

# Convergence Rate Determination for Gradient-based Adaptive Estimators\*

ROBERT R. BITMEAD,† BRIAN D. O. ANDERSON† and TUNG SANG NG‡

*In order to evaluate the performance of adaptive estimators it is necessary to quantify their convergence rate. For small adaptation gains stochastic averaging principles may be applied to determine this convergence rate as a function of the regression vector process.*

**Key Words**—Adaptive systems; averaging theory; convergence; parameter estimation; recursive algorithms; stochastic systems.

**Abstract**—A convergence rate estimate is derived for the homogeneous gradient-based adaptive linear estimator algorithm. This estimate involves the eigenvalues of the regression vector covariance matrix, yielding a useful measure for the choice of input signals for adaptive parameter estimation. The connection between this criterion and those more familiar from nonadaptive system identification is made and comparisons are drawn between the two areas.

## 1. INTRODUCTION

THE STANDARD adaptive linear parameter estimation problem is to take two time series,  $y_k$  of scalars and  $X_k$  of  $N$ -vectors and to attempt to fit the linear model  $\hat{y}_k = X_k' \theta$ , where  $\theta$  is an  $N$ -vector parameter, to minimize a criterion such as  $E(y_k - \hat{y}_k)^2$ . Many recursive algorithms for performing the adaptive estimation of this parameter have the form

$$\hat{\theta}_{k+1} = \hat{\theta}_k + \mu b_k(y_k - X_k' \hat{\theta}_k) \quad (1.1)$$

where  $\hat{\theta}_k$  is the parameter estimate and  $b_k$  is a vector function of the  $X_k$  data sequence. The simplest and most representative form of these algorithms is the LMS scheme where  $b_k = X_k$ . We shall examine this algorithm in particular here as the analysis carries over quite directly to many other choices of  $b_k$ , one of which will be instanced in the sequel.

One of the analytical problems of examining adaptive linear parameter estimation algorithms is that, in applications involving noises, time variations

or other additive extraneous signals, the parameter estimates do not necessarily converge to a constant in a deterministic or a probabilistic sense. Rather, subject to certain stationarity assumptions, the estimates will converge in a distributional manner (Bitmead, 1983). The development of performance measures for these algorithms often consists of determining or estimating parameters of this limiting distribution such as certain moments. Further, one is also frequently concerned with the quantitative dependence of these measures on the designer-adjustable parameters of the estimation scheme. There appears to be two distinct methodologies for proceeding with this performance analysis. Each approach has its own advantages and, thankfully, produces consistent results with the alternative method.

The first technique, taking a cue from identification theory, examines the adaptive estimation algorithm as a forced linear difference equation for the parameter estimate, e.g.

$$\hat{\theta}_{k+1} = \hat{\theta}_k + \mu X_k(y_k - X_k' \hat{\theta}_k) \quad (1.2)$$

or, in terms of parameter error  $\tilde{\theta}_k$ ,

$$\tilde{\theta}_{k+1} = (I - \mu X_k X_k') \tilde{\theta}_k + \mu X_k n_k \quad (1.3)$$

where the signal  $n_k$  embodies all the additive unmodelled disturbances. This approach has been adopted by Abu El Ata (1982), Farden (1981), Farden *et al.* (1980), Kim and Davisson (1975) and Macchi and Eweda (1983) and, while being technically very sound, does rely very heavily on stationarity and mixing assumptions on the underlying processes. The extension of these results to nonstationary applications requires some compromise because of the violation of the stationarity assumption.

The alternative approach is to consider initially

\* Received 14 December 1984; revised 18 September 1985. The original version of this paper was presented at the IFAC World Congress on a Bridge Between Control Science and Technology, in Budapest, Hungary, July 1984. The Published Proceedings of this IFAC Meeting may be ordered from Pergamon Press Limited, Headington Hill Hall, Oxford OX3 0BW, U.K. This paper was recommended for publication in revised form by Associate Editor R. Vinter under the direction of Editor P. C. Parks. This work was supported by the Radio Research Board of Australia.

† Department of Systems Engineering, Research School of Physical Sciences, Australian National University, Canberra, Australia.

‡ Department of Electrical Engineering, University of Wollongong, Wollongong, NSW, Australia.

only the homogeneous part of (1.3), to prove its exponential asymptotic stability and then to infer performance properties of the forced system. The extension of these results to nonstationary systems is relatively straightforward. However, there is frequently a need for the introduction of further conditions and assumptions to make the inference above. This method has been adopted in Bitmead and Anderson (1980), Johnson and Anderson (1981), Mendel (1973), Sondhi and Mitra (1976), Widrow *et al.* (1976) and Weiss and Mitra (1979), and will be further applied here.

In studying the properties of nonhomogeneous algorithms like (1.3) with the aid of the homogeneous algorithm

$$v_{k+1} = (I - \mu X_k X_k^T) v_k \tag{1.4}$$

it is desirable primarily to ensure the convergence of  $v_k$  to zero and then, secondarily, to determine or approximately quantify the convergence rate. All the available general nonstationary performance measure approximations rely (some implicitly) on knowing the exponential convergence rate of the homogeneous system and utilising bounded input/bounded output properties of linear systems.

Here we shall derive results for the case of stationary ergodic  $\{X_k\}$  which demonstrate that the exponential convergence rate of (1.4) is quantified by the second moment of the  $X_k$  process. Conditions are then derived which relate this quantification to the spectrum of  $X_k$ . Having determined a convergence rate in terms of the  $X_k$  process we then turn to consider these conditions in equation error and output error schemes where  $X_k$  is composed of system inputs and outputs. The experiment design problem is considered where one asks the question: given that we may choose our process inputs  $u_k$  which then determine the process outputs  $y_k$  and these jointly determine  $X_k$ , what is a sensible choice of  $u_k$  to guarantee a good convergence rate of an adaptive parameter estimator? The response to this enquiry is related to the equivalent problem in nonadaptive system identification.

2. CONVERGENCE RATE QUANTIFICATION

If the homogeneous LMS algorithm (1.4) with  $\{X_k\}$  stationary and ergodic is iterated over  $m$  time steps then we have

$$v_{n+m} = \prod_{i=n}^{n+m-1} (I - \mu X_i X_i^T) v_n \tag{2.1}$$

By stationarity we need only focus on  $n = 1$  and, writing  $A_i = X_i X_i^T$  and expanding, this yields.

$$\begin{aligned} \prod_{i=1}^m (I - \mu A_i) &= I - \mu \sum_{i=1}^m A_i + \mu^2 \sum_{i=1}^m \sum_{j=i+1}^m A_j A_i \\ &\quad - \mu^3 \sum_{i=1}^m \sum_{j=i+1}^m \sum_{k=j+1}^m A_k A_j A_i \\ &\quad + \dots + (-\mu)^m A_m A_{m-1} \dots A_1 \tag{2.2} \\ &= I - \mu m \bar{A} + (\mu m) \frac{1}{m} \sum_{i=1}^m (\bar{A} - A_i) \\ &\quad + (\mu m)^2 \frac{1}{m^2} \sum_{i=1}^m \sum_{j=i+1}^m A_j A_i \\ &\quad + \dots + (-\mu m)^m \frac{1}{m^m} A_m A_{m-1} \dots A_1 \tag{2.3} \end{aligned}$$

where  $\bar{A} = E(A_i) = E[X_i X_i^T]$ . We may now rewrite (2.1) using  $z_k = v_{mk+1}$  as

$$z_{k+1} = (I - \mu m \bar{A} + Q_k) z_k \tag{2.4}$$

where  $Q_k$  consists of the additional terms in the right hand side of (2.3).

The procedure taken will be to determine the behaviour of (2.4) by invoking the following theorem. Specifically, we are considering the convergence of  $v_k$  as  $k \rightarrow \infty$  and will study this by demonstrating (for a particular fixed  $m$ ) the convergence of  $z_k$  as  $k \rightarrow \infty$ .

*Theorem 1* (Bitmead and Anderson, 1981). The linear difference equation

$$w_{k+1} = \alpha_k w_k$$

with ergodic coefficient matrix  $\alpha_k$  will be exponentially asymptotically stable if  $\beta = E\|\alpha_k\| < 1$ . Further, this exponential rate is at least as fast as  $\beta^k$ .

The application of this result to (2.4) will require the bounding of the effects of the random term  $Q_k$ , and this is carried out by considering  $\mu$  to be small but not vanishing and  $m$  to be large but finite. Additional simple arguments are then necessary to show that properties of  $z_k$  are reflected in properties of  $v_k$ .

We may use the following result to bound some of these terms.

*Lemma 1.* Given a positive integer  $M$ , suppose that the moments of  $A_i$  up to the  $M$ th exist, i.e. there exist constants  $M, \delta < \infty$  such that

$$E\|A_i^M\| \leq M \delta^k \tag{2.5}$$

for all  $k \leq m$ . Then the  $l$ th term in (2.2) satisfies

$$E \left\| \frac{1}{(m)^i} \sum_{i=1}^m \sum_{j=i+1}^m \cdots \sum_{h=k+1}^m A_h A_k \cdots A_j A_i \right\| \leq \frac{1}{i!} M \delta^i. \tag{2.6}$$

Consequently, we have in (2.4)

$$E \|Q_k\| \leq E \|(\mu m) \frac{1}{m} \sum_{i=1}^m (\bar{A} - A_i)\| + M \sum_{i=2}^m \frac{(\mu m \delta)^i}{i!}. \tag{2.7}$$

The proof of this lemma, presented in the Appendix, involves a simple application of the triangle and Hölder inequalities.

Now it remains to bound the first term on the right-hand side of (2.7) and for this we invoke the following result presented in the Appendix as an extension of the work of Ibragimov and Linnik (1971).

*Lemma 2.* Let  $A_i$  be stationary with mean  $\bar{A}$  and let  $f^{ij}(\lambda)$  be the spectral density of the  $i$ - $j$  component of  $A_i - \bar{A}$  which we assume to exist and to be twice differentiable at  $\lambda = 0$ . Let

$$S_m = \frac{1}{m} \sum_{i=1}^m (A_i - \bar{A}). \tag{2.8}$$

Then

$$E \left[ \|S_m^{ij}\|^2 \right] = \frac{2\pi f^{ij}(0)}{m} + O \left( m^{-3/2} \left[ \frac{d^2}{d\lambda^2} f^{ij}(\lambda) \right]_{\lambda=0} \right) + O(m^{-7/4}) \tag{2.9}$$

and thus by Hölder's inequality,

$$E \left\| \frac{1}{m} \sum_{i=1}^m (\bar{A} - A_i) \right\| = O(m^{-1/2}). \tag{2.10}$$

The results of Lemmas 1 and 2 may now be combined to yield the main theorem.

*Theorem 2.* Subject to the hypotheses of Lemmas 1 and 2 we may always find a sufficiently small but constant value of the gain  $\mu$  so that the homogeneous adaptive estimation algorithm (1.4) converges to zero exponentially fast with a rate arbitrarily close to  $(1 - \mu\alpha)^k$  where  $\alpha$  is the minimum eigenvalue of  $\bar{A} = E[X_k X_k^T]$ .

*Proof.* Examining (2.4) in the light of Theorem 1 we see that the convergence rate of the homogeneous algorithm (1.4) is determined by

$$E \|I - \mu m \bar{A} + Q_k\| \leq \|I - \mu m \bar{A}\| + E \|Q_k\| = 1 - \mu m \lambda_{\min}(\bar{A}) + E \|Q_k\| \tag{2.11}$$

(assuming  $\mu m \lambda_{\max}(\bar{A}) < 1$ ). Now observe that

$$\begin{aligned} \sum_{i=2}^m \frac{(\mu m \delta)^i}{i!} &= (\mu m \delta)^2 \sum_{i=0}^{m-2} \frac{(\mu m \delta)^i}{(i+2)!} \\ &\leq (\mu m \delta)^2 \sum_{i=0}^{m-2} \frac{(\mu m \delta)^i}{i!} \\ &\leq (\mu m \delta)^2 e^{\mu m \delta} \\ &\leq (\mu m \delta)^2 e \end{aligned} \tag{2.12}$$

provided that  $\mu m \delta < 1$ . Consequently, using (2.10) and (2.12) in (2.7) we obtain for some constants  $K_1$  and  $K_2$

$$E \|Q_k\| \leq \mu m^{1/2} K_1 + (\mu m \delta)^2 K_2. \tag{2.13}$$

We now choose positive constants  $\alpha_1, \alpha_2$  such that  $\alpha_1 + \alpha_2 < 1$  and select  $m$  sufficiently large that

$$m^{1/2} K_1 \leq \alpha_1 m \lambda_{\min}(\bar{A}) \tag{2.14}$$

and then choose  $\mu$  sufficiently small that

$$(\mu m \delta)^2 K_2 \leq \alpha_2 \mu m \lambda_{\min}(\bar{A}) \tag{2.15}$$

$$\mu m \lambda_{\max}(\bar{A}) < 1. \tag{2.16}$$

These choices then ensure that

$$E \|I - \mu m \bar{A} + Q_k\| \leq 1 - (1 - \alpha_1 - \alpha_2) \mu m \lambda_{\min}(\bar{A}). \tag{2.17}$$

The choice of the constants  $\alpha_1, \alpha_2, m$  and  $\mu$  has been involved above. If  $M$  and  $\delta$  depend on  $m$  (as will normally be the case) then  $m$  must be chosen first to satisfy (2.14) for the chosen  $\alpha_1$ . Subsequently,  $\mu$  is chosen satisfying (2.15) and (2.16), where, provided we are prepared to take  $\mu$  very small, it is apparent that  $\alpha_2$  can be made arbitrarily small. In particular, by allowing  $\mu$  very small we can choose  $(1 - \alpha_1 - \alpha_2)$  arbitrarily close to one.

The theorem then follows by relating  $v_k$  in (1.4) to  $z_k$  in (2.4), to show that the convergence rate is arbitrarily close to  $(1 - \mu m \alpha)^k/m$ . Application of the binomial theorem then establishes the theorem.

To maximize the evaluated convergence rate by choice of  $\alpha_1$  and  $\alpha_2$  above, one may solve (2.14) replaced by equality for  $m$  and (2.15) replaced by equality for  $\mu$  and substitute into (2.17). The optimizing choice is  $\alpha_1 = 0.5, \alpha_2 = 0.25$ .

In the next section we shall discuss the implications of Theorem 2 in the problem of input selection for adaptive estimation, i.e. experiment design. However, our task now will be to relate the convergence conditions of the theorem statement to the more usual requirements of covariance decay (Abu El Ata, 1982; Farden *et al.*, 1980; Macchi and

Eweda, 1983) or mixing (Bitmead and Anderson, 1980) of the  $X_k$ -process.

Lemma 1 simply imposes requirements for the  $X_k$ -process to be bounded or at least have finite moments of sufficiently large orders. These restrictions only relate to the distribution of  $X_k$ . Lemma 2 however imposes indirectly a covariance decay condition. To see this consider

$$\begin{aligned} \left. \frac{d^2 f(\lambda)}{d\lambda^2} \right|_{\lambda=0} &= \frac{d^2}{d\lambda^2} \frac{1}{2\pi} \sum_{k=-\infty}^{\infty} e^{-jk\lambda} R_k \Big|_{\lambda=0} \\ &= \frac{-1}{2\pi} \sum_{k=-\infty}^{\infty} k^2 R_k \end{aligned} \quad (2.18)$$

where  $R_k$  is the autocorrelation function of  $A_i$ . The existence of the second derivative of  $f(\lambda)$  is thus a requirement that the sum in (2.18) be well defined. For this to occur it is sufficient that

$$\sum_{k=-\infty}^{\infty} k^2 |R_k| < \infty$$

so that the differentiability restriction may be interpreted as a covariance decay condition on  $A_i$ . Since  $A_i = X_i X_i'$  we in turn have a type of mixing requirement on  $X_i$ . This is entirely consistent with the previously known convergence requirements for these algorithms.

### 3. EXPERIMENT DESIGN FOR ADAPTIVE ESTIMATION

In practice, the algorithm gain  $\mu$  is typically chosen to be very small (usually  $< 0.001$  divided by the input signal power). The reasons for this choice of  $\mu$  are twofold. Firstly, it may be demonstrated (Abu El Ata, 1982; Farden *et al.*, 1980; Macchi and Eweda, 1983) that for stationary systems the asymptotic parameter error variance is linear in  $\mu$  so that time-invariant steady state performance requirements dictate that  $\mu$  should be small in order that the parameter estimates be smooth. Secondly, in order that the adaptive estimation problem with slowly time-varying parameters be well defined one requires that the estimator time constants should be sufficiently longer than the system time constants, but also sufficiently short compared to the time constants of parameter variation. The choice of  $\mu$  reflects the assumed relationships between these latter time constants and compromises estimate smoothness against parameter tracking. [To see the relationship between system time constants and estimator gain selection recall that the conditions of the previous section involved  $\mu m$  being small, where  $m$  reflected the convergence rate of the Cesaro sums in (2.8) and hence reflected the system time constants. The choice of  $m$  is thus related to the covariance of the input process.] The condition number of  $\bar{A}$  also plays

an important role in determining the maximum allowable value of  $\mu$ . Convergence requirements in the independent case demand  $\mu \lambda_{\max}(\bar{A}) < 1$  (Widrow *et al.*, 1976) while the convergence rate analysis of the previous section would favour  $\mu \lambda_{\min}(\bar{A})$  close to one. It is well known that a poor condition number severely affects the performance of gradient-based minimization methods (Luenberger, 1973) and this is clearly reflected here in these constraints on  $\mu$ .

Given that  $\mu$  is chosen small on the grounds above and that Theorem 2 applies, we are justified in attempting to choose an input process to maximize the convergence rate. The question of adjusting  $\mu$  is not addressed in this section but rather the questions of input selection after a particular choice of  $\mu$ . In order to maximize this convergence rate it is necessary to maximize the minimum eigenvalue of  $\bar{A}$  subject to input power constraints and we shall next address certain problems in achieving this.

Before proceeding we note that altering  $\bar{A}$  through manipulation of input signals can also lead to alteration of the error terms  $Q_k$  in (2.4) and this may complicate matters. However, as was indicated in the proof of Theorem 2, if  $\mu m$  is guaranteed small then we may neglect the role of this  $Q_k$  variation in the analysis.

There is a correspondence here between the notions of input choice for the maximization of the convergence rate of an adaptive estimator and of input choice for achieving the best efficiency of a parameter estimator. This latter problem is that of experiment design and is much studied in the area of system identification (Goodwin and Payne, 1977).

There is a natural performance benchmark in experiment design, the Cramer-Rao lower bound on the covariance of the estimates, and optimal experiment design in this context is directed towards ensuring that, since the Cramer-Rao bound is given by the inverse of the Fisher information matrix  $M$ , the input process spectrum extremizes a suitable scalar function of  $M$  in an attempt to minimize the parameter error covariance. (This function above typically involves an expectation taken over the *a priori* parameter distribution.)

For finite impulse response or moving average systems one typically has

$$X_k = (u_k, u_{k-1} \dots u_{k-N})' \quad (3.1)$$

where  $u_k$  is the scalar system input process. This is the well known adaptive transversal filter, and the best choice of input  $u_k$  to maximize convergence rate is an independent and identically distributed sequence. Performance in this case has been closely studied (Widrow *et al.*, 1976). When dealing with infinite impulse response or autoregressive moving average systems,

$$X_k = (y_{k-1}, y_{k-2} \dots y_{k-N} u_k, \dots u_{k-N})', \quad (3.2)$$

so that the actual system parameters have an implicit effect in determining the best input sequence  $u_k$  for adaptive parameter estimation.

The Cramer–Rao lower bound is met if  $E(\hat{\theta}\hat{\theta}') = M$  where, for moving average systems,  $M$  consists of the covariance matrix of the input process scaled by the measurement noise variance, while for autoregressive moving average systems it comprises the scaled covariance of a vector composed of input and output processes. In our notation here

$$M = E(X_k X_k') = \bar{A} \quad (3.3)$$

and the correspondence between the adaptive and nonadaptive schemes is clear, although it should be stated here that the estimate of convergence rate in the adaptive case is not necessarily tight and may be conservative. In spite of this disclaiming rider it is apparent that, for a given adaptive estimation problem with a particular system and small  $\mu$ , the estimate convergence rate will be enhanced by maximizing the minimum eigenvalue of  $\bar{A}$ .

What is important in this context is that the adaptive experiment design problem of maximizing convergence rate has a very similar formulation to that of system identification experiment design thus allowing the application of well-studied input selection procedures from this area.

There are several common scalar performance measures of  $M$  (Goodwin and Payne, 1977; Qureshi and Ng, 1982) and each appears to have some advantageous properties and some negative ones. Maximizing the trace of  $M$  has been used but this can lead to parameter nonidentifiability (Grewal and Glover, 1975). Maximizing  $\det M$  (or more commonly minimizing  $-\log \det M$ ) has better features as well as independence from parameter scaling and an interesting frequency domain signal to noise ratio interpretation (Qureshi and Ng, 1982). However, because  $\det M = 0$  if and only if  $M$  is singular, this criterion can be optimized by identifying perfectly the projection of  $\theta$  in any one particular direction, regardless of estimator performance in other directions. To this extent, this criterion does permit the trading off of identifier performance in one direction to the detriment of other directions.

The performance measure being advanced here is to maximize the minimum eigenvalue of  $M$ , which is equivalent to minimizing the maximum eigenvalue of the parameter error covariance. Thus the criterion attempts to control mean squared parameter error in every direction. It does suffer from being scaling dependent and analytically difficult to work with, however.

It certainly is not true that adaptive estimators of the form of (1.1) are efficient in the statistical sense

but rather that convergence rate improvement for these adaptive schemes is mathematically similar to nonadaptive experiment design. Thus the parameter error covariance is not equal to the Cramer–Rao lower bound of  $M^{-1}$ . In the case of independent  $X_k$  it is readily shown that the asymptotic parameter error covariance is given to first order in  $\mu$  by  $\frac{\mu\sigma^2}{2}$ , where  $\sigma^2$  is the variance of the Wiener residuals  $y_k - X_k'\theta$ , which is not the same as  $M^{-1}$ .

An alternative algorithm may be constructed which is closer in spirit to the asymptotically efficient methods used in nonadaptive systems identification. This alternative algorithm is based on least squares estimation:

$$\hat{\theta}_{k+1} = \hat{\theta}_k + \mu M^{-1} X_k (y_k - X_k' \hat{\theta}_k) \quad (3.4)$$

or

$$\tilde{\theta}_{k+1} = (I - \mu M^{-1} X_k X_k') \tilde{\theta}_k + \mu M^{-1} X_k n_k \quad (3.5)$$

where  $n_k$  is the Wiener residual. This algorithm may be written

$$M^{1/2} \tilde{\theta}_{k+1} = (I - \mu M^{-1/2} X_k X_k' M^{-1/2}) M^{1/2} \tilde{\theta}_k + \mu M^{-1/2} X_k n_k \quad (3.6)$$

for which one sees that the convergence rate will be linear in  $\mu$  and independent of  $M$ , while the asymptotic error covariance may be approximated (assuming independent  $X_k$  and to first order in  $\mu$ ) by the solution  $P$  to

$$P = (I - \mu I) P (I - \mu I) + \mu^2 M^{-1} \sigma^2 \quad (3.7)$$

or  $P = \mu \frac{\sigma^2}{2} M^{-1}$ . This modified algorithm (3.4) again demonstrates the connection with the experiment design problem as well as the trade-off between adaptive estimator convergence rate and parameter error covariance.

One criticism frequently levelled at experiment design is that, as the optimal design depends on unknown system parameters, there is a circularity in its justification. This criticism is frequently answered by using a Bayesian approach and averaging over a prior distribution; however in the general area of adaptive estimation (in contrast to nonadaptive system identification) where tracking of slowly varying parameters is often desired, one frequently will have a reasonable knowledge of the system parameter values thus allowing effective "causal" experiment design.

For parameter estimation problems for autoregressive moving average systems, particularly in output error or adaptive control where the regressor

(3.2) is a function of the parameter estimate, the convergence rates given by the preceding theory are only locally applicable in the neighbourhood of some parameter value admitting a stationary linearized description. Consequently, in these instances our results pertain to the local behaviour of the adaptive estimator rather than to its global (nonlinear) transient properties. This limitation, however, fits in with the idea of slow adaptation (small  $\mu$ ) and the dynamics of tracking slowly time-varying parameters.

4. CONCLUSIONS

We have presented a method for deriving a lower bound on the exponential convergence rate of the homogeneous part of gradient-based adaptive parameter estimators. The convergence rate is dependent on the minimum eigenvalue of the covariance matrix of the  $X_k$  process and, as a result, if the input signal is to be chosen to enhance the convergence rate, a useful criterion of performance is the magnitude of this smallest eigenvalue.

The comparison was drawn between this performance measure for adaptive estimators and the experiment design criteria for nonadaptive system identification. The important point here was to show that while these two problems have seemingly disparate origins there are considerable similarities and the input selection rules from identification may be carried over essentially intact to adaptive estimation.

We have concentrated specifically on the LMS gradient-based adaptive estimation algorithm because it provides performance criteria which are usually representative of the class of such schemes as was evidenced in Section 3. Equally, the exact analysis of the nonhomogeneous equations was not presented with dependent inputs, since the effect of dependence appears most critically in the convergence rate of the homogeneous equation, and the nonhomogeneous error covariances can be related to those in the independent  $X_k$  theory.

The importance of the results presented here is that they allow the approximate quantification of convergence rate and, more importantly, they allow the development of useful input selection procedures to maximize this rate. This represents a significant extension from the usual "persistence of excitation" (Bitmead and Anderson, 1980) requirements and would be readily interfaced with other restricted input designs for adaptive systems (Ioannou and Kokotovic, 1983). Our theory derives convergence rates for small values of  $\mu$  determined by covariance decay values—recall that  $\mu m$  was required to be small. This means that the adaptive estimation difference equation (1.4) is driven by a small "fast" term. An analogy can be drawn between these results and results of Khas'minskii (1980) and Arnold *et al.*

(1983) on the principle of averaging and its application to the stability of differential equations depending on a small parameter. Again there are similarities in the results of this theory and the diffusion approximation work of Kushner and Huang (1981) and Benveniste and Ruget (1982). But there is a crucial difference in that the results here pertain to small but not infinitesimal values of  $\mu$  while the diffusion results are asymptotic in character as  $\mu \rightarrow 0$ . The apparent consistency of the results produced by these diverse methods lends credence to their validity.

APPENDIX

Proof of Lemma 1

The difficult component here is the demonstration of (2.6). Using the Hölder inequality repeatedly we have

$$E\|A_h A_k \dots A_l\| \leq E^{1/l} \|A_h\|^l E^{1/l} \|A_k\|^l \dots E^{1/l} \|A_l\|^l \leq M \delta^l \quad \text{from (2.5)}$$

Thus the term on the left-hand side of (2.6) is bounded above by

$$\begin{aligned} \frac{1}{m^l} \sum_{i=1}^m \sum_{j=i+1}^m \sum_{h=k+1}^m M \delta^l &\leq \frac{M \delta^l}{m^l} \int_0^m \int_i^m \dots \int_k^m dh . dk \dots di \\ &= \frac{M \delta^l}{m^l} \int_0^m \int_0^i \dots \int_0^k dh . dk \dots di \\ &= \frac{M \delta^l}{l!} \end{aligned}$$

Proof of Lemma 2

The proof of Theorem 18.2.1 of Ibragimov and Linnik (1971) shows that if a stationary scalar sequence  $X_i$  is given, and one defines

$$T_m = \sum_{i=1}^m X_i$$

and if  $f(\lambda)$ , the spectral density function of  $X$ , is continuous at  $\lambda = 0$ , then

$$|\text{var } T_m - 2\pi f(0)m| \leq 2\pi m \max_{|\lambda| \leq m^{-1/4}} |f(\lambda) - f(0)| + O(m^{-1/2}). \quad \text{(A.1)}$$

Now if  $f(\lambda)$  has first and second derivatives at  $\lambda = 0$  then, since  $f'(0) = 0$ , by Taylor's theorem

$$|f(\lambda) - f(0)| = \frac{\lambda^2}{2} f''(0) + O(\lambda^3)$$

and (A.1) yields

$$\frac{\text{var } T_m}{m^2} = \frac{2\pi f(0)}{m} + O[\pi m^{-3/2} f''(0)] + O(m^{-7/4}) + O(m^{-5/2})$$

and the Lemma follows by identifying  $T_m/m$  with  $S_m^{ij}$ ,  $X_i$  with the  $i$ - $j$  component of  $A_i - \bar{A}$ , and  $f(\lambda)$  with  $f^{ij}(\lambda)$ .

REFERENCES

Abu El Ata, S. (1982). Asymptotic behaviour in an adaptive estimation algorithm with application to M-dependent data. *IEEE Trans. Aut. Control*, AC-27, 1255-1257.  
 Arnold, L., H. Crauerl and V. Wihstutz (1982). Stabilization of linear systems by noise. *SIAM J. Control Optimiz.*, 21, 451-461.  
 Benveniste, A. and G. Ruget (1982). A measure of the tracking capability of recursive stochastic algorithms with constant gains. *IEEE Trans. Aut. Control*, AC-27, 639-649.

- Bitmead, R. R. (1981). Convergence properties of LMS adaptive estimators with unbounded dependent inputs. *Proc. 20th IEEE Conf. Decision and Control*, San Diego, CA, pp. 607–612.
- Bitmead, R. R. (1983). Convergence in distribution of LMS-type adaptive parameter estimates. *IEEE Trans. Aut. Control*, **AC-28**, 54–60.
- Bitmead, R. R. and B. D. O. Anderson (1980). Performance of adaptive estimation algorithms in dependent random environments. *IEEE Trans. Aut. Control*, **AC-25**, 788–794.
- Bitmead, R. R. and B. D. O. Anderson (1981). Adaptive frequency sampling filters. *IEEE Trans. Circuits Syst.*, **CAS-28**, 524–534.
- Farden, D. C. (1981). Tracking properties of adaptive signal processing algorithms. *IEEE Trans. Acoust., Speech, Signal Processing*, **ASSP-29**, 439–446.
- Farden, D. C., J. C. Goding Jr. and K. Sayood (1980). On the 'desired behaviour' of adaptive signal processing algorithms. *Proc. 1979 IEEE Int. Conf. Acoust., Speech, Signal Processing*, pp. 466–496.
- Goodwin, G. C. and R. L. Payne (1977). *Dynamic System Identification: Experiment Design and Data Analysis*. Academic Press, New York.
- Grewal, M. S. and K. Glover (1975). Relationships between identifiability and input selection. *Proc. IEEE Conf. Decision and Control*, Phoenix AZ, pp. 526–528.
- Ibragimov, I. A. and Yu. V. Linnik (1971). *Independent and Stationary Sequences of Random Variables*. Wolters-Noordhoff, Groningen.
- Ioannou, P. A. and P. V. Kokotovic (1983). *Adaptive Systems with Reduced Models*. Springer, Berlin.
- Johnson Jr., C. R. and B. D. O. Anderson (1981). Sufficient excitation and stable reduced-order adaptive IIR filtering. *IEEE Trans. Acoust., Speech, Signal Processing*, **ASSP-29**, 1212–1215.
- Jones, S. K. (1973). Adaptive linear estimation for stationary M-dependent processes. Ph.D. dissertation, Dept. Elec. Eng., Southern Methodist University, Dallas, TX.
- Khas'minskii, R. Z. (1980). *Stochastic Stability of Differential Equations*. Sijthoff & Noordhoff, Alphen aan den Rijn.
- Kim, J.-K. and L. D. Davisson (1975). Adaptive linear estimation for stationary M-dependent processes. *IEEE Trans. Inform. Theory*, **IT-21**, 23–31.
- Kushner, H. J. and H. Huang (1981). Asymptotic properties of stochastic approximations with constant coefficients. *SIAM J. Control Optimiz.*, **19**, 87–105.
- Luenberger, D. G. (1973). *Introduction to Linear and Nonlinear Programming*. Addison-Wesley, Reading, MA.
- Macchi, O. and E. Eweda (1983). Second order convergence analysis of stochastic adaptive linear filtering. *IEEE Trans. Aut. Control*, **AC-28**, 76–85.
- Mendel, J. M. (1973). *Discrete Techniques of Parameter Estimation: The Equation Error Formulation*. Marcel Dekker, New York.
- Qureshi, Z. H. and T. S. Ng (1976). Parameter estimation: the  $D_s$ -optimality case. *SIAM J. Control Optimiz.*, **20**, 713–721.
- Sondhi, M. M. and D. Mitra (1976). New results on the performance of a well-known class of adaptive filters. *Proc. IEEE*, **64**, 1583–1597.
- Weiss, A. and D. Mitra (1979). Digital adaptive filters: conditions for convergence, rates of convergence, effects of noise and errors arising from the implementation. *IEEE Trans. Inform. Theory*, **IT-25**, 637–652.
- Widrow, B., J. M. McCool, M. G. Larimore and C. R. Johnson Jr. (1976). Stationary and nonstationary learning characteristics of the LMS adaptive filter. *Proc. IEEE*, **64**, 1151–1162.