# XIX. Communications Systems Research: Propagation Studies
## TELECOMMUNICATIONS DIVISION

## A. Two Stochastic Approximation Procedures for Identifying Linear Systems, *J. K. Holmes*

### 1. Introduction

Many important problems in communications research can be posed as a problem of system identification. That is, given an input signal record and an output signal record, find an equivalent system that fits these data. An important example occurs in the following communication problem where it is desired to study the effects of the atmosphere on signal transmission. A transmitter transmits directly to a tracking station via cables or a microwave link, and to a spacecraft which retransmits to the tracking station. The signal, as it passes through, for example, the Solar Corona to and from the spacecraft, undergoes distortion. One method of studying the distortion is to characterize it as a finite memory linear system. The identification could be derived from the input (the directly transmitted component) with the proper time delay, and the output, which is the retransmitted, distorted signal. Naturally, noise would, in general, corrupt both the input and output signal measurements.

This article, then, considers the basic problem of identifying linear systems from noisy input–output measurements. Related work is listed in Refs. 1–5.

Let $N$ denote the set of natural numbers and let $n$ be in $N$. Denote by $\mathcal{S}_L(\ell)$ the class of linear, finite memory, time invariant, time discrete systems. Then, for systems $s \in \mathcal{S}_L(\ell)$, we can relate the input and output random sequences by

$$v_n = \sum_{j=0}^{\ell-1} \phi_j\, u_{n-j}, \qquad \ell < \infty \qquad (1)$$

The vector $\varphi = (\phi_0, \cdots, \phi_{\ell-1})$ defines the system when it is in the class $\mathcal{S}_L(\ell)$. Further, we assume that the observable sequences $v'_n$ and $u'_n$ are defined as $v'_n = v_n + \delta_n$ and $u'_n = u_n + \epsilon_n$, where $\epsilon_n$ and $\delta_n$ are mutually and individually independent noise sequences with zero means and respective variances $\sigma_\epsilon^2$ and $\sigma_\delta^2$. Also, $u_n$ is assumed to be independent of $\epsilon_n$ and $\delta_n$. This article, then, is concerned with the following problem: From the noise-corrupted input and output measurements, estimate the unknown system (i.e., $\varphi$) via nonparametric methods. See Fig. 1 for a block diagram.
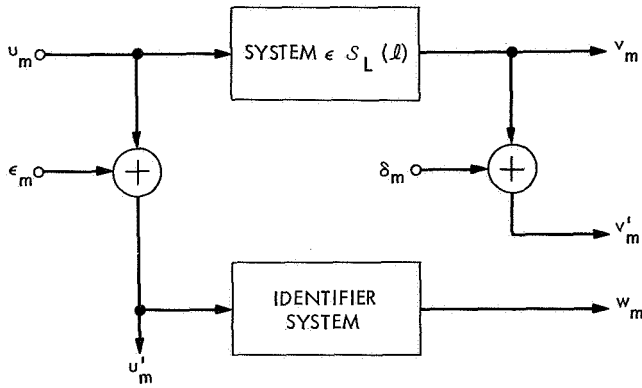
**Fig. 1. Block diagram of available measurements $u'_m$ and $v'_m$ and conceptual location of identifier system**

## 2. Development of the First Algorithm

Let the estimator of the sequence $v_m(\varphi)$ be

$$w_m(z) = \sum_{j=0}^{l-1} z_j \, u'_{m-j}$$

where $m = ln$ and the row vector $z \in E^l$. Now, we shall derive a specific sequential estimate of $\varphi$ denoted by $x^m$. Define an error measure in the following way:

$$M(z) = E\left[(w_m(z) - v'_m)^2 \,|\, z\right] \qquad (2)$$

The procedure to be used to estimate $\varphi$ will be to recursively determine $z$ in such a manner that $M(z)$ is a minimum. Specifically, we consider a Kiefer–Wolfowitz stochastic approximation method.

Let

$$\rho_m = (w_m(x^m) - v'_m)^2$$

and define

$$D_i^m = \rho_m(x^m + c_m e_i) - \rho_m(x^m - c_m e_i),$$
$$i = 0, \cdots, l-1$$

where $c_m$ satisfies

$$c_m \geqq 0 \qquad \text{and} \qquad \lim_{m \to \infty} c_m = 0$$

and the $e_i$ are the orthonormal unit vectors $\in E^l$, viz.,

$$e_0 = (1, 0, \cdots, 0) \qquad \text{and} \qquad e_{l-1} = (0, 0, \cdots, 1)$$

Then, recursively define $x^m$ by

$$x^{m+l} = x^m - a_m \frac{D^m}{c_m} \text{ for } m = l, 2l, 3l, \cdots \qquad (3)$$

where

$$D^m = (D_0^m, D_1^m, \cdots, D_{l-1}^m)$$

and $a_m$ satisfies

$$a_m \geqq 0 \qquad \text{and} \qquad \lim_{m \to \infty} a_m = 0$$

By using the definition of $D^m$ and $\rho_m$, one obtains the following scalar algorithm for each component of the estimate:

$$x_i^{m+l} = x_i^m - 4a_m u'_{m-i}\left[\sum_{j=0}^{l-1} x_j^m \, u'_{m-j} - v'_m\right],$$
$$i = 0, \cdots, l-1$$

or in vector form, designating $(u'_m)^T$ as the transpose of $u'_m$, we have

$$x^{m+l} = x^m - 4a_m x^m (u'_m)^T u'_m + 4a_m v'_m u'_m \qquad (4)$$

which is independent of $c_m$, and where

$$u'_m = (u'_m, u'_{m-1}, \cdots, u'_{m-l+1})$$

## 3. The Estimate Error

Based on minimizing the mean square error between the sequence $v_m(\varphi)$ and $w_m(x^m)$, where $x^m$ is defined by Eq. (3), we show that $x^m$ does not converge to $\varphi$ unless the input noise sequence $\epsilon_n$ is identically zero. More precisely, we have our first result.

***Theorem 1.*** If the conditions

(a) $a_m \geqq 0$, $\sum_{1}^{\infty} a_m = \infty$, and $\sum_{1}^{\infty} a_m^2 < \infty$.

(b) $yR(y)^T > 0$, $\forall y \neq 0$, and $E\left[u_m^T(u_m)\,|\,x^m\right] = R$.

(c) $E\left[(u'_m)^4\,|\,x^m\right] < c_1 < \infty$, and $E\left[u_m^2 \delta_m^2\,|\,x^m\right] < c_2 < \infty$.

(d) The random sequence $u_m$ and both noise sequences $\delta_m$ and $\epsilon_m$ are time-stationary.

(e) $E[\delta_m] = E[\epsilon_m] = 0, E[\delta_m \delta_n]$
$= \sigma_\delta^2 \delta_{mn}$, and $E[\epsilon_m \epsilon_n] = \sigma_\epsilon^2 \delta_{mn}$.

(f) $E[\delta_m \epsilon_n] = E[u_m \delta_n] = E[\delta_m \epsilon_n] = 0$.

(g) Unknown system $s \in \mathcal{S}_L(\ell)$.

are met, then the sequence $x^m$, defined by Eq. (3), converges in mean square to the vector

$$\theta = [R + \sigma_\epsilon^2 I]^{-1} R \phi$$

**Proof.** Equation (4) may be expanded to

$$x^{m+l} = x^m - 4a_m u'_m (u'_m)^T x^m + 4a_m [\delta_m u'_m + u'_m (u_m)^T \phi]$$

(5)

where in Eq. (5) and throughout the rest of this proof $x^m$, $u'_m$, $u_m$, and $\phi$ are now defined as column vectors instead of row vectors. Using $E[(\cdot)] = EE[(\cdot)|x^n]$ on both sides of Eq. (5) results in

$$E[x^{m+l}] = E[x^m] - 4a_m[R + \sigma_\epsilon^2 I] E[x^m] + 4a_m R\phi$$

(6)

Clearly, a solution to Eq. (6) is given by

$$E[x^m] = [R + \sigma_\epsilon^2 I]^{-1} R\phi$$

(7)

We shall now show that

$$\lim_{m \to \infty} E[\|x^m - \theta\|^2] = 0$$

Subtracting $\theta$ from both sides of Eq. (5) and rearranging yields

$$(x^{m+l} - \theta) = (x^m - \theta) - 4a_m u'_m (u'_m)^T (x^m - \theta)$$
$$+ 4a_m [\delta_m u'_m + u'_m u_m^T \phi - u'_m (u'_m)^T \theta]$$

Forming the norm squares of both sides and averaging, we obtain

$$b_{m+l} = b_m - 8a_m E \langle x^m - \theta, u'_m (u'_m)^T (x^m - \theta) \rangle$$
$$+ 16a_m^2 E[\|u'_m (u'_m)^T (x^m - \theta)\|^2]$$
$$+ 16a_m^2 E[\|\delta_m u'_m + u'_m u_m^T \phi - u'_m (u'_m)^T \theta\|^2]$$
$$+ 8a_m E \langle x^m - \theta, \delta_m u'_m + u'_m u_m^T \phi - u'_m (u'_m)^T \theta \rangle$$
$$- 32a_m^2 E \langle u'_m (u'_m)^T (x^m - \theta), \delta_m u'_m$$
$$+ u'_m u_m^T \phi - u'_m (u'_m)^T \theta \rangle$$

(8)

where we have let

$$b_m = E[\|x^m - \theta\|^2]$$

By Condition (b) we have for the second term in Eq. (8)

$$\sigma_\epsilon^2 b_m < \lambda_1 b_m \leq \langle x^m - \theta, [R + \sigma_\epsilon^2 I](x^m - \theta) \rangle$$

(9)

where $\lambda_1$ is the minimum eigenvalue of $[R]$. The third term can be bounded, using the Schwartz inequality, by

$$16a_m^2 b_m k_1$$

(10)

where $k_1 < \infty$ by Condition (c). The fourth term can be bounded by the following quantity

$$32a_m^2 E[\|\delta_m u'_m\|^2] + 32a_m^2 E[\|u'_m u_m^T \phi\|^2]$$
$$+ 32a_m^2 E[\|u'_m (u'_m)^T \theta\|^2]$$

All these terms are similarly bounded by

$$k_2 a_m^2$$

(11)

where $k_2 < \infty$. Now, the fifth term can be reduced to

$$8a_m E \langle x^m - \theta, R\phi - [R + \sigma_\epsilon^2 I] \theta \rangle = 0$$

since $\theta = [R + \sigma_\epsilon^2 I]^{-1} R\phi$. The last term is bounded by

---

$$32a_m^2 \{E[\|u'_m (u'_m)^T (x^m - \theta)\|^2] E[\|u'_m u_m^T \phi - u'_m (u'_m)^T \theta\|^2]\}^{1/2}$$
$$\leq 32a_m^2 \{E[\|x^m - \theta\|^2 E[\|u'_m (u'_m)^T\|^2 |x^m]] E[\|u'_m u_m^T \phi - u'_m (u'_m)^T \theta\|^2]\}^{1/2}$$

(12)

---

where $\|A\|$, with $A$ a square matrix, is the usual euclidean $\ell_2$ norm; i.e.,

$$\|A\| = (\sum_{i,j} a_{ij}^2)^{1/2}$$

Equation (12) can be shown to be bounded by

$$32a_m^2 (b_m k_3)^{1/2} \leq a_m^2 (1 + b_m) k_4$$

(13)

for $k_4 < \infty$. Hence, Eq. (8) yields

$$b_{m+l} \leqq b_m \left[ 1 - 8a_m \left( \lambda_1 - k_5 a_m \right) \right] + a_m^2 k_6$$

$$\leqq b_m \left( 1 - 4a_m \lambda_1 \right) + a_m^2 k_6 \qquad (14)$$

for $m$ sufficiently large and $k_5$ and $k_6$ finite. By applying an interesting application of Kronecker's theorem (e.g., Ref. 5) or Lemma I of Ref. (3), we have that $b_m \to 0$. The theorem follows directly.

To conclude the statement made in the first paragraph of this subsection, we have:

**Lemma 1.** If the conditions of Theorem 1 are satisfied and if $\varphi \neq 0$, then $x^m$ is a consistent estimator of $\varphi$ if and only if $\sigma_\epsilon^2 = 0$.

**Proof.** The Lemma follows from Theorem 1 and the fact that $[R + \sigma_\epsilon^2 I]^{-1} R \varphi$ is equal to $\varphi$ if and only if $\sigma_\epsilon^2 = 0$, since by Condition (b) $R$ is positive definite.

The fact that $x^m$ converges to $\theta \neq \varphi$, based on minimizing $M(x^m)$, is not so surprising since Eq. (2) can be written as

$$M(z) = E \left[ \left( \sum_{j=0}^{l-1} (z_j - \phi_j) u_{m-j} \right)^2 \Big| z \right] + \sigma_\epsilon^2 \|z\|^2 + \sigma_\delta^2 \qquad (15)$$

which reduces to

$$M(z) = (z - \theta) [R + \sigma_\epsilon^2 I] (z - \theta)^T$$
$$+ (\theta - \varphi) R (\theta - \varphi)^T + \sigma_\epsilon^2 \|\theta\|^2 + \sigma_\delta^2 \qquad (16)$$

where

$$\theta = [R + \sigma_\epsilon^2 I]^{-1} R \varphi$$

Since $R$ is positive definite, we have that the minimum occurs at $z = \theta$. Hence, we *cannot* exactly identify an unknown system, when the input noise variance is non-zero, without additional knowledge of the input noise statistics.

### 4. Second Algorithm

A modification of the original algorithm, defining a new sequence $y^m$, has the property that if $\sigma_\epsilon^2$ is known then the modified algorithm leads to the result that $y_m \to \varphi$ in mean square.

**Theorem 2.** If the conditions of Theorem 1 are fulfilled and if $y^m$ is defined recursively by

$$y^{m+l} = y^m - 4a_m \left[ u_m' (u_m')^T - \sigma_\epsilon^2 I \right] y^m + 4a_m u_m' v_m' \qquad (17)$$

then $y^m$ converges in mean square to $\varphi$.

**Proof.** The proof follows the lines of Theorem 1 and will not be indicated here.

### 5. Convergence to a Subspace

Up to this point, the condition that $R$ be positive definite was required when $\sigma_\epsilon^2 = 0$ to prove convergence (in mean square) to a unique point. It is to be noted that if $\sigma_\epsilon^2 > 0$ then it is not necessary to assume that $R$ is positive definite. However, if it turns out that $R$ is of rank $r < \ell$, then it can be proven that $x^n$ converges to a subspace of $E^l$. Obviously, this is a very weak form of convergence. However, looking at it philosophically, we cannot hope to do any better since if $x^n \in \eta$, the null space of $R$, $M(z)$ will be minimum and that is the best we can do using the mean square error criterion. This then leads us to the statement of Theorem 3.

**Theorem 3.** If the Conditions (a), (c)–(g), $0 < \text{rank } R = r < \ell$, and $\sigma_\epsilon^2 = 0$ are satisfied, then $R(x^n - \varphi)$ converges in mean square to zero.

**Proof.** Starting with Eq. (5), setting $\epsilon_m$ equal to zero, subtracting $\varphi$ from both sides of the equation, and premultiplying by $R$ yields

$$R (x^{m+l} - \varphi) = R (x^m - \varphi) - 4a_m R u_m u_m^T (x^m - \varphi)$$
$$+ 4a_m \delta_m R u_m$$

Forming the norm square of both sides and averaging, we obtain

$$d_{m+l} = d_m - 8a_m E \langle R (x^m - \varphi), R u_m u_m^T (x^m - \varphi) \rangle$$
$$+ 16a_m^2 E \left[ \| R u_m u_m^T (x^m - \varphi) \|^2 \right]$$
$$+ 16a_m^2 E \left[ \delta_m^2 \| R u_m \|^2 \right]$$
$$+ 8a_m E \langle R (x^m - \varphi), \delta_m R u_m \rangle$$
$$- 32a_m^2 E \langle R u_m u_m^T (x^m - \varphi), \delta_m R u_m \rangle \qquad (18)$$

where we have let

$$d_m = E \left[ \| R (x^m - \varphi) \|^2 \right]$$

Now, even though it is no longer true that

$$E \langle \mathbf{x}^m - \varphi, R (\mathbf{x}^m - \varphi) \rangle \gneqq \lambda_1 E [\|\mathbf{x}^n - \varphi\|^2]$$

it is true that

$$E \langle R (\mathbf{x}^m - \varphi), RR (\mathbf{x}^m - \varphi) \rangle$$
$$\gneqq \lambda_\rho E [\|R (\mathbf{x}^m - \varphi)\|^2] = \lambda_\rho \, d_m \qquad (19)$$

where $\lambda_\rho$ is the least nonzero eigenvalue of $R$.

The third term can be bounded by

$$16 a_m^2 \, d_m \, K_7 \qquad (20)$$

where $K_7 < \infty$ by Condition (c). The fourth term can be bounded by

$$16 K_8 \, a_m^2 \qquad (21)$$

The fifth and sixth terms are zero since $E [\delta_m] = 0$.

Hence, we have that

$$d_{m+I} \leqq d_m [1 - 8 a_m (\lambda_\rho - K_7 \, a_m)] + 16 K_8 \, a_m^2$$

or

$$d_{m+I} \leqq d_m (1 - 4 a_m \lambda_\rho) + 16 K_8 \, a_m^2 \qquad (22)$$

for $m$ sufficiently large and $K_7$ and $K_8$ finite. As before, by application of Kronecker's theorem to Eq. (22), we have that $d_m \to 0$. The theorem follows directly. It is not hard to show from Eq. (15) that $M (\mathbf{z})$ is minimized for any $\mathbf{z} \in \eta$.

## 6. An Example

A simple example was devised to compare the first with the second algorithm. The "unknown system" was programmed to be of the form

$$v_m = \sum_{J=0}^{9} \phi_j \, u_{m-j}, \qquad \phi_j = \exp \left( - \frac{j}{2} \right); j = 0, \cdots, 9 \qquad (23)$$

and the identifier system was programmed to be of the form

$$w_m = \sum_{J=0}^{9} x_j^m \, u_{m-j}'$$

The relative mean square error was used to measure the performance of the algorithm and, for our example,

is defined by

$$\text{rmse} = \frac{\displaystyle\sum_{i=0}^{9} (x_i - \phi_i)^2}{\displaystyle\sum_{i=0}^{9} \phi_i^2} \qquad (24)$$

where the true system values are given by $\phi_i$ and the estimate of the system by $x_i$.

The respective elements of the correlation matrix were simulated to satisfy

$$E [u_m u_n] = \delta_{mn}, \qquad E [\epsilon_m \epsilon_n] = K \delta_{mn},$$
$$E [\delta_m \delta_n] = 0 \qquad (25)$$

where $u_n$ and $\epsilon_n$ were programmed to simulate independent gaussian random sequences. In Eq. (25), $\delta_{mn}$ is the Kronecker delta and $K$ was chosen to be either (a) $K = 1$ or (b) $K = \frac{1}{4}$, corresponding to an observational signal-to-noise ratio $(\sigma_u^2/\sigma_\epsilon^2)$ of 0 and 6 dB, respectively. Forty thousand samples were used to obtain the system estimates.

Figure 2 illustrates the plot of the rmse as a function of the time index $m$. As can be seen from Fig. 2, the first algorithm, for both signal-to-noise ratios, approached the theoretical limit defined by the result of Theorem 1 and Eq. (24). Furthermore, the second algorithm continued to decrease, on the average, as $m$ increased, as claimed by Theorem 2.

## 7. Conclusions

Two sequential identification procedures were presented for the identification of a time invariant, linear system in which no knowledge of the dynamics of the system were known prior to the identification, with the exception that the memory of the system was required to be finite.

The first algorithm for identification required only a mild condition on the covariance function of the input random process and no knowledge of the input or output measurement noise statistics other than that they have finite variances. It was shown that the estimate of the system converged in mean square; however, a bias developed that was due to the input noise and consequently prevented the estimate from being consistent. This error or lack of complete identification, without further knowledge of the statistics of the noises, is an example of what has been called the "structural regression paradox" in the statistical literature.
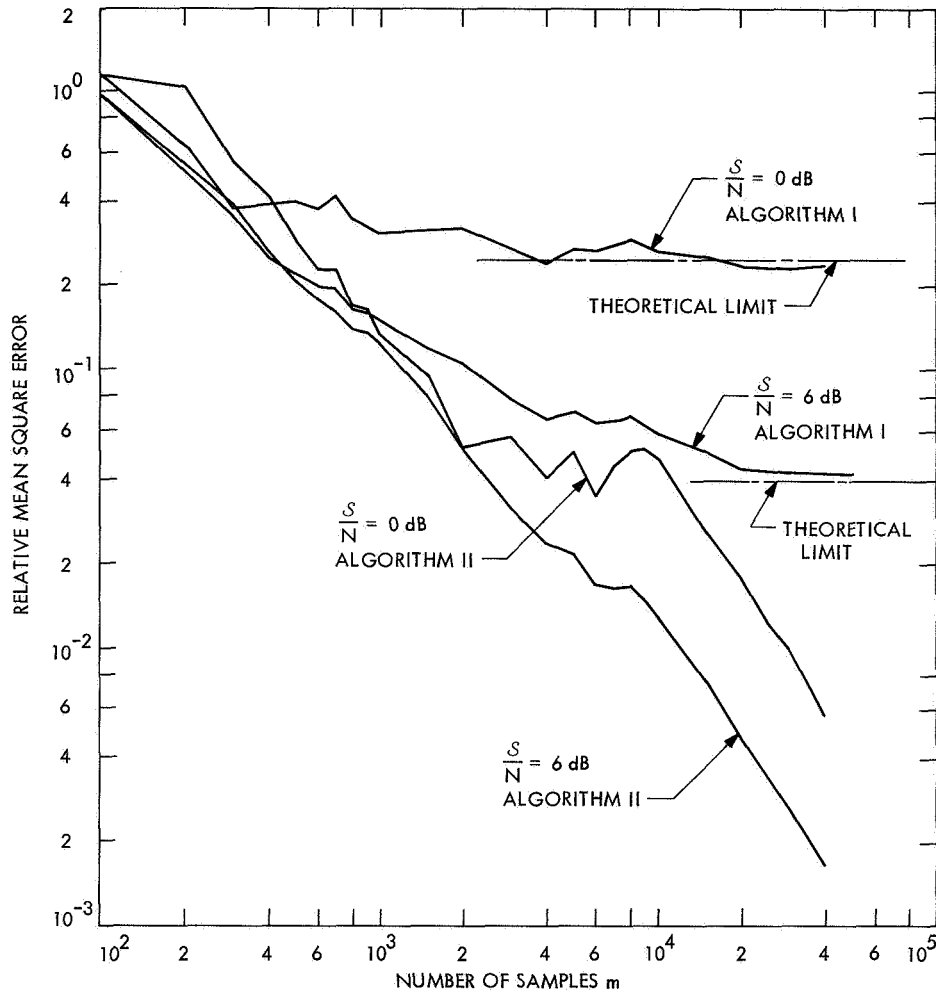
RELATIVE MEAN SQUARE ERROR

$\frac{S}{N} = 0$ dB
ALGORITHM I

THEORETICAL LIMIT

$\frac{S}{N} = 6$ dB
ALGORITHM I

$\frac{S}{N} = 0$ dB
ALGORITHM II

THEORETICAL
LIMIT

$\frac{S}{N} = 6$ dB
ALGORITHM II

NUMBER OF SAMPLES m

**Fig. 2. Relative mean square error as a function of number of samples m and theoretical limit of error for first algorithm**

The second algorithm was derived based on the additional assumption that the input measurement noise variance was known. With this additional knowledge, it was shown that the estimate of the unknown system converged in mean square to the unknown system. Hence, if the input noise variance is known, this algorithm can be used to obtain a consistent estimate of the unknown system.

Theorem 3 showed that as long as $r = $ rank $R$ was greater than zero the first algorithm would reduce $M(\mathbf{z})$ to a minimum; however, there could be an uncountable number of values of $\mathbf{z}$ that achieved this minimum.

The computer simulations agreed very well with the theory, indicating that the algorithms are useful and practical methods for identification of linear systems.

### References

1. Kushner, H. J., "A Simple Iterative Procedure for the Identification of the Unknown Parameters of a Linear Time Varying Discrete System," *J. Basic Eng.*, pp. 227–235, June 1963.

2. Papers presented at the IFAC Symposium on Identification in Automatic Control Systems, Prague, Czechoslovakia, 1967.

3. Saridis, G. N., and Stein, G., "Stochastic Approximation Algorithms for Linear Discrete-Time System Identification," *National Electronics Conference*, pp. 45–50, 1967.

4. Sakrison, D. J., "Application of Stochastic Approximation Methods to System Optimization," Technical Report 391. Massachusetts Institute of Technology Research Laboratory of Electronics, Cambridge, Mass., 1962.

5. Balakrishnan, A. V., "Determination of Non Linear Systems from Input-Output Data," paper presented at the 54th Meeting of the Princeton University Conference on Identification Problems in Communication and Control Systems, Princeton, N.J., 1963.