \equiv 1$ is optimal, but of course our proof assumes that $xg(x)$ is bounded, so it would need to be modified for this choice of g . For this g and its optimal D , i.e. $D = (c(EZ_1 Z_1'))^{-1}$, $S = \dagger$ so that $n^{1/2} \tilde{\theta}_n$ has asymptotic variance-covariance

$$c^{-2} \dagger^{-1} = (E\psi^2(u_1)/(E\dot{\psi}(u_1))^2) (EZ_1 Z_1')^{-1}.$$

To show that this g is optimal, we must show that for any real-valued g ,

$$(4.6) \quad (Eg(Z_1)Z_1 Z_1')^{-1} (Eg^2(Z_1)Z_1 Z_1') (Eg(Z_1)Z_1 Z_1')^{-1} - (EZ_1 Z_1')^{-1}$$

is p.s.d., but (4.6) is the variance-covariance matrix of

$$[g(Z_1)(Eg(Z_1)Z_1 Z_1')^{-1} - (EZ_1 Z_1')^{-1}]Z_1.$$

The use of $g(x) = 1$ allows outlying Z_n to have great influence on the estimation process, which could be disastrous if any of these Z_n have actually been observed with gross error, or if in some other way the pair (Y_n, Z_n) is an outlier. The optimal choice of g subject to the bound

$$||xg(x)|| \leq B \text{ for all } x,$$

where B is a fixed constant, has been considered by Hampel (1978) in a (slightly) different context, bounded influence linear regression.

Regardless of the value of g , we see from (4.5) that ψ should be chosen to minimize $E(\psi^2(u_1))/(E\dot{\psi}(u_1))^2$. If F has a density f

satisfying the typical regularity conditions of maximum likelihood estimation, then of course $\psi = \dot{f}/f$ is optimal. If one believes that f is only "close" to some known f_0 , then one should choose ψ robustly, and for this purpose the reader is referred to Huber (1964, 1981) for an introduction to robust estimation.

RECURSIVE NONLINEAR REGRESSION.

Now we consider the nonlinear regression model

$$Y_i = F(v_i, r) + e_i ,$$

where F is a known function from $\mathbb{R}^q \times \mathbb{R}^p$ to \mathbb{R} , v_i is a known vector of covariates, r is an unknown parameter vector, and e_i is a random noise term. Usually r is estimated using least squares, that is, by finding \tilde{r} which minimizes

$$\sum_{i=1}^n (Y_i - F(v_i, \tilde{r}))^2 .$$

However, just as with our last example, we may wish to use recursively defined estimators. We will study a particular recursive estimation sequence introduced by Albert and Gardner (1967), and we will show that its asymptotic distribution is the same as least squares, at least under certain regularity conditions. (Both the least squares estimator and Albert and Gardner's recursive estimator can be robustified as in the previous example.)

Define $h(v, w)$ to be the gradient $(\partial/\partial w)F(v, w)$, and $h_n(w) = h(v_n, w)$. Let \tilde{H}_0 be a p.d. $p \times p$ matrix, let \tilde{r}_1 be in \mathbb{R}^p , and define \tilde{H}_n and \tilde{r}_n by

$$\tilde{H}_n = n^{-1}(\tilde{H}_0 + \sum_{i=1}^n h_i(\eta_i)h_i'(\eta_i)) ,$$

where $\eta_i = \tilde{r}_\ell$ for some $\ell \leq i$, and

$$(4.7) \quad \tilde{r}_{n+1} = \tilde{r}_n - n^{-1}[F(v_n, \tilde{r}_n) - Y_n]\tilde{H}_n^{-1} h_n(\tilde{r}_n) .$$

Notice that

$$(4.8) \quad \tilde{H}_{n+1} = \tilde{H}_n - n^{-1}[\tilde{H}_n - h_n(\eta_n)h_n'(\eta_n)] ,$$

and if we define $B_n = n\tilde{H}_n$, then

$$B_{n+1}^{-1} = B_n^{-1} - \frac{B_n^{-1} h_{n+1}(\eta_{n+1})h_{n+1}'(\eta_{n+1})B_n^{-1}}{1 + h_{n+1}'(\eta_{n+1})B_n^{-1} h_{n+1}(\eta_{n+1})}$$

(Albert and Gardner (1967, page 111)). Consequently, $B_n^{-1} = n^{-1}\tilde{H}_n^{-1}$ can be calculated recursively, and the inversion of \tilde{H}_n is not necessary at each stage of the recursion. However, (4.8) is useful to our study of asymptotic properties of \tilde{r}_n and \tilde{H}_n .

First, we list the assumptions we will use. These assumptions are not the weakest possible, since we intend only to illustrate the content of our general results.

Let col be the function mapping the set of $p \times p$ matrices into \mathbb{R}^{p^2} by stacking columns one below the other.

C1. The functions D_1, D_2, D_3 , and D_4 map \mathbb{R}^p into \mathbb{R}^+ , $\mathbb{R}^{p \times p}$, $\mathbb{R}^{p \times p^2}$, and $\mathbb{R}^{p^2 \times p}$, respectively. Let H be a p.d. $p \times p$ matrix. The following hold for all v, w , and p.d. matrices M :

$$(4.9) \quad \begin{aligned} & || [F(v,w) - F(v,r)] M^{-1} h(v,w) - H^{-1} h(v,w) h'(v,w) (w-r) || \\ & \leq KD_1(v) (||w - r||^2 + ||M - H||^2) , \end{aligned}$$

$$(4.10) \quad \begin{aligned} & || M^{-1} h(v,w) - H^{-1} h(v,r) - D_2(v) (w-r) - D_3(v) \text{col}(M-H) || \\ & \leq KD_1(v) (||w - r||^2 + ||M - H||^2) , \end{aligned}$$

$$(4.11) \quad \begin{aligned} & || \text{col}[h(v,w)h'(v,w) - h(v,r)h'(v,r)] - D_4(v) (w-r) ||^2 \\ & \leq KD_1(v) ||w - r||^2 . \end{aligned}$$

C2. Let $\xi'_n = (v'_n, e_n)$. Suppose $R > 2$ and $\{\xi_n\}$ is strong mixing with mixing numbers $\{\alpha_n\}$ which are of size $-R(R-2)^{-1}$. McLeish (1975) defines size; here we mention McLeish's remark that $\{\alpha_n\}$ is of size $-q$ if $\{\alpha_n\}$ is monotonely decreasing and $\sum_{n=1}^{\infty} \alpha_n^\phi < \infty$ for some $\phi < 1/q$. For all n , the marginal distribution of ξ_n satisfies

$$E(e_n | v_n) = 0 ,$$

and

$$(4.12) \quad E h_n(r) h'_n(r) = H \text{ for all } n .$$

Moreover,

$$\begin{aligned} & \sup_n E \{ D_1^R(v_n) + ||D_2(v_n)||^{2R} + ||D_3(v_n)||^{2R} + ||D_4(v_n)||^R \\ & \quad + ||e_n||^{2R} + ||h_n(r)||^{2R} \} < \infty , \end{aligned}$$

and there exists a $p^2 \times p$ matrix D_4 such that

$$E D_4(v_n) = D_4 \text{ for all } n .$$

C3. $\tilde{r}_n \rightarrow r$ and $\tilde{H}_n \rightarrow H$.

REMARKS. Equations (4.9) to (4.11) merely state that certain Taylor series expansions are valid. Assumption C3 can be verified using the results of Ruppert (1981b) to prove that $\tilde{r}_n \rightarrow r$, and then using (4.12) and a continuity assumption on h to show that $\tilde{H}_n \rightarrow H$.

Using C2 and (2.3) of McLeish (1975) with $p = 2$, one sees that $\{D_1(v_n)\}$, $\{e_n D_2(v_n)\}$, $\{e_n D_3(v_n)\}$, $\{||e_n D_2(v_n)||\}$, $\{||e_n D_3(v_n)||\}$, $\{D_4(v_n)\}$, and $\{h_n(r)h'_n(r)\}$, when centered at expectations, are each mixingales of size $-1/2$. Moreover, $E(e_n D_2(v_n)) = E(e_n D_3(v_n)) = 0$.

Therefore, the following series converge by Corollary 1.8 of McLeish:

$$\sum_1^{\infty} n^{-1} e_n D_k(v_n) \text{ for } k = 2, 3 ,$$

$$\sum_1^{\infty} n^{-1} (||e_n D_k(v_n)|| - E||e_n D_k(v_n)||) \text{ for } k = 2, 3 ,$$

$$\sum_1^{\infty} n^{-1} (D_1(v_n) - ED_1(v_n))$$

$$\sum_1^{\infty} n^{-1} (D_4(v_n) - D_4) ,$$

$$\sum_1^{\infty} n^{-1} \{||D_4(v_n) - D_4|| - E||D_4(v_n) - D_4||\} ,$$

$$\sum_1^{\infty} n^{-(1/2 + \delta)} e_n h_n(r) \text{ for all } \delta > 0 ,$$

and

$$\sum_1^{\infty} n^{-(1/2 + \delta)} (h_n(r)h'_n(r) - H) \text{ for all } \delta > 0 .$$

Define

$$A_n = \begin{pmatrix} H^{-1} h_n(r) h_n'(r) + e_n D_2(v_n) & e_n D_3(v_n) \\ D_4(v_n) & I_{p^2} \end{pmatrix}.$$

By C1 and C3,

$$\begin{pmatrix} \tilde{r}_{n+1} - r \\ \text{col}(\tilde{H}_{n+1} - H) \end{pmatrix} = (I_{p(p+1)} - n^{-1} A_n) \begin{pmatrix} \tilde{r}_n - r \\ \text{col}(\tilde{H}_n - H) \end{pmatrix} - n^{-1} \begin{pmatrix} e_n H^{-1} h_n(r) \\ \text{col}[H - h_n(r) h_n'(r)] \end{pmatrix} + v_n,$$

where $\|v_n\| = o(D_1(v_n)(\|\tilde{r}_n - r\|^2 + \|\tilde{H}_n - H\|^2))$. Therefore, assumptions A1 to A7 hold with

$$\tilde{\theta}_n = \begin{pmatrix} \tilde{r}_n \\ \text{col} \tilde{H}_n \end{pmatrix}, \quad A = \begin{pmatrix} I_p & 0 \\ D_4 & I_{p^2} \end{pmatrix},$$

and $W_n = e_n H^{-1} h_n(r)$. It is easy to see that all eigenvalues of A are 1. Moreover, the $p \times p$ matrix in the top left-hand corner of n^{A-I} is I_p . Therefore, we have shown that C1 to C3 imply the existence of an $\varepsilon > 0$ such that

$$n^{\frac{1}{2}}(\tilde{r}_n - r) = -n^{-\frac{1}{2}} H^{-1} \sum_{k=1}^n e_k h_k(r) + o(n^{-\varepsilon}).$$

REFERENCES

ALBERT, A.E. and GARDNER, L.A., JR. (1967). *Stochastic Approximation and Nonlinear Regression*. The M.I.T. Press, Cambridge, Mass.

- BLUM, JULIUS R. (1954). Multidimensional stochastic approximation methods. *Ann. Math. Statist.* 25 737-744.
- BROWN, B.M. (1971). Martingale central limit theorems. *Ann. Math. Statist.* 42 59-66.
- CAMPBELL, KATHERINE (1979). Stochastic approximation procedures for mixing stochastic processes. Ph.D. Dissertation, The University of New Mexico, Albuquerque, N.M.
- CARROLL, RAYMOND J. and RUPPERT, DAVID (1980). Weak convergence of bounded influence estimates with applications. *Institute of Statistics Mimeo Series #1285*, University of North Carolina, Chapel Hill, N.C.
- CHUNG, K.L. (1954). On a stochastic approximation method. *Ann. Math. Statist.* 25 463-483.
- DERMAN, C. and SACKS, J. (1959). On Dvoretzky's stochastic approximation theorem. *Ann. Math. Statist.* 30 601-605.
- DOOB, JOSEPH L. (1953). *Stochastic Processes*. John Wiley and Sons, New York.
- FABIAN, VÁCLAV (1968). On asymptotic normality in stochastic approximation. *Ann. Math. Statist.* 39 1327-1332.
- HAMPEL, FRANK R. (1978). Optimally bounding the gross-error-sensitivity and the influence of position in factor space. *Proceedings of the American Statistical Association*, Statistical Computing Section.
- HIRSCH, MORRIS and SMALE, STEPHEN (1974). *Differential Equations, Dynamical Systems, and Linear Algebra*. Academic Press, New York.
- HUBER, PETER J. (1964). Robust estimation of a location parameter. *Ann. Math. Statist.* 35 73-101.

- HUBER, PETER J. (1981). *Robust Statistics*. John Wiley and Sons, New York.
- KIEFER, JACK and WOLFOWITZ, JACOB (1952). Stochastic estimation of the maximum of a regression function. *Ann. Math. Statist.* 23 462-466.
- KUSHNER, H.J. and CLARK, D.S. (1978). *Stochastic Approximation Methods for Constrained and Unconstrained Systems*. Springer-Verlag, New York.
- LJUNG, L. (1978). Strong convergence of a stochastic approximation algorithm. *Ann. Statist.* 6 680-696.
- MCLEISH, D.L. (1975). A maximal inequality and dependent strong laws. *Ann. Prob.* 3 829-839.
- PHAM, TUAN D. and TRAN, LANH T. (1980). The strong mixing property of the autoregressive moving average time series model. Technical Report, Department of Mathematics, Indiana University.
- ROBBINS, H. and MONRO, S. (1951). A stochastic approximation method. *Ann. Math. Statist.* 22 400-407.
- ROBBINS, HERBERT and SIEGMUND, DAVID (1971). A convergence theorem for non negative almost supermartingales and some applications. In *Optimizing Methods in Statistics* (J.S. Rustagi, ed.) 233-257. Academic Press, New York.
- RUPPERT, DAVID (1978). Almost sure approximations to the Robbins-Monro and Kiefer-Wolfowitz processes with dependent noise. *Institute of Statistics Mimeo Series #1203*, University of North Carolina, Chapel Hill, N.C. (To appear in *Ann. Prob.*)
- RUPPERT, D. (1979). A new dynamic stochastic approximation procedure. *Ann. Statist.* 7 1179-1195.

- RUPPERT, D. (1981a). Stochastic approximation of an implicitly defined function. To appear in *Ann. Statist.*
- RUPPERT, DAVID (1981b). Convergence of recursive estimators with applications to nonlinear regression. *Institute of Statistics Mimeo Series #1333*, University of North Carolina, Chapel Hill, N.C.
- SACKS, JEROME (1958). Asymptotic distributions of stochastic approximation procedures. *Ann. Math. Statist.* 29 373-405.
- SIELKEN, ROBERT L., JR. (1973). Stopping times for stochastic approximation procedures. *Z. Wahrscheinlichkeitstheorie und verw. Gebiete* 26 67-75.