211

# On the Statistical Efficiency of the LMS Algorithm with Nonstationary Inputs

BERNARD WIDROW, FELLOW, IEEE, AND EUGENE WALACH

*Abstract*—A fundamental relationship exists between the quality of an adaptive solution and the amount of data used in obtaining it. Quality is defined here in terms of "misadjustment," the ratio of the excess mean square error (mse) in an adaptive solution to the minimum possible mse. The higher the misadjustment, the lower the quality is. The quality of the exact least squares solution is compared with the quality of the solutions obtained by the orthogonalized and the conventional least mean square (LMS) algorithms with stationary and nonstationary input data. When adapting with noisy observations, a filter trained with a finite data sample using an exact least squares algorithms will have a misadjustment given by

$$M = \frac{n}{N} = \frac{\text{number of weights}}{\text{number of training samples}}.$$

If the same adaptive filter were trained with a steady flow of data using an ideal "orthogonalized LMS" algorithm, the misadjustment would be

$$M = \frac{n}{4\tau_{mse}} = \frac{\text{number of weights}}{\text{number of training samples}}.$$

Thus, for a given time constant $\tau_{mse}$ of the learning process, the ideal orthogonalized LMS algorithm will have about as low a misadjustment as can be achieved, since this algorithm performs essentially as an exact least squares algorithm with exponential data weighting. It is well known that when rapid convergence with stationary data is required, exact least squares algorithms can in certain cases outperform the conventional Widrow–Hoff LMS algorithm. It is shown here, however, that for an important class of nonstationary problems, the misadjustment of conventional LMS is the same as that of orthogonalized LMS, which in the stationary case is shown to perform essentially as an exact least squares algorithm.
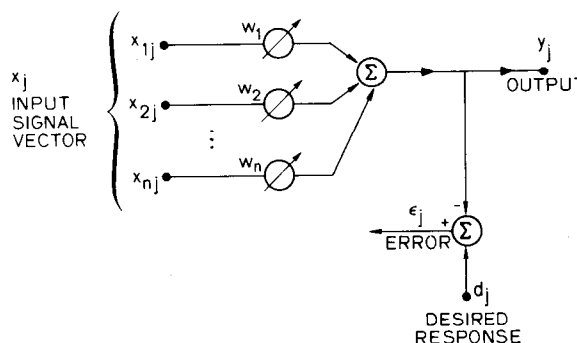
## I. INTRODUCTION

THE BASIC component of most adaptive filtering and signal processing systems is the adaptive linear combiner [1]–[4] shown in Fig. 1. An output signal is formed which is the weighted sum of a set of input signals. The output would be a simple linear combination of the inputs only if the weights were fixed. In actual practice, the weights are adjusted or adapted purposefully; the resulting weight values are signal dependent. This process causes the system behavior during adaptation to differ significantly from that of a linear system. However, when the adaptive process converges and the weights settle to essentially fixed values with only minor random fluctuations about the equilibrium solution, the converged system exhibits essentially linear behavior.
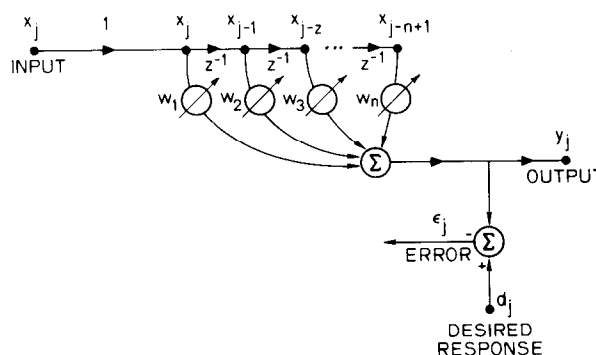
(a)



(b)

Fig. 1. Adaptive linear combiner and its application in an adaptive filter.

Adaptive linear combiners have been successfully used in the modeling of unknown systems [2], linear prediction [2], [5], adaptive noise cancelling [4], adaptive antenna systems [3], channel equalization systems for high-speed digital communications [6]–[9], echo cancellation [10]–[12], systems for instantaneous frequency estimation [13], receivers of narrow-band signals buried in noise (the "adaptive line enhancer") [4], [14], [15], adaptive control systems [16]–[18], and in many other applications.

In Fig. 1(a), the interpretation of the input signal vector $X_j = (x_{1j}, \cdots, x_{nj})^T$, and the desired response $d_j$ might vary depending on how the adaptive linear combiner is used. In Fig. 1(b), an application to adaptive filtering is shown. In turn, an application of adaptive filtering to plant modeling or system identification is shown in Fig. 2. Here, we can view the desired response $d_j$ as a linear combination of the last $n$ samples of the input signal vector, corrupted by a certain independent zero-mean plant noise $n_j$. Our
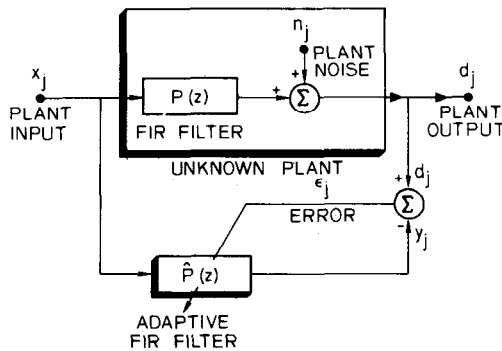
Fig. 2. Adaptive plant identification.

aim in this application is to estimate an unknown plant (represented by its transfer function $P(z) = w_1^* + \cdots + w_n^* z^{-n+1}$) through the minimization of the output error $\epsilon_j$ in the mean-square sense. For purposes of analysis, we consider the plant to be a transversal finite impulse response (FIR) filter. Referring to Figs. 1(a) and (b), the input signal vector at the $j$th sampling instant is designated by

$$X_j = (x_j, \cdots, x_{j-n+1})^T, \qquad (1)$$

and the set of the weights of adaptive transversal filter is designated by

$$W = (w_1, \cdots, w_n)^T. \qquad (2)$$

The $j$th output signal is

$$y_j = \sum_{i=1}^{n} w_i x_{j-i+1} = W^T X_j = X_j^T W. \qquad (3)$$

The input signals and desired response are assumed to be stationary ergodic processes. Denoting the desired response as $d_j$, the error at the $j$th time is

$$\epsilon_j = d_j - y_j = d_j - W^T X_j = d_j - X_j^T W. \qquad (4)$$

The square of this error is

$$\epsilon_j^2 = d_j^2 - 2d_j X_j^T W + W^T X_j X_j^T W. \qquad (5)$$

The mean-square error (mse) $\xi$, the expected value of $\epsilon_j^2$, is

$$\xi \triangleq E\left[\epsilon_j^2\right]$$
$$= E\left[d_j^2\right] - 2E\left[d_j X_j^T\right] W + W^T E\left[X_j X_j^T\right] W$$
$$= E\left[d_j^2\right] - 2P^T W + W^T R W, \qquad (6)$$

where the cross-correlation vector between the input signals and the desired response is defined as

$$E\left[d_j X_j\right] = E\begin{bmatrix} d_j x_j \\ \vdots \\ d_j x_{j-n+1} \end{bmatrix} \triangleq P, \qquad (7)$$

and the input autocorrelation matrix $R$ is defined as

$$E\left[X_j X_j^T\right] = E\begin{bmatrix} x_j x_j & \cdots & x_j x_{j-n+1} \\ x_{j-1} x_j & \cdots & x_{j-1} x_{j-n+1} \\ \vdots & & \vdots \\ x_{j-n+1} x_j & \cdots & x_{j-n+1} x_{j-n+1} \end{bmatrix} \triangleq R. \qquad (8)$$

It can be observed from (6) that with stationary ergodic inputs the mean square error performance function is a quadratic function of the weights, a paraboloidal "bowl" that has a unique minimal point for

$$W = W^* = R^{-1} P. \qquad (9)$$

In practice, of course, we do not know the exact statistics $R$ and $P$. One way of finding an estimate of the optimal weight vector $W^*$ would be to estimate $R$ and $P$ for the given input and given desired response. This approach would lead to what is called an exact least mean square solution. This approach is optimal in a sense that the sum of square errors will be minimal for the given data sample. However, such solutions are generally somewhat complex from the computational point of view [19]–[21]. On the other hand one can use one of the simpler gradient search algorithms such as the least mean square (LMS) steepest descent algorithm of Widrow and Hoff. However, this algorithm is sometimes associated with a certain deterioration in performance with problems for which there exists great spread between the eigenvalues of the autocorrelation matrix $R$ (see, for instance, [20]–[21]).

In order to establish a bridge between the LMS and the exact least squares approaches mentioned above, we will introduce an idealized orthogonalized LMS algorithm. For the implementation of this algorithm, we will have to assume perfect knowledge of the autocorrelation matrix $R$. Naturally that means that this idealized algorithm cannot be used in practice. However, its performance provides a convenient theoretical benchmark for the sake of comparison.[1]

In the next section we will analyze briefly the performance of the exact least squares solution when the weights are obtained with a finite data sample. Then in Sections III and IV we will analyze the idealized LMS algorithm and show, at least heurstically, that its performance is equivalent to that of an exact least squares algorithm. Based on this heuristic argument, we will view the idealized orthogonalized LMS process as an "optimal" gradient search algorithm.

In Section V we will define a class of nonstationary problems: problems in which an unknown plant $P(z)$ varies in a certain random way. Once again, the adaptive filter performs a modeling task. For this class of frequently encountered problems, we will analyze and compare the performances of the orthogonal LMS and the conventional steepest descent LMS algorithms. We will show that both perform equivalently (in the mean) for this class of nonstationary problems.

## II. Misadjustment of Finite - Data Exact Least - Squares Processes

Suppose that the adaptive linear combiner in Fig. 2 is fed $N$ independent zero-mean $n \times 1$ data vectors

---

[1] It should be noted that there are numerous algorithms in the literature which recursively estimate the autocorrelation matrix $R$ and use this estimation for the purpose of orthogonalization of the input data (see, for instance, [28]–[32]). These algorithms asymptotically converge to the idealized algorithm discussed here.
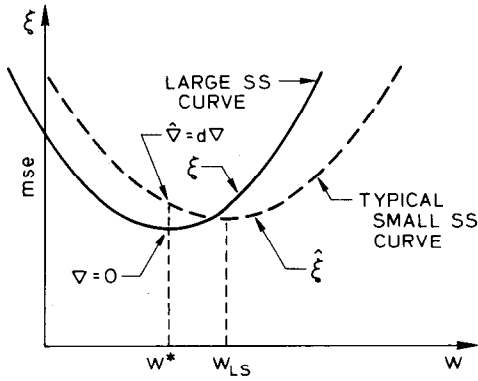
Fig. 3. Small and large sample-size mean square error curves.

$X_1, X_2, \cdots, X_N$ and their respective scalar desired responses $d_1, d_2, \cdots, d_N$, all drawn from a stationary ergodic process. Keeping the weights fixed, a set of $N$ error equations can be written as

$$\epsilon_i = d_i - X_i^T W, \qquad i = 1, 2, \cdots, N. \tag{10}$$

The objective is to find a weight vector that minimizes the sum of the squares of the error values based on the finite sample of $N$ items of data.

Eq. (10) can be written in matrix form as

$$\epsilon = D - \chi W, \tag{11}$$

where $\chi$ is an $N \times n$ rectangular matrix

$$\chi \triangleq [X_1, X_2, \cdots, X_N]^T, \tag{12}$$

where $\epsilon$ is an $N$ element error vector

$$\epsilon \triangleq [\epsilon_1, \epsilon_2, \cdots, \epsilon_N]^T, \tag{13}$$

and where $D$ is an $N$ element vector of desired responses

$$D \triangleq [d_1, d_2, \cdots, d_N]^T. \tag{14}$$

A unique solution of (11), a weight vector $W$ that would bring $\epsilon$ to zero, would exist only if $\chi$ is square and nonsingular. However, the case of greatest interest is that of $N \gg n$. As such, (11) would typically be overconstrained and one would generally seek a best least squares solution. The sum of the squares of the errors is

$$\epsilon^T \epsilon = D^T D + W^T \chi^T \chi W - 2 D^T \chi W. \tag{15}$$

This sum multiplied by $1/N$ is an estimate $\hat{\xi}$ of the mse $\xi$. Thus

$$\hat{\xi} = \frac{1}{N} \epsilon^T \epsilon$$

and

$$\lim_{N \to \infty} \hat{\xi} = \xi. \tag{16}$$

Note that $\hat{\xi}$ is a quadratic function of the weights, the parameters of the quadratic form being related to properties of the $N$ data samples. $(\chi^T \chi)$ is square and is assumed to be positive definite. $\hat{\xi}$ is a small sample-size mse function, while $\xi$ is the large sample-size "true" mse function. Fig. 3 shows a comparative sketch of these functions. Many small sample-size curves are possible, but there is

only one large sample-size curve. The large sample-size curve is the average of the many small sample-size curves.

The minimum of the small sample-size function can be found by differentiating (15), and the result is

$$W_{LS} = (\chi^T \chi)^{-1} \chi^T D. \tag{17}$$

This is the exact least squares solution for the given data sample. The Wiener solution $W^*$ is the expected value of $W_{LS}$.[2]

Each of the small sample-size curves of Fig. 3 is an ensemble member. Let the ensemble be constructed in the following manner. Assume that the vectors $X_1, X_2, \cdots, X_N$ are the same for all ensemble members, but that the associated desired responses $d_1, d_2, \cdots, d_N$ differ from one ensemble member to another because of the stochastic character of plant noise (refer to Fig. 2). Over this ensemble therefore, the $\chi$ matrix is constant while the desired response vector $D$ is stochastic. In order to evaluate the excess mean square error due to adaptation with the finite amount of data available, we have to find

$$\xi_{\text{excess}} = \frac{1}{N} E[V^T \chi^T \chi V], \tag{18}$$

where

$$V \triangleq W_{LS} - W^*. \tag{19}$$

Expectation is taken over the above-described ensemble. Eq. (18) can be written as

$$\xi_{\text{excess}} = \frac{1}{N} \text{tr}\left( E[V^T \chi^T \chi V] \right) = \frac{1}{N} E\left[ \text{tr}\left( V^T \chi^T \chi V \right) \right]$$
$$= \frac{1}{N} E\left[ \text{tr}\left( V V^T \chi^T \chi \right) \right] = \frac{1}{N} \text{tr}\left( E[V V^T] \cdot \chi^T \chi \right). \tag{20}$$

The covariance matrix of the weight error vector $V$ is known to be (see, for instance, [22])[3]

$$E(V V^T) = [\chi^T \chi]^{-1} \cdot \xi_{\min}, \tag{21}$$

where $\xi_{\min}$ is the minimum mean square error, the minimum of the true mse function.

Substitution of (21) into (20) yields

$$\xi_{\text{excess}} = \frac{n}{N} \cdot \xi_{\min}. \tag{22}$$

It is important to note that this formula does not depend on $\chi$. The above described ensemble can be generalized to an ensemble of ensembles, each having its own $\chi$, without changing formula (22). Hence this formula is valid for a very wide class of inputs.

It is useful to consider a dimensionless ratio between the excess mean square error and the minimum mean square error. This ratio is called in the literature (see, for instance, [1], [2], [4]) the misadjustment $M$. For the exact least square

---

[2] Note that validity of this statement is contingent on our earlier assumption that the plant $P(z)$ is FIR and both plant and plant model have the same order $n$.

[3] For the sake of simplicity we assumed here that the plant noise $n_j$ in Fig. 2 is white and that the adaptive filter has enough weights to match the unknown plant.

solution we can find the misadjustment (due to the finite length of input data) from (22) as

$$M = \frac{n}{N} = \frac{(\text{number of weights})}{(\text{number of independent training samples})}.$$

$$(23)$$

This misadjustment formula was first presented without detailed proof by Widrow and Hoff [1] in 1960, and has been used for many years in pattern recognition studies. For small values of $M$ (less than 25 percent), it has proven an excellent approximation. A formula similar to (23), based on somewhat different assumptions, was derived by Davisson [23] in 1970.

### III. STOCHASTIC GRADIENT SEARCH BY STEEPEST DESCENT

Gradient methods are commonly used to adjust adaptive parameters in order to search the quadratic mean square error performance function for its minimum. Most widely used is the method of steepest descent. With this method, a sequence of changes is made in the weight vector along the direction of the negative gradient. Thus the next weight vector $W_{j+1}$ is made equal to the present weight vector $W_j$ plus a change proportional to the negative gradient at the $j$th iteration:

$$W_{j+1} = W_j + \mu(-\nabla_j).$$

$$(24)$$

The parameter $\mu$ controls stability and rate of convergence. An "instantaneous gradient" $\hat{\nabla}$, an estimate of the true gradient $\nabla_j$, can be found by differentiation of (5) with respects to $W$:

$$\hat{\nabla}_j = -2\epsilon_j X_j.$$

$$(25)$$

Using the instantaneous gradient in place of the true gradient in (24) yields the LMS algorithm of Widrow and Hoff,

$$W_{j+1} = W_j + 2\mu\epsilon_j X_j.$$

$$(26)$$

The behavior of algorithm (26) has been analyzed extensively in the literature (see, for instance, [2], [3], [4], [18], [24], [25]). It was proved in [2], [4] that if the adaptation constant $\mu$ was chosen such that

$$0 < \mu < \frac{1}{\text{tr}(R)},$$

$$(27)$$

then the adaptive weights will relax from their initial condition to the Wiener solution[4] $W^*$. That means that the weight error vector

$$V_j = W_j - W^*$$

$$(28)$$

will converge to zero in the mean. The relaxation process will be governed by the relation

$$E[V_{j+1}] = (I - 2\mu R)E[V_j].$$

$$(29)$$

[4] In [2]–[4] it was proved that (27) is sufficient for the convergence in the mean. However it can be shown that this condition is both sufficient and necessary for the convergence of the variance. See also [33]–[35].

Therefore, there will be $n$ different modes of convergence corresponding to $n$ eigenvalues of the autocorrelation matrix $R$. Using normal decomposition of the matrix $R$:

$$R = Q\Lambda Q^T$$

$$(30)$$

$$QQ^T = I, \quad \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix},$$

$$(31)$$

we can find the corresponding $n$ time constants $\tau_i$ of the weight relaxation process as

$$\mu\lambda_i \ll 1, \quad \tau_i = \frac{1}{2\mu\lambda_i}, \quad 1 \le i \le n.$$

$$(32a)$$

As the weights relax toward the Wiener solution, the mean square error, a quadratic function of the weights, undergoes a geometric progression toward $\xi_{\min}$. The "learning curve" is a plot of mse versus number of adaptation cycles. The natural modes of the learning curve have time constants half as large as the corresponding time constants of the weights [2]–[4]. Accordingly, the mse (associated with the error $E[W_j] - W^*$) learning curve time constants are

$$\tau_{i_{\text{mse}}} = \frac{1}{4\mu\lambda_i}, \quad 1 \le i \le n.$$

$$(32b)$$

After convergence takes place, there remains noise in the weights due to the noise in the estimation (25) of the gradient. An approximate value of the covariance of the weight noise, valid for small $\mu$, was derived in [4, app. D]:

$$E[V_j V_j^T] = \mu \cdot \xi_{\min} \cdot I.$$

$$(33)$$

The noise in the weights will cause an excess error in the system output (in addition to the Wiener error):

$$\xi_{\text{excess}} = E\left[(V_j^T X_j)^2\right]$$

$$= E\left[\text{tr}(V_j V_j^T X_j X_j^T)\right].$$

$$(34)$$

Assuming, as has been done before [2]–[4], that $V_j$ and $X_j$ are independent, expression (33) can be substituted into (34) to obtain

$$\xi_{\text{excess}} = \mu\xi_{\min} n E\left[x_j^2\right] = \mu \text{tr}(R)\xi_{\min}.$$

$$(35)$$

Therefore, we can compute the misadjustment, defined as a ratio between the excess mean square and the minimum mean square error:

$$M = \mu \text{tr}(R).$$

$$(36)$$

The adaptation constant $\mu$ should be kept low in order to keep the misadjustment low. However, low $\mu$ is associated with slow adaptation in accordance with (32).

Expressions (27)–(36) illustrate the potential vulnerability of the steepest descent algorithm. The speed of convergence will depend on the choice of initial conditions. In the worst case, the convergence will be dominated by the lowest eigenvalue

$$\lambda_{\min} = \min(\lambda_1, \cdots, \lambda_n).$$

$$(37)$$

This implies that even if we choose the maximal value allowable for the adaptation constant $\mu$ (due to the stability

constraint (27)), the slowest time constant for the weights would be

$$\tau \geq \frac{\mathrm{tr}(R)}{2\lambda_{\min}}. \tag{38}$$

For the class of problems for which there exists a great spread of eigenvalues of the autocorrelation matrix $R$, this number will be high resulting in long convergence times (at least in the worst case).

## IV. Gradient Search by the Orthogonalized LMS Algorithm

In order to eliminate a potential deficiency of the Widrow–Hoff algorithm due to eigenvalue spread of the input autocorrelation matrix, we can prefilter the input signal in such a way that it will become orthogonal. Such a process would require perfect knowledge of the autocorrelation matrix $R$. Hence such an algorithm is not able to be implemented practically, although it is important from a theoretical point of view.

The block diagram of the idealized orthogonal LMS algorithm used in a modeling process is presented in the Fig. 4. In Fig. 4(a), each input signal vector to the adaptive process is preprocessed by the orthogonalization matrix

$$R^{-1/2} = Q \begin{pmatrix} \lambda_1^{-1/2} & & 0 \\ & \ddots & \\ 0 & & \lambda_n^{-1/2} \end{pmatrix} Q^T = Q\Lambda^{-1/2}Q^T. \tag{39}$$

Refer next to Fig. 4(b). Clearly this system is equivalent to that of Fig. 4(a) because the matrices $R^{1/2}$ and $R^{-1/2}$ incorporated into the unknown plant cancel each other. Denoting

$$X_j^\dagger \triangleq R^{-1/2}X_j, \tag{40a}$$

$$W_j^\dagger \triangleq R^{1/2}W_j, \tag{40b}$$

and

$$W^{\dagger*} \triangleq R^{1/2}W^*, \tag{41}$$

we can transform the system of Fig. 4(b) to that of Fig. 4(c). The adaptive process of Fig. 4(c) is, in turn, equivalent to the conventional LMS adaptive process of Fig. 1 and 2 with input signal $X_j^\dagger$ guaranteed to be orthogonalized. The plant impulse response is now modified to be $W^{\dagger*}$. From the systems equivalences of Fig. 4, we conclude that the orthogonalized LMS algorithm of Fig. 4(a) will perform exactly as the conventional steepest descent algorithm of Widrow and Hoff fed by an input signal with autocorrelation matrix

$$R^\dagger \triangleq E\left[X_j^\dagger\left(X_j^\dagger\right)^T\right] = R^{-1/2}E\left[X_jX_j^T\right]R^{-1/2} = I. \tag{42}$$

Therefore we can, evaluate the performance of the idealized orthogonalized LMS algorithm, using the well-known expressions (25)–(36).

The adaptation rule of the orthogonalized algorithm will be

$$W_{j+1}^\dagger = W_j^\dagger + 2\mu^\dagger\epsilon_j R^{-1/2}X_j = W_j^\dagger + 2\mu^\dagger\epsilon_j X_j^\dagger. \tag{43}$$

According to (27), for a choice of adaptation constant in the range

$$0 < \mu^\dagger < \frac{1}{\mathrm{tr}(R^\dagger)} = \frac{1}{n}, \tag{44}$$

the algorithm will converge to the optimal (Wiener) solution $W^{\dagger*}$. Since all the eigenvalues of $R^\dagger$ are equal to unity there will be only one mode of convergence and only one time constant. According to (32a) and (32b),

$$\tau = \frac{1}{2\mu^\dagger} \tag{45a}$$

$$\tau_{\mathrm{mse}} = \frac{1}{4\mu^\dagger}. \tag{45b}$$

These time constants will hold for an arbitrary choice of initial conditions. Therefore an idealized algorithm would indeed solve the problem of worst case slow convergence of the conventional steepest descent algorithm.

Next we will analyze the orthogonalized LMS algorithm to determine how close it brings us to the performance of the exact least mean square algorithm described in Section II.

According to (44) and (45a) the minimal possible time constant for the weights which is consistent with algorithm stability is

$$\tau_{\min} = \frac{n}{2}. \tag{46}$$

Completion of the convergence process will require several time constants or, in other words, a number of input samples which will be of the order of magnitude of $n$. This is consistent with the fact that the exact least mean square algorithm receives at least $n$ samples in order to compute an estimation of $n$ weights of the unknown plant.

After convergence takes place, the weights of the model remain noisy (due to the noise in the instantaneous estimation of the gradient). We will denote by $V_j^\dagger$ the weight error vector

$$V_j^\dagger = W_j^\dagger - W^{\dagger*} = W_j^\dagger - R^{1/2}W^*. \tag{47}$$

According to (33), the weight error covariance matrix is

$$E\left[V_j^\dagger\left(V_j^\dagger\right)^T\right] = \mu^\dagger \cdot \xi_{\min} \cdot I. \tag{48}$$

The noise in the adaptive weights will cause misadjustment which can be found from (36) as the ratio of excess mse to minimum mse:

$$M^\dagger = \frac{\left(\begin{array}{c}\text{excess mean}\\\text{square error}\end{array}\right)}{\left(\begin{array}{c}\text{Wiener error}\\\text{power}\end{array}\right)} = \mu^\dagger\,\mathrm{tr}(R^\dagger) = \mu^\dagger \cdot n. \tag{49}$$

Substituting (45b) into (49) yields

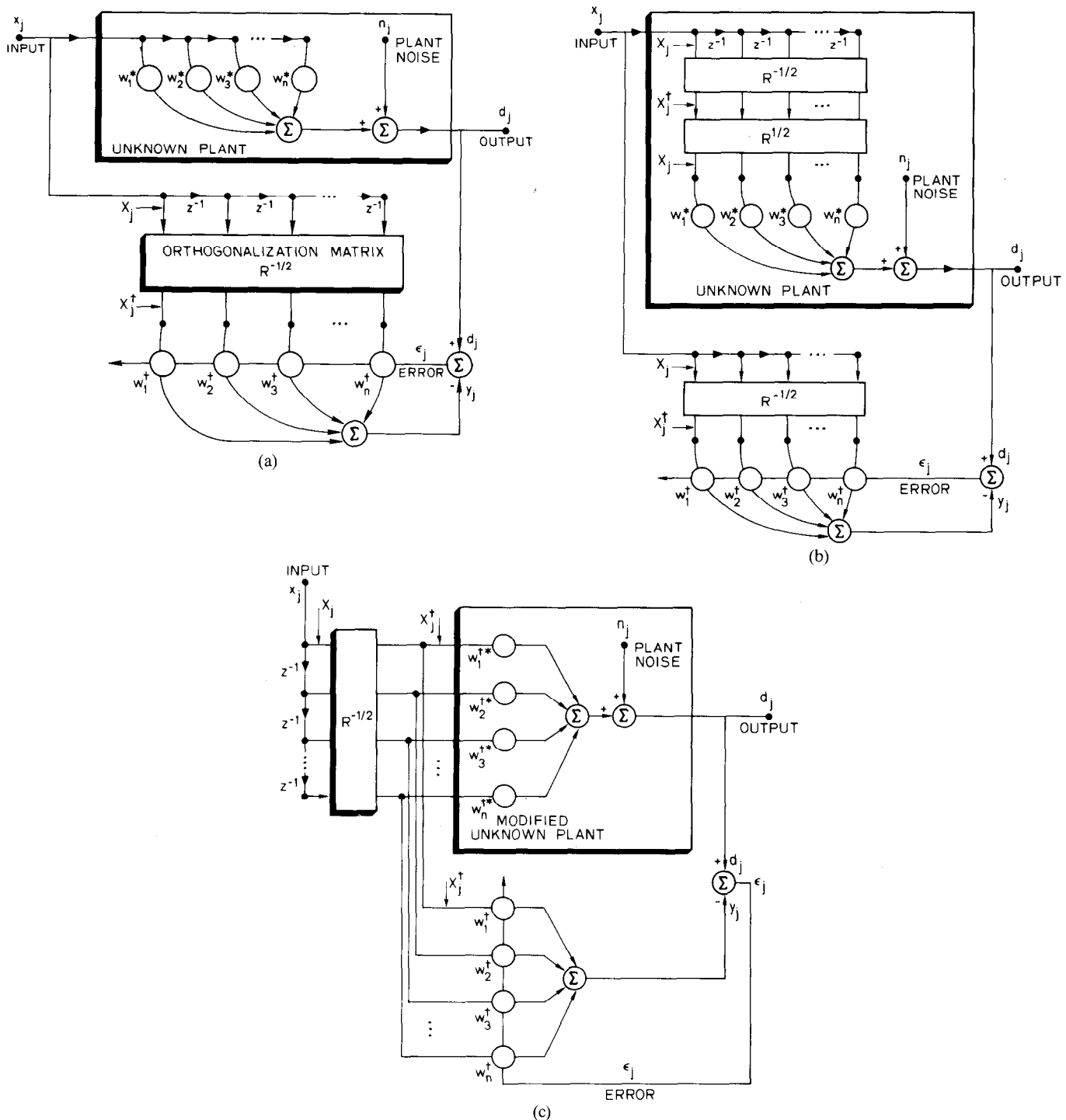$$M^\dagger = \frac{n}{4\tau_{\mathrm{mse}}}. \tag{50}$$

Fig. 4.  Orthogonalized LMS algorithm applied to plant identification. (a) Adaptation with orthogonalized inputs. (b) Adaptation with transformed plant inputs. (c) Modeling modified plant.

The question is, how does this compare with the misadjustment $M = n/N$ of the exact least mean squares algorithm found in Section II. Since we wish to compare the misadjustment of an exact least squares process that adapts with a finite equally weighted data sample to the orthogonalized LMS algorithm which weights its input data exponentially, we are to some extent comparing "apples with oranges." However, at least from a heuristic point of view, both (23) and (50) are equivalent.

The orthogonalized LMS algorithm exponentially weights its input data over time as it establishes its weight values. An exponential averaging window moves with time. The settling time of the adaptive process is of the order of four time constants of the mse learning curve. At any moment, the weights are determined by adaptation taken place over the last four time constants worth of data. Thus, in a steady flow situation, the training data "consumed" at any time by the orthogonalized LMS algorithm is essentially
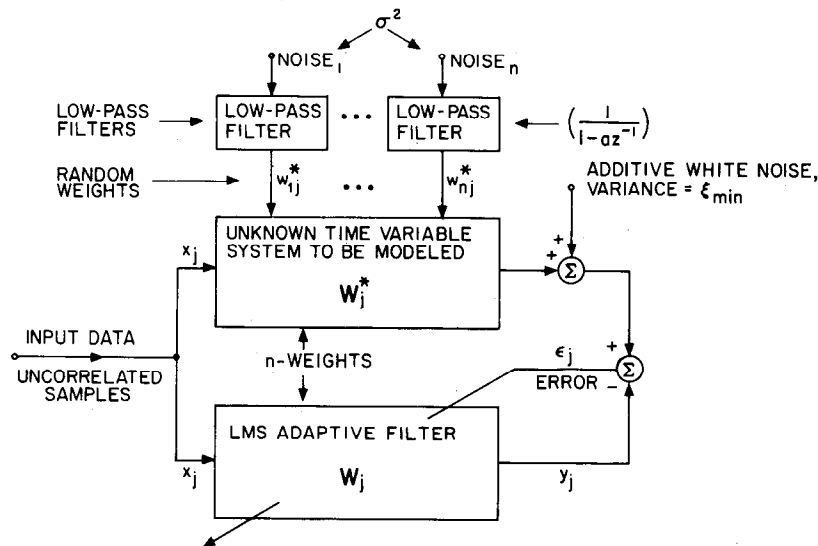
Fig. 5. Modeling an unknown time-variable system.

the most recent $4\tau_{mse}$ samples. The misadjustment of the orthogonalized LMS can therefore be expressed as

$$M^\dagger = \frac{n}{4\tau_{mse}} = \frac{\text{(number of weights)}}{\text{(number of independent training samples)}}.$$

$$(51)$$

It is clear from this expression that the orthogonalized LMS algorithm uses its training data about as efficiently as an exact least squares process.

## V. ORTHOGONALIZED LMS VERSUS CONVENTIONAL LMS: NONSTATIONARY INPUTS

Filtering nonstationary signals is a major area of application for adaptive systems. When the statistical character of an input signal changes gradually, randomly, and unpredictably, a filtering system that can automatically optimize its input–output response in accord with the requirements of the input signal could yield superior performance relative to that of a fixed, nonadaptive system. The performance of the conventional steepest descent LMS algorithm is compared here with orthogonalized steepest descent LMS (which, as demonstrated in the previous section, possesses certain optimality qualities), when both algorithms are used to adapt transversal filters with nonstationary inputs. The nonstationary situations to be studied are highly simplified, but they retain the essence of the problem that is common to more complicated and realistic situations.

The example considered here involves modeling or identifying an unknown time-variable plant, assumed to be transversal, of length $n$, whose weights (impulse response values) undergo independent stationary ergodic first-order Markov processes, as indicated in Fig. 5. The plant input signal $x_j$ is assumed to be stationary, ergodic, and, in general, colored. Additive plant output noise, assumed to be stationary, of mean zero, and of variance $\xi_{min}$, prevents a perfect match between the unknown system and the adaptive system. The minimum mse is, therefore, $\xi_{min}$, achieved whenever the $n$ weights of the adaptive filter $W_j$

match the corresponding $n$ weights of the unknown plant. The latter are at every instant the optimal values for the corresponding weights of the adaptive filter and are designated $W_j^*$, the subscript $j$ indicating that the unknown "target" to be tracked is time variable. It is tempting to call the target $W_j^*$ a time variable Wiener solution, but this would be improper since the nonstationary nature of the problem is beyond classical Wiener theory.

According to the scheme of Fig. 5, minimizing mse causes the adaptive weight vector $W_j$ to attempt to best match the unknown $W_j^*$ on a continual basis. The $R$ matrix, dependent only on the statistics of $x_j$, is constant even as $W_j^*$ varies. The desired response of the adaptive filter $d_j$ is nonstationary, containing the output of a time-variable system. The minimum mse $\xi_{min}$ is constant. Thus the mse function, a quadratic bowl, varies in position while its eigenvalues, eigenvectors, and $\xi_{min}$ remain constant. The adaptive process has the task of tracking the bottom of a randomly moving bowl in the presence of gradient noise.

In order to study this form of nonstationary adaptation both analytically and by computer simulation, a model comprising an ensemble of nonstationary adaptive processes has been defined and constructed as illustrated in Fig. 6. The ensemble of unknown filters to be modeled are all identical copies of the unknown plant and have the same time-variable weight vector $W_j^*$. Each ensemble member has its own independent stationary input signal going to both the unknown system and the corresponding adaptive system. The effect of output noise in the unknown systems is obtained by the addition of independent noise of variance $\xi_{min}$. All of the adaptive filters are assumed to start with the same initial weight vector $W_o$, each develops its own weight vector over time in attempting to pursue the Markovian target $W_j^*$.

The class of problems defined above was analyzed in [26], using the conventional steepest descent algorithm as an adaptation rule. Presently our aim will be to modify this line of reasoning in order to evaluate the nonstationary
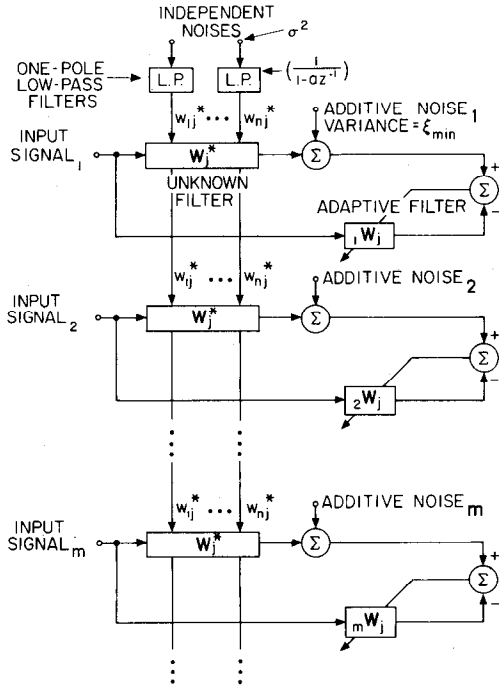
Fig. 6.   Ensemble of nonstationary adaptive processes.

performance of the orthogonalized LMS algorithm. We use once more the equivalent representations of Fig. 4. Clearly the adaptive model will attempt to track the moving target

$$W_j^{\dagger *} \triangleq R^{1/2} W_j^*. \tag{52}$$

For a given adaptive filter, the weight vector tracking error at the $j$th instant is $(W_j^\dagger - W_j^{\dagger *})$. This error is due to both the effects of gradient noise and weight vector lag, and may be expressed as

(weight vector error)$_j$

$$= \left( W_j^\dagger - W_j^{\dagger *} \right)$$

$$= \underbrace{\left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)}_{\substack{\text{weight vector} \\ \text{noise}}} + \underbrace{\left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right)}_{\substack{\text{weight vector} \\ \text{lag}}}. \tag{53}$$

The expectations are averages over the ensemble. Time averages are impossible to use with nonstationary statistics.

The components of error are identified in (53). Any difference between the ensemble mean of the adaptive weight vectors and the target value $W_j^{\dagger *}$ is due to lag in the adaptive process, i.e., due to the adapting weight vectors lagging behind the moving target $W_j^{\dagger *}$. The deviation of the individual adaptive weight vectors about their ensemble mean is due to gradient noise.

Weight vector error causes an excess mse. The ensemble average excess mse at the $j$th instant is

$$\begin{pmatrix} \text{ensemble} \\ \text{avg excess} \\ \text{mse} \end{pmatrix}_j = E\left[ \left( W_j^\dagger - W_j^{\dagger *} \right)^T R^\dagger \left( W_j^\dagger - W_j^{\dagger *} \right) \right]. \tag{54a}$$

The average excess mse at time $j$ is

$$\xi_{\text{excess}} = \begin{pmatrix} \text{avg excess} \\ \text{mse} \end{pmatrix}$$

$$= \left\langle E\left[ \left( W_j^\dagger - W_j^{\dagger *} \right)^T R^\dagger \left( W_j^\dagger - W_j^{\dagger *} \right) \right] \right\rangle. \tag{54b}$$

Using (53), this can be expanded as follows:

$$\begin{pmatrix} \text{avg excess} \\ \text{mse} \end{pmatrix}$$

$$= \left\langle E\left[ \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)^T R^\dagger \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right) \right] \right.$$

$$+ E\left[ \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right)^T R^\dagger \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right) \right]$$

$$\left. + 2E\left[ \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)^T R^\dagger \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right) \right] \right\rangle. \tag{55}$$

Expanding the last term of (55) and simplifying since $W_j^\dagger$ is constant over the ensemble,

$$\left\langle 2E\left[ W_j^{\dagger T} R^\dagger E\left[ W_j^\dagger \right] \right] - W_j^{\dagger T} R^\dagger W_j^{\dagger *} \right.$$

$$\left. - E\left[ W_j^\dagger \right]^T R^\dagger E\left[ W_j^\dagger \right] + E\left[ W_j^\dagger \right]^T R^\dagger W_j^{\dagger *} \right\rangle$$

$$= \left\langle 2E\left[ W_j^\dagger \right]^T R^\dagger E\left[ W_j^\dagger \right] - E\left[ W_j^\dagger \right]^T R^\dagger E\left[ W_j^\dagger \right] \right.$$

$$\left. - E\left[ W_j^\dagger \right]^T R^\dagger W_j^{\dagger *} + E\left[ W_j^\dagger \right]^T R^\dagger W_j^{\dagger *} \right\rangle$$

$$= 0. \tag{56}$$

Therefore, (55) becomes

$$\begin{pmatrix} \text{avg excess} \\ \text{mse} \end{pmatrix}$$

$$= \left\langle E\left[ \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)^T R^\dagger \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right) \right] \right.$$

$$\left. + E\left[ \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right)^T R^\dagger \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right) \right] \right\rangle. \tag{57}$$

The average excess mse is thus a sum of components due to both gradient noise and lag:

$$\begin{pmatrix} \text{avg excess} \\ \text{mse due to lag} \end{pmatrix}$$

$$= \left\langle E\left[ \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right)^T R^\dagger \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right) \right] \right\rangle$$

$$= \left\langle E\left[ \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right)^T \left( E\left[ W_j^\dagger \right] - W_j^{\dagger *} \right) \right] \right\rangle, \tag{58}$$

$$\begin{pmatrix} \text{avg excess} \\ \text{mse due to} \\ \text{gradient noise} \end{pmatrix}_j$$

$$= E\left[ \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)^T R^\dagger \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right) \right]$$

$$= E\left[ \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right)^T \left( W_j^\dagger - E\left[ W_j^\dagger \right] \right) \right]. \tag{59}$$

Using (48) can can evaluate the expression (59):[5]

$$\left(\begin{array}{c} \text{avg excess} \\ \text{mse due to} \\ \text{gradient noise} \end{array}\right)_j = \mu^\dagger n \xi_{\min} = \frac{n \xi_{\min}}{2\tau}. \qquad (60)$$

The next step is an evaluation of (58), the excess mse due to lag. Statistical knowledge of $(E[W_j^\dagger] - W_j^{\dagger *})$ will be required. In finding lag effects, we may eliminate gradient noise from consideration so that $E[W_j^\dagger] = W_j^\dagger$. Knowledge of $(W_j^\dagger - W_j^{\dagger *})$ will be sufficient.

Without gradient noise, orthogonalized LMS is represented by (43) as

$$W_{j+1}^\dagger = W_j^\dagger + 2\mu^\dagger E\left[\epsilon_j X_j^\dagger\right]$$

$$= W_j^\dagger + 2\mu^\dagger E\left[\left(X_j^\dagger\right)^T W_j^{\dagger *} - \left(X_j^\dagger\right)^T W_j^\dagger\right) X_j^\dagger\right]$$

$$= \left(1 - 2\mu^\dagger\right) W_j^\dagger + 2\mu^\dagger W_j^{\dagger *}. \qquad (61)$$

Substitution of (52) into (61) yields

$$W_{j+1}^\dagger - \left(1 - 2\mu^\dagger\right) W_j^\dagger = 2\mu^\dagger W_j^{\dagger *} = 2\mu^\dagger R^{1/2} W_j^*. \qquad (62)$$

The random vector $W_j^*$ is the driving function for this vector difference equation. Notice that there is no cross-coupling from any one coordinate to any other. Furthermore, all components of $W_j^*$ have been assumed to be stationary, ergodic, independent of each other, first-order Markov, and they have all been assumed to have the same variances and the same autocorrelation functions. Therefore, (62) may be treated as an array of $n$ independent first-order linear difference equations.

We will next derive an expression for the covariance of $(W_j^\dagger - W_j^{\dagger *})$ resulting from the random driving function $W_j^*$. Taking $z$ transforms of both sides of (62) yields

$$zW^\dagger(z) - \left(1 - 2\mu\right) W^\dagger(z) = 2\mu R^{1/2} W^*(z) \qquad (63)$$

or, equivalently,

$$W^\dagger(z) = \left[\frac{2\mu z^{-1}}{1 - (1 - 2\mu) z^{-1}}\right] R^{1/2} W^*(z) \qquad (64)$$

and

$$W^\dagger(z) - W^{\dagger *}(z) = \left[\frac{z^{-1} - 1}{1 - (1 - 2\mu) z^{-1}}\right] R^{1/2} W^*(z)$$

$$\triangleq R^{1/2} A(z). \qquad (65)$$

The vector $A(z)$, conveniently defined in (65), will be used below.

The random vector $W_j^*$ is assumed to be first-order Markov. A model for the generation of $W_j^*$ is shown in Fig. 5. Each of its components is generated by passing uncorrelated ("white") samples of variance $\sigma^2$ through a one-pole discrete filter of transfer function $1/(1 - az^{-1})$. This transfer function is in cascade with the above described transfer function. Thus, the transfer function from

[5] Strictly speaking, in the nonstationary case, $\xi_{\min}$ includes the power of the plant noise and the additional error power due to lag. For the sake of simplicity, we assume here that the adaptive process tracks well. Fluctuations in the plant parameters and error due to lag is assumed to be negligible relative to the $\xi_{\min}$ of stationary plant.

uncorrelated samples to the components of $A(z)$ is

$$T(z) = \left(\frac{1}{1 - az^{-1}}\right)\left(\frac{z^{-1} - 1}{1 - (1 - 2\mu) z^{-1}}\right), \qquad (66)$$

where $a$ is the geometric ratio of the nonstationarity. The variance of each of the independent components of $A(z)$ is equal to $\sigma^2$ multiplied by $S_t \triangleq$ sum of the squares of the impulses of the impulse response corresponding to the transfer function $T(z)$ given by (66). Transforming to the time domain and summing squares yields

$$S_t = \left(\frac{1}{a - 1 + 2\mu}\right)^2\left(\frac{1 - a}{1 + a} + \frac{\mu}{1 - \mu} - \frac{4(1 - a)\mu}{1 - a + 2\mu a}\right). \qquad (67)$$

The covariance of $A(z)$ may be expressed as

$$\text{cov}\left[A(z)\right] = S_t \sigma^2 I. \qquad (68)$$

Hence, the covariance of $(W_j^\dagger - W_j^{\dagger *})$ may be found as

$$\text{cov}\left[\left(W_j^\dagger - W_j^{\dagger *}\right)\right] = R^{1/2}(\text{cov}\left[A(z)\right]) R^{1/2} = S_t \sigma^2 R. \qquad (69)$$

Substituting (69) into (58) yields

$$\left(\begin{array}{c} \text{avg excess} \\ \text{mse due to lag} \end{array}\right) = S_t \sigma^2 \, \text{tr}\,(R). \qquad (70)$$

In order to use (70), we have to evaluate the value of the constant $S_t$ in expression (67). In practice it might be difficult, but in certain special cases of interest this task might be simplified. If the adaptive process is slow relative to the changes in the plant, we can assume $\mu^\dagger = 0$. Then

$$S_t = \frac{1}{1 - a^2}. \qquad (71)$$

A more common case occurs when the adaptive process is rapid relative to the time variation of the plant. As such,

$$1 - a \ll \mu^\dagger \ll 1, \qquad (72)$$

and therefore

$$S_t \approx \frac{1}{4\mu^\dagger} = \frac{\tau}{2}. \qquad (73)$$

Expression (73) can be used as an excellent approximation of $S_t$. Substitution of (73) into (70) yields

$$\left(\begin{array}{c} \text{avg excess} \\ \text{mse due to lag} \end{array}\right) = \frac{\tau \sigma^2}{2} \, \text{tr}\,(R) = \frac{\sigma^2}{4\mu^\dagger} \, \text{tr}\,(R). \qquad (74)$$

Substitution of (60) and (74) into (57) yields

$$\left(\begin{array}{c} \text{avg excess} \\ \text{mse} \end{array}\right) = \mu^\dagger n \xi_{\min} + \frac{\sigma^2}{4\mu^\dagger} \, \text{tr}\,(R)$$

$$= \frac{n \xi_{\min}}{2\tau} + \frac{\tau \sigma^2}{2} \, \text{tr}\,(R). \qquad (75)$$

Normalizing with respect to $\xi_{\min}$, we can compute the net misadjustment due to weight noise and weight lag as

$$M_{\text{sum}}^\dagger = \frac{n}{2\tau} + \frac{\tau \sigma^2}{2 \xi_{\min}} \, \text{tr}\,(R)$$

$$= \mu^\dagger n + \frac{\sigma^2}{4\mu^\dagger \xi_{\min}} \, \text{tr}\,(R). \qquad (76)$$

The optimal choice of $\mu^\dagger$ that minimizes $M_{\text{sum}}^\dagger$ is obtained by differentiating (76) with respect to $\mu^\dagger$, and setting the derivative to zero. The result is that the adaptive process is optimized when the component of $M_{\text{sum}}^\dagger$ due to the noise in the gradient is equal to the component of $M_{\text{sum}}^\dagger$ due to the lag:

$$\mu^{\dagger *} n = \frac{\sigma^2}{4\mu^{\dagger *}\xi_{\min}} \text{tr}(R). \tag{77}$$

Hence the optimal value of $\mu^\dagger$ is[6]

$$\mu^{\dagger *} = \sqrt{\frac{\sigma^2 \text{tr}(R)}{4n\xi_{\min}}}. \tag{78}$$

Substitution of (78) into (76) yields the minimal misadjustment

$$M_{\text{sum}}^{\dagger *} = \sqrt{\frac{\sigma^2 n \text{tr}(R)}{\xi_{\min}}}. \tag{79}$$

With *a priori* knowledge of the input signal power (to get $\text{tr}(R)$), and with knowledge of $\xi_{\min}$ and $\sigma^2$, these formulas could be used to set $\mu$ and predict misadjustment. Without such knowledge, one could slew the value of $\mu$ seeking the lowest system mse.

These results apply to orthogonalized LMS. An analogous problem was of considered in [26], where the performance of stochastic gradient LMS algorithm of Widrow and Hoff in a nonstationary environment was evaluated. Using as a model the same configuration presented in Figs. 5 and 6, the excess mean square error was computed. The approach used was similar to the one described above, i.e., the system error was divided into two uncorrelated parts: error due to the gradient noise and error due to the lag in the adaptation process. The overall misadjustment $M_{\text{sum}}$ equals the sum of misadjustments caused by each one of these errors:

$$M_{\text{sum}} = \mu \text{tr}(R) + \frac{1}{\mu} \frac{n\sigma^2}{4\xi_{\min}}. \tag{80}$$

Optimizing the choice of $\mu$ results in the minimum $M_{\text{sum}}$ when the two right-hand terms are equal,

$$\mu^* \text{tr}(R) = \frac{1}{\mu^*} \frac{N\sigma^2}{4\xi_{\min}}. \tag{81}$$

Hence for the steepest descent LMS algorithm on optimal choice of adaptation constant would be

$$\mu^* = \sqrt{\frac{n\sigma^2}{4\xi_{\min} \text{tr}(R)}} = \frac{n}{\text{tr}(R)}\mu^{\dagger *}. \tag{82}$$

In other words, on optimal value $\mu^*$ for the adaptation constant of the steepest descent LMS algorithm is equal to the optimal value of the adaptation constant $\mu^{\dagger *}$ for the

orthogonalized algorithm normalized by the factor of input signal power.

For the steepest descent LMS algorithm, the misadjustment will be minimal for $\mu = \mu^*$. Using (82) and (80) we have

$$M_{\text{sum}}^* = \sqrt{\frac{n\sigma^2 \text{tr}(R)}{\xi_{\min}}}. \tag{83}$$

Comparing (83) with (79), we note that

$$M_{\text{sum}}^* = M_{\text{sum}}^{\dagger *}. \tag{84}$$

The optimized misadjustment of the orthogonalized LMS algorithm turns out to be equal to the optimized misadjustment of the original Widrow–Hoff steepest descent LMS algorithm, regardless of the eigenvalue spread. The important conclusion is that, for the assumed form of nonstationary input, the two algorithms give identical mean square error performance although the adaptive steps of the two algorithms generally differ in detail.

## VI. CONCLUSION

Using an exact least squares algorithm to determine the $n$ weights of an adaptive filter from $N$ independent samples of data, we have shown the misadjustment to be $M = n/N$.

We have devised an ideal form of the steady flow LMS algorithm called orthogonalized LMS. To implement this algorithm, one would need perfect knowledge of the input covariance matrix $R$. Since this would not generally be known, the orthogonalized LMS algorithm is primarily of theoretical interest. We have shown that with a stationary input, this algorithm develops a misadjustment equal to $M^\dagger = n/4\tau_{\text{mse}}$. Comparing this result with the misadjustment formula for exact least squares, one sees great similarity. The amounts of data at any time being used to determine the weights is of the order of four time constants of the learning curve, the "settling time" of the adaptive process. One concludes that for a given speed of adaptation, the misadjustment of the orthogonalized LMS algorithm is about as low as can be obtained by any algorithm, and that this algorithm uses its input data about as efficiently as an exact least squares process.

Next, we have analyzed the behavior of the orthogonalized LMS algorithm in a nonstationary plant identification application. Comparing the results of this study with those obtained by Widrow and McCool for the same application but using the conventional steepest descent LMS algorithm, it was determined that the mean square error performance of the conventional LMS algorithm is on the average identical to that of the orthogonalized LMS algorithm, regardless of eigenvalue spread. For the nonstationary problem considered, we conclude that the conventional LMS algorithm performs as efficiently as exact least squares.

It is not always true that the performance of the LMS algorithm with nonstationary inputs will be independent of the eigenvalue spread. As a matter of fact, a form of nonstationarity studied by Macchi and Eweda [27] was

---

[6]When optimizing $\mu^\dagger$ is accord with the above procedure, one must keep this parameter within the stable range (44). Formula (78) applies as long as $\mu^{\dagger *}$ is in the stable range. Otherwise, $\mu^\dagger$ is set to its maximal stable value. The minimal misadjustment is determined either by (79), or by (76) when $\mu^\dagger$ is limited by stability considerations. Low plant noise would be the cause of $\mu^{\dagger *}$ exceeding the stable limit.

such that the LMS algorithm did show sensitivity to eigenvalue spread and the performance of an orthogonalizing algorithm for this case could have been better than that of conventional LMS. However, from the analysis conducted above, one can conclude that the steepest descent LMS algorithm of Widrow and Hoff, devised in 1960, not only is the simplest algorithm but is a highly efficient one for use with a variety of nonstationary inputs.

## ACKNOWLEDGMENT

We would like to thank John M. Cioffi, John M. McCool, and Timothy Saxe for many useful discussions.

## REFERENCES

[1] B. Widrow and M. E. Hoff, "Adaptive switching circuits," in *IRE WESCON Conv. Rec.*, pt. 4, pp. 96–104, 1960.

[2] B. Widrow, "Adaptive filters," in *Aspects of Network and System Theory*, R. Kalman and N. DeClaris, Eds. New York: Holt, Rinehart, and Winston, 1971, pp. 563–587.

[3] B. Widrow, P. Mantey, L. Griffiths, and B. Goode, "Adaptive antenna systems," *Proc. IEEE*, vol. 55, no. 12, pp. 2143–2159, Dec. 1967.

[4] B. Widrow, J. R. Glover, Jr., J. M. McCool, J. Kaunitz, C. S. Williams, R. H. Hearn, J. R. Zeidler, E. Dong, Jr., and R. C. Goodlin, "Adaptive noise cancelling: Principles and applications," *Proc. IEEE*, vol. 63, no. 12, pp. 1692–1716, Dec. 1975.

[5] J. Makhoul, "Linear prediction: A tutorial review," *Proc. IEEE*, vol. 63, no. 4, pp. 561–580, Apr. 1975.

[6] R. W. Lucky, "Automatic equalization for digital communication," *Bell Syst. Tech. J.*, vol. 44, pp. 547–588, Apr. 1965.

[7] R. W. Lucky, J. Salz, and E. J. Weldon, Jr., *Principles of Data Communication*. New York: McGraw Hill, 1968, ch. 6.

[8] R. D. Gitlin, E. Y. Ho, and J. E. Mazo, "Passband equalization of differentially phase-modulated data signals," *Bell Syst. Tech. J.*, vol. 52, no. 2, pp. 219–238, Feb. 1973.

[9] S. Qureshi, "Adaptive equalization—A comprehensive review," *IEEE Commun. Magazine*, pp. 9–16, Mar. 1982.

[10] M. M. Sondhi and A. J. Presti, "A self-adaptive echo canceller," *Bell Syst. Tech. J.*, vol. 46, no. 3, pp. 497–511, Mar. 1967.

[11] D. D. Falconer, K. H. Mueller, and S. B. Weinstein, "Simultaneous two-way data transmission over a two-wire circuit," presented at Proc. NTC, Dallas, TX, Dec. 1976.

[12] D. L. Duttweiler, "A twelve-channel digital echo canceller," *IEEE Trans. Comm.*, pp. 647–563, May 1978.

[13] L. J. Griffiths, "Rapid measurement of instantaneous frequency," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-23, pp. 209–222, Apr. 1975.

[14] J. R. Zeidler, E. H. Satorius, D. M. Chabries, and H. T. Wexler, "Adaptive enhancement of multiple sinusoids in uncorrelated noise," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-26, pp. 240–254, June 1978.

[15] J. T. Rickard, J. R. Zeidler, M. J. Dentino, and M. Shensa, "A performance analysis of adaptive line enhancer-augmented spectral detectors," *IEEE Trans. Circuits Syst.*, vol. CAS-28, no. 6, pp. 534–541, June 1981.

[16] B. Widrow, J. M. McCool, and B. P. Medoff, "Adaptive control by inverse modeling," presented at the Twelfth Asilomar Conference on Circuits, Systems, and Computers, Nov. 1978.

[17] B. Widrow, D. Schur, and S. Shaffer, "On adaptive inverse control," presented at the Fifteenth Asilomar Conf. Circuits, Syst., Comput., pp. 185–189, Nov. 1981.

[18] B. D. O. Anderson and R. M. Johnstone, "Convergence results for Widrow's adaptive controller," presented at the IFAC Conf. System Identification, 1981.

[19] D. T. L. Lee, M. Morf, and B. Friedlander, "Recursive least squares ladder estimation algorithms," *IEEE Trans. Circuits Syst.*, vol. CAS-28, no. 6, pp. 467–481, June 1981.

[20] L. J. Griffiths, "A continuously adaptive filter implemented as a lattice structure," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, Hartford, CT, May 1977, pp. 683–686.

[21] B. Friedlander, "Lattice filters for adaptive processing," *Proc. IEEE*, vol. 70, no. 8, pp. 829–867, Aug. 1982.

[22] G. C. Goodwin and R. L. Payne, "Dynamic system identification: Experiment design and data analysis," New York: Academic, 1977.

[23] L. D. Davisson, "Steady-state error in adaptive mean-square minimization," *IEEE Trans. Inform. Theory*, vol. IT-16, pp. 382–385, July 1970.

[24] R. R. Bitmead, "Convergence in distribution of LMS-type adaptive parameter estimates," *IEEE Trans. Automat. Contr.*, vol. AC-28, no. 1, pp. 54–60, Jan. 1983.

[25] O. Macchi and E. Eweda, "Second-order convergence analysis of stochastic adaptive linear filtering," *IEEE Trans. Automat. Contr.*, vol. AC-28, no. 1, pp. 76–85, Jan. 1983.

[26] B. Widrow, J. M. 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, pp. 1151–1162, Aug. 1976.

[27] O. Macchi and E. Eweda, "Tracking properties of adaptive nonstationary filtering," submitted for publication.

[28] R. W. Chang, "A new equalizer structure for fast start-up digital communication," *Bell Syst. Tech. J.*, vol. 50, no. 6, pp. 1969–2014, July–Aug. 1971.

[29] K. H. Mueller, "A new fast converging mean-square algorithm for adaptive equalizers with partial-response signaling," *Bell Syst. Tech. J.*, vol. 54, no. 1, pp. 143–153, Jan. 1975.

[30] D. Godard, "Channel equalization using a Kalman filter for fast data transmission," *IBM J. Res. Develop.*, pp. 267–273, May 1974.

[31] R. D. Gitlin and F. R. Magee, Jr., "Self-orthogonalizing adaptive equalization algorithms," *IEEE Trans. Commun.*, vol. COM-25, no. 7, pp. 666–672, July 1977.

[32] H. Sari, "Simplified algorithms for adaptive channel equalization," *Philips J. Res.*, vol. 37, pp. 56–77, 1982.

[33] L. L. Horowitz and K. D. Senne, "Performance advantage of complex LMS for controlling narrow-band adaptive arrays," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 562–576, June 1981.

[34] L. J. Griffiths, "Rapid measurements of digital instantaneous frequency," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-23, pp. 207–222, Apr. 1975.

[35] A. Nehorai and D. Malah, "On the stability and performance of adaptive line enhancer," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Processing*, 1980, pp. 478–481.