

Robert R. Bitmead  
 Department of Electrical Engineering  
 University of Newcastle  
 New South Wales, 2308  
 Australia

Brian D.O. Anderson  
 Department of Electrical Engineering  
 University of Newcastle  
 New South Wales, 2308  
 Australia

Abstract

A stochastic algorithm, familiar from adaptive estimation, is introduced and its homogeneous part is shown to be exponentially convergent for a wide class of inputs, which need not be stationary. The implications of this convergence rate for the non-homogeneous algorithm in practical situations are qualitatively examined and a possible approach to improving performance in use is suggested.

1. Introduction

In problems of parameter estimation for stochastic dynamic systems, iterative algorithms have found great application and, indeed, are at the heart of any recursive (on-line) procedure for the adaption of initial estimates to 'true' parameter values and the tracking of time-varying parameters. Their effectiveness in a broad variety of practical situations is well known and the volume of statistical and engineering literature regarding their convergence is substantial.

We consider here a variant of the familiar stochastic approximation scheme. It is well known that this scheme converges with probability one and in mean square for statistically stationary processes provided the convergence parameters tend to zero at specified rates. For the case of slowly time-varying systems it is often desirable to have these parameters tend to a small constant, say  $\mu > 0$ , in order that changes with time may be tracked. We concentrate on the performance of algorithms with a constant  $\mu$ .

As a background for our theoretical study we introduce the following filtering problem.

A time series  $\{x_k\}$  and a time series  $\{y_k\}$  are known to be related in some way, and we seek to model their relation via a moving-average equation of the type

$$y_k = w_0 x_k + w_1 x_{k-1} + \dots + w_{N-1} x_{k-N+1}$$

$$= (x_k \ x_{k-1} \ \dots \ x_{k-N+1}) (w_0 w_1 \ \dots \ w_{N-1})^T$$

$$\hat{y}_k = \hat{x}_k^T \hat{w}_k^* \quad (1)$$

The value of  $w^*$  is unknown, and in fact the real relation between  $\{x_k\}$  and  $\{y_k\}$  may not be of this form. A value of  $w^*$  is sought; just what the significance of  $w^*$  might be in case the real situation cannot be modelled as above is a question we leave for the moment in abeyance, except to comment that (1) would be in some way an approximation, or even best approximation to the true situation.

Now suppose that at time  $k$  an estimate  $w_k$  of  $w^*$  is available. Defining an error sequence  $e_k$  by

$$e_k = y_k - \hat{x}_k^T w_k$$

[which is  $\hat{x}_k^T (w^* - w_k)$  in case (1) is exact] we adjust  $w_k$  by some function of  $e_k$  to obtain an updated estimate  $w_{k+1}$ . The adjustment procedure is so designed that in case (1) is an exact model of reality,  $w_k \rightarrow w^*$ , and the convergence properties of the algorithm are fairly easily established. We are however interested in how these algorithms will behave in case (1) is not an accurate description of reality, because for example  $N$  is chosen too small, or  $w^*$  is (slowly) time-varying.

The behaviour of the algorithms has been considered by many authors, theoretically and experimentally. For example, for the LMS algorithm

$$w_{k+1} = w_k + \mu e_k x_k \quad (2)$$

there have been several results proved regarding convergence when  $N$  is too small. Widrow et al [1] prove many useful results on the convergence of  $w_k$  to the Wiener solution in the case that  $\{x_k\}$  is an independently and identically distributed process\* and, although this is a very restrictive assumption on  $\{x_k\}$ , they quantify several rules of thumb.

\* Since  $x_k = [x_k \ x_{k-1} \ \dots \ x_{k-N+1}]^T$ , it is apparently impossible when  $N > 1$  to secure independence of  $x_k, x_p$  for all  $k \neq p$ . However, work of [2] shows how independence of  $\{x_k\}$  can be relaxed to independence of  $\{x_k\}$ .

†Work supported by Australian Research Grants Committee.

Kim and Davisson [3] have proved that for a stationary M-dependent process† {x<sub>k</sub>} the mean squared error between w<sub>k</sub> and the Wiener solution can be made arbitrarily close to zero by choosing μ small enough. And finally Daniell [4] has proved similar results to [3] under assumptions of ergodicity of second moments, boundedness of conditional fourth moments, and asymptotic independence of {x<sub>k</sub>}. In [3,4] the performance of the LMS scheme under time varying conditions is only inferred from convergence, although the analysis of [1] presents a detailed investigation of nonstationary behaviour under the assumption of independent {x<sub>k</sub>}. This investigation shows the LMS algorithm to be exponentially convergent in the mean and the subsequent results follow largely from this.

We follow here a similar pattern: we present the particular algorithm

$$w_{k+1} = w_k + \mu e_k \frac{x_k}{\|x_k\|^2} \quad (3)$$

and demonstrate in Section II that, for certain broad classes of {x<sub>k</sub>} which encompass dependent sequences, (3) is exponentially convergent to w with probability one in the case that {y<sub>k</sub>} is actually a moving average of order N of {x<sub>k</sub>}. Then in Section III we investigate the consequences of exponential convergence of the homogeneous part of (3) in the situations of time variation and of too small an N. Here we utilise the bounded input/bounded output property of an exponentially stable system with a driving term. Section IV contains conclusions and directions for future research.

The study of deterministic adaption algorithms with exponential rates of convergence has been considered by Morgan and Narendra, by Kriesslmeier and by Anderson in [5], [6] and [7,8] and the necessary and sufficient conditions for this rate can be seen to correspond closely to the sufficient conditions that are derived in Section II for the exponential convergence of the stochastic algorithm. Both latter authors have used the deterministic algorithms to devise exponentially convergent adaptive identifiers and observers.

By exponential convergence of a random variable {z<sub>k</sub>} to zero we mean that the related random variable {(1+β)<sup>k</sup>z<sub>k</sub>} for some β > 0 converges to zero as k → ∞. In this case we say that z<sub>k</sub> converges exponentially fast with exponent less than or equal to -ln(1+β). Kushner [9] has a definition of exponential convergence of a random variable which is equivalent to demanding that {(1+β)<sup>k</sup>z<sub>k</sub>} be a positive supermartingale.

Nagumo and Noda [2] examine (3) for the homogeneous case where y<sub>k</sub> = x<sub>k</sub><sup>\*</sup>.

† M-dependence is defined as follows. Let A and B be two sets of integers with min A - max B > M. Then {x<sub>α</sub> | α ∈ A} and {x<sub>β</sub> | β ∈ B} are independent.

$$w_{k+1} - w^* = w_k - w^* - \mu \frac{x_k x_k^*}{\|x_k\|^2} (w_k - w^*) \quad (4)$$

or, writing v<sub>k</sub> = w<sub>k</sub> - w<sub>k</sub><sup>\*</sup>,

$$v_{k+1} = (1 - \mu \frac{x_k x_k^*}{\|x_k\|^2}) v_k \quad (5)$$

They present general results on the convergence of (5) with probability one. They also prove exponential convergence when {x<sub>k</sub>} is i.i.d. but do not proceed past this stage.

Polyak [10,11] examines the convergence and convergence rate of stochastic algorithms from a general viewpoint and is able to prove exponential rates although his first assumption on the input process represents a benign form of independence. Again, he establishes convergence rates and convergence in various probabilistic senses but does not consider their implications in these papers.

Finally, Sondhi and Mitra [12] have demonstrated exponential bounds on the performance of a continuous-time version of (5) subject to a mixing condition on the input (this condition is very similar to the deterministic conditions of [5,8]). They present an excellent analysis of the effect of exponential convergence of (5) upon the behaviour of (3) in real situations. The main drawback of their analysis, however, is that their mixing condition is extremely severe and, contrary to their claims, one cannot guarantee that it is satisfied by almost all sample paths of a stationary and ergodic input process with positive definite covariance matrix. We present a much less restrictive condition for exponential convergence.

## II. Main Result

In this section we examine the convergence rate of the algorithm (5) of Section I. We demonstrate exponential convergence of (5) under an ergodicity assumption - both upper and lower bounds on the exponent are given. Then the pivotal steps of the proof are examined and an extension of the result is derived where the ergodicity requirement is relaxed.

**Theorem 1:** Let there be given an ergodic sequence {x<sub>k</sub>} and let the parameter N-vector sequence {v<sub>k</sub>} be derived from the algorithm

$$v_{k+1} = \begin{cases} v_k & \text{if } x_k = 0 - \text{an event of zero probability} \\ (1 - \mu \frac{x_k x_k^*}{\|x_k\|^2}) v_k & \text{if } x_k \neq 0 \end{cases}$$

with initial value v<sub>0</sub> (|v<sub>0</sub>| < ∞) and μ ∈ (0,2).

Here, x<sub>k</sub> = [x<sub>k</sub> x<sub>k-1</sub> ... x<sub>k-N+1</sub>]<sup>T</sup>. Then if E[x<sub>k</sub>x<sub>k</sub><sup>\*</sup>] is positive definite

$$|v_0| \lambda^k \leq |v_k| \leq |v_0| \bar{\lambda}^k$$

with probability 1 as  $k \rightarrow \infty$  for  $1 > \bar{\lambda} \geq \lambda > 0$

Here  $\lambda = |1 - \mu|$  and

$$\bar{\lambda} = \min_{n > (\mu N)^{-1}} (1 - \frac{\bar{\beta}}{n})^{\frac{1}{n}} \quad \text{where}$$

$$\bar{\beta}_n = \frac{1 - (1 - \mu)^2}{(1 + \mu)^2} E \left[ \lambda \min_{j=0}^{n-1} \begin{matrix} x_j x_{j+1} \\ x_j x_{j+1} \end{matrix} \right]$$

and  $\lambda_{\min}(\cdot)$  is the minimum eigenvalue.

**Proof:** See Appendix. The minimization carried out for  $\bar{\lambda}$  is not necessary to prove exponential convergence but simply allows us to find the least  $\bar{\lambda}$ . Actually, if  $\{x_k\}$  is ergodic and full rank, then there exists a number  $p$  such that  $0 < \beta < 1$  for all  $n > p$  and any of these  $n$  will do.

Nagumo and Noda [2] have demonstrated an exponential convergence rate for this algorithm in the case that  $\{x_k\}$  is a zero mean independent, identically distributed process. Clearly, in this case we may take  $n = N$  for  $\mu$  big enough in order that  $\bar{\beta}$  be strictly greater than zero. Also, if the algorithm is used with an  $M$ -dependent process  $\{x_k\}$  (see [3]) then exponential convergence is again easily demonstrated by considering the expectation over an interval of  $n \geq M + N + 1$ , again for sufficiently large  $\mu$ . The principal tactic of the proofs of [2], [3] and here, to demonstrate convergence, is to avoid the explicit introduction of conditional expectations and to then devise representations that rely upon ordinary expectations. This clearly makes the results more applicable because they only require knowledge of certain moments of the  $\{x_k\}$  whereas conditional expectations, being random variables themselves, need to have a distribution supposed on the  $\{x_k\}$  before any thoughts of analysis can be entertained. Perhaps the major drawback of the work by Danielli [4] is his introduced bounds on conditional fourth moments which often can not be assumed a priori.

We believe that this is one of the most useful general convergence results on this algorithm to date because it makes no assumptions on conditional expectations or on absolute coherence lengths of the input sequence. However, while ergodicity is a handy supposition for the proof of convergence rates it is clearly not necessary, and we are able to produce extensions of theorem 1 by replacing the ergodicity requirement with information about the coherence of the input.

The two pivotal steps of the proof of theorem 1 are, firstly, that the matrix  $\begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \\ \vdots & \vdots \\ x_{N-1} & x_N \end{bmatrix}$  has one eigenvalue of 1 and  $N-1$  of zero - this is used in line (16) to bound  $\begin{bmatrix} x_j & x_{j+1} \\ x_{j+1} & x_{j+2} \end{bmatrix}$  by  $\begin{bmatrix} u_j & u_{j+1} \\ u_{j+1} & u_{j+2} \end{bmatrix}$  - and, secondly, that under the ergodicity assumption the random variables  $\frac{1}{j} \sum_{m=0}^j \beta_{km}$  converge to the

expected value with probability one as  $j \rightarrow \infty$ .

The following result of Cramér and Leadbetter [13] (adapted for discrete-time systems) is of interest in that it gives us an alternate condition for a random variable to have a "time-varying ergodicity in the mean", in the sense that time averages tend to ensemble averages. More precisely, we have:

**Lemma 1:** [13, p.94]

Let  $\{z_k\}$  be a sequence of random variables with means  $\bar{z}_k$  and with covariance satisfying

$$|E\{(z_k - \bar{z}_k)(z_{k+\tau} - \bar{z}_{k+\tau})\}| \leq K \frac{k^\alpha + (k+\tau)^\alpha}{1+\tau}, \quad 0 \leq 2\alpha < \beta < 1 \quad (6)$$

Then  $\frac{1}{m} \sum_{i=1}^m z_i \rightarrow \frac{1}{m} \sum_{i=1}^m \bar{z}_i$  with probability one and

in mean square as  $m \rightarrow \infty$ . This lemma has an obvious application to our current investigation.

**Corollary 1:** If the  $\{x_k\}$  process is loosely correlated (at least up to fourth moments) so that the  $\{\beta_{kj}\}$  sequence satisfies (6), then the algorithm (5) will converge exponentially fast with probability one and in mean square.

This corollary represents a crucial observation because it ensures that with certain nonstationary or nonergodic inputs the validity of our proof of exponential convergence of the homogeneous part of (3) remains.

It should be noted here that the ergodicity assumption of theorem 1 and the loose correlation assumption of corollary 1, together with the assumption of full rank covariances, will imply that always there exists a finite  $k$  which has  $\bar{\beta}_k$  bounded away from zero. That is, the ergodicity assumption performs a dual role in the proof of exponential convergence.

Finally, we remark that, as the boundedness of eigenvalues of the updating matrix appears to be a critical observation for the proof of exponential convergence, it might be expected that the results above could be applicable to a wider class of algorithm than (3) and, indeed, that exponential convergence may hold for any linear decrescent homogeneous algorithm such as the Inverse State scheme of Kumar, Moore and Evans [14]. This hypothesis has not yet been proved, but the similarity between experimental results of different linear algorithms would suggest some underlying unifying structure.

### III. Performance Implications of Exponential Convergence of the Homogeneous Algorithm

As has been noted in [1] and [12], the exponential convergence of the homogeneous algorithm (5) allows us to quantify the behaviour of the algorithm (3). There are three main perturbations to (5) which occur in its use in (3). We assume their independence and treat each singly. This is usually justified.

(a) Time Variation in the 'true' Parameter vector  $w^*$ . We suppose here that  $w^*$  takes on

different values at different times - denote these  $\{w_k^*\}$ . Then we may write (5) as

$$v_{k+1} = [I - \mu \frac{x_k^T x_k}{x_k^T x_k}] v_k - [w_{k+1}^* - w_k^*] \quad (7)$$

where now  $v_k = w_k - w_k^*$ . This equation will have a bounded solution for  $\|v_k\|$  as  $k \rightarrow \infty$ , provided  $\|w_{k-1}^* - w_k^*\| \leq w_\Delta^*$  for some fixed  $w_\Delta^*$  and the conditions for exponential convergence are satisfied. Suppose that  $\|v_k\|$  converges to zero in (5) faster than  $(1-\beta)^k$ . Then as  $k \rightarrow \infty$  in (7)

$$\|v_k\| \leq \sum_{i=1}^k \|w_{i-1}^* - w_i^*\| (1-\beta)^{k-i} \|v_0\| \doteq \frac{w_\Delta^*}{\beta} \|v_0\|$$

As the choice of possible  $\beta$  is affected by the value of  $\mu$ , it is clear that the misalignment between  $w_k$  and  $w_k^*$  as  $k \rightarrow \infty$  is  $\mu$ -dependent. One may see from the expression for  $\beta$  from theorem 1, as  $\mu$  increases from zero,  $\beta$  increases to a certain maximum value and then decreases again. Subject to the constraint  $\mu \in (0,1]$ , it is straightforward to maximize  $\beta$ , in which case the bound on the error due to time-variation in  $w_k^*$  is minimized.

If  $w^*$  suffers a jump change in some or all of its elements, at one time only, then we have exponentially fast reconvergence to the new value. This property is of great importance for situations such as fault detection.

The requirement of bounded variations in  $w^*$  is hardly restrictive as the types of variation most commonly met with in practice are (i) step changes, e.g. due to component failure or the detection of a new signal source etc., (ii) slow periodic changes e.g. due to diurnal changes in environment such as with communication channels, (iii) small random variations. These variations may include noise in the parameters, which need not be zero mean, i.e. slow drift may also occur. In situations (i) and (ii) above, particularly for the fault detection problem, it would be desirable to have rapid convergence, while for (iii) this need not be so.

#### (b) Measurement Noise in $\{y_k\}$ .

Suppose that the sequence  $\{y_k\}$  is not measurable but the sequence  $\{z_k\}$  where  $z_k = y_k + n_k$  is measurable. Here  $\{n_k\}$  is some noise sequence assumed zero mean, white and independent from  $\{x_k\}$ . Then (3) becomes

$$v_{k+1} = (I - \mu \frac{x_k^T x_k}{x_k^T x_k}) v_k + \mu \frac{x_k}{x_k^T x_k} n_k \quad (8)$$

and under the assumption of exponential convergence in mean square of (5), which occurs provided the conditions of corollary 1 or theorem 1 are satisfied, it appears that a relation similar to the following must hold:

$$E \|v_{k+1}\|^2 \leq K(1-\beta)^{2k} E \|v_0\|^2 + \sum_{j=1}^k \mu K(1-\beta)^{2(k-j)} E \left[ \left( \frac{x_j^T x_j}{x_j^T x_j} \right)^{-1} \right] E \left[ n_j^2 \right] \quad (9)$$

So from (8) and (9) we can see that if  $n_k$  has bounded variance &  $E \left[ \left( \frac{x_j^T x_j}{x_j^T x_j} \right)^{-1} \right]$  is bounded then

$v_{k+1}$  has bounded variance. The bound on this variance is clearly  $\mu$ -dependent (and  $\beta$ -dependent). Moreover, the dependence is different to that applying to the analysis tied to (7).

#### (c) Truncation Errors.

In fitting an  $N^{\text{th}}$ -order moving average to  $\{y_k\}$  we may be neglecting the effect on  $y_k$  of values of  $\{x_k\}$  before  $x_{k-N+1}$ . That is, we may really have  $y_k = w_{k-N}^T x_k$  for  $w$  and  $x$  of dimension  $M > N$  ( $M$  not necessarily finite). Write this as  $y_k = w^1 T x_k^1 + w^2 T x_k^2$  where  $w^1$  has dimension  $N$  and  $x_k^1 = (x_k, x_{k-1}, \dots, x_{k-N+1})^T$  and  $w^2 T x_k^2$  represents the tail effects of  $w$  and  $x_k^2 = (x_{k-N}, x_{k-N-1}, \dots)$  on  $y_k$ .

Then (3) becomes

$$v_{k+1} = \left[ I - \mu \frac{x_k^T x_k}{x_k^T x_k} \right] v_k + \mu \frac{x_k^1}{x_k^T x_k} x_k^2 T w^2 \quad (10)$$

The additive term in (10) will clearly produce a misadjustment in  $v_{k+1}$  from 0 or 'bias'. This bias is dependent upon  $\mu, w$  and the statistics of the  $\{x_k\}$  and a thorough analysis will be presented elsewhere. However, with knowledge of the statistics of  $\{x_k\}$  and of the proportion of energy in the tail of  $w$  a bound should be derivable and the dependence of this bound on  $\mu$  determined. The bias here is of a different character to that arising in (a) and (b). It may even be helpful, in that it yields an  $N$ -vector  $w^*$  which defines a more accurate approximation of  $[w^1 T w^2 T]^T$  than does  $w^1$ .

The point to be made from (a), (b) and (c) above is that, in practice, the choice of the adaption constant  $\mu$  is critical to the ultimate performance of the scheme. This has long been realized and the approach we are heading towards is one in which the choice of  $\mu$  is set up as an optimization problem where the cost function takes into account tracking ability, error variance and bias.

This optimization approach to identification schemes using stochastic-approximation-like algorithms seems to us better motivated and suited to usage than the classic formulas of Robbins and Munro [15] and Kiefer and Wolfowitz [16] which sacrifice performance over finite intervals for strong convergence with zero error in the limit.

#### IV. Conclusion

We have presented a familiar algorithm from adaptive estimation and established that, provided the input sequence is ergodic or loosely correlated, the homogeneous part of the algorithm converges to zero exponentially fast. The convergence rate depends on the expected value of the minimum eigenvalue of a sum of random matrices. For full rank input sequences the convergence rate may be bounded away from being arbitrarily slow.

We then gave a mainly qualitative examination of the practical implications of exponential convergence when driving terms were added to the homogeneous algorithm. If these terms are bounded then the bounded input/bounded output property of an

exponentially stable system is used to infer boundedness of the error. Similarly, subject to independence assumptions, bounded variance inputs produce bounded variance errors. The performance of the algorithm was shown to be dependent upon  $\mu$ , the adaption constant, as are the effects of time variation of the parameter, measurement noise, and error involved in truncation of the parameter vector. These effects are not all the same, and the fact they are different offers scope for optimally selecting  $\mu$ .

The results of section III only start to quantify some heuristic notions about the performance of the adaptive algorithm; there is clearly still a need to develop general rules of thumb, applicable in a wide variety of situations, that will aid in choosing good values of  $\mu$  in practice. Among these rules could be a simple method of associating performance with  $\mu$  and the spectrum of  $\{x_k\}$  if the input were stationary.

Before these rules could be devised, it would be necessary to provide a more thorough quantitative treatment of the propagation of errors and particularly of the dependence of bias upon  $\mu$ .

Apart from extending the results as above we believe that the general method of establishing exponential convergence as in theorem 1 should be applicable to prove exponential convergence rates for other algorithms arising in identification.

## V. Appendix - Proof of Theorem 1

### (i) Upper Exponential Bound

We carry through the proof assuming  $\mu \in (0,1]$  rather than  $\mu \in (0,2)$  as it allows the notation to be more concise. The extension to  $\mu \in (0,2)$  is simple and requires just a few modulus signs to be added at various points.

From (5) we have the following results:

Let

$$u_k = v_k - v_0 = (v_k - v_{k-1}) + (v_{k-1} - v_{k-2}) + \dots + (v_1 - v_0) \\ = \sum_{j=0}^{k-1} \mu \frac{x_j x_j^T}{x_j^T x_j} v_j \quad (11)$$

By an extended application of Cauchy's inequality we have  $(\sum_{i=0}^{k-1} a_i)^T (\sum_{i=0}^{k-1} a_i) \leq k \sum_{i=0}^{k-1} a_i^T a_i$  which yields

$$u_k^T u_k \leq k \mu^2 \sum_{j=0}^{k-1} \frac{(x_j^T v_j)^2}{x_j^T x_j} \quad (12)$$

upon application to the inner product of (11) with itself.

The evolution equation for  $v_k$  also yields

$$v_{j-1}^T v_{j-1} - v_j^T v_j = [1 - (1-\mu)^2] \frac{(v_{j-1}^T x_{j-1})^2}{x_{j-1}^T x_{j-1}}$$

and this shows that  $v_{j-1}^T v_{j-1}$  is decrescent, a fact

used below. Also

$$v_0^T v_0 - v_{k-k}^T v_{k-k} = [1 - (1-\mu)^2] \sum_{j=0}^{k-1} \frac{(v_j^T x_j)^2}{x_j^T x_j} \quad (13)$$

and (12) implies

$$v_0^T v_0 - v_{k-k}^T v_{k-k} \geq [1 - (1-\mu)^2] \frac{u_k^T u_k}{k \mu^2} \quad (14)$$

Equation (14) will be used further below. Meanwhile, rewrite (13) using  $u_k = v_k - v_0$ .

$$v_0^T v_0 - v_{k-k}^T v_{k-k} = [1 - (1-\mu)^2] \sum_{j=0}^{k-1} \frac{[(u_j + v_0)^T x_j]^2}{x_j^T x_j}$$

so that, using  $\|a+b\| \geq \|a\| - \|b\|$  with (a)  $= u_j^T x_j$  and (b)  $= v_0^T x_j$ ,

$$(v_0^T v_0 - v_{k-k}^T v_{k-k})^{1/2} \geq [1 - (1-\mu)^2]^{1/2} \left\{ \left[ \sum_{j=0}^{k-1} \frac{(u_j^T x_j)^2}{x_j^T x_j} \right]^{1/2} + \left[ \sum_{j=0}^{k-1} \frac{(v_0^T x_j)^2}{x_j^T x_j} \right]^{1/2} \right\} \quad (15)$$

We next find a bound on the first sum of (15)

$$\sum_{j=0}^{k-1} \frac{(u_j^T x_j)^2}{x_j^T x_j} \leq \sum_{j=0}^{k-1} u_j^T u_j \quad (16)$$

$$\leq \sum_{j=0}^{k-1} \frac{v_0^T v_0 - v_{j-j}^T v_{j-j}}{[1 - (1-\mu)^2]} j \mu^2, \text{ using (14)}$$

$$\leq \frac{(v_0^T v_0 - v_{k-k}^T v_{k-k}) k^2 \mu^2}{[1 - (1-\mu)^2]}, \text{ because } v_{j-j}^T v_{j-j} \text{ is decrescent.}$$

Now (15) yields

$$(v_0^T v_0 - v_{k-k}^T v_{k-k})^{1/2} \geq -(v_0^T v_0 - v_{k-k}^T v_{k-k})^{1/2} k \mu \\ + [1 - (1-\mu)^2]^{1/2} \left[ v_0^T \sum_{j=0}^{k-1} \frac{x_j x_j^T}{x_j^T x_j} v_0 \right]^{1/2} \\ \geq -(v_0^T v_0 - v_{k-k}^T v_{k-k})^{1/2} k \mu \\ + [1 - (1-\mu)^2]^{1/2} (v_0^T v_0)^{1/2} \left[ \lambda_{\min} \sum_{j=0}^{k-1} \frac{x_j x_j^T}{x_j^T x_j} \right]^{1/2}$$

or

$$\frac{v_0^T v_0 - v_{k-k}^T v_{k-k}}{v_0^T v_0} \geq \frac{1 - (1-\mu)^2}{(1+k\mu)^2} \lambda_{\min} \left[ \sum_{j=0}^{k-1} \frac{x_j x_j^T}{x_j^T x_j} \right] \quad (17)$$

As  $\frac{1 - (1-\mu)^2}{(1+k\mu)^2} \geq \frac{1}{k^2 \mu^2}$  and  $\lambda_{\min} \left( \sum_{j=0}^{k-1} \frac{x_j x_j^T}{x_j^T x_j} \right) \leq \frac{k}{N}$  (using

the facts that  $\lambda_{\min}(A) \leq \frac{\text{trace } A}{\dim A}$ ,  $\text{trace } A+B = \text{trace } A + \text{trace } B$  and  $\text{trace} \frac{x_j x_j^T}{x_j^T x_j} = 1$ ,

provided  $kN > 1$ , the right hand side of (17), call it  $\delta_k$  say, is less than one and greater than or equal to zero.

The parameters  $\mu$  and  $N$  are of course fixed.

Let us also temporarily fix  $k$  so that  $k\mu N > 1$ . We consider the sequence  $v_0, v_{k-k}, v_{2k-2k}, \dots, v_{mk-mk}, \dots$  for interger  $m$ .

We have established above that  $v_{k-k} \leq (1-\beta_k) v_0$  and more generally we have

$$v_{mk-mk} \leq (1-\beta_{mk}) v_{(m-1)k} v_{(m-1)k}$$

Setting  $\tau_{m+1} = (1-\beta_{mk}) \tau_m$  and  $\tau_0 = v_0$  we observe that  $0 \leq v_{mk-mk} \leq \tau_m$  and prove an exponential convergence to zero of  $\tau_m$ .

Clearly 
$$\tau_m = \left[ \prod_{i=0}^{m-1} (1-\beta_{ik}) \right] \tau_0$$
 and taking logarithms

$$\ln \tau_m - \ln \tau_0 = \sum_{i=0}^{m-1} \ln(1-\beta_{ik}) \quad (18)$$

The requirement that  $\tau_m$  converge to zero regardless of initial conditions is that the right hand side of (18) diverge to minus infinity. Furthermore, if this term diverges faster than  $m \ln(1-\alpha)$  for  $\alpha(0,1)$  then  $\tau_m$  converges exponentially fast to zero with exponent less than  $\ln(1-\alpha)$ . These requirements are clearly satisfied provided  $\{\beta_{ik}\}$  is bounded away from zero (infinitely) often enough.

We now make the observation that if  $\{x_i\}$  is ergodic then so is  $\{\beta_{ik}\}$  and also  $\{\ln(1-\beta_{ik})\}$ . And, if  $\{\ln(1-\beta_{ik})\}$  is ergodic, then

$$m^{-1} \sum_{i=0}^{m-1} \ln(1-\beta_{ik}) \rightarrow E[\ln(1-\beta_{mk})]$$

with probability one as  $m \rightarrow \infty$ .

Applying Jensen's inequality

$$E[\ln(1-\beta_{mk})] \leq \ln(1-E(\beta_{mk})) = \ln(1-\bar{\beta}_k)$$

Further, if  $x_k$  is ergodic and has full rank covariance then  $\frac{1}{l} \sum_{i=0}^{l-1} \frac{x_i x_i^T}{x_i x_i^T}$  converges to some full rank matrix as  $l \rightarrow \infty$  so that there always exists some integer  $p$  such that, provided  $k > p$ ,  $\bar{\beta}_k > 0$ .

Thus  $\tau_m$  converges to zero with probability one and at an exponential rate faster than  $(1-\bar{\beta}_k)^m$ . Now, although the  $\{\tau_m\}$  sequence only captures every  $k$ th element of  $v_{l-l}^m$ , the fact that  $v_{l-l}^m$  is decrescent shows that  $v_{l-l}^m$  converges exponentially fast to zero with rate faster than

$$(1-\bar{\beta}_k)^{\frac{m}{k}}$$

(ii) Lower Exponential Bound

From (5) we have, since the eigenvalue of minimum magnitude of

$$(I - \mu \frac{x_k x_k^T}{x_k x_k^T})$$
 has magnitude  $|1-\mu|$ , that

$$\|v_{k+1}\| \geq |1-\mu| \|v_k\|$$

whence  $\|v_{k+1}\| \geq |1-\mu|^{k+1} \|v_0\|$ . Taking  $\underline{\lambda} = |1-\mu|$

we have the final part of the theorem, provided  $0 < \mu < 2$ .

VVV

REFERENCES

- [1] B. Widrow, J. McCool, M.G. Larimore and C.R. Johnson Jr., "Stationary and nonstationary learning characteristics of the LMS adaptive filter", *Proc. IEEE*, vol.64, No.8, August 1976, pp. 1151-1162.
- [2] J-I. Nagumo & A. Noda, "A Learning method for system identification", *IEEE Trans.Auto.Control* Vol.AC-12, No.3, June 1967, pp.282-287.
- [3] J-K. Kim & L.D. Davisson, "Adaptive Linear Estimation for stationary M-dependent processes", *IEEE Trans. Inf. Theory*, Vol.IT-21, No.1, Jan. 1975, pp. 23-31.
- [4] T.P. Daniell, "Adaptive estimation with mutually correlated training sequences", *IEEE Trans. Sy.Sci. & Cyb.*, Vol.SSC-6, No.1, January 1970, pp.12-19.
- [5] A.P. Morgan & K.S. Narendra, "On the Uniform Asymptotic Stability of Certain Linear Nonautonomous Differential Equations", *SIAM J.Cntrl. & Optimization*, Vol.15, No.1, Jan.1977, pp.5-24.
- [6] G. Kreisselmeier, "Adaptive observers with exponential rate of convergence", *IEEE Trans.Auto Cntrl.*, Vol.AC-22, No.1, Feb.1977, pp.2-8.
- [7] B.D.O. Anderson, "Exponential stability of linear equations arising in adaptive identification", *IEEE Trans.Auto.Cntrl.*, Vol.AC-22, No.1, February 1977, pp. 83-88.
- [8] B.D.O. Anderson, "Adaptive identification of multiple-input, multiple-output plants", *Proc. 1974 IEEE Conf. Dec. & Cntrl.*, pp.273-281.
- [9] H.J. Kushner, *Introduction to Stochastic Control*, Holt, Rinehart & Winston Inc, 1971.
- [10] B.T. Polyak, "Convergence & convergence rate of iterative stochastic algorithms, I. General Case", *Avt. i Telemekh.*, No.12, December 1976 pp. 1858-1868.
- [11] B.T. Polyak, "Convergence & convergence rate of iterative stochastic algorithms, II. The Linear case", *Avt. i Telemekh.*, No.4, April 1977, pp. 537-542.
- [12] M.M. Sondhi & D. Mitra, "New results on the performance of a well-known class of adaptive filters", *Proc.IEEE*, Vol.64, No.11, November 1976, pp. 1583-1597.
- [13] H. Cramér & M.R. Leadbetter, *Stationary and Related Stochastic Processes*, John Wiley and Sons, Inc., 1967
- [14] R. Kumar, J.B. Moore & R.J. Evans, "Inverse state stochastic approximation recursions", Preprint Elec. Eng. Dept. Univ. of Newcastle, Australia, submitted for publication.
- [15] H. Robbins & S. Munro, "A Stochastic approximation method", *Ann.Math.Stat.*, Vol. 22, 1951, pp. 400-407.
- [16] J. Kiefer & J. Wolfowitz, "Stochastic Estimation of the maximum of a Regression Function", *Ann.Math.Stat.* Vol.23, 1952, pp. 462-466.