

Mathematics Notes

Note 66

May 1980

Matrix Methods For Determining System
Poles From Transient Response

T. L. Henderson
University of Kentucky
Department of Electrical Engineering
Lexington, Kentucky

Abstract

The problem treated is that of identifying the poles of a finite order system by observing its transient decay after cessation of input, for a limited time, using (possibly) multiple observation points and experimental repetition. Various approaches are studied, having the common characteristic that a homogeneous matrix equation must be solved. Several techniques that have been given scant attention in the literature are consolidated into the treatment, together with new results including an analytical treatment of the consequences of assuming an excessively high system order, derivation of a statistically unbiased estimate for an intermediate parameter in the solution, new theorems on error effects, a recipe for effective use of the singular value decomposition, a new method for suppression of extraneous poles, an elucidating derivation and extension of the method of Jain, a new form of the problem wherein the system poles are eigenvalues, and a study of the relationship between various pole identification methods.

ACKNOWLEDGMENTS

The author gratefully acknowledges the support of W. Pearson, who first introduced him to the Prony method, and fruitful discussions with a colleague, H. Yeh, and two students, J. Auton and P. Wigler. This work was partially supported by the Office of Naval Research under Contract No. N0014-77-C-0362. Mrs. Betty Bradshaw and Miss Gail McAlister did the manuscript preparation.

INTRODUCTION

One often seeks to determine the poles of a system by observing its natural response after cessation of input. Several examples can be cited: (i) Acoustic transducers [1,2] and electromagnetic antennae [3] are lately being tested by applying a pulse input, (ii) Aircraft and other military hardware are being subjected to electromagnetic pulse (EMP) tests [4], (iii) Engines and other cast metal objects are often tested by direct mechanical impact [5], (iv) Rooms are excited acoustically and decay characteristics are recorded [6]. In these and other cases the response after cessation of input stimulus may be expressed as

$$y_t = \sum_k \alpha_k \exp(s_k t) \quad (1)$$

The α_k 's depend upon the excitation, location of the sensor that extracts y_t , and selection of time origin, but the s-poles, s_k , are inherent characteristics of the system and remain invariant so long as the parameters of the applicable wave equation and its post-excitation boundary conditions are not disturbed. Sampling at interval T , starting at t_0 , gives a sequence $\{y(n)\}$:

$$y(n) = \sum_k r_k z_k^n u(n) \quad \text{for } -\infty < n < +\infty \quad (2)$$

where $z_k \triangleq \exp(s_k T)$, $r_k \triangleq \alpha_k \exp(s_k t_0)$ and the unit step sequence, $u(n)$, has

been used to extend $\{y(n)\}$ with zeroes for negative n . Taking Z-transforms gives

$$Y(z)/z = \sum_k r_k / (z - z_k) \quad (3)$$

so that $\{z_k\}$ and $\{r_k\}$ are the poles and residues of $Y(z)/z$. This paper concerns the problem of identifying the "system poles" $\{z_k\}$ through observation of $\{y(n)\}$ sequences. Then the s-poles can be determined from

$$s_k = T^{-1} \log z_k = (\log |z_k| + j[\arg(z_k) + 2\pi\ell]) / T \quad (4)$$

The ambiguity due to aliasing is expressed by ℓ . If π/T is known to be greater than the imaginary part of every s-pole then $\ell = 0$. Even when this is not true, if one can repeat the identification process with a slightly different sampling interval, say $T + dT$, giving displaced system poles $z_k + dz_k$, then since $dz_k = s_k z_k dT$ it follows that $s_k = dz_k / (z_k dT)$. This computation only has to be accurate enough to resolve the uncertainty in ℓ . Some of the pole identification methods to be discussed here use decimated subsequences and determine only the q^{th} -power of each system pole. This contributes an additional "decimation aliasing" expressed by ℓ in the formula

$$z_k = [z_k^q]^{1/q} = |z_k^q|^{1/q} \exp\{j(\arg[z_k^q] + 2\pi\ell)/q\} \quad (5)$$

But this ambiguity can be resolved in a similar manner. If z_k^{q+1} can also be determined then $z_k = z_k^{q+1} / z_k^q$, which should be accurate enough to resolve $\arg(z_k)$.

In all that follows we shall assume that the system poles are non-zero, distinct, and K in number. Throughout this paper we shall use the notation $K' \triangleq K + 1$. Equation (3) can then be expressed with a ratio of polynomials,

$$Y(z)/z = \frac{\beta_{K-1}z^{K-1} + \dots + \beta_1z + \beta_0}{\theta_Kz^K + \theta_{K-1}z^{K-1} + \dots + \theta_1z + \theta_0} \quad (6)$$

where the system poles are the roots of the denominator polynomial. Our failure to normalize either polynomial enables us to scale the coefficients of the denominator freely, with the numerator coefficients then being uniquely determined. Due to the assumptions above, θ_K and θ_0 and at least one of the β 's must be non-zero, and the polynomials can have no roots in common. We shall define the $K \times 1$ vector, $\theta \triangleq (\theta_0, \theta_1, \dots, \theta_K)^T$. Furthermore the symbol Z will be used to denote the "power vector" $(z^0, z^1, z^2, \dots)^T$ where z is a generic complex variable. The number of elements in the vector must be inferred from the context of its use. If z is explicitly z_i , the i^{th} system pole, then its power vector will be denoted Z_i . (Other than this Z function, all vectors and matrices in this paper will be real. Note also that the first element of a vector will always be denoted by a "0" subscript.) Using this notation the pole polynomial is simply $\theta^T Z$, and the system poles satisfy $\theta^T Z_i = 0$. Clearing the denominator and taking inverse transforms in Eq. (6) gives

$$\theta_0 y(n) + \theta_1 y(n+1) + \dots + \theta_K y(n+K) = \beta_0 \delta(n+1) + \dots + \beta_{K-1} \delta(n+K) \quad (7)$$

in particular, for $n \geq 0$,

$$(y(n), y(n+1), \dots, y(n+K)) \cdot \theta \equiv 0 \quad (8)$$

Since in Eq. (8) one can solve for $y(n+K)$ as a linear combination of the preceding $y(i)$'s, the set $\{y(n)\}$ forms an autoregressive (AR) sequence. Another interpretation is that when the sequence is passed through a finite impulse response (FIR) filter, represented by Eq. (8), whose transfer

function has zeroes that coincide with the system poles, then its output will be zero after enough time has elapsed to fill the delay line with signal. But the geometric interpretation is that the vector θ , which is an invariant parameter of the system and unique to within a scalar multiple, is orthogonal to any post-excitation output sequence of length K' regardless of the system input stimulus and sensor location. A variety of such subsequences can be extracted from a single long output sequence. But if the system can be repeatedly excited in a variety of ways, if the sensor location can be varied, or if multiple sensors can be used to record response data, then an enormous amount of information can be gotten on what θ is orthogonal to. This should enable one to accurately determine the direction of θ , which is all that matters since its magnitude is arbitrary.

Unfortunately the problem is confounded by the presence of noise in the recorded data or in numerical computations. The θ vector must be determined very precisely to ensure that the roots of $\theta^T Z$ accurately estimate the system poles. The selection of T relative to a given s-pole's frequency and decay time can greatly influence the results. If T is too short then all of the elements of an output subsequence of length K' have about the same numeric value, so that every subsequence "points" approximately in the direction represented by the single vector $(1, 1, \dots, 1)^T$. Being repeatedly told that θ is orthogonal to this vector is not much help. Consider the case of a system whose natural reverberant response is a simple undamped (or very lightly damped) oscillation at an unknown frequency ω_0 , so that there are two s-poles, $s_{1,2} = \pm j\omega_0$, so $z_{1,2} = \exp(\pm j\omega_0 T)$. Then the pole polynomial is $(z^2 - 2\cos(\omega_0 T)z + 1)$ so $\theta^T = (1, -2\cos(\omega_0 T), 1)^T$. Since this vector is orthogonal to every subsequence of length 3 we have $y(n) - 2\cos(\omega_0 T)y(n+1) + y(n+2) \equiv 0$, so ω_0 can be determined from

$\omega_0 = T^{-1} \arccos[(y(n)+y(n+2))/2y(n+1)]$. Clearly if T is so short that $y(n) \approx y(n+1) \approx y(n+2)$ then the accuracy would be very poor. A crude sensitivity analysis of this formula suggests that ω_0 can most accurately be determined when the argument of the arc cosine is near zero, which implies $T \approx \frac{1}{2}(2\pi/\omega_0)$ or some odd multiple thereof. This leads to the conjecture that, even for systems with many poles, z_k might be most accurately estimated if T is one half the oscillation period of that pole, provided the decay time of the pole is long enough that it rings loudly throughout the sampling process. In any case some control over the sampling rate is obviously desirable, but this may be achieved simply by using decimation subsequences, assuming that T is not too long an interval to begin with. Decimation is equivalent to multiplying T by an integer q (the decimation "epoch") and shifting t_0 by some integral multiple of T , and leads directly to a modified version of Eq. (8):

$$(y(n), y(n+q), y(n+2q), \dots, y(n+Kq)) \psi \equiv 0 \quad (9)$$

where $\psi \triangleq (\psi_0, \psi_1, \dots, \psi_K)^T$ is the vector of coefficients for a polynomial $\psi^T Z$ whose roots are $\{z_k^q\}$. Thus subsequences of length K' whose elements are spaced q samples apart are orthogonal to ψ , and one can determine the system poles from ψ provided the decimation aliasing ambiguity can be resolved. If T was rather short to begin with, then these decimation subsequences may point in much more varied and useful directions than if the sequence had not been decimated. Note that in decimating the data one does not need to discard anything, since many subsequences can be staggered along the original data sequence. In virtually all that follows, the techniques proposed for estimating θ can also be used to estimate ψ provided one uses decimation subsequences.

The pole identification problem has been addressed by many investigators, and the number of references to it in the literature is almost unbounded. Entire areas of control theory and speech processing have been devoted to it, but often with the following differences: (1) The system is assumed to be persistently excited, perhaps by a pulse train or by noise, (2) The system input is measured and incorporated into the analysis. However, much of the literature does pertain to our problem. Nevertheless, we shall break with tradition by presenting our results first and then discussing their relation to those of other researchers.

Finding θ when shown what it is orthogonal to is equivalent to solving a homogeneous matrix equation $Ax=0$, where the rows of the A matrix consist of output data subsequences. Because the homogeneous problem is lightly treated in most textbooks, a section of our paper summarizes some of the relevant techniques. This is followed by some autoregression matrix terminology, theorems, and algorithms developed by the author especially for this application, although their simplicity suggests that they may have been discovered in some form by mathematicians long ago. With this preparation, various approaches to solving the central problem are discussed. Connections with the work of Jain [7] are considered and a simpler and more elucidating pathway to his results is found, providing a generalization of the method. Some topics from the literature are discussed and several other results are presented before posing unanswered questions.

Our treatment departs from most of the literature in that it stresses the use of matrix methods and vector space geometry rather than "sequential" concepts, primarily because we permit the use of several output sequences derived from possibly different input excitations or sensors, and it strives to avoid ad hoc and asymmetrical treatments of the data that have been used

in the past. In particular we avoid the artificial conversion of what is fundamentally a homogeneous problem into an inhomogeneous one. We present for the first time a rigorous theoretical foundation for the popular practice of using "extra-wide" data matrices, i.e. of pretending there are more poles than actually exist.

SOLVING THE HOMOGENEOUS MATRIX EQUATION

For an $M \times N$ matrix A define the nullspace, $\phi(A)$, and rowspace, $\phi^\perp(A)$, as $\{x: Ax = 0\}$ and $\{x: x = A^T y\}$. If the rank of A is K , then $\phi^\perp(A) \subset \mathbb{R}^N$ and $\phi^\perp(A^T) \subset \mathbb{R}^M$, both being subspaces of dimension K . Also $\phi(A) \subset \mathbb{R}^N$, and it is a subspace of dimension $N-K$. Indeed $\phi(A)$ and $\phi^\perp(A)$ are orthogonal complements of \mathbb{R}^N , so that any $x \in \mathbb{R}^N$ can be expressed as an orthogonal sum, $x = x_A + x_{\bar{A}}$, where x_A is the projection onto the rowspace of A and $x_{\bar{A}}$ is the projection onto the nullspace. Given a matrix E of the same shape as A we shall denote by E_A and $E_{\bar{A}}$ the matrices obtained by decomposing the rows of E similarly; thus $E = E_A + E_{\bar{A}}$, $A_A = A$, and $A_{\bar{A}} = 0$. Moreover, occasionally we shall use the notation " $[x]_{\text{unit}}$ " to denote the unit vector in the direction of x . $A^T A$ and AA^T also have rank K and $\phi(A^T A) = \phi(A)$. The homogeneous equation $Ax=0$ is thus solved by any x in $\phi(A^T A)$. If the rank of A is $N-1$, then the solution is unique to within a scalar multiple, since $\phi(A)$ is of unit dimension. If $\text{rank}(A)$ is N , then no solution exists.

If A is a square $N \times N$ matrix, then $A \times \text{adj}(A) = \det(A) \times I$, recalling that $\text{adj}(A)$ is the matrix obtained by replacing each element of A by its cofactor and then transposing. In particular if $\text{rank}(A)$ is N then the matrix is non-singular and the inverse can be defined as $A^{-1} = \text{adj}(A)/\det(A)$. However if $\text{rank}(A) < N$, i.e. A is singular, then $\det(A) = 0$ so that $A \times \text{adj}(A) = 0$, implying that every column of $\text{adj}(A)$ solves the homogeneous equation $Ax=0$.

Unfortunately $\text{adj}(A)$ is identically zero whenever $\text{rank}(A) < N-1$, so a non-trivial solution is provided only if $\text{rank}(A)$ is $N-1$, in which case all the columns of $\text{adj}(A)$ are collinear and at least one is non-zero. Thus the solution to $Ax=0$ can be taken as any non-trivial column of $\text{adj}(A)$ or alternatively as $x = \text{adj}(A) \times x_0$ for any $N \times 1$ vector x_0 so long as it is chosen to avoid a trivial solution.

Given an $M \times N$ matrix of rank K , the traditional approach to solving $Ax=0$ is to proceed as follows: First discard rows of A so as to form a rank-preserving $K \times N$ matrix B ; since $\phi(A) = \phi(B)$, it follows that $Bx=0$ has the same solutions. Now delete columns of B to leave a $K \times K$ nonsingular matrix B_1 , and move the deleted columns to the right side of the equation together with the corresponding elements of x to produce the inhomogeneous equation

$B_1 x_1 = -B_0 x_0$, where B_0 is $K \times (N-K)$ and x has been separated into the $K \times 1$ and $(N-K) \times 1$ vectors x_1 and x_0 . If, as is usually the case, the rank of A can be preserved by simply deleting its rightmost columns, then one can simply partition $B = [B_1; B_0]$, and $x^T = [x_1^T; x_0^T]$ to achieve the desired result. Now $x_1 = -B_1^{-1} B_0 x_0$ so that the general solution is

$$x = Cx_0 \triangleq \begin{bmatrix} -B_1^{-1} B_0 \\ I \end{bmatrix} x_0 \quad (10)$$

where I is $(N-K) \times (N-K)$ and x_0 is an arbitrary $(N-K) \times 1$ vector, thus giving a solution subspace of dimension $N-K$. If something more complicated than a simple partitioning is used to form the matrices B_1 and B_0 , then Eq. (10) is still valid with the rows of C interchanged appropriately.

This traditional method is fine if the data (i.e. A) are noiseless and computations are exact. Otherwise it is worthwhile to search for better alternatives. In particular if $M \gg K$ then it is wasteful to discard so

many rows to get B, since each row provides some information on what θ is orthogonal to. It would be better to solve $Px=0$, where $P = A^T A$ so that $\phi(P) = \phi(A)$; then only a few rows of P will have to be discarded to get B, and every element of A contributes something to the solution. (An alternative would be to fabricate rows of B as averages of those of A, perhaps averaging out some of the noise in the process. But this could as easily average out some of the signal, and it would be hard to guarantee preservation of rank. This alternative will not be explored further.) Typically there are many rank-preserving ways to discard rows of P, and to choose which columns to transfer to the other side of the equation. These choices are arbitrary and may influence the accuracy of the result; indeed they control the degree to which each of the elements of A contribute to the final result. This arbitrary and asymmetrical use of data is bothersome and might produce statistical bias.

In the special case where $\text{rank}(A)$ is $N-1$, the adjoint solution can postpone or remove the arbitrary asymmetry. The direct solution to $A^T Ax=0$ is just $x = \text{adj}(A^T A)x_0$, where x_0 is an arbitrary $N \times 1$ vector. Even the arbitrariness of x_0 can be symmetrically removed. Since $Q \triangleq \text{adj}(A^T A)$ has collinear columns and is symmetric, its elements obey $q_{ij}^2 = q_{ii}q_{jj}$ and the j^{th} column can be expressed as $|q_{jj}| \times [(\pm)_1 |q_{11}|^{\frac{1}{2}}, (\pm)_2 |q_{22}|^{\frac{1}{2}}, \dots, (\pm)_N |q_{NN}|^{\frac{1}{2}}]^T$, where $(\pm)_i \triangleq \text{sign}(q_{ij})$. Thus the solution can be expressed formally as $x^T = [\Delta_1, \Delta_2, \dots, \Delta_N]^T$, where Δ_i is the square root of the magnitude of the i^{th} diagonal cofactor of $A^T A$, with the proper sign affixed. The signs can be determined by examining any column of $\text{adj}(A^T A)$.

The adjoint solution enables the development of an unbiased estimate, based on the following theorem:

§ 1 THEOREM: Let A_1, A_2 be statistically independent, random matrices,

whose elements are also statistically independent. Then $E\{\text{adj}(A_1^T A_2)\} = \text{adj}(E\{A_1^T\}E\{A_2\})$, where E is the expectation operator. \square

Thus if A_1 and A_2 are separate measurements of the same A corrupted by additive, zero-mean noise that is independent from element to element, then the estimate $\hat{x} = \text{adj}(A_1^T A_2)x_0$ is an unbiased estimate of a true solution x .

The singular value decomposition (SVD) can always be used to obtain a solution; indeed its applicability overlaps that of the adjoint solution. The SVD of an $M \times N$ matrix A is expressed by $U^T A V = S$ and $A = U S V^T$ where S is an $M \times N$ diagonal matrix of elements $s_1 \geq s_2 \geq \dots \geq s_{K+1} = s_{K+2} = \dots = s_N = 0$, where rank (A) is K . The "singular values" s_i are the square roots of eigenvalues of the non-negative definite matrix $A^T A$, whose eigenvectors also form the columns of the $N \times N$ orthogonal matrix V . The eigenvectors of $A A^T$ form U . (Incidentally $A^T A$ and $A A^T$ agree as to their non-zero eigenvalues.)

The notation $sv_i(A)$ will denote the i^{th} singular value of any matrix A .

The singular values can be used to bound $\|Ax\|$ for any x as follows [8]:

$$\|x\| \times sv_1(A) \geq \|Ax\| \geq \|x\| \times sv_N(A) \text{ and } \|Ax_A\| \geq \|x_A\| \times sv_K(A)$$

where rank (A) is K . Clearly $\|A\| \equiv sv_1(A)$. (Note: $\|x\|$ denotes the Euclidean norm of the vector, and $\|A\|$ is defined as the maximum of $\|Ax\|$ for all unit vectors x .) The U and V matrices are not quite unique; the

direction of any column vector can be reversed, and corresponding to a multiple eigenvalue (including the zero eigenvalues) any orthonormal basis

for the eigenspace can be used. The homogeneous equation $Ax=0$ now takes

the form $U S V^T x = 0$ or simply $Sy=0$ where $y \triangleq V^T x$, and we have premultiplied

the equation by $U^T = U^{-1}$. But due to the diagonal form of S , the solution is

immediate: $y = [0, 0, \dots, 0; y_0^T]^T$ where y_0 is an arbitrary $(N-K) \times 1$ vector.

Thus $x = V y = V_0 y_0$, where V_0 consists of the rightmost $N-K$ columns of V .

But these columns are just the basis vectors for the eigenspace of the zero

eigenvalue of $A^T A$, i.e. basis vectors for $\phi(A)$, so it seems we have only reiterated that $Ax=0$ implies $x \in \phi(A)$. But the superiority of the SVD method emanates from the fact that, due to noise or computing inaccuracy, $A^T A$ will not be exactly singular; the SVD method will reveal the extremely small eigenvalues that arise, and the associated columns of V should approximately span $\phi(A)$. Furthermore the SVD can be achieved by a very orderly procedure. An extremely well documented algorithm and FORTRAN program appear in the textbook of Lawson and Hanson [9] under the label SVDRS. (Note: In their notation one should set $BA=0$ so that array B is not referenced.)

When the diagonal matrix S resulting from the SVD algorithm is replaced by \hat{S} , wherein all but the first k elements (i.e. the k largest) have been forced to zero, and then used to compute $\hat{A}=U\hat{S}V^T$, the result is an optimum approximant to A of rank k , denoted $\hat{A}(k;SVD)$. It has been shown [9] to be closest as measured by either the ordinary matrix norm, $\|A-\hat{A}\|$, or the Frobenius norm, $\|A-\hat{A}\|_F$, where $\|A\|_F \triangleq (\sum a_{ij}^2)^{1/2}$. Thus if A has been corrupted by noise then a logical estimate of the nullspace is $\hat{\phi}(A) = \phi(\hat{A}(K;SVD))$, where K is the "true" rank A would have if it were uncorrupted by noise. This estimate of the nullspace thus provides a general solution to $Ax=0$. Using the approximant to solve nonhomogeneous equations is a common practice. The justification for our approach to solving the homogeneous equation is strengthened by the following theorem.

§2 THEOREM: If $\tilde{A}=A+E$, where nothing is known of the $M \times N$ matrix A except that it is of rank K , and the elements of E are zero-mean, independent, Gaussian random variables having equal variance, then $\hat{\tilde{A}}(K;SVD)$ and its nullspace are maximum likelihood estimates of A and $\phi(A)$, respectively. Moreover, if it is known only that A is singular, with $\text{rank}(A)$ unknown, then the maximum likelihood estimate is obtained with $K=N-1$. If on the other hand K is

random with known probability distribution P_K , then choose K to minimize $\|\tilde{A} - \hat{A}(K; \text{SVD})\|_F^2 - 2\sigma^2 \log P_K$, where σ^2 is the noise variance. \square

If only a single solution to $Ax=0$ is desired rather than the general (nullspace) solution, then one can simply point x in the direction of the rightmost column vector of V , since it is the basis vector of $\phi(A)$ that has possibly been perturbed the least. But that is the same as picking x to be the eigenvector corresponding to the smallest eigenvalue of $\tilde{A}^T \tilde{A}$; i.e. picking x to minimize $(x^T \tilde{A}^T \tilde{A} x) / \|x\|^2$ or $\|\tilde{A}x\| / \|x\|$. When normalized to $\|x\|=1$ we shall denote this estimate by $\hat{x}[\text{norm}]$, where "norm" stands for both "normalized" and "norm-minimizing", but it is only unique to within a direction reversal. Since the addition of noise can be expected to destroy all singularity of A , and even eliminate any multiple eigenvalues, then $\tilde{A}^T \tilde{A}$ will be nonsingular and its smallest eigenvalue is the largest of $(\tilde{A}^T \tilde{A})^{-1}$. Matrix iteration, $x_{i+1} = [(\tilde{A}^T \tilde{A})^{-1} x_i]_{\text{unit}}$, will converge to $\hat{x}[\text{norm}]$. However since $\tilde{A}^T \tilde{A}$ is nonsingular only because of noise, it may be poorly conditioned and difficult to invert accurately, so the method must be used cautiously. Of course direct matrix inversion can be avoided by numerically solving the equation $\tilde{A}^T \tilde{A} x_{i+1} = x_i$ with subsequent normalization for each iteration. It is interesting to note that since $(\tilde{A}^T \tilde{A})^{-1} = \Delta \times \text{adj}(\tilde{A}^T \tilde{A})$, where the scalar Δ is just the determinant, the adjoint solution discussed previously may be regarded as one step in the matrix iteration process* (except for the trivialities of normalization and possible direction reversal) from an arbitrary starting point x_0 . Indeed an "improved" adjoint solution may be put forth as $[\text{adj}(\tilde{A}^T \tilde{A})]^\ell x_0$ for any integer ℓ . In the absence of noise it will still

* A connection between the adjoint solution and the traditional least-squares solution of the inhomogeneous equation will be given in a later section.

give the same solution when rank(A) is N-1, and in the presence of noise it tends toward the direction of \hat{x} [norm]. Moreover by using several separate measurements of A in conformance with Theorem §1 one can improve the unbiased estimate in a similar manner.

When A is corrupted by noise one expects that \hat{x} will depart from $\phi(A)$. The following theorem provides an upper bound for that error.

§ 3 THEOREM: If the unit vector \hat{x} is an estimated solution to $Ax=0$ constructed to be within the nullspace of $\hat{A}(K;SVD)$ where $\tilde{A}=A+E$ and rank(A) = K, then the error \hat{x}_A is bounded according to $\|\hat{x}_A\| \leq b_1, b_2, b_3, b_4$; and if $\|E\| \ll sv_K(A)$ then the denominator of each bound is approximately $sv_K(\tilde{A})$, where $b_1 \triangleq \{\|\tilde{A}\hat{x}\| + \|E\|\}/sv_K(A)$; $b_2 \triangleq 2\|E\|/sv_K(A)$; $b_3 \triangleq \{\|\tilde{A}\hat{x}\| + \|E_A\|\}/sv_K(A+E_A)$; $b_4 \triangleq 2\|E_A\|/sv_K(A+E_A)$. Moreover $b_1 \leq b_2 \leq 2b_1$ and $b_3 \leq b_4 \leq 2b_3$. \square

Discussion: Note that the theorem includes $\hat{x} = \hat{x}$ [norm] as a special case. \hat{x}_A is the "correct" portion of \hat{x} (i.e. the component actually in $\phi(A)$) and \hat{x}_A is the error. Since \hat{x} is a unit vector, $\|\hat{x}_A\|$ is the sine of the angle between \hat{x} and $\phi(A)$. (The angle between two unit vectors is the arc cosine of their inner product; the angle between a vector and a subspace is then defined as the smallest such angle, always taken positive.) Note that $\|E\|$ is the square root of the largest eigenvalue of $E^T E$, with a similar result holding for E_A . Also $(sv_K(A))^2$ is the smallest non-zero eigenvalue of $A^T A$. Clearly if $\|E\| \ll sv_K(A)$ then $b_2 \ll 1$ so the error component is small, but the b_4 bound shows that error occurs only when some portion of E "complies" with the row space of A. Observation of a single noisy matrix \tilde{A} will tell us little of the error, E; thus knowledge of $\|E\|$ or $\|E_A\|$ must be obtained a priori, e.g. in statistical form. The singular values of \tilde{A} can be used a posteriori to approximate the denominators of the bound s. Since bounds b_2 and b_4 are almost as tight as b_1 and b_3 , use of $\|\tilde{A}\hat{x}\|$ a posteriori does

little to refine our error estimate. The distinction between $||\tilde{A}\hat{x}||$ and $||x_A||$ is important: The "residual error", i.e. $||\tilde{A}\hat{x}|| = sv_N(\tilde{A})$, measures the degree to which the noisy matrix \tilde{A} refuses to admit a homogeneous solution. On the other hand $||\hat{x}_A||$ measures the actual error in the estimate \hat{x} . Although scaling a row of A does not affect the solution of the ideal equation $Ax=0$, scaling of a row of \tilde{A} may influence the error in \hat{x} since it can affect $sv_K(\tilde{A})$.

Even when \hat{x} is truly a solution of $Ax=0$, the "residual" vector $\tilde{A}\hat{x}$ cannot be expected to be a null vector. In some problems the statistics of the residual vector may be known, at least approximately. Indeed suppose that it is zero-mean with positive definite covariance matrix Λ . Then instead of choosing \hat{x} to minimize $||\tilde{A}\hat{x}||$ it is more equitable to minimize the weighted norm $||\tilde{A}\hat{x}||_{\Lambda^{-1}} = (x^T \tilde{A}^T \Lambda^{-1} \tilde{A} x)^{1/2}$. We shall denote by $\hat{x}[\text{norm}:\Lambda^{-1}]$ the unit vector that minimizes this norm. It is the weakest eigenvector of $\tilde{A}^T \Lambda^{-1} \tilde{A}$. The SVD method can also be modified to incorporate this covariance weighting. One simply does the SVD of $\Lambda^{-1/2} \tilde{A}$ instead of \tilde{A} . The matrices A and $A_w \stackrel{\Delta}{=} \Lambda^{-1/2} \tilde{A}$ have the same nullspace and $\hat{A}_w(k; \text{SVD})$ is the best k -rank approximant to A_w , minimizing $||A_w - \hat{A}_w|| = ||A - \Lambda^{1/2} \hat{A}_w||_{\Lambda^{-1}}$. So $\Lambda^{1/2} \hat{A}_w$ is the best weighted approximant to A . Therefore the corresponding weighted estimate of $\phi(A)$ is $\phi(\Lambda^{1/2} \hat{A}_w) = \phi(\hat{A}_w)$. All of the results given in this section hold when the weighted norm is substituted for the ordinary norm, with appropriate interpretations and adjustments.

AUTOREGRESSION MATRICES, GENERATORS, AND SUBSPACES

§ 4 DEFINITIONS: For any $K \times 1$ vector θ , the "polynomial produced by θ " will refer to $\theta^T Z$, where Z is the complex power vector defined previously. Any vector θ is said to be a "generator" if its first and last elements are both

non-zero and the roots of the "generator polynomial" $\theta^T Z$ are distinct. Generators that differ by a scalar multiple are considered equivalent. (Note: the interpretation of the elements of θ as coefficients of the system pole polynomial will be ignored until the next section.) Given any generator θ and a positive integer ℓ define the $\ell \times N$ matrix $G(\ell; \theta)$, where $N = K' - 1 + \ell = K + \ell$, as:

$$G(\ell; \theta) = \begin{bmatrix} \theta_0, \theta_1, \dots, \theta_K, 0, 0, \dots, 0 \\ 0, \theta_0, \theta_1, \dots, \theta_K, 0, \dots, 0 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ 0, 0, \dots, 0, \theta_0, \theta_1, \dots, \theta_K \end{bmatrix} \quad (11)$$

Now define the " ℓ^{th} autoregression nullspace generated by θ ", denoted $\Omega_\ell(\theta)$, as the row space of G . Then it is easy to prove the following theorems:

§ 5 THEOREM: Rank(G) is ℓ , so that $\Omega_\ell(\theta)$ is a subspace of R^N having dimension ℓ , where $N=K+\ell$. Moreover the generator of the subspace Ω_ℓ is unique to within a scalar multiple. \square

§ 6 THEOREM: For every $N \times 1$ vector x in $\Omega_\ell(\theta)$, the roots of the polynomial $x^T Z$ constitute a superset of the roots of the generator polynomial $\theta^T Z$. The "extraneous roots" are $\ell-1$ in number, and any complex extraneous roots must occur in complex conjugate pairs. Indeed an x in Ω_ℓ can be found to produce any specified set of extraneous roots. \square

§ 7 THEOREM: The subspace $\Omega_\ell(\theta)$ can also be defined as that subspace of R^N which is orthogonal to each of the complex power vectors Z_i , of the K roots of the generator polynomial z_i , with $N=K+\ell$. \square

§ 8 DEFINITION: A sequence $\{y(n)\}$ is said to be " θ -autoregressive" (θ -AR), where θ is $K' \times 1$, if every subsequence of length K' forms a vector orthogonal

to θ . This implies that for $n \geq K$, $y(n) = \theta_K^{-1} \times (\theta_{K-1}y(n-1) + \theta_{K-2}y(n-2) + \dots + \theta_0y(n-K))$. Thus the first K members of the sequence are arbitrary and will be termed the "seed", while the remaining members are determined by a recursion formula. Clearly if any leading group of members is deleted, the remaining sequence is still θ -AR.

§ 9 THEOREM: The elements of an $N \times 1$ vector y form a θ -AR sequence if and only if $y \in \Omega_\ell(\theta)$ where $\ell = N-K = N-(K'-1)$ and θ is $K' \times 1$. \square

§ 10 DEFINITION: An $M \times N$ matrix A is said to be θ -AR if each row constitutes a θ -AR sequence. Clearly the definition makes no sense unless $N \geq K'$. The matrix is termed "minimal-width" if $N=K'$ and "extra-wide" if $N > K'$. If $\text{rank}(A)$ is $K(=K'-1)$, which according to the next theorem is the most it can be, then the θ -AR matrix A is said to be of "sufficient rank".

§ 11 THEOREM: If the $M \times N$ matrix A is θ -AR where θ is $K' \times 1$, then $\text{rank}(A) \leq K=K'-1$. Further, if A is of sufficient rank then $\phi(A) \equiv \Omega_\ell(\theta)$ where $\ell=N-K$, and any solution of $Ax=0$ produces a polynomial of degree $N-1$ whose roots constitute a superset of the roots of the generator polynomial. \square

§ 12 MORE DEFINITIONS AND DISCUSSION: Clearly a θ -AR matrix is the continuation eastward of its "seed submatrix" consisting of the leftmost K columns. If it happens that A^T is θ -AR as well, then the seed submatrix may instead be regarded as the $K \times K$ matrix in the northwest corner, since from it the entire matrix can be propagated autoregressively. A sufficient (but not necessary) condition that a square seed matrix can propagate such a matrix is that the seed be symmetric. A special case is the "Hankel" θ -AR matrix, so named because the $M \times N$ matrix A is formed from a single θ -AR sequence $\{h(n)\}$ of length $M+N+1$ as $A_{ij} = h(i+j-2)$. The Hankel matrix has an interesting property: For any $N \times 1$ vector x , $Ax = G(M;x)h$ where h is the vector whose elements constitute $\{h(n)\}$. This property is just a consequence of the fact that

discrete convolution is commutative. The numbering pattern for a 8×6 Hankel matrix is shown in Figure 1. For this case there is no need to define the kernel as a submatrix: the kernel simply consists of the first K matrix elements encountered along either edge starting from the northwest corner, and the entire matrix is propagated from it via the θ -AR property. Any submatrix of contiguous rows and columns remains a Hankel θ -AR matrix, and A and A^T are both θ -AR.

The central problem considered in this paper may be stated succinctly as follows: Given a θ -AR matrix A or a noisy version thereof, \tilde{A} , where the $K' \times 1$ generator θ is unknown, find the roots of the generator polynomial. Finding θ is normally an intermediate step, albeit one we would prefer to bypass since the roots may be poorly conditioned with respect to the polynomial coefficients. Various paths to the solution are possible, given a matrix A of sufficient rank and at least minimal-width. From Theorem §11 $\phi(A) = \Omega_\lambda(\theta)$, and any x solving $Ax=0$ produces a polynomial whose roots include the desired roots of the generator polynomial. If the extraneous roots can somehow be identified then the problem is solved. If instead we find ℓ linearly independent solutions of $Ax=0$, i.e. a basis for $\Omega_\lambda(\theta)$, then θ can be determined by a process described below without having to identify the extraneous poles. If a noisy matrix \tilde{A} is used then the method used to solve $Ax=0$ becomes significant. The estimate $\hat{x}[\text{norm}]$ can be found as the weakest eigenvector of $\tilde{A}^T \tilde{A}$ and used to produce a polynomial whose extraneous roots may now be regarded as "noise poles". The more traditional methods of solving the equation may also be used; they may entail less computation but perhaps more error. If the A matrix is of minimal width then any solution of $Ax=0$ is collinear with θ and there are no extraneous roots. However the roots may be perturbed greatly by noise in A .

Suppose that a set of ℓ basis vectors g_1, g_2, \dots, g_ℓ has been found for

$\phi(A) = \Omega_\ell(\theta)$; i.e. the g_i 's are linearly independent solutions of $Ax=0$. In practice the g_i 's may be obtained as the rightmost columns of the V matrix in the SVD expansion of A, or by some cruder technique. For example if some quick method is used to solve the noisy equation $\tilde{A}x=0$ then repetition of the solution after tinkering with A (i.e. inserting new data, scaling rows, etc.) would likely result in linearly independent solutions. Regardless of how obtained, the solution vectors g_i can be used to find θ through the following algorithm:

§ 13 ALGORITHM: First define G_1 as the $\ell \times N$ matrix whose rows are $g_1^T, g_2^T, \dots, g_\ell^T$, so that the rowspace of G_1 is $\phi(A) = \Omega_\ell(\theta)$, thus identical to the rowspace of $G(\ell; \theta)$ defined in Eq. (11). Then starting with the topmost row use Gaussian elimination to form the upper-trapezoidal matrix G_2 as exemplified in Figure 2(b). Then repeat the process starting with the bottom row and eliminating terms on the right to form the upper-parallelogramic matrix G_3 exemplified in Figure 2(c). Assuming the rank of G_1 has been preserved the rowspace of G_3 is the same as that of $G(\ell; \theta)$, so the non-zero portion of each row of G_3 is a replica of θ^T to within a scalar multiple (If a noisy matrix \tilde{A} is used then the rows are only estimates of θ .) Now define G_4 as the $\ell \times K'$ matrix formed by "straightening out" the non-zero parallelogram of G_4 and discarding the zeroes. In the absence of noise or computational error the rows are collinear and $\text{rank}(G_4)$ is 1. If noise is present then each row is an estimate of θ^T . Logically a "best" estimate of θ^T can be obtained by using the SVD to find the best unit-rank approximant to G_4 , whose collinear rows then produce identical estimates of θ^T . But this is completely equivalent to estimating θ as the strongest eigenvector of $G_4^T G_4$, a task that can be performed very easily with matrix iteration.

The procedure presented above can be modified so as to constitute a

stepwise reduction of $\Omega_\ell(\theta)$ to $\Omega_{\ell-1}(\theta)$, ultimately arriving at $\Omega_1(\theta)$ whose single vector direction is θ .

§ 14 THEOREM: Let the $N \times 1$ vectors g_1, g_2, \dots, g_ℓ form a basis of $\Omega_\ell(\theta)$, of which at least one must then have a non-zero first element; assume it is g_1 . Furthermore assume the g_i 's are scaled so that the first element of each is either 1 or 0. Then for $i=2$ to ℓ define $a_i \triangleq g_i - g_1$ if the first element of g_i is 1, otherwise $a_i \triangleq g_i$. Then the set $\{g_1, a_2, a_3, \dots, a_\ell\}$ is still a basis of $\Omega_\ell(\theta)$, and the set of contracted vectors $\{\bar{a}_2, \bar{a}_3, \dots, \bar{a}_\ell\}$ obtained by discarding the first element of each a_i forms a basis of $\Omega_{\ell-1}(\theta)$. \square

§ 15 COROLLARY: The theorem remains valid if the words "last element" are substituted for "first element" at every occurrence. \square

It should be remembered that $\Omega_\ell \subset \mathbb{R}^N$ while $\Omega_{\ell-1} \subset \mathbb{R}^{N-1}$. If in the process of stepwise reduction one uses Theorem §14 repeatedly the final result will be a single vector equivalent to the bottom row of the matrix G_4 obtained with Algorithm §13. But if one switches at some point to using the corollary then the result will correspond to one of the other rows of G_4 .

Now we present a formula for the sensitivity of the roots of the generator polynomial, when estimated as the non-extraneous roots of the polynomial produced by \hat{x} , where \hat{x} is some approximate solution to $Ax=0$. Recall that if \hat{x} is a unit vector then its error component is measured by $\|\hat{x}_A\|$.

§ 16 THEOREM: Suppose $\|\hat{x}\|=1$ and \hat{x} lies approximately within $\phi(A)$, i.e. $\|\hat{x}_A\| \ll \|\hat{x}_A^-\|$. Then the non-extraneous roots of the polynomial $x^T Z$ depart from the true roots z_1, z_2, \dots, z_K of the generator polynomial according to the first-order formula:

$$\frac{dz_i}{z_i} = \frac{-\hat{x}_A^T Z_i}{\hat{x}_A^T D Z_i}$$

where D is a diagonal matrix, $D_{11}=0$, $D_{kk}=1/k$ for $k=2$ to N . \square

If in particular \hat{x} is $\hat{x}[\text{norm}]$, then the above result can be combined with Theorem §3 to produce an approximate error bound in terms of $\|E_{\tilde{A}}\|$ and the non-extraneous roots, z_i , of the polynomial it produces. (The bound is approximate because the sensitivity formula is correct only to first order, and estimated roots are used.)

$$\left| \frac{dz_i}{z_i} \right| \lesssim \frac{2 \|E_{\tilde{A}}\| \times \left[\sum_{k=0}^K |z_1^k|^2 \right]^{1/2}}{sv_K(\tilde{A}) \times |x^T D z_1|} \quad (12)$$

In the problem of finding θ given a noisy version θ -AR matrix A, it might appear that Theorem §2 should always apply, and that $\hat{\phi} = \phi(\hat{A}(K;SVD))$ is a maximum likelihood estimate of $\phi(A)$. However, there are two critical assumptions in the hypothesis of Theorem §2: (1) that the noise corrupts the elements independently, and (2) that absolutely nothing of A is known except for its rank. Clearly if A is known only to be of minimal-width and sufficient rank with θ unknown, then we know only that its rank is one less than its width. (If any $M \times N$ matrix A has $\text{rank}(A) = N-1$, then it is automatically θ -AR for the vector θ that spans $\phi(A)$.) So if assumption (1) is satisfied then Theorem §2 does indeed apply. On the other hand if A is extra-wide and θ -AR then we know more than just its rank, so assumption (2) is violated. If A is known to have been formed as a Hankel matrix then both assumptions are violated. (But if by some fortuitous circumstance $\hat{A}(K;SVD)$ does indeed conform to all of our prior knowledge of A, then the conclusion of Theorem §2 should still apply.) These facts might lead one to avoid using either extra-wide or Hankel matrices. However Hankel matrices appear to be very economical since each new element of data permits the inclusion of another row. Furthermore several investigators have proven that if one uses an extra-wide Hankel matrix and

traditional methods of solving $Ax=0$, then the non-extraneous roots of the resulting polynomial often estimate the true roots much more accurately than if a minimal-width matrix had been used [10,11,12]. Clearly there remain many unanswered questions regarding the efficacy of the various alternatives.

POLE IDENTIFICATION

To apply the methods developed above to the pole identification problem posed in the introduction one has merely to observe that every post-excitation output sequence that can be elicited from a system of order K is θ -AR, where the generator θ is the $K \times 1$ vector of the pole polynomial's coefficients. Thus every output subsequence of length $N=K+l$ will be orthogonal to $\Omega_l(\theta)$, which we shall now denote simply as Ω_l , the " l^{th} polespace", to concede its relationship to the system poles. Of course Ω_1 is just the space spanned by θ itself. Thus any $M \times N$ data matrix A whose rows are output sequences will obey $\phi(A) = \phi(A^T A) = \Omega_l$, where $l = N-K$, provided enough rows have been used to give the matrix sufficient rank (i.e. $\text{rank}(A) = K$). This will hold true even if the various rows were obtained as a result of separate tests, with different input excitation, or different sensor placement. Any solution of $Ax=0$ will produce a polynomial whose roots include the system poles together with $N-K'$ extraneous ones. If decimation sequences of epoch q are used, then the appropriate subspace is $\Omega_l(\psi)$, which we shall denote as Ω_l^q . The $K' \times 1$ vector ψ is formed of the coefficients of the polynomial whose roots are the q^{th} powers of the system poles.

§ 17 PRONY'S METHOD [14]: Given a single output sequence of length $2K$, form an $M \times N$ minimal-width Hankel matrix A , where $M=K$ and $N=K'$. Then solving $Ax=0$ gives the generator θ . Traditionally this is done by converting to a non-homogeneous equation, which in this case is entirely equivalent to forcing $\theta_K=1$. Since the equation is not overspecified it always has an exact

solution, and the only error is due to noise in A. Experience proves that even a small amount of noise can be devastating [14].

§ 18 LEAST-SQUARES PRONY METHOD [10,11,12]: Given an output sequence of length greater than $2K$ one can still form an $M \times N$ minimal, width Hankel matrix A, where $N=K'$ but now $M > K$. The equation $Ax=0$ is now overspecified. Following the traditional approach discussed earlier one seeks to determine the null-space $\phi(A) = \phi(A^T A)$ by solving $A^T Ax=0$. Even this equation is overspecified since $A^T A$ is $K' \times K'$ but $\text{rank}(A^T A) = K$. If the rightmost column of A can be partitioned off without altering the rank, i.e. $A = [A_+; a]$ where "a" is the rightmost column of A and the submatrix A_+ is nonsingular, then the equation $A^T Ax=0$ can be expressed as

$$\begin{bmatrix} A_+^T A_+ & A_+^T a \\ \hline a^T A_+ & a^T a \end{bmatrix} \cdot x = 0.$$

The solution, normalized so that the last element of x is one, can be expressed immediately as

$$x = \begin{bmatrix} -(A_+^T A_+)^{-1} A_+^T a \\ -\frac{1}{a^T a} \end{bmatrix} \quad (13)$$

Since A is a minimal width data matrix the solution, x, given by Eq. (12) should be the generator vector θ , since the solution of $Ax=0$, or equivalently $A^T Ax=0$, is unique to within a scalar multiple. Thus the same answer would result if one used the adjoint solution $x = \text{adj}(A^T A)x_0$, regardless of the choice of x_0 . On the other hand if a corrupted version of A is used, namely $\tilde{A} = A+E$ where the noise matrix E increases the rank so that $\text{rank}(\tilde{A})$ is K' , then one expects that the adjoint solution would produce a result different from that of Eq. (13). However the following theorem

establishes a connection.

§ 19 THEOREM: Let $A = [A_+; a]$ where $\text{rank}(A) = \text{rank}(A_+) = K$ and A_+ is non-singular and let $\tilde{A} = A+E = [\tilde{A}_+; \tilde{a}]$ where $\text{rank}(\tilde{A})$ is K' . Then

$$\text{adj}(\tilde{A}^T \tilde{A}) x_0 = \begin{bmatrix} -(\tilde{A}_+^T \tilde{A}_+)^{-1} \tilde{A}_+^T \tilde{a} \\ -\frac{1}{\alpha} \end{bmatrix}$$

provided

$$x_0 = [0, 0, \dots, 0, \alpha]^T$$

$$\alpha \triangleq \frac{\tilde{a}^T \{I - \tilde{A}_+ (\tilde{A}_+^T \tilde{A}_+)^{-1} \tilde{A}_+^T\} \tilde{a}}{\det(\tilde{A}^T \tilde{A})}. \quad \square$$

Many users of this least-squares Prony Method have demonstrated quite clearly that when one "requests" more poles than actually exist, i.e. when one constructs A to be extra-wide, the resulting non-extraneous poles estimate the system poles with greater accuracy [10,11,12]. This fact persists even when data are generated artificially to ensure that only a finite number of poles are really present. Theoretically $(A_+^T A_+)^{-1}$ should not even exist, since $A_+^T A_+$ is then an $(N-1) \times (N-1)$ matrix of rank $K < N-1$. Clearly it is possible to carry out the computation of Eq. (13) only because of noise inherently present in A or introduced by imperfect computation of $A_+^T A_+$. Indeed in the presence of noise we may regard Eq. (13) as a special case of the adjoint solution through Theorem §19. Thus when an extra-wide data matrix is used the least-squares Prony Method succeeds only because of noise, and Eq. (13) produces an (approximate) solution to the equation $A^T A x = 0$. Moreover the fact that the roots of the resulting polynomial $x^T Z$ include the true system poles as a subset has previously been observed only experimentally; it has remained unjustified by theory until now (i.e. in our Theorem §6). But the fact that poles can be estimated more accurately by using an extra-wide matrix remains to be justified by theory. (The reader may argue this point

with references to "generalized least squares", noise modelling, decorrelation of residuals and the like, but we shall contend in a later section that such an explanation has serious logical gaps, at least as it pertains to our problem.)

The methods described in preceding sections of this paper provide many alternate avenues toward the solution. Specific approaches, results, and adaptations are discussed below:

§ 20 CONSTRUCTING DATA MATRICES FROM A SINGLE OUTPUT SEQUENCE $\{y(n):n=0,1,2, \dots\}$: The construction of a Hankel matrix as described under the least squares Prony method is an obvious course of action. However in that case the noise contributions are not statistically independent from element to element in the matrix, since the noises appear repeatedly along with the elements of $\{y(n)\}$. But it is possible to construct an $M \times N$ data matrix in which the noises are independent by using a "sliding grid" of decimated subsequences, for example $A_{ij} = y_r$ where $r = (i-1) + (j-1)q$ and $q (> M)$ is the decimation epoch. An advantage of the noise independence is that it is possible to apply Theorems §1 and §2. But in applying Theorem §1 to obtain an unbiased estimate of θ it is necessary to have two replicates of the same sequence $\{y(n)\}$, each with different (and totally independent) noise components, to construct A_1 and A_2 . Also note that Theorem §2 is of somewhat limited value, since the K -rank approximant $\hat{A}(K;SVD)$ probably will not strictly possess the "sliding grid" characteristic (i.e. its rows cannot be reassembled into a single autoregressive sequence), and thus $\hat{\theta}$ is not the true maximum likelihood estimate. Whether or not the construction of matrices with independent noise components has other benefits beyond application of Theorems §1 and §2, or perhaps even has drawbacks, has not been demonstrated.

§ 21 USING ROWS OF A DERIVED FROM DIFFERENT EXCITATIONS OR SENSORS: A

single output sequence $\{y(n)\}$ may reveal some poles only weakly, and provide only a brief glimpse of poles having a rapid decay rate (i.e. small $|z_i|$). The advantages of obtaining separate output sequences, using totally different system excitations (to within the limits of one's ability to control the excitation at all), and then using them to construct more rows of A, has been largely ignored in the literature. Further, the use of independent rows may enable use of Theorems §1 and §2 to obtain "nice" estimates of θ .

§ 22 USING "BETTER" SOLUTIONS TO $Ax=0$. Since investigations of the past have usually converted $Ax=0$ to a nonhomogeneous problem, the cost-effectiveness of using the other solutions we have discussed is worth exploring. Indeed the norm minimizing solution \hat{x} [norm] is not that much more difficult to obtain. The least-squares Prony solution of Eq. (13) is usually obtained by solving $(A_+^T A_+)x_+ = A_+^T a$, where x_+ is all of x except for its last element, rather than actually inverting the $(N-1) \times (N-1)$ matrix. However if an extra-wide matrix A is being used (which is almost always the case) then the matrices $(\tilde{A}_+^T \tilde{A}_+)$ and $(\tilde{A}^T \tilde{A})$ are nonsingular only because of noise, and matrix iteration with $(\tilde{A}^T \tilde{A})^{-1}$ can find \tilde{x} [norm]. But this can be done by successively solving $\tilde{A}^T \tilde{A} x_{i+1} = x_i$, with occasional renormalization. Indeed the solution for x found from Eq. (13) could be used as the initial guess for \hat{x} . Thus instead of solving a single $(M-1) \times (M-1)$ nonhomogeneous matrix equation we solve an $M \times M$ equation repeatedly; hopefully a few repetitions will suffice. Of course with A extra-wide, \hat{x} 's polynomial will have some extraneous poles to sort out.

§ 23 EXPUNGING EXTRANEIOUS POLES: When an extra-wide A is being used, Algorithm §13 or Theorem §14 provide a means for eliminating the extraneous poles if K linearly independent solutions of $Ax=0$ are first found. These independent solutions can be obtained by use of SVD or perhaps simply by repeating the solution after some tinkering is performed on A (for example

by the introduction of some new data containing new noise). Whether the roots of the resulting θ -polynomial will be more accurate estimates of the true poles than are the non-extraneous roots of the larger \hat{x} -polynomial remains to be determined, but at least this approach can reveal which roots of the \hat{x} -polynomial are extraneous in a fashion reminiscent of the de-aliasing procedures described earlier.

§ 24 PREFILTERING OF DATA: It may be desirable to prefilter a data sequence $\{y(n)\}$ by some simple digital filter, for example to improve the signal-to-noise ratio. This may be done to enhance the accuracy in determining a particular pole or set of poles, as by the use of a bandpass filter. After such filtering the data sequence has the filter poles incorporated into its generator polynomial. This increase in the number of poles must be taken into account when constructing the A matrix. These new poles are known, and ought to be forced into the solution somehow (see below). Moreover if the use of a particular prefilter enables us to accurately determine some subset of the poles, then when we use another prefilter to enhance estimation of other poles we should force our previously estimated poles into the solution.

§ 25 FORCING KNOWN POLES: The use of prefilters is not the only impetus for wanting to force known poles. The choice of the sampling interval T or decimation epoch "q" may be tailored to accurate estimation of particular poles, and once determined these estimates should be forced into subsequent analyses done with different T or q. Forcing poles is not difficult. Assuming for simplicity that A is a matrix of ordinary (i.e. not decimated) data, then $\phi(A) = \Omega_\ell$. And the power vector Z_i of each system pole z_i is orthogonal to Ω_ℓ by Theorem §7. But if a z_i is known then we need only solve $Ax=0$ subject to $Z_i^T x=0$. If the pole z_i is real, then $Z_i^T = [z_i^0, z_i^1, \dots, z_i^N]$ can merely be adjoined to A as a new row with a large scaling constant C. The larger C,

the more strongly the pole z_i is forced into the solution of $A_f x = 0$, where A_f is the new matrix formed by addition of this new row. If the pole z_i is complex then the same procedure could be used, but this has the unfortunate effect of making A complex. A better approach is afforded by realizing that complex poles occur only in complex conjugate pairs, so that instead of adjoining Z_i^T and $Z_{i+1}^T = [(z_i^*)^0, (z_i^*)^1, \dots, (z_i^*)^N]$ one adjoins $a_1^T \triangleq C/2 [Z_i^T + (Z_i^*)^T]$ and $a_2^T \triangleq -\frac{1}{2} jC [Z_i^T - (Z_i^*)^T]$ which are both purely-real row vectors, and can be expressed as

$$a_1^T = C[1, r \cos \omega, r^2 \cos 2\omega, \dots, r^N \cos N\omega],$$

$$\text{and } a_2^T = C[1, r \sin \omega, r^2 \sin 2\omega, \dots, r^N \sin N\omega]$$

with $r \triangleq |z_i|$ and $\omega \triangleq \arg(z_i)$.

Thus for each pole forced, one more row is added to A . The choice of the scaling factor C is ad hoc, and some adjustment may be necessary.

§ 26 USE OF DECIMATED DATA. Everything that has been said concerning the linkage between A and θ when A is not composed of decimated data carries over to a linkage between A and ψ when A is decimated, where ψ defines a polynomial whose roots are the q^{th} powers of the system poles. Obvious modifications are needed in a few places.

THE METHOD OF JAIN AND ITS EXTENSIONS

Transformation of a complex variable according to some formula $\zeta = f(z)$ is often employed in polynomial root solving programs, and the resulting roots are then transformed back into the z -space. Previously we have seen that if one transforms the pole polynomial produced by θ to a new polynomial whose roots are the q^{th} powers of the true poles, the new ψ vector is related to decimated output sequences in the same way that θ is related to undecimated sequences $\{y(n)\}$. Suppose instead that the pole polynomials were modified by a simple linear transformation of the variable z . Would the new polynomial

coefficient vector be related to some other modification of the data sequence $\{y(n)\}$? The answer is yes. Suppose the transformation is $\zeta = [1-c_1z]/c_2$ so $z = [1-c_2\zeta]/c_1$, converting the θ -polynomial to a μ -polynomial:

$$(\theta_0 + \theta_1 z + \dots + \theta_K z^K) = (\mu_0 + \mu_1 \zeta + \dots + \mu_K \zeta^K) \quad (14)$$

where $\zeta = [1-c_1z]/c_2$. Solving the μ -polynomial and then transforming the roots will give the roots of the θ -polynomial. Applying Eq. (14) to Eq. (6) and clearing the denominators gives

$$\{\mu_0 + \mu_1 [(1-c_1z)/c_2] + \dots + \mu_K [(1-c_1z)/c_2]^K\} Y(z) = \beta_{K-1} z^K + \dots + \beta_1 z^2 + \beta_0 z,$$

where $Y(z)$ is the Z-transform of $\{y(n)\}$ and the β_i 's are the coefficients in the numerator polynomial of $Y(z)/z$. Defining $H(z) \triangleq c_2/(1-c_1z)$ and multiplying both sides by $H(z)^K$ gives

$$\mu_0 H^K(z) Y(z) + \mu_1 H^{K-1}(z) Y(z) + \dots + \mu_K Y(z) = H^K(z) \{\beta_{K-1} z^K + \dots + \beta_1 z^2 + \beta_0 z\} \quad (15)$$

The assumption that $H(z)$ is the Z-transform of some filter impulse response $\{h(n)\}$ whose region of convergence is compatible with that of $\{y(n)\}$ permits us to take the inverse transform of Eq. (15) to get

$$\mu_0 \eta_K(n) + \mu_1 \eta_{K-1}(n) + \dots + \mu_K \eta_0(n) = \{h(n)\} *^K \{\beta_{K-1} \delta(n+K) + \dots + \beta_0 \delta(n+1)\} \quad (16)$$

where $\{h(n)\} *^K$: represents convolution (filtering) by $h(n)$ for a total of K times, and $\{\eta_0(n)\} \triangleq \{y(n)\}$; $\{\eta_1(n)\} \triangleq \{h(n)\} * \{\eta_0(n)\}$; etc., i.e. the $\{\eta_i(n)\}$ sequence is the result of filtering the data sequence i times. But what kind of filter is represented by $H(z)$? Actually $H(z)$ has two inverse Z-transforms, one causal and one not, corresponding to the difference equations:

$$\text{Causal: } x_{\text{out}}(n+1) = c_1^{-1} x_{\text{out}}(n) - (c_2/c_1) x_{\text{in}}(n) \quad (17)$$

$$\text{Non Causal: } x_{\text{out}}(n) = c_1 x_{\text{out}}(n+1) + c_2 x_{\text{in}}(n) \quad (18)$$

We shall choose the non-causal filter, for whose Z-transform the region of convergence is $|z| < |c_1|^{-1}$, so that to be compatible with the region of convergence for $\{y(n)\}$ one must have $|c_1|^{-1} > \max_k \{|z_k|\}$ where the z_i 's are the system poles. The filter has the property that inputs arriving at $n \leq 0$ will contribute to the output only for $n \leq 0$. Since the δ -sequences on the right hand side of Eq. (16) are already zero for $n > 0$, K-fold convolution by $\{h(n)\}$ will preserve this property, so for $n \geq 0$ Eq. (16) gives $\mu_0 \eta_0(n) + \mu_1 \eta_1(n) + \dots + \mu_K \eta_K(n) = 0$ for $n \geq 0$, i.e. the sequences $\eta_i(n)$ are linearly dependent through the vector $\mu \triangleq [\mu_0, \mu_1, \dots, \mu_K]$. Furthermore recall that $\{\eta_0(n)\} \equiv \{y(n)\}$, and, using the difference equation for $H(z)$, $\eta_{i+1}(n) = c_1 \cdot \eta_{i+1}(n+1) + c_2 \eta_i(n)$ for $i, n \geq 0$.

Thus given a $\{y(n)\}$ sequence generated by K poles, we can prefilter it K times using the first order non-causal filter described in Eq. (18) to get the $\eta_i(n)$ sequences, and then construct an $M \times K'$ data matrix A, where $A_{ij} = \eta_{j-1}(i-1)$ and then solve $A\mu = 0$ or $(A^T A)\mu = 0$ for the $K' \times 1$ vector μ . Then the corresponding μ -polynomial can be solved for its roots, which are then subjected to a simple linear transformation to get the system poles.

Jain's method [7] is a special case of this procedure wherein $c_1 = c_2 = 1$ so that the filtering is a simple reverse-time integration of a discrete sort, and his A matrix is essentially infinite in height (i.e. $M \gg 1$). This gives a solution

$$\mu = \text{adj}(A^T A) \mu_0$$

where μ_0 is an arbitrary $K' \times 1$ vector. Note that with $P \triangleq (A^T A)$ we have

$$P_{ij} = \sum_{\ell=1}^M (A)_{\ell i} (A)_{\ell j} = \sum_{n=0}^{M-1} \eta_{i-1}(n) \eta_{j-1}(n), \text{ where the summation is extended}$$

to ∞ when $M \rightarrow \infty$ to be in conformance with Jain's results, wherein our P matrix corresponds to Jain's "Grammian" matrix. Moreover Jain uses the diagonal elements of $\text{adj}[P]$ so that the μ -vector is $[(\pm)_1 |\Delta_{11}|^{1/2}, \dots, (\pm)_K |\Delta_{K'K'}|^{1/2}] = \mu^T$ where Δ_{ii} are the diagonal cofactors of R, and the $(\pm)_j$ notation has been discussed earlier. The potentially annoying $(\pm)_j$ signs might better be avoided by simply using one of the other ways of constructing the solution of $P\mu = 0$. (See our section entitled "SOLVING THE HOMOGENEOUS MATRIX EQUATION.")

THE POLES AS EIGENVALUES

§ 27 THEOREM: Given an $M \times K'$ data matrix A, whose rows are post-excitation output sequences from a system having K poles, the poles z_i are the eigenvalues λ in the following generalized eigenvalue problem:

$$(A_+^T A_-)x = \lambda (A_+^T A_+)x \quad (19)$$

where A_+ is A with its rightmost column removed, and A_- is A with its leftmost column removed. \square

The usefulness of this result is difficult to determine. At least it provides a restatement of the problem in which the poles are eigenvalues, and in which the data appear in fairly simple form with no matrix inverses, offering the hope that a method can be employed to directly determine the eigenvalues without having to compute polynomial coefficients, hence avoiding what is generally considered to be a poorly posed problem. Perhaps as better numerical algorithms become available for the $Ax = \lambda Bx$ eigenvalue problem, one of them can successfully be applied.

It is interesting to note that the eigenvalue problem can also be written as

$$A_+^T (A_- - \lambda A_+)x = 0. \quad (20)$$

Equation (20) has a solution only if the matrix $A_+^T(A_- - \lambda A_+)$ is singular, which must occur whenever λ is a system pole. The similarity of this to Jain's "pencil of functions" concept is noteworthy, particularly in view of the following observation: The i^{th} row of A is a data sequence $\{d_i(n) : n = 0, \dots, K\}$ so that the i^{th} row of $(A_- - \lambda A_+)$ is $\{d_i(n+1) - \lambda d_i(n) : i = 0, \dots, (K-1)\}$.

OTHER ESTIMATION CRITERIA, ITERATIVE METHODS, AND ASYMPTOTIC ERROR

For a single post-excitation output sequence $\{y(n)\}$, Eq. (6) shows that the Z-transform, $Y(z)$, can be expressed as a function of the parameter vectors $\theta^T = [\theta_0, \theta_1, \dots, \theta_K]$ and $\beta^T = [\beta_0, \beta_1, \dots, \beta_{K-1}]$ that consist of the coefficients of the pole polynomial and numerator polynomial. Equation (6) can be expressed as

$$Y(z) = \frac{\beta_{K-1} + \beta_{K-2}z^{-1} + \dots + \beta_0 z^{-(K-1)}}{\theta_K + \theta_{K-1}z^{-1} + \dots + \theta_0 z^{-K}} \quad (21)$$

We can define $\beta\uparrow$ and $\theta\uparrow$ as the finite impulse response (FIR) operators (Filters) whose transfer functions are the numerator and denominator polynomials of Equation (21), and $\theta\downarrow$ as the recursive infinite impulse response (IIR) operator that is the reciprocal of $\theta\uparrow$. To illustrate, $\{a(n)\} = \beta\uparrow\{b(n)\}$ means $a(n) = \beta_{K-1}b(n) + \beta_{K-2}b(n-1) + \dots + \beta_0 b(n-K+1)$ and $\{a(n)\} = \theta\downarrow\{b(n)\}$ means $a(n) = [b(n) - \theta_{K-1}a(n-1) - \dots - \theta_0 a(n-K)]/\theta_K$ for $-\infty < n < +\infty$. With this notation $\theta\uparrow\theta\downarrow$ is the unit operator and the sequence $\{y(n)\}$ may be expressed by inverting the Z-transforms in Equation (21) as

$$\{y(n)\} = \theta\downarrow\beta\uparrow \{\delta(n)\} \quad (22)$$

where $\{\delta(n)\}$ is the impulse sequence. If the data sequence is corrupted by noise to give $\{\tilde{y}(n)\} = \{y(n) + e(n)\}$ where the $e(n)$'s are zero mean,

independent Gaussian random variables with common variance σ^2 , then the log-likelihood function for $\{\tilde{y}(n)\}$ given β and θ is

$$L = -\frac{1}{2\sigma^2} \sum_{n \geq 0} (\tilde{y}(n) - y(n))^2 + \text{const.}$$

where $\{y(n)\}$ is defined in terms of the parameter vectors β and θ by Equation (22). Thus the maximum likelihood estimate of θ, β is obtained by choosing them to minimize $J_1 \triangleq \sum_{n \geq 0} (\tilde{y}(n) - y(n))^2$, i.e. to minimize the mean square error in fitting $\{y(n)\}$ to the data sequence $\{\tilde{y}(n)\}$. If the noise components $\{e(n)\}$ are vanishingly small then the minimum value is $J_1=0$. In that case $\{\tilde{y}(n) - y(n)\} \equiv 0$ and this null sequence could be operated upon by $\theta \uparrow$ to give $\theta \uparrow \{\tilde{y}(n)\} - \theta \uparrow \{y(n)\} \equiv \theta \uparrow \{\tilde{y}(n)\} - \beta \uparrow \{\delta(n)\} \equiv 0$. Thus the same result would have been achieved by minimizing $J_2 \triangleq \sum_{n \geq 0} r^2(n)$ where the sequence $\{r(n)\} = \theta \uparrow \{\tilde{y}(n)\} - \beta \uparrow \{\delta(n)\}$. But J_2 can be decomposed as $J_2 = J_{20} + J_{21}$, where $J_{20} \triangleq \sum_{n=0}^{K-1} r^2(n)$ and $J_{21} \triangleq \sum_{n \geq K} r^2(n)$. Furthermore the sequence $\beta \uparrow \{\delta(n)\}$ is identically zero for $n \geq K$, so J_{21} depends only upon θ . Indeed J_{20} can always be minimized to zero by setting $\beta_{K-1} = \theta_K \tilde{y}(0)$; $\beta_{K-2} = \theta_K \tilde{y}(1) + \theta_{K-1} \tilde{y}(0)$; etc. ... ; and ultimately $\beta_0 = \theta_K \tilde{y}(K-1) + \theta_{K-1} \tilde{y}(K-2) + \dots + \theta_1 \tilde{y}(0)$.

Thus the estimate of θ is obtained by choosing it to minimize J_{21} , i.e. to minimize the mean square output of the $\theta \uparrow$ FIR filter for $n \geq K$ when the input is the data sequence $\{\tilde{y}(n)\}$. But if one constructs an $M \times K'$ Hankel data matrix \tilde{A} according to $\tilde{A}_{ij} = \tilde{y}(i+j-2)$ then $J_{21} = \theta^T \tilde{A}^T \tilde{A} \theta$ so that the optimum θ is simply the $\hat{x}[\text{norm}]$ solution of $\tilde{A}x=0$ discussed in earlier sections of this paper.

However minimizing J_1 is not really equivalent to minimizing J_2 unless both can actually be minimized to zero, i.e. the noise is vanishingly small. Otherwise the $\hat{x}[\text{norm}]$ solution does not minimize J_1 , and is therefore not

the maximum likelihood estimate. Direct minimization of J_1 is difficult, being a highly nonlinear problem. Steiglitz [18] has very neatly described the "iterative-prefiltering" procedure for minimizing J_1 by choosing the FIR operators θ^\dagger and β^\dagger to minimize

$$J_3 = \sum_{n>0} (\theta^\dagger \hat{\theta}^\dagger \{\tilde{y}(n)\} - \beta^\dagger \hat{\theta}^\dagger \{\delta(n)\})^2$$

where $\hat{\theta}^\dagger$ is the IIR operator defined as the reciprocal operator to the FIR operator $\hat{\theta}^\dagger$ resulting from the previous step in the iteration process. (The procedure can be started by taking $\hat{\theta}$ as the \hat{x} [norm] solution.) Each time the minimization is done the resulting θ is used as $\hat{\theta}$ in the next step. If the procedure converges then $\theta = \hat{\theta}$ and $J_3 = J_1$ so the resulting θ is a maximum likelihood estimate.

A different iteration method can be used to "improve" the estimate of θ beyond that of simply minimizing $J_{21} = \theta^T \tilde{A}^T \tilde{A} \theta = \|\tilde{A}\theta\|^2$, although it is based on heuristic arguments. Even when θ is truly correct the residual vector $\tilde{A}\theta$ cannot be expected to possess uniform statistical variance in its elements. Due to the Hankel matrix form of \tilde{A} the residual vector can be expressed as $\tilde{A}\theta = G(M;\theta)\tilde{y}$ where G is the matrix defined in Equation (11) and \tilde{y} is $(M+K) \times 1$ the vector whose elements form the sequence $\{\tilde{y}(n): n=0, \dots, (M+K-1)\}$. But \tilde{y} can be expressed as $\tilde{y} = y + e$ where y and e are the vectors of the uncorrupted and noise components respectively, and if θ is the true solution then $Gy = 0$ since every data subsequence of length K is orthogonal to θ . Hence $\tilde{A}\theta = G(M;\theta)e$ and is therefore a zero-mean, Gaussian random vector with co-variance matrix $\Lambda = \mathcal{E}[\tilde{A}\theta(\tilde{A}\theta)^T] = \mathcal{E}[G(M;\theta)ee^T G^T(M;\theta)]$, or $\Lambda = \sigma^2 GG^T$ since $\mathcal{E}[ee^T] = \sigma^2 I$ and the G matrix is not random. If Λ were known a priori then instead of minimizing $J_{21} = \|\tilde{A}\theta\|^2$ one might prefer to minimize the weighted quadratic form $J_4 \triangleq \|\tilde{A}\theta\|_{\Lambda^{-1}}^2 = \theta^T \tilde{A}^T \Lambda^{-1} \tilde{A} \theta$. Unfortunately the

matrix $\Lambda = \sigma^2 GG^T$ cannot be computed without knowing θ , but an iterative procedure can be employed in which each new estimate of θ is used to compute G and estimate Λ for the next step. This procedure is a modification of that used in Refs. [19] and [31].

In the control theory literature the J_2 criteria is the one most often used for minimization [15,16]. However in that context the problem is usually complicated by the presence of a persistently exciting input to the system. Furthermore there is much emphasis placed on the problem of "bias" in the estimate of θ , but it is not statistical bias of the type dealt with in our Theorem §1 and the discussion following. Rather it pertains to asymptotic bias in the estimate of θ as the observation interval of the system output $\{y(n)\}$ becomes infinite. Indeed in that context the asymptotic bias is connected to statistical correlation in the residuals [15], a problem that can be alleviated by the use of pre-whitening filters, "generalized-least-squares", or often simply by supposing that the system is of higher order [16]. By these approaches one can theoretically produce an estimate of θ that converges to the actual θ (not just its maximum likelihood estimate) as the interval of observation becomes infinite. In view of the last approach one might conclude that the use of extra-wide data matrices in our problem, since it is equivalent to supposing a higher system order, is therefore justified by the experience of researchers in the control theory area. However in our problem there is no persistently exciting input, and the "signal" portion of a noise-corrupted data sequence $\{\tilde{y}(n)\} = \{y(n)\} + \{e(n)\}$ will eventually decay to insignificance, leaving only the noise, and we cannot expect the estimate of θ to converge to the true value. Indeed it will almost certainly begin to diverge as soon as the signal component decays to the point of becoming lost in the noise.

But there is a line of heuristic reasoning that lends some relevance to the asymptotic behavior described above. Suppose the signal component of $\{\tilde{y}(n)\}$ does strongly persist for a long enough time that the statistical law of large numbers can be applied to computations involving the noisy data; i.e. the observation interval is infinite with respect to the noise but finite with respect to the signal. Largeness of the interval of observation corresponds to largeness of M in the $M \times N$ Hankel data matrix \tilde{A} . The nature of the "quasi-asymptotic error" in estimating θ using \tilde{A} can be determined by studying the matrix $\tilde{A}^T \tilde{A} = (A+E)^T (A+E) = A^T A + 2(A^T E)_s + E^T E$, where E is the Hankel matrix of error components, and $(\)_s$ denotes the symmetric part of a matrix. Every time M is increased, new rows are added to the A and E matrices. This means that the elements of the $N \times N$ matrix $A^T A$ will continue to reflect these additional summed components for as long as the signal persists. However the matrix E is a Hankel matrix derived from a sequence of independent, zero mean Gaussian variables of variance σ^2 , and it immediately follows that $E^T E$ tends asymptotically toward $M \cdot I \cdot \sigma^2$ where I is the identity matrix. Moreover $(A^T E)_s$ is a linear combination of the independent, zero-mean noise variables, and can be expected to converge to its mean (zero) by the law of large numbers; i.e., it becomes insignificant compared with $A^T A$ and $E^T E$. (This is a broad conclusion which would require an unwieldy set of assumptions and pre-conditions to attain mathematical rigor). The conclusion can be stated succinctly: For $M \gg 1$, $\tilde{A}^T \tilde{A} \approx A^T A + M \cdot \sigma^2 \cdot I$ where I is an $N \times N$ identity matrix. Moreover if the rightmost column of \tilde{A} is isolated by partitioning it as $\tilde{A} = [\tilde{A}_+; \tilde{a}]$, then by a similar line of reasoning $\tilde{A}_+^T \tilde{A}_+ \approx A_+^T A_+ + M \cdot \sigma^2 \cdot I$ where I is now an $(N-1) \times (N-1)$ identity matrix. With these approximations it is possible to estimate the quasi-asymptotic error in determining θ given an $M \times N$ data matrix A with $M \gg 1$.

First consider the least-squares Prony solution. In this case A is $M \times K'$ and the solution $x^T = [x_+^T; 1]$ estimates θ^T , where from Eq. (13), $x_+ = -(A_+^T A_+)^{-1} A_+^T a$ which means that x_+ is the solution of $\tilde{A}_+^T \tilde{A}_+ x_+ = -\tilde{A}_+^T a$. If one uses the approximations just derived, the result is $\tilde{A}_+^T \tilde{A}_+ x_+ + M\sigma^2 x_+ \approx -\tilde{A}_+^T a$ or $\tilde{A}_+^T \tilde{A}_+ x_+ \approx -\tilde{A}_+^T a + \gamma$ where $\gamma = -M\sigma^2 x_+ - A_+^T e - E_+^T a - E_+^T e$. The result is clearly a perturbation in the solution, a fact that has been explored in more detail by Kay [17]. The behavior of the adjoint solution, $x = \text{adj}(\tilde{A}^T \tilde{A}) x_0$, in light of the approximation $\tilde{A}^T \tilde{A} \approx A^T A + M\sigma^2 I$ is somewhat similar, but it is more instructive to view it as one step in the matrix iteration process from x_0 toward $\hat{x}[\text{norm}]$, whose asymptotic behavior is discussed below.

Fortunately the solution $\hat{x}[\text{norm}]$, wherein \hat{x} is chosen to minimize $\|\tilde{A}\hat{x}\|$ subject to $\|\hat{x}\| = 1$, has no quasi-asymptotic error since $\|\tilde{A}\hat{x}\|^2 = \hat{x}^T \tilde{A}^T \tilde{A} \hat{x} \approx \hat{x}^T (A^T A + M\sigma^2 I) \hat{x} = \|\tilde{A}\hat{x}\|^2 + M\sigma^2 \|\hat{x}\|^2$ which obviously leads to the same solution as if there had been no noise. Of course we do not mean to suggest that $\hat{x}[\text{norm}]$ is absolutely errorless, since our argument is fundamentally limited by the accuracy of the approximation $\tilde{A}^T \tilde{A} \approx A^T A + M\sigma^2 I$. Nevertheless it is clear that $\hat{x}[\text{norm}]$ is free of the "asymptotic bias" that is given so much attention in the control theory literature, and this is perhaps the strongest argument in its favor. Before continuing to the next section we remind the reader that the above discussion pertains strictly to estimating θ from a single data sequence $\{y(n)\}$.

PREVIOUS WORK AND UNANSWERED QUESTIONS

Since the problem treated in this paper is in most respects an extension of that studied by Prony in the eighteenth century, not surprisingly there exists more literature than can be referenced here. Nevertheless, we shall

attempt to reference some of the more relevant results, particularly those with which some of our readers may not be aware.

An interesting treatment of the least-squares Prony method, in which the polynomial of the system poles is expressed directly as a determinant, is given by Ellington et al. [20]. Demonstrations of the noise sensitivity problem have been presented by Hildebrand [14] and some analytical work along those lines has been done by Dudley [21]. A very interesting study of the behavior of "noise poles" as a function of signal to noise ratio has been recently published by Kay [17], including a theoretical analysis that may be considered applicable to our problem when the observation interval is infinite with respect to the noise but finite with respect to survival of the signal, and an "autocorrelation" approach is appropriate. Accuracy of time domain and spectral domain reconstructions of noisy signals using Prony-type estimates have been studied by Spitznogle and Quazi [22] and Beatty and George [1], and the latter paper provides a rare look at the use of decimated data sequences.

The use of the SVD decomposition has been explored by Holt and Antill [23], but peculiarly enough they have applied it only to the non-homogeneous version of the least-squares Prony solution, in which the problem becomes poorly conditioned if extra wide matrices are used. Indeed, they use the SVD to "recondition" the problem, rather than as a direct solution; thus their adjusted matrix must still be inverted. Earlier approaches using Householder triangularization have been reported by others, in particular Van Blaricum and Mitra [32].

Price [28] has approached the problem from an eigenvector viewpoint, although without reference to SVD. Moreover, he has addressed the problem of forcing known poles by a transformation of the space rather than through augmentation of the data matrix as we have done.

An apparently different approach to the problem of exponential representations of signals has been used by Jain with proven success [7,24], but our derivation and extension has demonstrated its kinship with other methods. Applications of adaptive filters inspired by Jain's method have been done by Auton [25]. Modification of the Prony method to represent a set of several waveforms with a common pole set has rarely been mentioned in the literature, although it was utilized by Young and Huggins [26].

There remain many unanswered questions of which we mention only a few: Does the advantage of using extra-wide data matrices persist if one uses the methods for solving $Ax=0$ that are emphasized in our paper (i.e., homogeneous methods, $\hat{x}[\text{norm}]$, adjoint solution)? Or when one uses a non-Hankel matrix? How cost-effective is the use of the unbiased version of the adjoint solution, when applicable? Following our derivation of Jain's method, what happens when one uses different transformations of the z variable? Would that approach lead to other, perhaps superior, prefilters?


```

0 1 2 3 4 5
1 2 3 4 5 6
2 3 4 5 6 7
3 4 5 6 7 8
4 5 6 7 8 9
5 6 7 8 9 10
6 7 8 9 10 11
7 8 9 10 11 12

```

Figure 1. Numbering pattern for an 8x6 Hankel matrix.

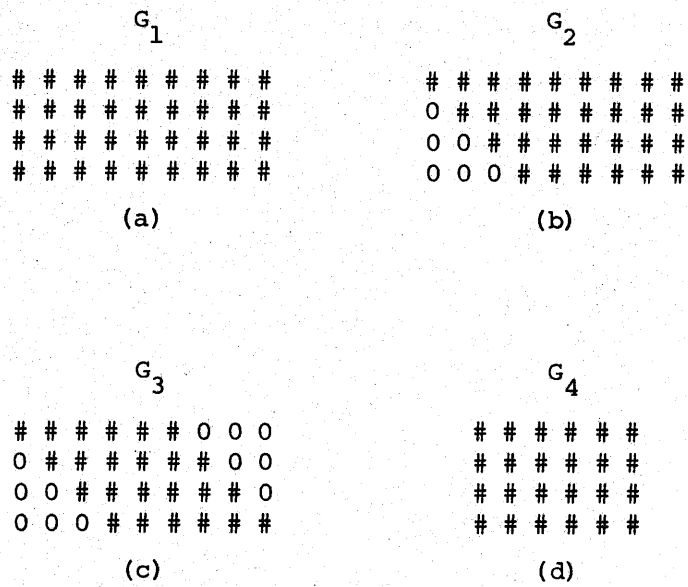


Figure 2. Successive states in the algorithm for suppressing extraneous poles.

APPENDIX

(THEOREM PROOFS AND PROOF OUTLINES)

THEOREM §1: If $P \stackrel{\Delta}{=} A_1^T A_2$ then $P_{ij} = \sum_k (A_1)_{ki} (A_2)_{kj}$ so that P_{ij} and P_{mn} are statistically independent unless $i = m$ or $j = n$. Each element of $\text{adj}(P)$ is a determinant of a "minor" submatrix of P with a row and column deleted. But the determinant of any matrix may be defined [27] as the sum of all possible (appropriately signed) products of elements of the matrix, wherein each product has exactly one representative of each column and each row of the matrix; i.e., no such product contains more than one element from either a single row or a single column. Since the latter property applies even to each minor submatrix of P , it follows that the expectation operator distributes on each element of $\text{adj}[P]$ first over the sum (due to linearity) and then onto the members within each product (since P_{ij} and P_{mn} have the statistical independence property described above), and ultimately distributes to the individual elements of A_1 and A_2 , thus proving the theorem.

THEOREM §2: In the first version of the theorem the unknown parameter A determines the joint probability density of the data matrix \tilde{A} , and due to the Gaussian assumption the log-likelihood function is $L = -(2\sigma^2)^{-1} \|\tilde{A} - A\|_F^2 - \frac{1}{2} MN \log(2\pi\sigma^2)$. The maximum likelihood choice for A , given \tilde{A} , is that which maximizes L within the known set of possible A 's (in this case the set of $M \times N$ matrices of rank K). Clearly $\hat{A}(K; \text{SVD})$ is that choice, since it minimizes the Frobenius norm. If K is also an unknown parameter, known only to be less than N (where $N \leq M$), then L is clearly maximized by using $\hat{A}(K; \text{SVD})$ with $K = N-1$, since \tilde{A} can always be approximated as well by matrices of

larger and larger rank. If some probability distribution for K is known a priori, say P_K , then the traditional modification to the maximum likelihood procedure is to choose A to maximize the weighted likelihood, which gives the log-likelihood function $L = -(2\sigma^2)^{-1} \|\tilde{A}-A\|_F^2 + \log P_K - \frac{1}{2} MN \log (2\pi\sigma^2)$. Since $\hat{A}(K;SVD)$ then maximizes L for each particular K , the optimum K is that which minimizes $\|\tilde{A}-\hat{A}(K;SVD)\|_F^2 - 2\sigma^2 \log P_K$. If one seeks to estimate $\phi(A)$ rather than A , the situation is more complicated. The nullspace, $\phi(A)$, is a peculiar sort of parameter, and is not sufficient to determine the probability distribution of \tilde{A} . It does however determine a family of possible distributions, and thus a family of log likelihood functions $L = -(2\sigma^2)^{-1} \|\tilde{A}-A\|_F^2 + \text{const.}$, parametrized by the matrix A (compatible with the specified nullspace). Clearly setting $A = \hat{A}(K;SVD)$ achieves a global maximum of L for matrices of rank K , and taking $\hat{\phi}(A) = \phi(\hat{A})$ then gives a nullspace estimate whose family of likelihood functions includes one that achieves the maximum. In another sense, estimating $\phi(A)$ is a partial estimation problem in which additional parameters (i.e., the elements of A) must be estimated incidentally.

THEOREM §3: To prove this theorem we first develop several lemmas.

Lemma 1: For any $M \times N$ matrix B and its K^{th} -rank approximant $C \stackrel{\Delta}{=} \hat{B}(K;SVD)$ where $K \leq \text{rank}(B)$, then for any $N \times 1$ vector x ,

$$(a) \quad \|Bx_C\| \geq sv_K(B) \|x_C\|, \text{ and}$$

$$(b) \quad \|Bx_C\| \leq sv_{K+1}(B) \|x_C\|.$$

Proof: Let the SVD of B be given as $B = USV^T$ and let S and V be partitioned as $S = [S_1; S_0]$, $V = [V_1; V_0]$ where in each case the leftmost K columns are isolated. Then $B = U(S_1 V_1^T + S_0 V_0^T)$, and the diagonal elements of S_1 are the

first K singular values of B , thus non-zero. The approximant C may be obtained by artificially setting S_0 to zero, which means that $\phi(C)$ is simply the space orthogonal to the columns of V_1 , i.e., spanned by the columns of V_0 . Thus $V_1^T x_C = 0 = V_0^T x_C$. Then since U and V are orthogonal transformation matrices,

$$\|Bx_C\| = \|U(S_1 V_1^T + S_0 V_0^T)x_C\| = \|S_1 V_1^T x_C\| \geq sv_K(S_1) \|V_1^T x_C\|$$

or since V_1 consists of orthonormal columns,

$$\|Bx_C\| \geq sv_K(B) \|x_C\|.$$

Similarly

$$\|Bx_C^-\| = \|U(S_1 V_1^T + S_0 V_0^T)x_C^-\| = \|S_0 V_0^T x_C^-\| \leq sv_1(S_0) \|V_0^T x_C^-\|$$

or simply

$$\|Bx_C^-\| \leq sv_{K+1}(B) \|x_C^-\|.$$

Lemma 2: For any $M \times N$ matrix A of rank K , and any $N \times 1$ vector x ,

$$\|x_A\| \leq \|Ax_A\| / sv_K(A)$$

Proof: Use Lemma 1(a) with $B \stackrel{\Delta}{=} A$ so that $C = A$ also; solve for $\|x_A\|$.

Lemma 3: With same hypothesis as Lemma 2, and another $N \times M$ matrix E , where $\text{rank}(A) = \text{rank}(A+E_A)$, then $\|x_A\| \leq \|(A+E_A)\| / sv_K(A+E_A)$.

Proof: Apply Lemma 2 with A replaced by $A+E_A$, noting that $x_A = x_{(A+E_A)}$.

Lemma 4: If \hat{x} is a unit vector lying in the nullspace of $\hat{A}(K; \text{SVD})$, where $\tilde{A} = A+E$ and A is an $M \times N$ matrix of rank K , then

$$\|\tilde{A}\hat{x}\| \leq sv_{K+1}(\tilde{A}) \leq \|E\|$$

Proof: Apply Lemma 1(b) with $B \stackrel{\Delta}{=} \tilde{A}$ and $x \stackrel{\Delta}{=} \hat{x}$. Then by assumption $\hat{x}_C^- = \hat{x}$ and $\|\hat{x}_C^-\| = 1$. To get the rightmost inequality we use a fundamental

property of singular values [9]:

$$sv_{K+1}(\tilde{A}) = sv_{K+1}(A+E) \leq sv_{K+1}(A) + sv_1(E) = 0 + \|E\|.$$

Lemma 5: Same as Lemma 4, ending with $\|E_A^-\|$ instead of $\|E\|$.

Proof: As before, but use

$$sv_{K+1}(\tilde{A}) = sv_{K+1}(A+E_A^+ + E_A^-) \leq sv_{K+1}(A+E_A^+) + sv_1(E_A^-) = 0 + \|E_A^-\|$$

Lemma 6: With A, E as before and x any unit vector, then

$$(a) \quad \|Ax_A\| \leq \|\tilde{A}x\| + \|E\|, \text{ and}$$

$$(b) \quad \|(A+E_A^-)x_A\| \leq \|\tilde{A}x\| + \|E_A^-\|$$

Proof: $Ax_A = Ax = (\tilde{A}-E)x = \tilde{A}x - Ex$, so

$$\|Ax_A\| \leq \|\tilde{A}x\| + \|Ex\| \leq \|\tilde{A}x\| + \|E\| \text{ since } \|x\| = 1.$$

Part (b) results similarly upon noting that

$$(A+E_A^-)x_A = (A+E_A^-)x = (\tilde{A}-E_A^-)x.$$

Lemma 7: If α and β are positive and $\alpha \leq \beta$, then $\alpha + \beta \leq 2\beta \leq 2(\alpha + \beta)$.

Proof: Obvious.

The theorem follows immediately, defining $b_1 \triangleq \{\|\tilde{A}x\| + \|E\|\}/sv_K(A)$ and $b_2 \triangleq 2\|E\|/sv_K(A)$ and using Lemmas 2, 6(a), 4, and 7 in sequence; and defining $b_3 \triangleq \{\|\tilde{A}x\| + \|E_A^-\|\}/sv_K(A+E_A^-)$ and $b_4 \triangleq 2\|E_A^-\|/sv_K(A+E_A^-)$ before applying Lemmas 3, 6(b), 5 and 7.

THEOREM §5: Since both ends of the θ vector are non-zero by assumption, no row of G can be expressed as a linear combination of the rows above, and all the rows are therefore linearly independent by induction. To prove uniqueness of the generator, suppose that there are two matrices G_1 and G_2 as in Eq. (11) that span the same Ω_ℓ ; i.e., they have the same rowspace. Then by

construction the top row of G_2 is non-zero only for its first K' elements, and it can be expressed as a linear combination of the rows of G_1 . But that combination cannot include the last row of G_1 since it would produce a non-zero rightmost element. Similarly each higher row of G_1 can be ruled out so that the top row of G_2 is a scalar multiple of that of G_1 ; i.e., the generator is unique to within a scalar multiple.

THEOREM §6: Although the dimension of the power vector Z has always been inferred from the context of its use, for the proof of this theorem we shall promote clarity by appending the dimension as a subscript in parentheses, e.g. $Z_{(N)}$. Thus if z is a root of the generator polynomial then $\theta^T Z_{(K)} = 0$, which clearly implies $GZ_{(N)} = 0$, where G is the matrix $G(\ell; \theta)$ of Eq. (11). But $x \in \Omega_\ell(\theta)$ implies x lies in the rowspace of G , i.e., $x = G^T \alpha$ for some $\ell \times 1$ vector $\alpha^T \triangleq [\alpha_0, \alpha_1, \dots]^T$. Then it is easily shown that $x^T Z_{(N)} = \alpha^T GZ_{(N)} = (\alpha^T Z_{(\ell)}) \times (\theta^T Z_{(K')})$. Clearly the roots of the generator polynomial are a subset of those of $x^T Z_{(N)}$, and the $\ell-1$ extraneous roots can be placed at will by selecting the elements of α as the coefficients of a polynomial whose roots are the desired extraneous set. Moreover since the polynomial coefficients of all three polynomials are real, all three root sets contain only conjugate pairs.

THEOREM §7: The N -dimensional power vectors of the (distinct) roots z_i of the generator polynomial constitute a set of K linearly independent vectors since they can be juxtaposed to form columns of a matrix for which the upper $K \times K$ submatrix is Vandermonde, and therefore nonsingular. Since each such power vector obviously lies within the nullspace of the G matrix of Eq. (11)

so that the set spans that nullspace, the K -dimensional subspace of R^N that is the orthogonal complement is identically $\Omega_\ell(\theta)$, the rowspace of G .

THEOREM §9: Obviously $y \perp \Omega_\ell(\theta)$ if and only if $Gy = 0$, where G is the matrix of Eq. (11), and the theorem follows immediately.

THEOREM §11: Clearly if A is θ -AR then $AG^T = 0$ where G is the matrix of Eq. (11). But $AG^T = 0$ implies [29] $\text{rank}(A) \leq N - \text{rank}(G^T) = N - \ell = K$. However $x \in \Omega_\ell(\theta)$ iff $x = G^T \alpha$ for some α , but then $Ax = AG^T \alpha = 0$, so $x \in \phi(A)$. If $\text{rank}(A) = K$ then $\phi(A)$ is a subspace of dimension $N - K = \ell$, the same as the dimension of $\Omega_\ell(\theta)$, so in that case $\phi(A)$ and $\Omega_\ell(\theta)$ are identical. Thus $Ax=0$ implies $x \in \Omega_\ell(\theta)$ and therefore the roots of $x^T Z$ constitute a superset of the roots of the generator polynomial by Theorem §6.

THEOREM §14: Since both ends of θ are non-zero, at least one member of any basis of $\Omega_\ell(\theta)$ must have its first element non-zero, and the scaling described in the theorem statement clearly does no harm. Furthermore $\{g_1; a_2, a_3, \dots\}$ is still a basis since g_1 could be added to the a 's to recover the original basis. Now if \bar{x} is an arbitrary vector in $\Omega_{\ell-1}(\theta)$ and we augment it by attaching an initial zero element to form a vector x , then it can be expressed as a linear combination of the bottom $\ell-1$ rows of $G(\ell; \theta)$, hence as a linear combination of the basis vectors $\{g_1; a_2, a_3, \dots\}$. Indeed g_1 can obviously be ruled out of the combination, and it follows easily that x lies in the space spanned by $\{\bar{a}_2, \bar{a}_3, \dots, \bar{a}_\ell\}$ as defined in the theorem. Since that set is of dimension $\ell-1$, it must constitute a basis of $\Omega_{\ell-1}(\theta)$. The corollary follows directly.

THEOREM §16: Perturbations in the (assumed distinct) roots of the polynomial $x^T Z$ due to perturbations in the x vector may be determined to first order by setting the total differential of the polynomial to zero at each root location: $(x+\delta x)^T Z + x^T (Z+\delta Z) = 0$ at $Z = Z_i$. But since $x^T Z_i = 0$ we have

$$(\delta x)^T Z_i = -x^T \delta Z_i = -x^T DZ_i z_i^{-1} dz_i$$

where the diagonal matrix is as defined in the theorem. Solving for dz_i/z_i and noting that $x \approx x+\delta x$ to first order gives

$$dz_i/z_i = -(\delta x)^T Z_i / (x+\delta x)^T DZ_i .$$

For the purpose of the theorem we set the nominal value, x , as the "correct" portion of the approximate solution \hat{x} , i.e., $x = \hat{x}_A^-$. Then clearly $\delta x = \hat{x}_A^+$, and $x+\delta x = \hat{x}$ so the theorem follows immediately.

THEOREM §19: The theorem follows directly from the fact that $\tilde{A}^T \tilde{A}$ is nonsingular, so that premultiplying the equation by $\tilde{A}^T \tilde{A}$ gives $\det(\tilde{A}^T \tilde{A}) x_0$ on the left-hand side, and $\tilde{A}^T \tilde{A}$ can be partitioned as

$$\tilde{A}^T \tilde{A} = \begin{bmatrix} \tilde{A}_+^T \tilde{A}_+ & | & \tilde{A}_+^T \tilde{a} \\ \hline \tilde{a}^T \tilde{A}_+ & | & \tilde{a}^T \tilde{a} \end{bmatrix}$$

The rest is simple algebra.

THEOREM §27: The least-squares Prony solution of Eq. (13) can be expressed as $y^T = [y_+^T; 1]^T$ where y_+ is a $K \times 1$ vector defined as $y_+ = -(\tilde{A}_+^T \tilde{A}_+)^{-1} \tilde{A}_+^T \tilde{a}$, and the system poles are the roots of the polynomial $y^T Z$. But the roots of any such normalized polynomial are identically the eigenvalues of the companion matrix [30]:

$$C(y_+) \triangleq \begin{bmatrix} 0, 0, \dots, 0 & | \\ \hline & | y_+ \\ I & | \end{bmatrix},$$

i.e., as the eigenvalues λ in $Cx = \lambda x$, where I is a $(K-1) \times (K-1)$ identity matrix. Premultiplying this equation by $A_+^T A_+$ and carefully carrying out the partitioned multiplication leads directly to the theorem.

REFERENCES

1. L. G. Beatty and J. D. George, "Use of the Complex Exponential Expansion as a Signal Representation for Underwater Acoustic Calibration," J. Acoust. Soc. Am. 63, 1782-1794 (1978).
2. W. M. Leach, R. W. Schafer, and T. P. Barnwell, "Time Domain Measurement of Loudspeaker Driver Parameters," IEEE Trans. Acoust., Speech, Signal Processing ASSP-27, 734-739 (1979).
3. C. E. Baum, "Emerging Technology for Transient and Broad-Band Analysis and Synthesis of Antennas and Scatterers," AFWL EMP Interaction Note 300, Air Force Weapons Laboratory, Kirtland AFB, New Mexico 87117 (1976).
4. Anon., "Electromagnetic Pulse Handbook for Missiles and Aircraft in Flight," AFWL EMP Interaction 1-1, Air Force Weapons Laboratory, Kirtland AFB, New Mexico 87117 (1972).
5. D. Brown, G. Carbon, and K. Ramsey, "Survey of Excitation Techniques Applicable to the Testing of Automotive Structures," available from the Soc. of Automotive Engineers as SAE paper 77029 (1977).
6. W. T. Chu, "Comparison of Reverberation Measurements Using Schroeder's Impulse Method and Decay-Curve Averaging Method," J. Acoust. Soc. Am. 63, 1444-1450 (1978).
7. V. K. Jain, "Representation of Sequences" IEEE Trans. Audio Electroacous. AU-19, 208-215 (1971).
8. G. W. Stewart, Introduction to Matrix Computations, Academic Press, 1973.
9. C. L. Lawson and R. J. Hanson, Solving Least Squares Problems, Prentice-Hall, 1974.
10. A. J. Poggio, M. L. Van Blaricum, E. K. Miller, and R. Mittra, "Evaluation of a Processing Technique for Transient Data," IEEE Trans. Antennas Propagat. AP-26, 165-173 (1978).
11. D. L. Lager, H. G. Hudson, and A. J. Poggio, "User's Manual for SEMPEX: A Computer Code for Extracting Complex Exponentials from a time Waveform," AFWL Mathematics Note 45, Air Force Weapons Laboratory, Kirtland AFB, New Mexico 87117 (1977).
12. T. Cordaro, "A Note on Representing a Transient Waveform by a Finite Sum of Complex Exponentials," AFWL Mathematics Note 46, Air Force Weapons Laboratory, Kirtland AFB, New Mexico 87117 (1977).
13. R. Prony, "Essai Experimental et Analytique," Paris Journal l'Ecole Poltechnique 1, 24-76 (1795).

14. F. B. Hildebrand, Introduction to Numerical Methods, McGraw-Hill (1956).
15. P. Eykhoff, System Identification, Wiley Interscience (1974).
16. T. C. Hsia, System Identification: Least-Squares Methods, D. C. Heath and Company, Lexington, Massachusetts (1977).
17. S. M. Kay, "The Effects of Noise on the Autoregressive Spectral Estimator," *IEEE Trans. Acoust., Speech, Signal Processing ASSP-27*, 478-485 (1979).
18. K. Steiglitz, "On the Simultaneous Estimation of Poles and Zeros in Speech Analysis," *IEEE Trans. Acoust., Speech, Signal Processing ASSP-25*, 229-234 (1977).
19. A. G. Evans and R. Fischl, "Optimal Least Squares Time-Domain Synthesis of Recursive Filters," *IEEE Trans. Audio Electroacoust. AU-21*, 61-65 (1973).
20. J. P. Ellington, H. McCallion, "The Determination of Control System Characteristics from a Transient Response," *Proc. Inst. Elec. Eng. 105* (British), 370-373 (1958).
21. D. G. Dudley, "Parametric Modeling of Transient Electromagnetic Systems," *Radio Science 14*, 387-396 (1979).
22. F. R. Spitznogle and A. H. Quazi, "Representation and Analysis of Time-Limited Signals Using a Complex Exponential Algorithm," *J. Acoust. Soc. Am. 47*, 1150-1155 (1970).
23. J. N. Holt and R. J. Antill, "Determining the Number of Terms in a Prony Algorithm Exponential Fit," *Mathematical Biosciences 36*, 319-332 (1979).
24. V. K. Jain, "Filter Analysis by Use of Pencil Functions: Part 1," *IEEE Trans. Circuits Systems CAS-21*, 574-579 (1974).
25. J. Auton, (Master's thesis in preparation), Department of Electrical Engineering, University of Kentucky, Lexington, Kentucky 40506 (1980).
26. T. Y. Young and W. H. Huggins, "On the Representation of Electrocardiograms," *IEEE Trans. Bio-Med. Electronics BME-10*, 86-95 (1963).
27. W. H. Greub, Linear Algebra, Springer-Verlag (1967).
28. H. J. Price, "An Improved Prony Algorithm for Exponential Analysis," AFWL Mathematics Note 59, Air Force Weapons Laboratory, Kirtland AFB (1978).
29. D. C. Murdoch, Linear Algebra for Undergraduates, John Wiley & Sons (1957).

30. C. R. Rao and S. K. Mitra, Generalized Inverse of Matrices and Its Applications, John Wiley & Sons (1971).
31. J. T. Cordaro, "Pole Measurements for the ATHAMAS Pipe Test," ACT Mathematics Note 56, Air Force Weapons Laboratory, Kirtland AFB, August 1977.
32. M. Van Blaricum and R. Mitra, "A Technique for Extracting the Poles and Residues of a System Directly from Its Transient Response," IEEE Trans. Antennas Propagat. AP-23, 777-781 (1975), also Interaction Notes, Note 245, Air Force Weapons Laboratory, Kirtland AFB, February 1975.