THE UNIVERSITY OF MICHIGAN
INDUSTRY PROGRAM OF THE COLLEGE OF ENGINEERING

LINEAR SYSTEM APPROXIMATION BY MEAN SQUARE ERROR MINIMIZATION
IN THE TIME DOMAIN

Elmer G. Gilbert

This dissertation was submitted in partial
fulfillment of the requirements for the degree of
Doctor of Philosophy in the University of Michigan,
1956.

January, 1957

IP-201

Биас
(
има
1834

## ACKNOWLEDGEMENTS

TABLE OF CONTENTS

TABLE OF CONTENTS (CONT'D)

## LIST OF TABLES

LIST OF FIGURES

# LIST OF FIGURES (CONT'D)

# LIST OF SYMBOLS

$a$ - lower limit of $x$ interval

$a_n$ - nth coefficient in approximating series

$a]$ - column of coefficients $a_1, - - - a_N$

$A$ - gain constant of system function $H*(s)$

$A_n$ - coefficients in partial fraction expansion of $H*(s)$

$b$ - upper limit of $x$ interval

$b_n$ - nth coefficient in constrained approximating series

$b]$ - column of coefficients $b_1, - - - b_N$

$B_m$ - coefficients in differential equation for Prony method

$c_{nm}$ - constant defined by weighted integral of $\theta_n(t)\theta_m(t)$

$c_{nm}^{-1}$ - elements of matrix $C^{-1}$

$C$ - $N$ by $N$ matrix with elements $c_{nm}$

$C^{-1}$ - inverse $C$ matrix

$d_n$ - constant defined by weighted integral of $f_o(t)\theta_n(t)$

$d]$ - column of coefficients $d_1, - - - d_N$

$\Delta d_n$ - computer error in $d_n$

$\Delta d]$ - column of coefficients $\Delta d_1, - - - \Delta d_N$

$e(t), E(s)$ - error between $f_o(t)$ and $f_o*(t)$

$e*(t)$ - approximation of $e(t)$

$e*_{max}$ - estimate of peak approximation error

$e_h(t), E_h(s)$ - error between $h(t)$ and $h*(t)$

$e_o(t)$ - Prony method error

$E$ - relative mean square error in $d_1, - - - d_N$

$f(t)$ - solution of Prony method differential equation

$f_i(t)$, $F_i(s)$ - prescribed input

$f_o(t)$, $F_o(s)$ - prescribed response

$f_o'(t)$ - first derivative of $f_o(t)$

$f_o^{(M)}(t)$ - Mth derivative of $f_o(t)$

$f_o*(t)$, $F_o*(s)$ - approximation of prescribed response

$f_{kj}$ - weighted integral of $f_o^{(k)}(t)f_o^{(j)}(t)$

$g(t)$ - x transformation function, function used in computer approximation

$g]$ - column used in derivation of $I_c$

$h(t)$, $H(s)$ - prescribed impulse response, prescribed system function

$h*(t)$, $H*(s)$ - approximation of prescribed impulse response, approximation of prescribed system function

$I$ - identity matrix

$I$ - weighted mean square error

$\Delta I$ - increase in I as a result of $\Delta d]$

$I_c$ - weighted mean square error for constrained approximation

$I_{max}$ - weighted integral of $f_o^2(t)$

$I_M$, $I_P$, $I_R$, $I_I$ - mean square error in the magnitude, angle, real part, and imaginary part of $E(j\omega)$

$I_o$ - weighted mean square error for Prony method

$I_{rel}$ - relative weighted mean square error $I/I_{max}$

$\Delta I_{rel}$ - increase in $I_{rel}$ as a result of $\Delta d]$

$j$ - the imaginary unit

$k$, $k_q$ - value of constraint condition

$k_1$, $k_N$ - smallest eigenvalue of $C^{-1}$, largest eigenvalue of $C^{-1}$

$k_{qn}$ - $K_q\{ \theta_n \}$

$k]$ - column of values $k_1, --- k_Q$

$K$ - $Q$ by $N$ matrix with elements $k_{qn}$

$K_n$ - normalization constant for $\theta_n$

$K_T$ - transposed $K$ matrix

$K\{\}$, $K_q\{\}$ - operator which measures constraint condition

$L_n(x)$ - Laquerre polynomial of order $n$

$M$ - number of poles in $H*(s)$

$n(t)$ - noise signal in optimum filter problem

$n_i(t)$ - undesired component of input

$n_o(t)$ - response component dependent on $n_i(t)$

$N$ - number of approximating functions

$P$ - number of zeros in $H*(s)$

$P_n(x)$ - functions orthogonal in $(a, b)$

$Q$ - number of constraints

$r$ - sensitivity of $\Delta I_{rel}$ to $E$

$r_{nm}$ - elements of the matrix $R$

$R$ - matrix which defines one set of approximating functions in terms of another set of approximating functions

$s$ - complex variable of Laplace transform

$s_n$ - pole of $H*(s)$

$\bar{s}_p$ - zero of $H*(s)$

$\tilde{s}_k$ - exponential constant in weight factor series

$\overset{*}{s}_k$ - exponential constant in predistortion series

$s_{nm}$ - pole in nth approximating function

$\bar{s}_{nm}$ - zero in nth approximating function

$S(t)$ — signal in detection problem

$t$ — time variable

$T$ — constant used in time scaling or translation of time function

$T_1$, $T_2$ — limits on time interval of approximation

$\overline{T}_1$, $\overline{T}_2$ — interval for determination of $e^*_{max}$

$u_0$ — unit impulse

$u_1$, $u_2$, — — — — derivatives of unit impulse

$u_{-1}$, $u_{-2}$, — — — — integrals of unit impulse

$w_k$ — coefficients in exponential weight factor series

$\overset{*}{w}_k$ — coefficients in predistortion series for $W^{-1/2}(t)$

$W(t)$ — weight factor

$W_0(t)$ — weight factor for Prony method

$W_p(t)$ — weight factor for functions $P_n(x)$

$W^{-1}_{max}$ — maximum value of $W^{-1}(t)$

$x$ — argument of orthogonal function $P_n(x)$

$x(t)$, $X(s)$ — response of computer circuits

$y(t)$, $Y(s)$ — input to computer circuits

$\alpha_n$ — negative real part of nth pole

$\beta_n$ — distance of nth pole from origin

$\Theta_n(t)$, $\Theta_n(s)$ — nth approximating function of $f^*_0(t)$

$\lambda_n$ — Lagrange multiplier

$\Lambda_n(s)$ — all pass function with poles of the nth orthonormal function

$\mu$ — exponential constant of weight factor

$\pi$ — 3.1416. . .

$\sigma$ — real part of s

$\sigma_n$ - real part of $s_n$

$\tau$ - variable of integration, argument of correlation functions

$\varphi_n(t)$, $\Phi_n(s)$ - nth approximating function of $\overset{*}{h}(t)$

$\psi_{ii}(\tau)$ - auto-correlation function of $f_i(t)$

$\psi_{oi}(\tau)$ - cross-correlation function of $f_o(t)$ and $f_i(t)$

$\psi_{fg}(\tau)$, $\Psi_{fg}(\omega)$ - cross-correlation function of $f(t)$ and $g(t)$

$\omega$ - imaginary part of $s$

$\omega_n$ - imaginary part of $s_n$

# I. INTRODUCTION

The synthesis of linear systems to have prescribed transient
response has become increasingly important in recent years. Present
applications in automatic control, electronic computation, data trans-
mission, noise filtering, and measurement of linear-system characteris-
tics are concerned with input functions which are not periodic in time.
It is therefore understandable that synthesis procedures which are
stated and carried out in terms of time response, like the ones pre-
sented here, are of considerable interest.

The synthesis problem is rarely solved without error. Lim-
itations such as those imposed by noise and linear-system physical
realizability assure this. Consequently, synthesis really involves
the solution of two problems: (1) the approximation problem, the
determination of an approximate system function which matches closely
a prescribed system function; (2) the realization problem, the con-
struction of a physical linear system which possesses the approximate
system function. This investigation is aimed primarily at solving the
first problem. The approximate system function is expanded in a series
of realizable approximating functions, and the coefficients in this
series are chosen to minimize a time weighted average of the squared
response error. The theoretical development and practical application
of this weighted mean square error approximation method are the pur-
pose of this dissertation.

Before discussing the advances made in mean square approxi-
mation it is relevant to review briefly the present state of linear-
system approximation for prescribed transient response. The greater
part of past work is recent and has been restricted primarily to the

determination of system functions suitable for realization in the form of finite, lumped-element electrical networks with fixed parameters.[1] A number of useful methods[2] have been proposed which produce time-response errors tending toward zero with increased network complexity. No optimum is obtained, however, in the sense that a measure of the error is truly minimized. Furthermore, the procedures are only practicable when the prescribed input is the impulse function and the prescribed response is an analytic expression. The mean square error approach has been investigated by relatively few workers [1,2,17,18,22,31]. By far the most general and complete treatment is the report by Kautz [17], which is devoted primarily to the approximation of impulsive responses with known Laplace transforms. Certain fairly broad classes of exponential time functions, orthonormal in the semi-infinite interval, are developed and employed in an approximating series. Although Kautz and others mention the approximation problem when the input function is arbitrary they do not present wholly satisfactory solutions.

Thus, present methods of approximation for prescribed transient response are limited mainly to the following areas: (1) approximation by systems which are finite, lumped, and fixed in time, (2) approximation of impulsive responses expressable in analytic or Laplace transform form, (3) means square approximation by certain classes of

_____

[1] Winkler [41] gives an excellent review of the approximation problem as applied to electrical networks for both time and frequency response. An extensive bibliography is also included.

[2] The following methods and corresponding references are noteworthy: (1) matching of time moments [15,34,37], (2) numerical calculation by time series [3,24,25], (3) Pade approximates and continued fraction expansions [17,26,29,36], (4) Prony's method [6,28] (5) rational fraction approximation along contour in complex plane [7,14], (6) Fourier series [35].

orthonormal exponential functions.

It is felt that the contributions of this investigation largely overcome these limitations. A summary of the more important results includes: (1) a general theory of constrained and unconstrained mean square approximation by linearly independent but not necessarily orthogonal functions; (2) a practical solution of the arbitrary input problem; (3) procedures for generating wider classes of orthonormal approximating functions, especially orthonormal exponential functions; (4) the application of analog computer techniques to the mean square error approximation problem and the realization problem; (5) methods for experimentally measuring linear-system characteristics, processing experimental data, and experimentally synthesizing optimum filters.

In order to best develop these results the text has been divided into two parts. The first part discusses the basic theory involved in the weighted mean square approximation of time invariant linear systems. It includes chapters on notation and preliminary assumptions, the meaning of the weighted mean square error criterion, constrained and unconstrained approximation by linearly independent approximating functions, the orthogonalization of linearly independent functions with special emphasis on complex exponential functions, and the choice of suitable approximating functions. In the second part, the theory of part one is applied to practical synthesis problems. In this part there are two chapters, the first, on the impulse response approximation problem, and the second, on the arbitrary input problem and the handling of experimental data. An appendix tabulates various families of exponential approximating functions and inverse matrices useful in analytic approximations.

PART I

BASIC THEORY

## II.  NOTATION AND PRELIMINARY ASSUMPTIONS

A resume of notation and preliminary assumptions is required before the detailed discussion of basic theory can begin.  This chapter will consider theory pertaining to linear-system description, a precise statement of the approximation problem and the mean square error approach, physical and mathematical limitations involved, and preliminary simplifications of the prescribed response.

### Linear System Theory[1]

The linear-system notation used is shown in Figure 2.1.  The input is $f_i(t)$, and the response is $f_o(t)$.  Laplace transforms of

**LINEAR  SYSTEM**



Figure 2.1  Linear System Notation

lower case letter time functions are given by corresponding upper case letter functions  of the complex variable $s = \sigma + j\omega$.  Thus

$$F_o(s) = \int_0^\infty e^{-st} f_o(t) \, dt \tag{2.1}$$

In what follows all such transforms are presumed to exist.  To avoid omission of the time functions for $t < 0$ the time origin is chosen so that the time functions are zero for $t < 0$.

---

[1]  A detailed discussion of linear-system theory is given by Gardner and Barnes [10].

The synthesis problem is essentially an input-response problem. It is therefore assumed that the linear system is initially at rest with the total response dependent entirely on the input. In this case,

$$F_o(s) = H(s)F_i(s). \tag{2.2}$$

The system function $H(s)$ may be defined either in terms of a transformed input and response by equation (2.2) or in terms of the transform of the impulse response $h(t)$,

$$H(s) = \int_o^\infty e^{-st} h(t)\ dt\ . \tag{2.3}$$

The impulse response or weighting function is important because it allows the response to an arbitrary input to be expressed by means of the superposition integrals,

$$f_o(t) = \int_{-\infty}^\infty f_i(\tau)h(t-\tau)\ d\tau \tag{2.4}$$

and

$$f_o(t) = \int_{-\infty}^\infty f_i(t-\tau)h(\tau)\ d\tau\ . \tag{2.5}$$

If the linear system is realizable it must not be a predictor, and $h(t) = 0$ for $t < 0$. This causes the upper limit in equation (2.4) to become t and the lower limit in equation (2.5) to become zero. The Laplace transform equation (2.2) and the superposition integrals will be used frequently in future sections.

The Approximation Problem

The approximation problem involves a prescribed system function and an approximation to it. Problem notation is shown in Figure 2.2. The prescribed system function is defined in terms of

**PRESCRIBED SYSTEM**



Figure 2.2   Approximation Problem Notation

a given input $f_i(t)$ and response $f_o(t)$ and is not necessarily physically realizable. In the special case where $f_i(t)$ is the unit impulse, the prescribed response $f_o(t)$ is simply the impulse response $h(t)$. The approximate system function (note that approximate functions are distinguished from their exact counterparts by an asterik) is subject to conditions of realizability and is chosen to make the approximation error $e(t)$(or $E(s)$ ) small.

Two methods of error evaluation may be considered. In a frequency domain approximation the prescribed input and response are stated in Laplace or Fourier transform notation so that a function of $E(s)$(or $E(j\omega)$ ) measures the approximation tolerance. In a time domain approximation the prescribed input and response are stated as time functions so that a function of $e(t)$ measures the approximation tolerance. When accurate duplication of transient response is of primary importance the second procedure is more direct and allows better control of approximation error.

The mean square error approach to the time domain approximation depends on two things: (1) a series expansion of the approximate weighting function, and (2) the minimization of the weighted mean square error in time response. The series expansion includes N predetermined approximating functions and is written as

$$h^*(t) = \sum_{n=1}^{N} a_n \varphi_n(t) \quad .$$
(2.6)

Consequently, the approximate system function is given by the Laplace transform of equation (2.6).

$$H^*(s) = \sum_{n=1}^{N} a_n \Phi_n(s) \quad .$$
(2.7)

Conditions which $\varphi_n(t)$ and $\Phi_n(s)$ must satisfy if $H^*(s)$ is to be realizable are considered in the next section. Factors controlling the choice of the functions themselves are examined in Chapter VI. The weighted mean square error is defined by

$$I = \int_{T_1}^{T_2} W(t) e(t)^2 dt = \int_{T_1}^{T_2} W(t) [f_0(t) - f_0^*(t)]^2 dt \quad .$$
(2.8)

where $T_1 \leq t \leq T_2$ is the interval of approximation and $W(t)$ is a positive, bounded[1] weight function. When I is minimized with respect to the coefficients $a_1, \text{------} a_N$, the approximation problem is considered solved.

Before I can be minimized it must be expressed in terms of the coefficients. This is possible through application of the superposition integral

---

[1] Occasionally, it may be desirable to let $W(t)$ be unbounded. This is permissible if the limitations of the next section are satisfied.

$$f_o^*(t) = \int_0^\infty f_o(t-\tau)h^*(\tau)d\tau = \sum_{n=1}^N a_n \int_0^\infty f_i(t-\tau)\varphi_n(\tau)\,d\tau \quad . \quad (2.9)$$

Notation is simplified by letting

$$\Theta_n(t) = \int_0^\infty f_i(t-\tau)\,\varphi_n(\tau)\,d\tau \quad , \quad (2.10)$$

then

$$f_o^*(t) = \sum_{n=1}^N a_n\Theta_n(t) \quad (2.11)$$

and

$$I = \int_{T_1}^{T_2} W(t)[f_o(t) - \sum_{n=1}^N a_n\Theta_n(t)]^2 dt \quad . \quad (2.12)$$

Detailed procedures for minimizing I are taken up in Chapter IV.

The equations of the previous paragraph can be understood more clearly by reference to Figure 2.3. The approximate linear system



Figure 2.3  Block Diagram Representation of Series Expansion of H*(s)

is divided into N parallel paths with respective system functions $\Phi_n(s)$, gains $a_n$, and outputs $a_n\theta_n(t)$ which are summed to form the approximate response $f_o^*(t)$. Such a representation is schematic; it does not mean that $H^*(s)$ must be realized in the same way.

### Physical and Mathematical Limitations

If the approximate linear system is to be realizable it must be non-predicting and stable.[1] Or equivalently, the following conditions hold: (1) $h^*(t) = 0$ for $t < 0$, (2) $\int_o^\infty |h^*(t)|dt$ is bounded or $H^*(s)$ is bounded and analytic for $\sigma \geq 0$ ($s = \sigma+j\omega$) and has an integrable $|\frac{dH}{ds}|^2$ on the entire $j\omega$ axis.[2] Clearly, the same two conditions must apply also to the functions of the series expansion, $\varphi_n(t)$ and $\Phi_n(s)$.

Conditions (2) means, among other things, that $H^*(s)$ is bounded at the point at infinity; i.e., the approximate system has finite gain at infinite frequency. This excludes certain ideal devices such as perfect differentiators. Actually, most practical systems have the even greater restriction of zero gain at infinite frequecy.

In the majority of useful approximation problems it is further specified that the approximate system be lumped and finite. Unless otherwise noted this assumption will be made in the following work. Mathematically, the lumped and finite condition demands that: (1) $H^*(s)$ be a real, rational function of s with a finite number of poles, (2) $h^*(t)$ be a real, finite sum of complex exponential functions with a possible impulse added at $t = 0$. Stating these conditions in equation form for a Mth order system yields

---

[1] A stable system is defined as one producing bounded responses for bounded inputs.

[2] These conditions are discussed more fully by James, Nichols, and Phillips [16] .

$$H^*(s) = A \frac{(s-\bar{s}_1)(s-\bar{s}_2)----(s-\bar{s}_P)}{(s-s_1)(s-s_2)----(s-s_M)} = A_o + \sum_{m=1}^{M} \frac{A_m}{s-s_m} \qquad (2.13)$$

and

$$h^*(t) = A_o u_o(t) + \sum_{m=1}^{M} A_m e^{s_m t} \qquad 1,2 \qquad (2.14)$$

where $\bar{s}_p$, $s_m$, and $A_m$, occur in conjugate pairs when complex, and

where $\sigma_m < 0$ ($s_m = \sigma_m + j\omega_m$) for m = 1,-----M. If $H^*(s)$ has zero

gain at infinite frequency, P < M and $A_o$ = 0. The series expansion

of the approximate system function given in equation (2.7) must contain

the same poles as $H^*(s)$. Since the approximating functions in the series

are predetermined, this means that the poles of the final approximation

are predetermined. The zeros of $H^*(s)$ depend, of course, on the coeffi-

cients $a_1$,-----$a_N$.

In addition to the above system limitations, solution of the

minimization problem requires I to be bounded (and hence, in this case

continuous) for finite $a_1$,-----$a_N$. From equation (2.12) it is evident

that this is assured if and only if all the integrals

$$\int_{T_1}^{T_2} W(t) f_o^2(t)\, dt, \quad \int_{T_1}^{T_2} W(t)\theta_1^2(t)\, dt, --- \int_{T_1}^{T_2} W(t)\theta_N^2(t)\, dt$$

$$(2.15)$$

are bounded. The first integral depends on the prescribed response.

If it is unbounded a preliminary simplification of the prescribed

---

1 An nth order pole in $H^*(s)$ will cause equation (2.13) to have the
terms $\frac{A_m}{s-s_m}$, $\frac{A_{m+1}}{(s-s_m)^2}$,----$\frac{A_{m+n-1}}{(s-s_m)^n}$ and equation (2.14) to have the

terms $A_m e^{s_m t}$, $A_{m+1} t e^{s_m t}$,---$\frac{A_{m+n-1}}{(n-1)!} t^{n-1} e^{s_m t}$ .

2 $u_o(t)$ is the unit impulse function.

response is necessary. Such simplications are described in the next section. The remaining integrals depend on the approximating functions $\varphi_n(t)$ and the prescribed input $f_i(t)$. For realizable, finite, lumped element systems the integrals always will converge provided the energy in $f_i(t)$ between $T_1$ and $T_2$ is finite, i.e.

$$\int_{T_1}^{T_2} f_i^2(t) \, dt < \infty \qquad (2.16)$$

For a bounded input and a finite interval $(T_1, T_2)$ equation (2.16) is obviously satisfied. In an infinite interval the entire input must have finite energy. Thus, the infinite interval may be used for pulse -like inputs but not for ever present random inputs.[1] Useful, unbounded inputs are the unit impulse $u_0(t)$ and its successive derivatives $u_1(t)$ $u_2(t)$,------. If they occur in $(T_1, T_2)$ the integrals will be unbounded unless $H^*(s)$ falls off at a sufficient rate as $s \to \infty$. The exact condition required is $M-P \geq n+1$ where $n$ is the subscript on the highest order impulse appearing in $(T_1, T_2)$; that is, $H^*(s)$ must have at least an $(n+1)$th order zero at $s=\infty$.

Last of all, the approximating series $f_0^*(t)$ must not contain any redundant terms which can be expressed as linear combinations of the other terms of the series. More precisely, the functions $\Theta_1(t),------\Theta_N(t)$ must be linearly independent or $f_0^* = \sum_{n=1}^{N} a_n \Theta_n(t)$ must be zero only when all the coefficients $a_1,------a_N$ are zero. But $f_0^*(t) = \int_0^\infty h^*(\tau) f_i(t-\tau) \, d\tau$ is identically zero only when $h^*(t)$ is zero (excluding the trivial case where $f_i = 0$). Thus, the linear independence of the system approximating functions $\varphi_1(t),------\varphi_N(t)$ is a sufficient condition for the linear independence of $\Theta_1(t),------\Theta_N(t)$. Therefore, it will be assumed that $\varphi_1(t),---\varphi_N(t)$

1 The input energy may even be infinite if $W(t)$ goes to zero at a sufficient rate as $|t| \to \infty$.

are always linearly independent. Since M poles can generate only M independent functions, this further requires that $N \leq M$.

## Preliminary Simplifications of the Prescribed Response

In order to satisfy the previously mentioned integral restriction on $f_o(t)$ and to reduce approximation complexity and error, it is desirable to simplify the prescribed response $f_o(t)$ before the actual approximation process begins. Possible simplifications include:

1. <u>Change of time scale</u>. The prescribed input and response are replaced by $f_i(\frac{t}{T})$ and $f_o(\frac{t}{T})$ and the approximation completed. The resulting system function is $H^*(Ts)$. Through a logical choice of T, notation and numerical work is much simplified.

2. <u>Extraction</u>. This technique is feasible when $h(t)$ is approximated directly ($f_o = h$) and is expressed in analytic form. Terms which are realizable (the impulse, real and complex exponentials) or better handled by direct approximation (the derivatives or integrals of the impulse) are subtracted from $h(t)$ and realized separately. The remaining part of $h(t)$ is then approximated by the mean square error method.

3. <u>Delay removal</u>. The prescribed response is replaced by $f_o(t+T)$ and the approximation completed. The resulting system function is $H^*(s) e^{Ts}$ expressed as a rational function of s. If a shift in the time origin of the response is unimportant the result may be used directly. $H^*(s)$ is realized by following $H^*(s) e^{Ts}$ by an ideal delay of T seconds. Delay removal allows more accurate approximation when the approximated system possesses a delay-like character.

4. <u>Integration removal</u>. The prescribed response is replaced by $f_0'(t)$ and the approximation completed. The resulting system function is $s\,H^*(s)$, and multiplication by $\frac{1}{s}$ achieves the desired approximation. Integration removal is important in approximating prescribed systems known to have an integrating behavior. The multiplication by $\frac{1}{s}$ tends to produce an $f_0^*(t)$ which is smoothed and has less approximation error ripple.

5. <u>Differentiation removal</u>. The procedure is similar to 4 except $\int_0^t f_0(t)\,dt$ is approximated and the multiplication is by s. The error in $f_0^*(t)$ has increased approximation ripple.

One or more of the above operations will usually make $\int_{T_1}^{T_2} Wf_0^2\,dt$ bounded; if not, the problem statement should be again inspected. Perhaps, an equivalent but more suitable prescribed input and response can be chosen. Subsequent chapters assume that the above conditions are satisfied and that the preliminary simplifications are completed.

## III. THE WEIGHTED MEAN SQUARE ERROR CRITERION

As seen in Chapter II the weighted mean square error criterion (abbreviated WME criterion) is fundamental to the approximation process. Therefore, it is important to investigate the criterion and the reasons for its choice. In addition to doing this, the following sections will examine the resulting errors in the frequency domain and in system function approximation.

### The WME Criterion

In order to solve the minimization problem it is necessary to define a suitable measure of the approximation error. Such a measure should meet the following requirements: (1) be zero for zero error, (2) be positive for non-zero error, (3) decrease with decreasing error, (4) permit mathematical solution of the minimization problem. While the first three conditions are easily satisfied, it is the last condition which really forces the choice of the WME criterion. All other proposed criteria fail on this count and lead invariably to trial and error solutions of the minimization problem [13]. The WME criterion has a number of other important characteristics.

A most important characteristic is the square weighting of error amplitude, which tends to reduce strongly large errors at the expense of increased small errors. The error magnitude criterion is better in this respect but does not allow mathematical solution of the minimization

$$I = \int_{T_1}^{T_2} W(t) \, |e(t)| \, dt \qquad (3.1)$$

problem. Fortunately, the deficiency is not severe in most applications where the prescribed response is fairly smooth. The square weighting may even be preferred when large peak errors are particularly undesirable.

Because of the heavy weight placed on large errors the unweighted mean square error criterion ($W(t) = 1$) has a tendency to cause approximation errors which oscillate symmetrically about zero with relatively constant peak amplitude. This property seems characteristic of most mean square approximations, especially when the approximated functions are continuous. The Fourier series expansion of a triangular wave is a good example, the peak error deviations being almost equal with slight increases at points of slope discontinuity. An exception occurs when the prescribed response and approximating functions are both small for any length of time, since then the error amplitude must also be small. But in time regions where both the prescribed response and the approximating functions have appreciable value, the error does tend to oscillate with approximately constant peak deviations.

The WME criterion permits the peak deviations to be varied in a prescribed way as a function of time and is therefore superior to the simpler, more commonly used, unweighted mean square error criterion.[1] The variation is achieved by making $W(t)$ large in the time regions where the error is to be made small. In the weighted error integral $W^{\frac{1}{2}}f_0 - W^{\frac{1}{2}}f_0^*$ is equivalent to $f_0 - f_0^*$ in the unweighted error integral and has as a result approximately constant oscillation peaks in regions where $W^{\frac{1}{2}}f_0$ and $W^{\frac{1}{2}}f_0^*$ have appreciable value. Thus, the envelope of the oscillating error is approximately proportional to $W^{-\frac{1}{2}}(t)$. While $W$ may be any bounded, positive function there are certain

---

1 A good example of WME criterion application is given by Westcott [40] where $W(t) = t$.

-19-

more desirable choices, discussed in part two, which reduce numerical
and analytic work in impulse response approximation.

Although it is impossible to calculate the error $e(t)$ without
evaluating the approximating series, it is possible to make a rough
estimate of its maximum value from the quanity I. From the discussion
of the previous paragraph the following somewhat crude approximation
to $e(t)$ seems reasonable:

$$e(t) \simeq e^{*}(t) = AW^{-\frac{1}{2}}(t)\sin \Omega_0 t, \quad \overline{T}_1 \leq t \leq \overline{T}_2 \qquad (3.2)$$

$$= 0, \quad t < T_1, \ t > T_2.$$

The interval $(\overline{T}_1, \overline{T}_2)$ is the region in which $W^{\frac{1}{2}}f_0(t)$ and $W^{\frac{1}{2}}f_0^{*}(t)$ have
appreciable value, and $\sin \Omega_0 t$ represents approximately the rather
oscillatory nature of $e(t)$ within the envelope $W^{-\frac{1}{2}}(t)$. Substituting
equation (3.2) in the WME integral gives

$$I \simeq \int_{\overline{T}_1}^{\overline{T}_2} W(t)[AW^{-\frac{1}{2}}(t) \sin \Omega_0 t]^2 dt = A^2 \int_{\overline{T}_1}^{\overline{T}_2} [\tfrac{1}{2} - \tfrac{1}{2}\sin 2\Omega_0 t]\, dt. \qquad (3.3)$$

Assuming an integral number of periods of $\sin 2\Omega_0 t$ in $(\overline{T}_1, \overline{T}_2)$,

$$I \simeq \frac{A^2}{2}(\overline{T}_2 - \overline{T}_1). \qquad (3.4)$$

Finally, eliminating A between equation (3.4) and equation (3.2) yields,

$$e_{max} \simeq e_{max}^{*} = \sqrt{\frac{2I}{\overline{T}_2 - \overline{T}_1}} \ W^{-\frac{1}{2}}(t) \sin \Omega_0 t \Big|_{max}$$

$$\simeq \sqrt{\frac{2IW^{-1}_{max}}{\overline{T}_2 - \overline{T}_1}}, \qquad (3.5)$$

the estimate of the maximum error deviation. Despite the rather gross
assumptions made in deriving equation (3.5), it has given in many cases
accuracies of better than 50%. Application of equation (3.5), including
the choice of $(\bar{T}_2, \bar{T}_1)$, is demonstrated in part two, Chapter VII.

### Frequency Domain Errors

Minimization of the WME integral brings about an approximation
in the frequency domain as well as in the time domain. To see the re-
lation between frequency domain errors and time domain errors, consider
first the infinite interval, unweighted mean square error integral,

$$I = \int_{-\infty}^{\infty} [f_0 - f_0^*]^2 \, dt \quad .$$  (3.6)

By means of the complex convolution integral[1] of the Laplace transform
I may be written as

$$I = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} [F_0(-s) - F_0^*(-s)][F_0(s) - F_0^*(s)] \, ds \quad .$$  (3.7)

Evaluating equation (3.7) on the imaginary axis (real frequency axis)
yields

$$I = \frac{1}{2\pi} \int_{-\infty}^{\infty} |F_0(j\omega) - F_0^*(j\omega)|^2 \, d\omega \quad .$$  (3.8)

Thus, the unweighted mean square approximation in the time domain is a
mean square error-magnitude approximation in the frequency domain.

But frequency domain approximation errors are more usually
measured in terms of magnitude error, phase error, real part error,
and imaginary part error. The mean square integrals of these errors
are respectively:

$$I_M = \frac{1}{2\pi} \int_{-\infty}^{\infty} [\,|F_0| - |F_0^*|\,]^2 d\omega,$$  (3.9)

---

[1] See Gardner and Barnes [10], page 275.

$$I_P = \frac{1}{2\pi} \int_{-\infty}^{\infty} [angF_o - angF_o{}^*]^2 d\omega , \qquad (3.10)$$

$$I_R = \frac{1}{2\pi} \int_{-\infty}^{\infty} [ReF_o - ReF_o{}^*]^2 d\omega , \qquad (3.11)$$

$$I_I = \frac{1}{2\pi} \int_{-\infty}^{\infty} [ImF_o - I_mF_o{}^*]^2 d\omega . \qquad (3.12)$$

Expansion of equation (3.8) into imaginary and real parts of $F_o$ and $F_o{}^*$ gives

$$I = I_R + I_I. \qquad (3.13)$$

so that $I \geq I_R$ and $I \geq I_I$. Reference to Figure 3.1 shows that $|F_o - F_o{}^*| \geq |F_o| - |F_o{}^*|$ so that $I \geq I_M$. It is seen that $I_M$, $I_R$, and $I_I$



Figure 3.1  Vector Diagram of $F_o$, $F_o{}^*$, and $F_o - F_o{}^*$

are not minimized by minimizing I.  They do converge toward zero, however, as I converges toward zero.

No similar bound on the mean square phase error $I_P$ is possible. When $|F_o - F_o{}^*|$, $|F_o|$ and $|F_o{}^*|$ are all small, then the angular error between $F_o$ and $F_o{}^*$ may be large though $|F_o - F_o{}^*|$ is small. Unfortunately, all three quantities often become small together as $\omega \to \infty$

and the minimization of I (a measure of $F_O - F_O^*$) allows large phase errors at high frequencies. This is what happens in practice. The phase of $F_O^*(j\omega)$ as $\omega \to \infty$ is dependent on the approximating functions used and not on the phase of $F_O(j\omega)$ as $\omega \to \infty$. This limitation of the mean square error approximation can be eliminated by constrained approximations of the type described in the next chapter. Good phase approximation, as well as good magnitude, real part, and imaginary part approximation, is then possible.

The introduction of the weight factor $W(t)$ and the finite interval $(T_2, T_1)$ makes the above relations much more complicated. One interpretation of the WME integral is to consider I as a weighted average of frequency domain errors arising from time errors which exist in small time intervals. Thus,

$$I \cong \sum_{n=1}^{N} W(T_1 + n\Delta t) \frac{1}{2\pi} \int_{-\infty}^{\infty} |F_{on} - F_{on}^*|^2 d\omega \qquad (3.14)$$

where $N \Delta t = T_2 - T_1$ and

$$F_{on}(s) = \int_{T_1+(n-1)\Delta t}^{T_1+n\Delta t} e^{-st} f_o(t) dt \quad . \qquad (3.15)$$

Though this representation is admittedly strained it does allow the previously developed results to be extended to the WME criterion.

### System Function Approximation Errors

It is useful to know how the prescribed impulse response function $h(t)$ is approximated. If $f_i(t)$ is the unit impulse, the WME criterion applies directly to the impulse response error $e_h(t) = h(t) - h^*(t)$. If $f_i(t)$ is an arbitrary input the approximation criterion which applies to $e_h(t)$ is not as simple.

Again, the initial development will be restricted to the infinite interval, unweighted mean square error criterion. Substituting $F_o = HF_i$ and $F_o^* = H^* F_i$ in equation (3.8) yields

$$I = \frac{1}{2\pi} \int_{-\infty}^{\infty} |F_i H - F_i H^*|^2 d\omega = \frac{1}{2\pi} \int_{-\infty}^{\infty} |F_i|^2 |H - H^*|^2 d\omega . \quad (3.16)$$

The effect of an arbitrary input is now clear. The mean square error in the $h(t)$ approximation is weighted in the frequency domain by $|F_i|^2$. The result seems reasonable, for the best approximation would be expected in the frequency range where the input amplitudes are greatest.

It should be noted that the weighting depends only on $|F_i|$ and not on the phase of $F_i$. Thus, when the unweighted mean square error criterion is used, the input need not be the prescribed input $f_i(t)$ but may be any time function which possesses the same spectral magnitude. For example, a random input with spectral energy $|F_i|^2 = \frac{1}{1+\omega^2}$ may be replaced by the exponential pulse

$$f_i(t) = e^{-t} , \quad t \geq 0$$
$$= 0 , \quad t < 0 .$$

which has the same $|F_i|$. The approximation independence of input phase does not indicate, though, that phase characteristics of the prescribed system function $H(s)$ are ignored. A change in the phase of $H(j\omega)$ does change $I$ as given by equation (3.16).

The extension of these relations to the WME integral requires techniques similar to those used in the previous section. A preferable alternative is to consider two block diagrams which are equivalent to the diagram of Figure 2.2. Figure 3.2 shows the first modification of Figure 2.2. The input is replaced by a unit impulse followed by a system function $F_i(j\omega)$ (not necessarily realizable). Since the system

Figure 3.2  Equivalent Block Diagram of Approximation Problem with Input Replaced by $F_i(j\omega)$ and $u_o(t)$

is linear from point A to point B, the order of the operations is immaterial, and $F_i(j\omega)$ may be shifted from A to B as in Figure 3.3. Figure 3.3 shows that the error $e_h$ in impulsive response approximation is first frequency weighted by $F_i(j\omega)$ and then squared and time weighted



Figure 3.3  Equivalent Block Diagram of Approximation Problem Showing the Weighting of $e_h$ by $F_i(j\omega)$

by $W(t)$. Since the phase of $F_i(j\omega)$ affects the time distribution of the filtered error $e(t)$, it must also affect the value of I. Thus, contrary to the unweighted approximation, the weighted approximation is sensitive to forcing function phase.

In summary, the following statements concerning the WME approximation can be made:

1. The WME criterion is the only practical error criterion permitting mathematical solution of the minimization problem.

2. Large approximation errors are weighted heavily compared with small approximation errors.

3. The approximation error tends to oscillate about zero and when $W^{\frac{1}{2}}f_0$ and $W^{\frac{1}{2}}f_{0_1}^*$ are reasonably large has an envelope roughly proportional to $W^{-\frac{1}{2}}(t)$.

4. Good time domain approximation assures good frequency domain approximation of the magnitude, the real part, and the imaginary part of the response $F_0(j\omega)$.

5. Good phase approximation of $F_0(j\omega)$ for small $F_0(j\omega)$ generally demands a constrained WME approximation.

6. For arbitrary inputs the impulse response error $e_h$ is weighted frequency-wise by $F_i(j\omega)$.

# IV. THE MINIMIZATION PROBLEM

The preceding sections have described in detail the approximation problem, mathematical and realizability restrictions, and properties of the WME approximation. This chapter will develop the general theory of WME approximation and derive the equations necessary to solve the minimization problem for simple and constrained approximations. Practical applications of the results are delayed until part two.

Before beginning it is advisable to review briefly the conditions specified in Chapter II. They are: (1) $H*(s)$ must be physically realizable, (2) the weighted square integrals of $f_o$, $\theta_1$, $---$ $\theta_N$ must be bounded (to assure the continuity of I for finite $a_1$, $---$ $a_N$), (3) the system approximating functions $\varphi_1$, $---$ $\varphi_N$ must be linearly independent (to prevent $\theta_1$, $---$ $\theta_N$ from being linearly dependent). The first section that follows is devoted to the simple or unconstrained WME approximation.

## The Unconstrained Approximation

The WME integral is given by

$$I = \int_{T_1}^{T_2} W\,[f_o - f_o*]^2\,dt = \int_{T_1}^{T_2} W\,[f_o - \sum_{n=1}^{N} a_n\,\theta_n]^2\,dt$$

$$= \int_{T_1}^{T_2} W\,f_o^2\,dt - 2\sum_{n=1}^{N} a_n \int_{T_1}^{T_2} W\,f_o\theta_n\,dt + \sum_{n=1}^{N}\sum_{m=1}^{N} a_n a_m \int_{T_1}^{T_2} W\,\theta_n\theta_m\,dt \;.$$

$$(4.1)$$

Since the integrals of equation (4.1) appear frequently it is desirable to introduce the following notation:

$$d_n = \int_{T_1}^{T_2} W \, f_o \, \Theta_n \, dt, \qquad n = 1, \; - \; - \; - \; N \; , \qquad (4.2)$$

$$c_{nm} = \int_{T_1}^{T_2} W \, \Theta_n \Theta_m \, dt = c_{mn}, \quad n = 1, \; - \; - \; - \; N, \; m = 1, \; - \; - \; - \; N. \qquad (4.3)$$

Substituting these expressions into equation (4.1) yields

$$I = \int_{T_1}^{T_2} W f_o^2 dt \; - \; 2 \sum_{n=1}^{N} a_n d_n \; + \; \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m c_{nm} \qquad . \qquad (4.4)$$

To obtain the **WME** approximation I must be minimized with respect to $a_1, \; - \; - \; - \; a_N$. Since I is continuous for finite $a_1, \; ---- \; a_N$ it possesses continuous derivatives with respect to the $a_n$ which may be set equal to zero to find the stationary point. Thus,

$$\frac{\partial I}{\partial a_n} = \; - \; 2 \, d_n \; + \; 2 \sum_{m=1}^{N} a_m c_{nm} = 0, \quad n = 1, \; - \; - \; - \; N. \qquad (4.5)$$

These N equations defining the stationary point coefficients may be expressed more simply in matrix notation by

$$d] = C \; a] \qquad . \qquad (4.6)$$

If the matrix C is non-singular (det C $\neq$ 0) a unique solution of equation (4.6) results. But the det C is the Gram[1] determinant of the functions $W^{1/2}\Theta_n$, which are linearly independent because the

---

[1] See Courant and Hilbert[8], page 61.

functions $\Theta_n$ are linearly independent. The determinant is therefore non-zero, and equation (4.6) does define a unique stationary point. Furthermore, the quadratic form

$$\sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m c_{nm}$$

is positive definite.[1] Hence I can be made arbitrarily large by increasing $|a_1|, - - - |a_N|$ indefinitely.[2] Since the stationary point is unique it must therefore be a minimum.

The value of the minimum may be calculated by expressing equation (4.4) in matrix notation and substituting the known value of d]. Thus

$$I = \int_{T_1}^{T_2} W f_o^2 \, dt - 2 \underline{a} \, d] + \underline{a} \, C \, a]$$

$$= \int_{T_1}^{T_2} W f_o^2 \, dt - \underline{a} \, d] \quad .$$

(4.7)

Other equally good expressions are

$$I = \int_{T_1}^{T_2} W f_o^2 \, dt - \underline{a} \, C \, a] \quad , \qquad (4.8)$$

$$I = \int_{T_1}^{T_2} W f_o^2 \, dt - \underline{d} \, C^{-1} \, d] \quad . \qquad (4.9)$$

---

1  Ibid.

2  For large $|a_n|$ the quadratic form in equation (4.4) increases much more rapidly than $\sum_{n=1}^{N} a_n d_n$.

Because the quadratic form $\underline{a}\,C\,a]$ is positive definite

equation (4.8) shows that

$$I \leq \int_{T_1}^{T_2} W\, f_o^2\, dt \quad .$$

(4.10)

It is therefore logical to define the worst possible WME as

$$I_{max} = \int_{T_1}^{T_2} W\, f_o^2\, dt$$

(4.11)

which exists when there are no terms in the approximating series. The

relative error $I/I_{max}$ is less than or equal to one and is particularly

useful because it provides a normalized indication of approximation ac-

curacy which is independent of the amplitude or time scaling of the

approximation problem.

To summarize, the solution of minimization problem is unique

and involves the following steps: (1) the evaluation of the $1/2\ N(N+3)$

integrals (equations 4.2 and 4.3) which are the elements of matrices

d] and C, (2) the inversion of the matrix C and the calculation of the

column of approximating coefficients from

$$a] = C^{-1}\, d] \quad ,$$

(4.12)

(3) the calculation of I from known a] and d] by equation (4.7).

Usually, the inversion of C entails the greatest amount of work, es-

pecially when C is large. It should be noted, however, that C depends

only on the approximating functions $\theta_1,\ -\ -\ -\ \theta_N$ and not on the

prescribed response $f_0$. Thus, a <u>single inverse matrix suffices for the approximation of any number of prescribed responses</u>.

It would be particularly fortunate if

$$c_{nm} = \int_{T_1}^{T_2} W\, \Theta_n \Theta_m\, dt = 1\ , \qquad n = m \qquad\qquad (4.13)$$
$$= 0\ , \qquad n \neq m$$

for then the coefficients would be given simply by

$$a_n = d_n = \int_{T_1}^{T_2} W\, f_0\, \Theta_n\, dt, \qquad n = 1,\ -\ -\ -\ N \qquad (4.14)$$

and the WME by

$$I = I_{max} - \underline{a}\,\ a] = I_{max} - \sum_{n=1}^{N} a_n \qquad\qquad . \qquad (4.15)$$

While such "orthonormal" function approximations are valuable they are not necessarily easier to implement since the generation of orthonormal families of functions may be as difficult as matrix inversion. Ortho-gonalization procedures are discussed fully in the next chapter. Relative merits of orthogonal and non-orthogonal approximation depend on application and are treated in part two.

Since the generation of families of orthonormal functions depends on the linear combination of linearly independent functions, it is important to investigate the approximation properties of approximating functions which are linear combinations of $\Theta_1,\ -\ -\ -\ \Theta_N$. The most general set of such functions $\overline{\Theta}_1,\ -\ -\ -\ \overline{\Theta}_N$ is given by the matrix equation

$$\overline{\theta}] = R \, \theta] \quad , \tag{4.16}$$

where the matrix R is any matrix whose determinant is non-zero. The latter requirement is necessary to assure the linear independence of the new functions $\overline{\theta}]$. By substituting $\overline{\theta}]$ in the previously developed equations it is readily shown that approximation by $\overline{\theta}]$ is exactly equivalent to approximation by $\theta]$. Hence, it makes no difference in the final approximation whether or not the approximating functions are first orthonormalized.

### The Constrained Approximation

Frequently it is necessary for the approximate response to satisfy exactly certain specific conditions which stem from important properties of the prescribed response or from the requirements of a particular application. Unfortunately, such specific conditions can only be approached, not equaled, by the finite, unconstrained approximation of the previous section. It is therefore desirable to formulate a WME approximation which is constrained.

Typical conditions which can be enforced by constrained approximation are: (1) the values of $f_0^*(t)$ or its derivatives at specified instants of time, (2) the values of $F_0^*(s)$ or its derivatives at specified points in the complex plane, and (3) the asymptotic behavior of $F_0^*(s)$ at $s = \infty$. The last condition is especially important when the asymptotic phase is to be exact, or when $\lim_{|s| \to \infty} s^n H^*(s) = \text{const}$ is fixed by practical realizability limitations. In any event, none of the above conditions are normally obtained with a simple WME approximation. They must be forced on

$f_o*(t)$ by restricting the way in which the approximating functions are combined.

Symbolically, a constraint condition may be written as

$$k = K \quad \{f_o*(t)\} \quad . \tag{4.17}$$

$K \{ \ \}$ is the operator which measures the desired condition, and $k$ is its specified value. For a particular set of approximating functions equation (4.17) is a relation between the approximation coefficients. To distinguish these constrained approximation coefficients from the unconstrained coefficients, they will be denoted by $b]$, i.e.,

$$f_o* = \sum_{n=1}^{N} b_n \theta_n (t) \quad . \tag{4.18}$$

The following examples show how constraints restrict the way in which the coefficients can vary:

1. The value of $f_o*(t)$ equals the value of $f_o(t)$ at $t = t_1$.

$$f_o(t_1) = k = f_o*(t_1) = \sum_{n=1}^{N} b_n \theta_n (t_1) \quad .$$

2. The area under $f_o*(t)$ equals one.

$$1 = k = \int_{T_1}^{T_2} f_o* \, dt = \sum_{n=1}^{N} b_n \int_{T_1}^{T_2} \theta_n (t) \, dt \quad .$$

3. The $\lim_{|s| \to \infty} s^3 F_o*(s) = $ const. This condition is restated as $f_o*(0) = 0$ and $f_o'*(0) = 0$.

$$0 = k_1 = f_o*(0) = \sum_{n=1}^{N} b_n \theta_n (0) \quad ,$$

$$0 = k_2 = f_o'*(0) = \sum_{n=1}^{N} b_n \theta_n' (0) \quad .$$

4. The mean square value of $f_o*$ equals one.

$$1 = k = \int_{T_1}^{T_2} f_o*^2 \, dt = \sum_{n=1}^{N} \sum_{m=1}^{N} b_n b_m \int_{T_1}^{T_2} \theta_n \theta_m \, dt \quad .$$

The first three conditions might be termed linear constraints in that

$$k = K \{f_o*(t)\} = \sum_{n=1}^{N} b_n K \{\theta_n(t)\} \quad . \qquad (4.19)$$

Linear constraints are by far the most common and permit a simple theoretical solution of the constrained WME approximation problem. The fourth constraint is nonlinear in the approximation coefficients and makes practical solution of the minimization problem very difficult. Constraints which are nonlinear are not considered further in this investigation.

In general, there may be more than one constraint equation. If there are Q of them, they may be ordered by subscripts as follows:

$$k_1 = K_1 \{f_o*\} = \sum_{n=1}^{N} b_n K_1 \{\theta_n\} = \sum_{n=1}^{N} b_n k_{1n}$$

$$k_2 = K_2\{f_o*\} = \sum_{n=1}^{N} b_n \; K_2\{\theta_n\} = \sum_{n=1}^{N} b_n \; k_{2n}$$

$$k_Q = K_Q\{f_o*\} = \sum_{n=1}^{N} b_n \; K_Q\{\theta_n\} = \sum_{n=1}^{N} b_n \; k_{Qn} \quad \cdot (4.20)$$

where

$$k_{qn} = K_q\{\theta_n\} \quad . \tag{4.21}$$

Equations (4.20) are written more compactly as

$$k] = K \; b] \quad . \tag{4.22}$$

The matrix $K$ has elements $k_{qn}$ given by (4.21) and has $Q$ rows and $N$ columns.

All the constraint equations in the set (4.20) must be consistent. If the approximating functions are such that $k_{q1}$, $k_{q2}$, - - - $k_{qN}$ are all zero, then the specified value $k_q$ must also be zero. In this case no limitation is imposed on b] and the whole equation is superfluous. Similarly, the rows of $K$ must not be linearly dependent. If they are, the same constraint condition holds at least twice, perhaps with different specified values.

The solution of the constrained approximation problem requires the minimization of the constrained WME integral

$$I_c = \int_{T_1}^{T_2} W(t) \; [f_o - \sum_{n=1}^{N} b_n \; \theta_n]^2 \; dt \tag{4.23}$$

under the supposition that the coefficients b] always obey the matrix constraint equation. There are several ways of approaching the minimization problem. The most direct method is to use equation (4.22) to eliminate Q of N coefficients in (4.23) and to vary the remaining N-Q coefficients to produce the minimum.[1] A more satisfactory procedure, which makes better use of the equations already developed for unconstrained approximation, is the method of Lagrange multipliers.[2]

The multipliers $\lambda_1, - - - \lambda_Q$ are introduced and the quantity

$$I_e = \int_{T_1}^{T_2} W[f_o - f_o*]^2 \, dt + \sum_{q=1}^{N} \lambda_q [ \quad K_q \{f_o*\} - k_q ] \qquad (4.24)$$

is minimized with respect to $b_1, - - - b_N$ and $\lambda_1, - - - \lambda_Q$. This is equivalent to minimizing equation (4.23) subject to equation (4.22). Setting

$$\frac{\partial I_c}{\partial \lambda_q} = 0, \ q = 1, - - - Q; \quad \frac{\partial I_c}{\partial b_n} = 0, \ n = 1, - - - N$$

yields equation (4.22) and

$$0 = -2 d_n + 2 \sum_{m=1}^{N} b_m c_{nm} + \sum_{q=1}^{N} \lambda_q k_{qn}, \ n = 1, - - - N , \qquad (4.25)$$

where $d_n$ and $c_{nm}$ are given by the same integral formulas as before. In matrix notation equation (4.25) is simply

$$d] = C \ b] + 1/2 \ K_T \lambda] \qquad . \qquad (4.26)$$

---

[1] Note that N must be greater than Q.

[2] Courant and Hilbert[8] page 164.

$K_T$ has N rows and Q columns and is obtained by transposing the rows and columns of K.

Since C has an inverse, equation (4.26) may be solved for b] to give

(4.27)

$$b] = C^{-1} d] - 1/2 \ C^{-1} K_T \lambda] = a] - 1/2 \ C^{-1} K_T \lambda]$$

where a] is the column of unconstrained approximation coefficients. To find b] the column $\lambda]$ must be obtained. This is done by substituting b] from equation (4.27) into equation (4.22),

$$k] = K b] = K a] - 1/2 K \ C^{-1} K_T \lambda], \quad (4.28)$$

and solving equation (4.28) for $\lambda]$

$$\lambda] = 2 \left\{ KC^{-1} K_T \right\}^{-1} \left\{ K a] - k] \right\} . \quad (4.29)$$

The matrix $K C^{-1} K_T$ has Q rows and columns and has an inverse if the rows of K are linearly independent, a condition which has already been assumed. Equation (4.27) and equation (4.29) define a unique stationary point which by the argument used earlier must be a minimum.

Before calculating the value of the minimum it is worthwhile to compare the solution steps for unconstrained and constrained approximations. The unconstrained approximation requires: (1) the evaluation of the integrals for $d_n$ and $c_{nm}$, (2) the inversion of C and the calculation of a]. In addition to these steps the constrained approximation necessitates: (3) the evaluation of the coefficients $k_{qn}$ from equation (4.21), (4) the calculation and inversion of

$K \ C^{-1} \ K_T$, and (5) the solution of the following equation for b]
(obtained by substituting $\lambda$] from equation (4.29) in equation (4.27)).

$$b] = a] - C^{-1} K_T \left\{ K \ C^{-1} \ K_T \right\}^{-1} \left\{ K \ a] - k] \right\}.$$
$$(4.30)$$

None of these additional equations are particularly complicated. The
inversion of $K \ C^{-1} \ K_T$ is usually not too difficult since the
number of constraints Q is generally small. Where different prescribed
responses are approximated with the same approximating functions and
constraints, only the columns a] in equation (4.30) are changed. Thus,
the constrained approximation is a relatively simple and straight-
forward extension of the unconstrained approximation.

The minimum value of $I_c$ is given by

$$I_c = I_{max} - 2 \ \underline{b} \ d] + \underline{b} \ C \ b].$$
$$(4.31)$$

To simplify the manipulations let

$$g] = 1/2 \ K_T \ \lambda] \ ,$$
$$(4.32)$$

then

$$C \ b] = d] - g]$$

and equation (4.31) becomes

$$I_c = I_{max} - 2 \ \underline{b} \ d] + \underline{b} \ \left\{ d] - g] \right\}$$

$$= I_{max} - \underline{b} \ \left\{ d] + g] \right\} \qquad .$$

$$I_c = I_{max} - \{ \underline{a} - \underline{g} \quad C^{-1} \} \{d] + g]\}$$

$$= I_{max} - \underline{a} \quad d] + \underline{g} \quad C^{-1} g] \quad .$$

$$I_c = I + \underline{g} \quad C^{-1} g] \tag{4.33}$$

$$I_c = I + 1/4 \underline{\lambda} \quad K \quad C^{-1} \quad K_T \lambda] \tag{4.34}$$

$$I_c = I + 1/2 \underline{\lambda} \left\{ K \; a] - k]\right\}, \tag{4.35}$$

where I is the WME for the unconstrained approximation.  Either equation (4.34) or equation (4.35) may be used to calculate $I_c$, the latter probably being the easiest since both $\lambda]$ and $\left\{ K \quad a] - k]\right\}$ are encountered in application of equation (4.30).

Since the quadratic form $\underline{g} \quad C^{-1} g]$ is positive definite the constrained error $I_c$ must be greater than or equal to the unconstrained error I.  The number $1/4 \underline{\lambda} \quad K \quad C^{-1} \quad K_T \lambda]$ (or $1/2 \underline{\lambda} \left\{ K \; a] - k]\right\}$) measures the difference between $I_c$ and I and thus the severity of the imposed constraints.  If the prescr.bed response and constraints are incompatible $I_c$ may be intolerably large even when N is arbitrarily large.

As in the unconstrained approximation the functions $\bar{\theta}]$, which are defined by equation (4.16) produce the same approximation as $\theta]$.  Hence the functions $\theta]$ may be orthogonalized and still produce the same final approximation.  Orthonormal function approximation causes little change in the derived formulas.  The matrices C and $C^{-1}$ are simply replaced by the identity matrix.

## Approximation by Constrained Functions

The method of constrained approximation just presented is important but does involve additional work over the unconstrained approximation. It would be useful, if possible, to work instead with a series of approximating functions, each constrained so that equation (4.30) would reduce to

$$b] = a]. \qquad (4.36)$$

Such a set of functions would make constrained approximation more practical.

Inspection of equation (4.30) shows that equation (4.36) holds if $K = 0$. This places a restriction on the type of constraints that can be used. They must be homogeneous, i.e. k] must equal zero. Fortunately, homogeneous constraints are relatively common. Examples include: (1) zeros of $f_o^*(t)$, (2) zeros of $F_o^*(s)$, (3) asymptotic behavior of $F_o^*(s)$ at $s = \infty$. It will now be shown that with homogeneous constraints it is always possible to construct a set of constrained functions.

The development depends on an alternative solution of the constrained WME approximation problem. The approximate response is given by

$$f_o^* = \sum_{n=1}^{N} b_n \theta_n(t) = \underline{\theta} \ b] \qquad (4.37)$$

where

$$Kb] = 0 \qquad (4.38)$$

By partitioning the Q by N matrix K it is possible to eliminate Q of the N coefficients in equation (4.37). Thus equation (4.38) becomes

$$Kb] = 0 = [K_{QQ} K_{Q(N-Q)}] \begin{array}{c} b]_Q \\ b]_{(N-Q)} \end{array}. \qquad (4.39)$$

$$0 = K_{QQ} b]_Q + K_{Q(N-Q)} b]_{(N-Q)}$$

The subscripts indicate the number of elements in matrices. For example,

$$K_{Q(N-Q)} = \begin{bmatrix} K_{1(Q+1)} ---K_{1N} \\ \\ \\ K_{Q(Q+1)} ---K_{QN} \end{bmatrix}, \qquad b]_Q = \begin{bmatrix} b_1 \\ \\ b_Q \end{bmatrix}.$$

The particular choice of partitioning in equation (4.38) makes no real difference since the ordering of the original functions $\theta$] is arbitrary. Since the rows of K are assumed to be linearly independent the ordering may be arranged so det $K_{QQ} \neq 0$ and

$$b]_Q = -K_{QQ}^{-1} K_{Q(N-Q)} b]_{(N-Q)} \qquad (4.40)$$

Substitution of equation (4.40) in equation (4.37) yields the result

$$f_o^* = \underline{\theta} \; b] = \underline{\theta}_Q \; b]_Q + \underline{\theta}_{(N-Q)} b]_{(N-Q)}$$

$$= \{ \underline{\theta}_{(N-Q)} - \underline{\theta}_Q K_{QQ}^{-1} K_{Q(N-Q)} \} b]_{(N-Q)} \qquad (4.41)$$

The minimization problem is solved by varying the remaining $N-Q$ coefficients of $b]_{(N-Q)}$ to make the WME integral a minimum. In terms of the $N-Q$ functions $\overline{\theta}_1, ----- \overline{\theta}_{N-Q}$,

$$\underline{\overline{\theta}} = \underline{\theta}_{(N-Q)} - \underline{\theta}_Q K_{QQ}^{-1} K_{Q(N-Q)} \qquad (4.42)$$

and the $N-Q$ coefficients $\overline{b}_1, ----- \overline{b}_{N-Q}$,

$$\overline{b}] = b]_{N-Q} \qquad (4.43)$$

the approximation is an unconstrained one yielding

$$\overline{b}] = \overline{C}^{-1} \overline{d}] = \overline{a}].$$

Thus the equivalent unconstrained functions $\underline{\bar{\Theta}}$ are a linearly independent[1] set satisfying the condition $\bar{K} = 0$.

Often, the family of constrained functions may be found more simply by inspection. Suppose $\Theta_1 = \dfrac{1}{s-s_1}$ , ----- $\Theta_N = \dfrac{1}{s-s_N}$ and the condition $s^3 F_o{}^*(s) \to \text{const.}$ as $s \to \infty$ is required. This causes the two constraints

$$\bar{\Theta}_n(o) = 0 , \quad \bar{\Theta}_n{}'(o) = 0.$$

An obvious set of N-2 constrained, linearly independent functions is simply

$$\bar{\Theta}_1 = \frac{1}{(s-s_1)(s-s_2)(s-s_3)}, \quad \bar{\Theta}_2 = \frac{1}{(s-s_1)(s-s_2)(s-s_4)},$$

$$- - - - - - - - \quad \bar{\Theta}_{N-2} = \frac{1}{(s-s_1)(s-s_2)(s-s_N)}.$$

This particular type of constrained function is very useful and its application in a practical problem will be demonstrated later.

While limited to homogeneous constraints, the constrained function approach is also valuable when additional non-homogeneous conditions are specified, since the number of Lagrange multipliers is then only the number of non-homogeneous conditions. Further convenience in the application of constrained functions would be obtained by their orthogonalization.

To review, WME approximation may be either constrained or unconstrained. Unconstrained approximations result in simple coefficient equations but do not allow approximation conditions to be specified precisely. Constrained approximations are more difficult to obtain

---

[1] They are linearly independent because each function contains separately one of the independent functions of $\underline{\Theta}$ $(N-Q)$.

but permit arbitrary linear constraints to be imposed in a straight-

forward manner.  And finally, homogeneous linear constraints may often

be implemented with less work by means of suitably constrained approxima-

ting functions.

# V.  ORTHONORMAL FUNCTIONS

The simplification in coefficient evaluation which results from orthonormal function approximation is especially important when many different approximations are to be made with the same set of approximating functions.  It is therefore appropriate to discuss various sets of orthonormal functions and the procedures for obtaining them by orthogonalization of linearly independent functions.

Two classes of orthogonalization procedures are considered in this chapter.  The first is applicable to any arbitrary set of linearly independent functions.  The second is restricted to linearly independent exponential functions[1] of the type which occur in impulse response approximation.  In this class simple formulas are available to define the orthonormal functions.

## Orthogonalization of Arbitrary Functions

Various orthogonalization procedures which are applicable to any set of linearly independent approximating functions may be derived, all producing different families of orthonormal functions.  The Schmidt process, which is developed in this section, is as practical as any, even though it demands considerable numerical effort.  Although the details of the process are well known [8], they are repeated here for convenience.

The desired orthonormal functions $\bar{\theta}_1, \text{-----} \bar{\theta}_N$ are expressed by linear combinations of the given linearly independent functions $\theta_1, \text{-----} \theta_N$.  To avoid complete repetition of the orthogonalization

---

[1] The exponential functions considered include real exponential functions and exponentially damped cosine and sine functions.

procedure when the number of functions N is changed, it is assumed that $\bar{\Theta}_1$ depends only on $\Theta_1$, $\bar{\Theta}_2$ depends only on $\Theta_1$ and $\Theta_2$, and so on. Thus

$$\bar{\Theta}_1 = r_{11}\Theta_1$$

$$\bar{\Theta}_2 = r_{21}\Theta_1 + r_{22}\Theta_2$$

$$\bar{\Theta}_3 = r_{31}\Theta_1 + r_{32}\Theta_2 + r_{33}\Theta_3$$

$$\bar{\Theta}_N = \sum_{n=1}^{N} r_{Nn}\,\Theta_n \tag{5.1}$$

The transformation coefficients $r_{nm}$ must be chosen so the orthonormal condition holds; i.e. so

$$\tag{5.2}$$

$$\bar{c}_{nm} = \int_{T_1}^{T_2} W(t)\,\bar{\Theta}_n(t)\,\bar{\Theta}_m(t)dt = (\bar{\Theta}_n,\bar{\Theta}_m) = 1,\ n=m$$
$$= 0,\ n\neq m \quad .$$

Applying equation (5.2) with $n = m = 1$ yields

$$(\bar{\Theta}_1,\bar{\Theta}_1) = r_{11}^2\,(\Theta_1,\Theta_1) = 1 \tag{5.3}$$

and

$$r_{11} = \frac{1}{\sqrt{(\Theta_1,\Theta_1)}} = \frac{1}{\sqrt{c_{11}}} \quad . \tag{5.4}$$

Considering next $\bar{\Theta}_2$, equation (5.2) must hold for $n = 2$ and $m = 1,2$. If $\bar{\Theta}_2$ is written as a linear combination of $\bar{\Theta}_1$ and $\Theta_2$, it is apparent that $(\bar{\Theta}_2,\bar{\Theta}_1) = 0$ requires

$$\bar{\Theta}_2 = r_{22}\,[\Theta_2 - (\Theta_2,\bar{\Theta}_1)\,\bar{\Theta}_1] \quad .$$

Taking $r_{22} = \dfrac{1}{\sqrt{\text{norm }[\Theta_2 - (\Theta_2,\,\bar{\Theta}_1)\,\bar{\Theta}_1]}} = \dfrac{1}{\sqrt{(\Theta_2 - (\Theta_2,\bar{\Theta}_1)\,\bar{\Theta}_1,\Theta_2-(\Theta_2,\bar{\Theta}_1)\bar{\Theta}_1)}}$

makes $(\bar{\Theta}_2,\bar{\Theta}_2) = 1$ and gives

$$\bar{\Theta}_2 = \frac{\Theta_2 - (\Theta_2,\bar{\Theta}_1)\,\bar{\Theta}_1}{\sqrt{(\Theta_2,\Theta_2) - 2(\Theta_2,\bar{\Theta}_1)^2 + (\Theta_2,\bar{\Theta}_1)^2}}$$

$$= \frac{\theta_2 - r_{11} c_{21} \bar{\theta}_1}{\sqrt{c_{22} - r_{11}^2 c_{21}^2}} \quad , \tag{5.5}$$

and consequently

$$r_{22} = \frac{1}{\sqrt{c_{22} - r_{11}^2 c_{21}^2}} \tag{5.6}$$

$$r_{21} = - r_{22} r_{11}^2 c_{21} \quad . \tag{5.7}$$

In a similar manner

$$\bar{\theta}_3 = \frac{\theta_3 - (\theta_3,\bar{\theta}_1) \bar{\theta}_1 - (\theta_3,\bar{\theta}_2) \bar{\theta}_2}{\sqrt{(\theta_3,\theta_3) - (\theta_3,\bar{\theta}_1)^2 - (\theta_3,\bar{\theta}_2)^2}}$$

$$= \frac{\theta_3 - r_{11} c_{31} \bar{\theta}_1 - (r_{21}c_{31}+r_{22}c_{32})\bar{\theta}_2}{\sqrt{c_{33} - r_{11}^2 c_{31}^2 - (r_{21}c_{31} + r_{22}c_{32})^2}} \tag{5.8}$$

and

$$\bar{\theta}_N = \frac{\theta_N - r_{11}c_{N1}\bar{\theta}_1 - (r_{21}c_{N1}+ r_{22}c_{N2})\bar{\theta}_2 - \cdots - \left(\sum_{n=1}^{N-1} r_{(N-1)n} c_{Nn}\right)\bar{\theta}_{N-1}}{\sqrt{c_{NN} - r_{11}^2 c_{N1}^2 - (r_{21}c_{N1}+ r_{22}c_{N2})^2 - \cdots - \left(\sum_{n=1}^{N-1} r_{(N-1)n} c_{Nn}\right)^2}} \tag{5.9}$$

from which $r_{31}, r_{32}, r_{33}$ and $r_{N1}, r_{N2}, \cdots r_{NN}$ may be calculated. All indicated divisions are possible since the numerator functions and the square roots of their norms are non-zero.[1] Furthermore, $r_{11}$ depends on $c_{11}$; $r_{21}, r_{22}$ depends on $r_{11}, c_{21}, c_{22}$; $r_{31}, r_{32}, r_{33}$ depend on $r_{11}, r_{21}, r_{22}, c_{31}, c_{32}, c_{33}$; and so on. Thus the orthogonalization of any set of linearly independent approximating functions $\theta_1, \cdots \theta_N$ is always possible in terms of the constants $c_{nm}$ of the set.

It is apparent from inspection of the above equations that application of the Schmidt process is not simple. Numerical work grows

---

1 This is a direct result of the linear independence of $\theta_1, \cdots \theta_N$.

very rapidly with the number of functions, making machine calculation almost mandatory. Because of this and other considerations, which are discussed later, it is sometimes preferable to invert C.

### Orthogonalization of Exponential Functions

Approximation by exponential functions occurs, among other times, when the prescribed input is an impulse or a sum of exponential functions. Since such approximations are common it is worthwhile to consider separately the orthogonalization of linearly independent exponential functions. Fortunately, these functions admit mathematical techniques which are far more general and tractable than the Schmidt process.

Restrictions.-- The prescribed input is assumed to begin at t = 0 so the approximating functions $\theta_1$,----- $\theta_N$ are zero for negative time and are linear combinations of functions of the form

$$t^p e^{\sigma t} \quad , \quad t^p e^{\sigma t} \cos \omega t, \quad t^p e^{\sigma t} \sin \omega t$$

$$(p = 0 \text{ or an integer})$$

for positive time. Since the interval of expansion is semi-infinite, i.e. $T_1 = 0$ and $T_2 = \infty$, the integrals $\int_0^\infty W(t) \, \theta_n(t)^2 dt$ must be finite. This is insured if all $\sigma$ are negative.[1] That is, the prescribed input must be damped and $\Phi_1$,--------$\Phi_N$ stable.

The Laplace transforms of the approximating functions must correspondingly be linear combinations of functions of the form

$$\frac{1}{(s-\sigma)^{p+1}} \quad , \quad \frac{1}{(s-\sigma-j\omega)^{p+1}} \quad , \quad \frac{1}{(s-\sigma+j\omega)^{p+1}} \quad .$$

---

[1] If $W(t)$ becomes exponentially unbounded as $t \to \infty$ the $\sigma$ must be more than negative. They must be sufficiently negative to cause the integrand to vanish as $t \to \infty$.

As expected, the poles of these functions are real (at $s = \sigma$) or complex conjugate pairs (at $s = \sigma \pm j\omega$) and lie in the left half plane. When $p \neq 0$ they are multiple.

Last of all, the functions $\theta_1, \text{-----} \theta_N$ (or $\Theta_1, \text{-----} \Theta_N$) must be linearly independent. This requires the total number of poles (counting multiple poles according to their order) to be greater than or equal to N.

<u>Orthogonalization by Transformation</u>--Limited classes of orthonormal exponential functions can be generated from known sets of orthogonal functions by transformation of variables.

For example, suppose the functions $P_n(x)$ are known to be orthonormal with respect to $W_p(x)$ in the interval (a,b). Then,

$$\int_a^b W_p(x)P_n(x)P_m(x) \, dx = 1 \quad , \quad m = n$$
$$= 0 \quad , \quad m \neq n \quad . \tag{5.10}$$

By letting

$$x = g(t) \tag{5.11}$$

such that

$$a = g(0) \tag{5.12}$$

$$b = g(\infty) \tag{5.13}$$

equation (5.10) becomes

$$\int_0^\infty W_p[g(t)] \, P_n[g(t)] \, P_m[g(t)] \, g'(t) \, dt = 1, \quad m = n \tag{5.14}$$
$$= 0, \quad m \neq n \quad .$$

Thus, the functions $\bar{\theta}_n(t) = P_n[g(t)]$ are orthonormal with respect to $W(t) = W_p(g)g'$ in the interval $(0, \infty)$. Or if a unity weight is desired, $\bar{\theta}_n(t) = P_n(g)\sqrt{W_p(g) \, g'}$ .

Relatively few transformations of the above type are workable. $P_n(x)$ and $g(t)$ must be chosen so that the $\bar{\Theta}_n(t)$ are of the proper exponential form. Kautz [17] has used $g(t) = 1 - 2e^{-2t}$ transform the x interval $(-1, 1)$. Of the various known functions orthogonal in $(-1,1)$ he took the Legendre and Tchebycheff polynomials. The resulting $\bar{\Theta}_n$ functions had poles at $s = -(2n-1)$. Other possible transformations would include $e^{-t} = \cos x$ (orthonormal cosine functions in $(0, \frac{\pi}{2})$) and $e^{-t} = x$ (orthogonal even Legendre functions in $(0, 1)$). Again, these transformations produce $\bar{\Theta}_n$ functions with poles spaced evenly along the negative real axis.[1]

An alternative approach to the WME approximation is proposed by Papoulis [31]. Instead of transforming the approximating functions to the variable t, he transforms the prescribed response to the variable x. Although the method is different the final approximations are much the same, since identical $P_n$ functions and transformations yield identical poles in $F_o{}^*(s)$.

Further elaboration on transformation procedures is possible; other types of orthogonal functions can be transformed. However, the results of the remaining part of the chapter are sufficiently general to make such investigation of little value.

Residue Statement of Orthonormal Condition--The orthogonalization of exponential functions by application of the theory of residues has many advantages. To develop this method the Laplace transform equivalent

---

1 Transformations such as $x = e^{-t}\cos t$ generate complex poles but are rather inefficient because relatively few orthonormal functions require many poles.

of the orthonormal condition must be considered. In what follows $W(t) = 1$.

Hence,

$$\bar{c}_{nm} = (\bar{\Theta}_n, \bar{\Theta}_m) = \int_0^\infty \bar{\Theta}_n \bar{\Theta}_m \, dt = 1 , \quad m = n$$

$$= 0 , \quad m \neq n . \tag{5.15}$$

By means of complex convolution integral[1] of the Laplace transform, equation (5.15) may be expressed by[2]

$$\bar{c}_{nm} = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \bar{\Theta}_n(s) \, \bar{\Theta}_m(-s) ds = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \bar{\Theta}_n(-s) \, \bar{\Theta}_m(s) \, ds \tag{5.16}$$

where c must be taken as zero if the integrand is always to be analytic on the path of integration.

The assumptions concerning $\Theta_n(s)$ assert that $\Theta_n(s) \rightarrow \dfrac{const.}{s^{p+1}}$ as $|s| \rightarrow \infty$ where $p \geq 0$. Therefore, $\bar{\Theta}_n(s) \, \bar{\Theta}_m(-s)$ must go to zero at least as rapidly as $\dfrac{const.}{s^2}$ as $|s| \rightarrow \infty$. Hence, equation (5.16) may be evaluated along the closed contour which traverses the $j\omega$ axis and encircles the entire right half plane at infinity. For the contour $C_0$ shown in Figure 5.1 $(R \rightarrow \infty)$ $\bar{c}_{mn}$ becomes simply

$$\bar{c}_{mn} = -\frac{1}{2\pi j} \int_{C_0} \bar{\Theta}_n(s) \, \bar{\Theta}_m(-s) \, ds , \tag{5.17}$$



Figure 5.1 The $C_0$ Contour in the s Plane

1 See Gardner and Barnes [10], page 275.

2 Note that $\displaystyle\int_0^\infty e^{-st} \bar{\Theta}_n(t) \, \bar{}(t) dt = \bar{c}_{nm}$ at $s = 0$.

since the contribution to the integral along the infinite semi-circle
is zero.

Finally, because $\overline{\Theta}_n(s) \ \overline{\Theta}_m(-s)$ is analytic along $C_0$, the residue
theorem may be applied to yield the desired result, $\overline{c}_{nm} = -\sum$ residues
of $\overline{\Theta}_n(s) \ \overline{\Theta}_m(-s)$ in the right half plane

$$= 1, \quad m = n$$
$$= 0, \quad m \neq n \ .$$

$$(5.18)$$

Orthogonalization by Means of the Residue Condition--The
steps taken in the orthogonalization of exponential functions by
application of the residue condition are similar to those of the
Schmidt process. They are: (1) the specification of the linearly
independent functions $\Theta_1, ----- \ \Theta_N$, i.e. the specification of the poles
of $\overline{\Theta}_1, ----- \ \overline{\Theta}_N$; (2) the linear transformation of these functions into
the orthonormal set. As before, complete repetition of the orthogon-
alization process for different N is avoided if $\overline{\Theta}_1$ depends only on $\Theta_1$,
$\overline{\Theta}_2$ depends only on $\Theta_1$ and $\Theta_2$, and so on. Or in terms of the frequency
domain functions, $\overline{\Theta}_1$ should contain only the first pole, $\overline{\Theta}_2$ should
contain only the first and second poles,......, and $\overline{\Theta}_N$ should contain
all the poles.

To begin, assume that the given poles all lie on the negative
real axis at s $=-\alpha_1, -\alpha_2, ----- \ -\alpha_N$ (the $\alpha_n$ are positive). The first
function is then

$$\overline{\Theta}_1 = \frac{K_1}{s+\alpha_1} \ , \qquad (5.19)$$

where $K_1$ is an arbitrary constant. Application of equation (5.18)
for n = m = 1 gives

$$-\sum \text{ residues of } \frac{K_1^2}{(s+\alpha_1)(-s+\alpha_1)} \text{ in the right half plane } = + \frac{K_1^2}{2\alpha_1} = 1.$$

$$(5.20)$$

Thus,

$$\bar{\Theta}_1 = \frac{\sqrt{2\alpha_1}}{s+\alpha_1} , \qquad (5.21)$$

The higher order functions require that equation (5.18) holds when $m \neq n$ (really $m < n$ since $\bar{c}_{mn} = \bar{c}_{nm}$ ). This is assured if $\bar{\Theta}_n(s)$ is constructed so $\bar{\Theta}_n(s) \bar{\Theta}_m(-s)$ has no poles in the right half plane for all $m$ less than $n$. The obvious step is to choose the zeros of $\bar{\Theta}_n(s)$ so they cancel the right half plane poles of $\bar{\Theta}_m(-s)$. But these poles occur at $s = -\alpha_1, ---- -\alpha_{n-1}$. Thus equation (5.18) is satisfied for $m \neq n$ when the zeros of $\bar{\Theta}_n(s)$ are taken at $s = \alpha_1, \alpha_2, ----- \alpha_{n-1}$.

Applying this result to $\bar{\Theta}_2$ yields

$$\bar{\Theta}_2 = \frac{K_2(s-\alpha_1)}{(s+\alpha_1)(s+\alpha_2)} , \qquad (5.22)$$

where $K_2$ is given by

$$-\sum \text{residues of } \bar{\Theta}_2(s)\bar{\Theta}_2(-s) \text{ in the right half plane } = \frac{K_2^2}{2\alpha_2} = 1 . \qquad (5.23)$$

Therefore,

$$\bar{\Theta}_2 = \frac{\sqrt{2\alpha_2} \ (s-\alpha_1)}{(s+\alpha_1)(s+\alpha_2)} . \qquad (5.24)$$

Extension of the process to higher order functions is clear, leading to the formula

$$\bar{\Theta}_n = \frac{\sqrt{2\alpha_n} \ (s-\alpha_1)(s-\alpha_2)-------(s-\alpha_{n-1})}{(s+\alpha_1)(s+\alpha_2)-----------(s+\alpha_n)} . \qquad (5.25)$$

When both real and complex poles are specified the same orthogonalization scheme is valid, although the formulas are somewhat more complicated. If the real poles are given by $s = -\alpha_n$ and the complex poles by $s = -\alpha_n \pm j \sqrt{\beta_n^2 - \alpha_n^2}$ ($-\alpha_n$ is the real part, $\beta_n$ is the distance from the origin); then the orthonormal functions will be proven to be

$$\bar{\Theta}_n = \Lambda\ (s)_{n-1}\ \frac{\sqrt{2\alpha_n}}{s+\alpha_n} \qquad\qquad (5.26)$$

$$\bar{\Theta}_n = \Lambda\ (s)_{n-1}\ \frac{\sqrt{2\alpha_n}(s+\beta_n)}{s^2+2\alpha_n s +\beta_n^2} \qquad\qquad (5.27)$$

$$\bar{\Theta}_{n+1} = \Lambda\ (s)_{n-1}\ \frac{\sqrt{2\alpha_n}(s-\beta_n)}{s^2+2\alpha_n s +\beta_n^2} \qquad . \qquad [1] \qquad (5.28)$$

Equation (5.26) holds if the nth pole is real, and equation (5.27) and equation (5.28) hold if the nth and (n+1)th poles are complex conjugates.[2]

$\Lambda_n(s)$ is the all-pass function of magnitude one which has its poles the first n poles of specified set, i.e.

$$\Lambda_n(s)\ =\ \frac{(s+s_1)(s+s_2)\ \text{-------}(s+s_n)}{(s-s_1)(s-s_2)\ \text{-------}(s-s_n)} \qquad . \qquad (5.29)$$

For example, the real pole functions of the above paragraph yield

$$\Lambda_n(s)\ =\ \frac{(s-\alpha_1)(s-\alpha_2)\ \text{-------}(s-\alpha_n)}{(s+\alpha_1)(s+\alpha_2)\ \text{-------}(s+\alpha_n)} \qquad .$$

To show that these functions all satisfy equation (5.18) it is only necessary to note the following: (1) $\Lambda(s)_n\ \Lambda(-s)_n\ = 1$, (2) the sum of the residues of $\dfrac{2\alpha_n}{(s+\alpha_n)(-s+\alpha_n)}$ or $\dfrac{2\alpha_n(s+\beta_n)(-s+\beta_n)}{(s^2+2\alpha_n s+\beta_n^2)(s^2-2\alpha_n s+\beta_n^2)}$ in the right half plane is -1, (3) $\Lambda_{n-1}(s)\bar{\Theta}_m(s)$ has no poles in the right half plane for m < n, (4) the sum of the residues of $\dfrac{2\alpha_n(s-\beta_n)(-s+\beta_n)}{(s^2+2\alpha_n s+\beta_n^2)(s^2-2\alpha_n s+\beta_n^2)}$ in the right half plane is zero.

Unfortunately, the determination of the corresponding time functions $\bar{\Theta}_1, \text{-----}\ \bar{\Theta}_N$ is not as simple. Each Laplace function

---

1 A somewhat more restricted version of these equations is given by Kautz [17].

2 Two equations are needed for complex poles if the corresponding time functions $\bar{\Theta}_n$ and $\bar{\Theta}_{n+1}$ are to be real

$\overline{\Theta}_n(s)$ must be expanded in a partial fraction series of $n$ terms, where each term has as an inverse transform an exponential time function. Thus, for first order real poles[1]

$$\overline{\Theta}_1 = r_{11} e^{-\alpha_1 t}$$

$$\overline{\Theta}_2 = r_{21} e^{-\alpha_1 t} + r_{22} e^{-\alpha_2 t}$$

$$\vdots$$

$$\overline{\Theta}_N = \sum_{n=1}^{N} r_{Nn} e^{-\alpha_n t} \quad , \quad (5.30)$$

and first order complex poles

$$\overline{\Theta}_1 = r_{11} e^{(-\alpha_1 + j\omega_1)t} + r_{12} e^{(-\alpha_1 - j\omega_1)t}$$

$$\overline{\Theta}_2 = r_{21} e^{(-\alpha_1 + j\omega_1)t} + r_{22} e^{(-\alpha_1 - j\omega_1)t}$$

$$\overline{\Theta}_3 = r_{31} e^{(-\alpha_1 + j\omega_1)t} + r_{32} e^{(-\alpha_1 - j\omega_1)t} + r_{33} e^{(-\alpha_3 + j\omega_3)t}$$

$$+ r_{34} e^{(-\alpha_3 - j\omega_3)t}$$

$$\vdots$$

$$(5.31)$$

With multiple poles the equations are the same except some of the exponential functions include multiplicative integral power of $t$.

The orthonormal functions produced by multiple poles at $s = -1$ are of particular interest. In this case,

---

[1] The constants $r_{mn}$ are tabulated by Kautz [17] for several real and complex pole sets.

$$\overline{\Theta}_n(s) \;=\; \sqrt{2}\; \frac{(s-1)^{n-1}}{(s+1)^n} \tag{5.32}$$

and the corresponding time functions are the well known orthonormal Laguerre functions

$$\overline{\Theta}_n(t) \;=\; \sqrt{2}\; e^{-t}\, L_{n-1}(2t) \qquad\cdot \tag{5.33}$$

The polynomials $L_n(x)$ are the Laguerre polynomials of order $n$ given by

$$L_n(x) \;=\; \frac{e^x}{n!}\, \frac{d^n}{dx^n}\, (x^n e^{-x}) \tag{5.34}$$

Historically, Laguerre functions were used by Lee [22] in the first published report on mean square approximation of linear system response. They also were discussed and applied by Wiener, and subsequently, by many others.

It should be noted that the functions $\overline{\Theta}_n$ are limited in their characteristics. They all act as $\dfrac{const.}{s}$ for large $s$. Hence, prescribed fall off rates on $F_o^{*}(s)$ different than $\dfrac{const.}{s}$ must be obtained by constrained approximation. An alternative approach is discussed in the next subsection where formulas are given for constrained orthonormal exponential functions.

One modification of the orthonormal function formulas is easily made at this time. The complex pole equations (5.27) and (5.28) may be replaced by

$$\overline{\Theta}_n \;=\; \Lambda_{n-1}(s)\; \frac{2\,\beta_n\,\sqrt{\alpha_n}}{s^2 + 2\alpha_n s + \beta_n^2} \tag{5.35}$$

$$\overline{\Theta}_{n+1} \;=\; \Lambda_{n-1}(s)\; \frac{2\,\sqrt{\alpha_n}\,s}{s^2 + 2\alpha_n s + \beta_n^2} \quad , \tag{5.36}$$

a substitution which is easily verified through equation (5.18). The simplicity of these expressions is useful when the functions must be realized by a physical system. Note that $\bar{\Theta}_n$ in equation (5.35) differs from the members of the earlier set in that $s^2\bar{\Theta}_n \to$ const. as $|s| \to \infty$.

Orthogonalization of Constrained Functions by Means of the Residue Condition.--Constrained exponential functions can be orthogonalized by techniques which are similar to those just developed. As a beginning, consider the functions which act as $\dfrac{\text{const.}}{s^K}$ as $|s| \to \infty$ (K-1 constraints).

Applying the condition that $\bar{\Theta}_n(s)\bar{\Theta}_m(-s)$ has no poles in the right half plane for $m < n$ yields

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \frac{K_n}{(s-s_{n1})(s-s_{n2})\text{-----}(s-s_{nK})} \qquad (5.37)$$

where $\Lambda_{n-1}$ is again an all pass function containing the poles of the first n-1 functions. The normalization constant depends on the type and number of poles. If K=2, real and complex poles give respectively

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \frac{\sqrt{2\alpha_{n1}\alpha_{n2}(\alpha_{n1} + \alpha_{n2})}}{(s+\alpha_{n1})(s+\alpha_{n2})} \qquad (5.38)$$

and

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \frac{2\sqrt{\alpha_n\beta_n}}{s^2 + 2\alpha_n s + \beta_n^2} . \qquad (5.39)$$

Since $\bar{\Theta}_n(s)$ must contain as many zeros as $\bar{\Theta}_{n-1}(s)$ has poles, it is clear that $\bar{\Theta}_n$ has K more poles than $\bar{\Theta}_{n-1}$. Thus, a set of N constrained orthonormal functions has a total of KN poles. Chapter IV indicates, however, that N + K - 1 poles (i.e. one additional independent function for each constraint) are sufficient. Consequently, equation (5.37) is rather extravagant in the use of poles. It is shown

later that this is a price which must be paid if the orthonormal

constrained functions are to be defined by a simple formula.

Arbitrary zeros are generated in the orthonormal functions

by letting

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \frac{K_n(s-\bar{s}_{n1})(s-\bar{s}_{n2})-------(s-\bar{s}_{nJ})}{(s-s_{n1})(s-s_{n2})---------(s-s_{nK})} ,$$

$$(5.40)$$

where $J < K$ so that $\bar{\Theta}_n(s) \to 0$ as $|s| \to \infty$. Each zero in equation

(5.40) is chosen to achieve a desired homogeneous constraint condition,

such as $\bar{\Theta}_n(t_1) = 0$ or $\bar{\Theta}_n(s_1) = 0$. To illustrate, consider the case

where $\bar{\Theta}_n(1) = 0$. Then,

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \frac{K_n(s-1)}{(s-s_{n1})(s-s_{n2})} .$$

$$(5.41)$$

If the poles are real and $\alpha_{n1} \neq \alpha_{n2}$

$$(5.42)$$

$$\bar{\Theta}_n = \Lambda_{n-1}(s) \sqrt{\frac{2\alpha_{n1}\alpha_{n2}(\alpha_{n1}+\alpha_{n2})(\alpha_{n1}-\alpha_{n2})}{\alpha_{n1}\alpha_{n2}^2 - \alpha_{n2}\alpha_{n1}^2 - \alpha_{n1}+\alpha_{n2}}} \frac{s-1}{(s+\alpha_{n1})(s+\alpha_{n2})} .$$

Although the normalization constants for the constrained functions

appear complicated, they are readily derived and present no real

difficulty. As before the equations are wasteful of poles since N

functions and K-1 constraints require KN poles.

The inefficiency of pole utilization in the above functions

is a direct result of the simplified orthogonalization procedure em-

ployed. The functions are constructed so that no poles of $\bar{\Theta}_n(s) \bar{\Theta}_m(-s)$

lie in the right half plane for $m \leq n$. This condition is a good deal

less general than that stated by equation (5.18) which applies to the

sum of residues in the right half plane.

To expand these remarks consider the residue orthogonalization of the functions $\frac{1}{(s+1)^2}$ , $\frac{1}{(s+1)^3}$ , etc.; that is, an orthogonalization under the constraint that $\lim_{|s| \to \infty} s \, \overline{\Theta}_n(s) = 0$. According to the theory of Chapter IV, N orthogonal functions should require (N+1) poles. The first function is simply

$$\overline{\Theta}_1 = \frac{K_1}{(s+1)^2} . \tag{5.43}$$

The second function has three poles and one zero $\overline{s}_{21}$ which must be chosen so

$$\overline{\Theta}_2(s)\overline{\Theta}_1(-s) = \frac{K_1 K_2 \,(s-\overline{s}_{21})}{(-s+1)^2(s+1)^3} \tag{5.44}$$

has a zero sum of residues in the right half plane. Thus

$$\overline{\Theta}_2 = \frac{K_2 \,(s-\frac{1}{3})}{(s+1)^3} . \tag{5.45}$$

Continuing, $\overline{\Theta}_3$ has four poles and two zeros, $\overline{s}_{31}$ and $\overline{s}_{32}$, and must satisfy equation (5.18) with n=3 and m=1,2. Solution of simultaneous algebraic equations in $\overline{s}_{31}$ and $\overline{s}_{32}$ gives

$$\overline{\Theta}_3 = \frac{K_3\left(s^2 -\frac{2}{3} s+\frac{1}{3}\right)}{(s+1)^4} . \tag{5.46}$$

Derivation of the higher order functions proceeds in a like manner; however, the numerical work multiplies rapidly, forcing a practical end to the calculations. Certainly, no simple formula is available for expressing the Nth function. The orthogonalization of constrained functions which have a minimum number of poles must depend, therefore, on a general orthogonalization method like the Schmidt process.

Summary--In conclusion, the available methods for orthogonalizing linearly independent exponential functions possess considerable variety, simplicity, and generality. Orthogonalization by transformation generates various exponential sets with different weight factors. The orthonormal residue condition allows arbitrary pole functions of both the unconstrained and constrained types to be written in simple Laplace transform formulas. All of these factors do much to make approximation by exponential functions practicable.

On the negative side, the transformation method is limited by the number of known, applicable transformations and orthonormal function sets. Furthermore, the residue theory fails when applied to constrained functions having a minimum number of poles or to weight functions $W(t)$ which are not a constant. Difficulties of this sort are to be expected in a process as complicated as orthogonalization. Part two will discuss various ways for minimizing some of these deficiencies in practical applications.

# VI.  CHOICE OF APPROXIMATING FUNCTIONS

Earlier chapters have assumed the predetermination of the approximating functions $\Theta_1$, - - - $\Theta_N$.  It is the purpose of this chapter to investigate the factors influencing the choice of these functions, or more explicitly, the choice of $\varphi_1$, - - - $\varphi_N$ or the poles of $H^*(s)$. Although an optimum solution of this pole location problem exists, it is numerically impractical and is usually replaced by less exact and more feasible techniques.  The following sections discuss, therefore, not only the pole optimization equations and their solution, but also the criticalness of pole locations and several of the more important approximate methods for pole determination.

## Optimum Choice of Approximating Functions

Two sets of parameters characterize the approximate linear system and affect the weighted mean square error I.  They are the approximation coefficients $a_1$, - - - $a_N$ and the poles $s_1$, - - - $s_M$.[1] Chapter III describes the minimization of I with respect to $a_1$, - - $a_N$. If optimum (i.e., least WME) poles positions are also to be found, additional partial derivatives of I must be set equal to zero.  Repeating equation (4.7)

$$I = I_{max} - 2 \lfloor a \rfloor d] + \lfloor a \rfloor C a]  \qquad (6.1)$$

and setting $\dfrac{\partial I}{\partial a_n} = 0$ for n = 1, - - - N yields the N equations obtained before,

$$d] = C a] \qquad . \qquad (6.2)$$

---

[1]  Note M > N.  See Chapter II.

But I also depends on d] and C which in turn depend on $\Theta_1, - - - \Theta_N$ and hence on $s_1, - - - s_M$. Thus setting $\frac{\partial I}{\partial s_m} = 0$ for m $=$ 1, - - - M yields the further equations

$$\frac{\partial I}{\partial s_m} = - 2 \underline{a} \left\{ \frac{\partial}{\partial s_m} d] \right\} + \underline{a} \left\{ \frac{\partial}{\partial s_m} C \right\} a] = 0.$$

$$(6.3)$$

Simultaneous solution of equations (6.2) and (6.3) for $a_1, - - - a_N$ and $s_1, - - - s_M$ gives the desired optimum.[1]

Difficulty in obtaining an actual numerical solution of these equations arises because $d_n$ and $c_{nm}$ are complicated transcendental functions of $s_1, - - - s_M$.[2] Derivation and trial and error solution of such equations for N > 2 becomes a formidable task, even when automatic computing machines are available. An added complication is the possibility of multiple stationary points, each of which must be determined if the true optimum is to be established.

Several alternative procedures for pole optimization exist. An iterative numerical method described by Aigrain and Williams [2] makes use of an intermediate approximation by a large number of Laguerre functions. It is, however, only applicable to the approximation of impulsive responses having known Laplace transforms. Another approach is trial and error approximation by orthonormal exponential functions. Here $I = I_{max} - \sum_{n=1}^{N} a_n^2$ is calculated repeatedly for different combinations of pole locations until the optimum solution is found. Later

---

[1] A similar set of equations written in Laplace transform variables was derived by Aigrain and Williams [1] for limited types of input functions.

[2] An exception occurs when $W(t)$ and $f_o(t)$ are exponential and $T_1 = 0$, $T_2 = \infty$. Then $d_n$ and $c_{nm}$ are rational functions of $s_1, - - - s_M$. The equations, however, are still difficult to solve.

on, it will be shown how this method can be practically implemented with a repetitive analog computer.

## Approximate Choice of Approximating Functions

Since none of the above solutions are particularly straight-forward it is important to question the criticalness of pole location. Does an approximate solution of the pole optimization problem appreciably alter approximation accuracy? For numerous practical examples[1] the answer is fortunately no. Pole positions can be changed rather drastically in certain directions with little effect on system response (provided, of course, that the approximation coefficients are always chosen for minimum error). However, in some cases the tolerance to pole postion shifts may be poor or the approximate poles may be poorly chosen. In these cases, the additional error can always be corrected by the addition of more terms (i.e., more poles) in the approximating series.

The Prony Method--The Prony method [6,28] is one of the better techniques for approximate choice of approximating functions. Though limited to exponential functions, it has considerable generality and good accuracy. It is based on the determination of a linear differential equation with constant coefficients of order M which best fits the pre-scribed response. The M linearly independent exponential functions which are solutions of the equation are used as source functions for $\theta_1, ----- \theta_N$. "Best fit" means that the coefficients $B_1, ----- B_M$ in

---

1 See the examples given by Kautz [17].

$$B_M f_o^{(M)}(t) + \cdots - B_1 f_o^{(1)}(t) + f_o(t) = e_o(t) \overset{\sim}{=} 0 \quad {}^1 \qquad (6.4)$$

are chosen so the WME

$$I_o = \int_{T_1}^{T_2} W_o(t) e_o^{\,2}(t) dt = \int_{T_1}^{T_2} W_o [\sum_{m=1}^{M} B_m f_o^{(m)}(t) + f_o(t)]^2 \, dt \qquad (6.5)$$

is a minimum. The exponential functions are the M linearly independent

solutions of the equation

$$B_M f^{(M)}(t) + \cdots - B_1 f^{(1)}(t) + f(t) = 0. \qquad (6.6)$$

Hence, the poles of the approximation $F_o^*(s)$ are the roots of the

equation

$$B_M s^M + \cdots - B_1 s + 1 = 0. \qquad (6.7)$$

It should be noted that the poles produced by the Prony me-

thod are in no way related to the poles defined by the optimization

equations of the previous section. Nothing is said relative to the co-

efficients of each exponential function; only the exponential decay

constants are considered. Of course, the poles determined will be

exact if $f_o(t)$ is itself a sum of M exponential functions. The weight

factor $W_o(t)$ has the same function as previously discussed weight

factors. It emphasizes the undesirability of error as a function of

time.

The constants $B_1$, $\cdots B_M$ are found by setting $\dfrac{\partial I_o}{\partial B_M} = 0$ for

$m = 1$, $\cdots M$. Defining

---

[1] $\quad \dfrac{d^M f_o(t)}{dt^M} = f_o^{(M)}(t)$

$$f_{kj} = \int_{T_1}^{T_2} W_o f_o{}^{(k)} f_o{}^{(j)} dt \qquad (6.8)$$

equation (6.5) becomes

$$I_o = \sum_{j=1}^{M} \sum_{k=1}^{M} B_j B_k f_{jk} + 2 \sum_{k=1}^{M} B_k f_{k0} + f_{00} \qquad (6.9)$$

and

$$\frac{\partial I_o}{\partial B_m} = 0 = 2 B_1 f_{m1} + 2 B_2 f_{m2} + - - - 2 B_m f_{mM} + 2 f_{m0}, \qquad (6.10)$$

$$m \quad 1, ------M .$$

Or in matrix notation equations (6.10) are simply

$$\begin{bmatrix} f_{11} --------- f_{1M} \\ \\ \\ f_{M1} --------- f_{MM} \end{bmatrix} \begin{bmatrix} B_1 \\ \\ \\ B_M \end{bmatrix} = - \begin{bmatrix} f_{10} \\ \\ \\ f_{M0} \end{bmatrix} . \qquad (6.11)$$

After the constants $B_1$, - - - $B_M$ are obtained by solution[1] of equation

(6.11) they are substituted in equation (6.7) which is then solved for

the M roots which are the approximation poles.

A number of observations concerning the Prony method can be

made:

   1.  It requires complicated calculations (the inversion of an

       Mth order matrix and the solution of an Mth degree

---

[1] Equation (6.11) has a solution. If the determinant of the matrix is zero then the functions $W^{1/2} f_o{}^{(M)}$, - - - $W^{1/2} f_o{}^{(1)}$ are not linearly independent and $f_o{}^{(M)}$ equals a linear combination of $f_o{}^{(M-1)}$, - - - $f_o{}^{(1)}$. But this means that $f_o$ satisfies a linear differential equation with constant coefficients of order less than M. A lower order matrix will thus yield the coefficients $B_m$ which define the exponential functions of which $f_o$ must be composed.

algebraic equation). The calculations are, however,
easily made by a high speed digital computer.

2. The poles obtained are not necessarily in the left half
plane. If unstable poles are near the imaginary axis,
they may be shifted into the left half plane with little
added error. If they are deep inside the right half
plane, the entire procedure must be repeated with more
favorable $W_o$ or $(T_1, T_2)$.

3. The constants $f_{kj}$ can be computed even if $f_o$ is not ex-
pressed in analytic form. This means the Prony method
is applicable when $f_o$ is given as experimental data.

4. Equation (6.4) may be integrated an arbitrary number of
times before minimization of $I_o$. This weights more
heavily the undesirability of approximation errors in
the low frequency region and allows some or all of the
derivatives of $f_o$ to be replaced by integrals of $f_o$.
The latter change is particularly important when the deri-
vatives of $f_o$ are difficult to obtain. This may be the
case when $f_o$ is obtained from experimental measurements.

The Prony method has been discussed in some detail because of
its close relation to WME approximation and generality. One other me-
thod which is perhaps less exact but simpler to apply will be des-
cribed now.

The Preliminary Approximation Method--In this method a pre-
liminary approximation is made by an approximation procedure different
from the WME approach, and the poles of the preliminary approximation
are used as the poles of the WME approximation. While many types of

preliminary approximations may be used, those which have reasonable

simplicity are preferred.  Depending on accuracy and computational

complexity, the approximation may require varying degrees of effort.

Kautz [17] , for example, discusses an accurate preliminary

approximation which employs a continued fraction expansion.  In his

method $F_0(s)$ is expanded in a power series about a point[2] in the com-

plex plane, and the leading terms of this power series are then dupli-

cated by a continued fraction expansion.  The poles of the terminated

expansion become, finally, the poles of the WME approximation.  The

overall work is somewhat less than that required by the Prony method

and is better suited to hand calculations.  While continued fraction

expansions have certain advantages (among them control of approximation

error in the time domain), other kinds of preliminary approximations

of similar complexity are undoubtedly as good.  Their study is, how-

ever, beyond the scope of this investigation.

A less precise method of preliminary approximation is to take

the poles of well known approximations which have the same character-

istics as the prescribed response.  For example, a WME low-pass filter

might use the poles of a Butterworth or Tchebycheff low-pass filter.

Or a prescribed weighting function $h(t)$ which is similar to an expo-

nentially damped periodic function might be approximated by a sum of

exponentially damped harmonic sine and cosine functions.  Though em-

pirical, this approach is simple and often very successful.  To aid

its application, several sets of functions with different approximation

---

[1]  See the references given in Chapter I.

[2]  The point is chosen to make approximation errors small in certain
time regions.

qualities are tabulated in the appendix. Their use is illustrated in the next chapter.

From the preceding discussion it is clear that the method of preliminary approximation does not lend itself to a brief and specific description. Further elaboration is therefore limited to the examples of part two.

To review, the optimum determination of approximating functions is very difficult and must usually be replaced by approximate methods. The Prony method and preliminary approximation method often offer satisfactory solutions when $f_o$ is to be approximated by a sum of exponential functions. In other cases little can be done to solve what is an extremely complicated problem, and good approximation accuracy must be achieved by the expedient of including many terms in the approximating series.

PART TWO

APPLICATIONS

## VII. IMPULSE RESPONSE APPROXIMATION

Since a large fraction of approximation problems are stated by prescribing the impulse response, it is desirable to consider separately and in detail the factors relevant to WME approximation of impulse responses. As noted earlier, the impulse response approximation problem is a special case of the more general problem in which $f_i(t) = u_o(t)$, $f_o(t) = h(t)$, and $\theta_n(t) = \varphi_n(t)$ $(n = 1, - - - N)$ and in which the approximating functions are linear combinations of exponential functions. This restriction on the form of the approximating functions permits not only the development of the orthonormal function formulas of Chapter V, but also the approximation procedures which follow. The description of these procedures is the purpose of this chapter and divides conveniently into two sections: (1) the analytic approximation of $h(t)$ $(f_o(t))$ from its Laplace transform $H(s)$ $(F_o(s))$, and (2) the analog computer approximation and simulation of $h(t)$. To illustrate further and to apply the results of these sections, a concluding section presents several examples.

### Laplace Transform Approximation Procedures

Since the approximating functions are linear combinations of exponential functions, it is possible to write

$$\theta_n = \sum_{m=1}^{M} r_{nm}\, e^{s_m t}$$

or equivalently

$$\theta] = R \begin{bmatrix} e^{s_1 t} \\ \vert \\ \vert \\ e^{s_M t} \end{bmatrix} \qquad (7.1)$$

where $s_1, ----- s_M$ are the poles of the approximation. As indicated be-fore, the determination of the approximation coefficients a] requires the calculation of the matrices C and d], and the inversion of C. Ex-pressing C in terms of equation (7.1) gives (assume for the time being that $W(t) = 1$)

$$C = \int_0^\infty \Theta] \underline{\Theta} \; dt = R \left\{ \int_0^\infty \begin{bmatrix} e^{s_1 t} \\ \vdots \\ e^{s_M t} \end{bmatrix} \begin{bmatrix} e^{s_1 t} & - - - & e^{s_M t} \end{bmatrix} dt \right\} R_T$$

$$= R \begin{bmatrix} \dfrac{1}{s_1 + s_1} & -------- & \dfrac{1}{s_1 + s_M} \\ & & \\ & & \\ \dfrac{1}{s_M + s_1} & -------- & \dfrac{1}{s_M + s_M} \end{bmatrix} R_T$$

$$(7.2)$$

Usually, the above formula is simplified by taking $R = I$. If ortho-normal functions are employed, equation (7.2) is not required since $C = I$.

In a similar manner,

$$d] = \int_0^\infty \dot{f}_0(t) \; \Theta] \; dt = R \int_0^\infty f_0(t) \begin{bmatrix} e^{s_1 t} \\ \vdots \\ e^{s_M t} \end{bmatrix} dt \qquad (7.3)$$

But $\int_0^\infty f_0(t) \; e^{s_m t} \; dt = F_0(-s_m)$, the Laplace transform of the prescribed response at $s = -s_m$. Hence,

$$d] = R \begin{bmatrix} F_0(-s_1) \\ \vdots \\ F_0(-s_M) \end{bmatrix} \qquad (7.4)$$

Because $F_O(s)$ is known to be regular in the right half plane, the numbers $F_O(-s_m)$ exist and yield by equation (7.4) the coefficients $d_1, - - - d_N$.

From the preceding paragraphs it is clear that there are two alternative approaches for analytic approximation of $f_O(t)$ from $F_O(s)$:

1. The orthogonal function approach. By applying the techniques of Chapter V, C = I, and the calculation and inversion of C is avoided. The determination of d] = a] must be preceded, however, by the evaluation of R which in turn requires a partial fraction expansion of each orthonormal function $\Theta_1, - - - \Theta_N$.[1]

2. The exponential function approach. The matrix R is made equal to the identity matrix I so that d] is simply

$$
\begin{bmatrix} F_O(-s_1) \\ \vdots \\ F_O(-s_M) \end{bmatrix} \quad \text{and} \quad a] = C^{-1} \begin{bmatrix} F_O(-s_1) \\ \vdots \\ F_O(-s_M) \end{bmatrix}. \quad \text{The}
$$

necessary calculation and inversion of C complicates this procedure.[2]

_____

[1] This is substantially the procedure proposed by Kautz (17).

[2] For a given set of poles the calculation of $C^{-1}$ in approach 2 can be made from R in approach 1. The orthonormal functions give

$$
f_O^* = \underline{a}, \Theta] = \underline{d}, \Theta] = \underline{F_O(-s_1) - - F_O(-s_M)}, R_T R \begin{bmatrix} e^{s_1 t} \\ \vdots \\ e^{s_M t} \end{bmatrix}
$$

while the exponential functions yield

$$
f_O^* = \underline{a}, \Theta] = \underline{d}, C^{-1} \Theta] = \underline{F_O(-s_1) - - F_O(-s_M)}, C^{-1} \begin{bmatrix} e^{s_1 t} \\ \vdots \\ e^{s_M t} \end{bmatrix}.
$$

By comparing these equations it is seen that $C^{-1} = R_T R$ .

Little difference in overall work exists between the two approaches if all the approximation details are worked out from the beginning. However, when the inverse matrix $C^{-1}$ is known prior to approximation, the second approach is especially convenient to use. To make such a technique possible, a tabulation of inverse matrices for certain pole combinations is given in the appendix. Slight variations in the above formulas are made to avoid the introduction of the complex time functions $e^{s_m t}$ and to simplify the form of approximations. A set of constrained approximating functions is also included. For these and other details see the appendix and the examples at the end of this chapter.

Although the above description of mean square approximation by means of the Laplace transform has been limited to simple approximations, the constraint equations and constrained functions of Chapters IV and V are also applicable. The only difference lies in the evaluation of the coefficients $d_1$, - - - $d_N$ by means of equation (7.4).

When the weight factor $W(t)$ is not a constant the theory is not as simply extended since the integral

$$\int_0^\infty W(t)\ e^{s_m t}\ f_o(t)\ dt$$

cannot, in general, be interpreted in terms of $F_o(s)$. Furthermore, as mentioned in Chapter V, the residue orthogonalization of exponential functions with respect to an arbitrary $W(t)$ is not practical.

However, one particular type of weight factor does permit the application of the Laplace transform procedure to non-orthogonal approximating functions. It is defined by the sum of exponential functions

$$W(t) = \sum_{k=1}^{K} w_k \, e^{\tilde{s}_k t} \qquad \cdot \quad 1 \qquad\qquad (7.5)$$

With this weight factor

$$C = \sum_{k=1}^{K} w_k \int_0^{\infty} e^{\tilde{s}_k t} \; \Theta] \;\; \underline{\Theta} \;\; dt = \sum_{k=1}^{K} w_k \quad R \begin{bmatrix} \dfrac{1}{s_1 + s_1 + \tilde{s}_k} & ---- & \dfrac{1}{s_1 + s_M + \tilde{s}_k} \\ \vdots & & \vdots \\ \dfrac{1}{s_M + s_1 + \tilde{s}_k} & ---- & \dfrac{1}{s_M + s_M + \tilde{s}_k} \end{bmatrix} R_T$$

and $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (7.6)$

$$d] = \sum_{k=1}^{K} w_k \int_0^{\infty} e^{\tilde{s}_k t} \, f_0(t)\Theta] \; dt = \sum_{k=1}^{K} w_k \quad R \begin{bmatrix} F_0(-s_1 - \tilde{s}_k) \\ \vdots \\ F_0(-s_M - \tilde{s}_k) \end{bmatrix}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (7.7)$$

Though the computational effort is evidently increased, the weight $W(t)$ given by equation (7.5) is quite flexible and can be made to cover a wide range of practical situations with relatively few terms. Again, much of the detailed work would be eliminated by a tabulation of inverse matrices for various pole combinations and weight factors.

Predistortion[2] of the prescribed response is another method for handling a somewhat more restricted type of weight factor. It has the advantage that it may be applied when the approximating functions are unweighted orthonormal functions. To begin, the WME integral

---

[1] The terms of equation (7.5) may also contain multiplicative integral powers of t. Equation (7.6) would then contain terms such as $\dfrac{1}{(s_m + s_n + \tilde{s}_k)^p}$ and equation (7.7) the derivatives of $F_0(s)$.

[2] The predistortion method is discussed by Kautz. However, some of the weight factors proposed by him are unbounded and lead to divergent coefficient integrals.

$$I = \int_0^\infty W \, [f_o - f_o{}^*]^2 \, dt$$

is written as an unweighted mean square error integral

$$I = \int_0^\infty [W^{1/2} f_o - W^{1/2} f_o{}^*]^2 \, dt$$

in which the predistorted prescribed response $W^{1/2} f_o$ is approximated by $W^{1/2} f_o{}^*$ a sum of unweighted approximating functions $a_n \Theta_n]$. It is clear then, that

$$d] = \int_0^\infty W^{1/2} f_o \Theta] \, dt \tag{7.8}$$

and
$$W^{1/2} f_o{}^* = a_n \Theta]$$

$$f_o{}^* = W^{-1/2} a_n \Theta] \qquad . \tag{7.9}$$

If the impulse response $h^*(t) = f_o{}^*(t)$ in equation (7.9) is to be a sum of complex exponentials ($a_n \Theta]$ is already such a sum) then $W^{-1/2}$ must also be a sum of complex exponentials, say

$$W^{-1/2} = \sum_{k=1}^{K} w_k{}^* e^{s_k{}^* t} \qquad . \tag{7.10}$$

Substituting equation (7.10) in equation (7.9) and taking the Laplace transform of the result yields the final approximation

$$F_o{}^*(s) = \sum_{k=1}^{K} w_k{}^* \sum_{n=1}^{N} a_n \Theta_n(s - s_k{}^*) \qquad . \tag{7.11}$$

Several disadvantages of the predistortion method are obvious: (1) the approximation has KM poles rather than M poles;[1] (2) the evaluation of the integrals,

$$\int_0^{\infty} \frac{f_o e^{s_m t}}{\sum\limits_{k=1}^{K} w_k^* e^{s_k^* t}} \, dt \quad ,$$

in equation (7.8) becomes difficult, eliminating one of the main advantages of the Laplace transform procedure; (3) the class of allowed weight factors is restricted. Consequently, $W^{-1/2}$ is usually limited to one or two real exponential functions. When $W^{-1/2} = e^{-1/2 \mu t}$ the formulas become particularly simple:

$$W = e^{\mu t} \tag{7.12}$$

$$d] = \begin{bmatrix} F_o(-s_1 - \mu/2) \\ \vdots \\ F_o(-s_M - \mu/2) \end{bmatrix} \tag{7.13}$$

$$F_o^*(s) = \lfloor a \rfloor \begin{bmatrix} \Theta_1(s + \mu/2) \\ \vdots \\ \Theta_N(s + \mu/2) \end{bmatrix} \quad . \tag{7.14}$$

It is this special case of predistortion that has the most practical value. More complicated weight factors are treated better with non-orthogonal functions and equations (7.6) and (7.7).

In summary, the unweighted mean square error approximation of impulse responses having known Laplace transforms is straightforward,

---

[1] A judicious choice of poles in $\Theta_1, - - - \Theta_N$ and of $s_1^*, - - - s_K^*$ can reduce considerably this number of poles.

requiring only the values of $F_o(s)$ at $s = -s_1, - - - -s_M$. Either ortho-
normal or independent exponential approximating functions can be used
and yield similar amounts of computational effort,except when tables
of inverse matrices are available. Then, the independent exponential
functions are preferred. Constrained approximations can be made simply,
but the introduction of arbitrary weight factors is more difficult and
must be limited primarily to non-orthogonal approximating functions and
the weight factor defined by equation (7.5).

### Analog Computer Approximation Procedures

The Laplace transform procedures of the previous section have
two primary disadvantages: (1) they are limited to prescribed response
functions having known Laplace transforms, and (2) they require lengthy
calculations to obtain the approximate function $f_o^*(t)$ and the error
$e = f_o - f_o^*$. Conventional digital or analog computing techniques
offer the most obvious and straightforward solution to these problems.
An alternative computer method, more closely related to the realiza-
tion aspects of linear system synthesis, is proposed in the present
section. It is based on the physical realization of the functions
$\Phi_1, - - - \Phi_N$ (i.e., $\Theta_1, - - - \Theta_N$) which make up $H^*(s)$ and depends on
these physical realizations to compute the approximation coefficients
and to simulate the approximate system. The following discussion of
this "analog" computer approximation procedure includes the mathemati-
cal basis for analog computation of the approximation coefficients,
the electronic differential analyzer circuits for realizing non-ortho-
gonal and orthonormal approximating functions, a computer method for
trial and error solution of the pole optimization problem, computer

evaluation of exponentially weighted coefficient integrals, and the

effects of computer inaccuracy.

Analog Computation of Approximation Coefficients--Figure 7.1

shows the coefficient computation process as applied to the integral

$$d_n = \int_0^\infty \Theta_n(t)f_o(t)dt \qquad .$$

A physical system which is synthesized to have a transfer function

$\Theta_n(s)$ receives as an input a forcing function which is made equal to
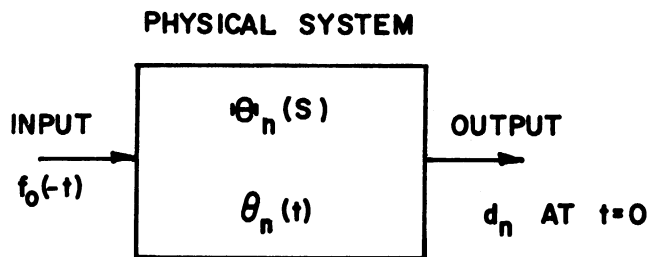
**PHYSICAL SYSTEM**



Figure 7.1  Physical System for Computation of the Integral $d_n$

$f_o(-t)$.  In practice $f_o(-t)$ would be obtained from $f_o(t)$ by some kind

of function generator, perhaps a tape recorder run backwards.  By

means of the superposition integral the output of the system is

$$\text{output} = \int_0^\infty \Theta_n(\tau)f_o(\tau-t)d\tau \qquad (7.15)$$

and at t = 0 is seen to be equal to the desired coefficient $d_n$.[1] If

need be, the coefficients $c_{nm}$ can be obtained in the same manner by

replacing $f_0(-t)$ with $\Theta_m(-t)$.

If the physical systems $\Theta_1$, - - - $\Theta_N$ were able to compute

just the coefficients and nothing else, the above scheme would have

little advantage over a direct analog computer integration of the for-

mulas for $d_1$, - - - $d_N$. However, the physical systems can also be com-

bined as in Figure 2.3 to simulate $H*(s)$ and to measure the approximation

error $e = f_0 - f_0*$ for arbitrary input functions. When the complexity

of the simulation is not too great it may even be used directly as

the final approximate system, bypassing the usual realization step of

the synthesis problem.

Differential Analyzer Circuits for Realization of Approxi-

mating Functions--Various physical embodiments of the functions $\Theta_1$,---$\Theta_N$

are possible, but those which employ passive or active electrical net-

works are favored because of their simplicity and inherent accuracy.

The physical systems which follow are examples of lumped[2] active net-

works and make use of electronic differential analyzer computing elements.[3]

---

[1] The use of the convolution property to compute integrals is well
known (see for example Laning and Battin [21]). McCool [27] dis-
cusses the application of a second order linear system with no
damping to the determination of Fourier series coefficients, and
Bose [5] utilizes the orthonormal Laguerre functions to find the co-
efficients $a_1$, - - - $a_N$ as functions of time in a non-linear system
representation. Only the last two references are concerned with
WME approximation.

[2] Nothing precludes the use of linear systems which are not lumped.

[3] Korn and Korn [19] discuss electronic differential analyzers and
explain the standard circuit symbology that follows.

Simpler passive networks can also be derived but lack convenience and flexibility.

Figure 7.2 shows the circuits for obtaining the commonly used functions $\dfrac{1}{s + \alpha}$ , $\dfrac{\beta}{s^2 + 2\alpha s + \beta^2}$ , and $\dfrac{s}{s^2 + 2\alpha s + \beta^2}$ . [1]



$$\frac{X}{Y} = -\frac{1}{s+\alpha}$$

$$\frac{X_1}{Y} = \frac{-s}{s^2 + 2\alpha\, s + \beta^2}$$

$$\frac{X_2}{Y} = \frac{\beta}{s^2 + 2\alpha s + \beta^2}$$

Figure 7.2  Analog Computer Circuits for Non-orthogonal Functions.

The pole positions are changed with the coefficient potentiometers which have the settings $\alpha$ and $\beta$.  As in Chapter V, $\alpha$ is the negative real part of the pole, and $\beta$ is the distance of the pole from the origin.

When the orthonormal functions of Chapter V are to be realized, the circuits become more complex and a cascade connection of individual elements, as in Figure 7.3, is preferred.  For the functions defined by equations (5.35) and (5.36) the ratio $\dfrac{\Lambda_n}{\Lambda_{n-1}}$ assumes either

[1] Otterman [30] employs circuits of this type for realization of Laguerre function approximations, but he computes the coefficients analytically.

the form $\dfrac{s - \alpha_n}{s + \alpha_n}$ or the form $\dfrac{s^2 - 2\alpha_n s + \beta_n^{\,2}}{s^2 + 2\alpha_n s + \beta_n^{\,2}}$ . Figure 7.4 gives the

circuits for realizing these functions along with taps for obtaining



Figure 7.3  Cascade Connection of All-Pass Elements for Realization
of Orthonormal Functions



$$\frac{X_1}{Y} = -\frac{1}{s + a_n} = -\frac{1}{\sqrt{2a_n}}\frac{\Theta_n}{\bigwedge_{n-1}}$$

$$\frac{X_2}{Y} = -\frac{s - a_n}{s + a_n} = -\frac{\bigwedge_n}{\bigwedge_{n-1}}$$

**(a)**



$$\frac{X_1}{Y} = \frac{-s}{s^2 + 2a_n s + \beta_n^2} = -\frac{1}{2\sqrt{a_n}}\frac{\Theta_{n+1}}{\bigwedge_{n-1}} \qquad \frac{X_2}{Y} = \frac{\beta_n}{s^2 + 2a_n s + \beta_n^2} = -\frac{1}{2\sqrt{a_n}}\frac{\Theta_n}{\bigwedge_{n-1}}$$

$$\frac{X_3}{Y} = -\frac{s^2 - 2a_n s + \beta_n^2}{s^2 + 2a_n s + \beta_n^2} = -\frac{\bigwedge_{n+1}}{\bigwedge_{n-1}}$$

**(b)**

Figure 7.4  Analog Computer Circuits for Orthonormal Functions
(a)  Real Pole.          (b)  Complex Poles.

$\Theta_n(s)$. When the individual elements are connected in cascade, as in Figure 7.3, most of the summing operations that take place in the end amplifiers of Figures 7.4a and 7.4b can be performed at integrator inputs, so that only one summing amplifier and N integrating amplifiers are needed for a total of N orthonormal functions.

It should be noted that the orthonormal functions generated by the preceding circuits require no preliminary mathematical calculations. They are produced entirely by the physical arrangement of the circuit compon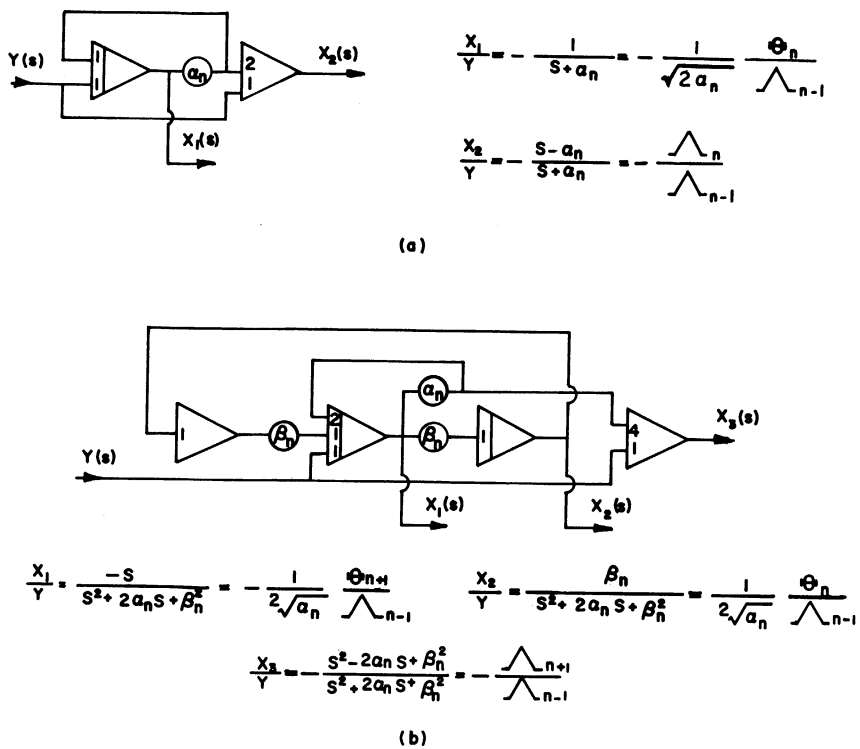ents and remain orthonormal even when the pole positions are varied by changing computer settings of $\alpha_n$ and $\beta_n$. A set of such circuits could be connected permanently in a practical, special purpose computer for orthonormal function approximation and simulation with arbitrary pole positions.

A further advantage of orthonormal function approximation is the simplicity of evaluating the approximation coefficients a] = d]. The matrix R is not required as in the Laplace transform procedure since $a_n = d_n$ is obtained directly from the system function $\Theta_n(s)$ and the forcing function $f_0(-t)$. This eliminates one of the main disadvantages of the orthonormal function expansions given in Chapter V. It also suggests a simplified computer method for optimum pole determination.

Computer Method for Pole Optimization--The computer block diagram for pole optimization is shown in Figure 7.5. A series of squaring devices and a summing amplifier are used in conjunction with the orthonormal function circuits to compute

$$\sum_{n=1}^{N} a_n^2$$

and consequently the mean square approximation error

$$I = I_{max} - \sum_{n=1}^{N} a_n^2 \quad .$$

Because the approximation error is so easily obtained, it is practical

to choose the optimum pole positions through trial and error adjustment

of the N computer settings $\alpha_n$ and $\beta_n$.[1] Suitable convergence schemes,



Figure 7.5  Computer Block Diagram for Pole Optimization

such as the method of steepest descent,would speed the process, and the

computer could be run repetitively at high speed so that changes in I

could be observed almost instantaneously.  The major disadvantage of

the method is the high computer accuracy which is demanded because I

becomes small as the square of approximation error.  For example, an

[1] Trial and error procedures for system approximation have been pro-
posed by many authors [9, 42].  As far as is known, none have used
orthonormal functions to automatically compute the zeros of H*(s)
and thus halve the number of unknowns.

average approximation error of 1 percent might require a computer error of .01 percent. Practically then, computer optimization must be limited to a reasonable number of poles, perhaps four or five. Fortunately, it is the simple approximation with relatively few poles where optimum pole determination is most important.

Computation of Weighted Integrals--Little elaboration on the previous remarks concerning the weight factor $W(t)$ and constrained approximations is needed here. The introduction of the weight factor $W(t)$ again offers the major difficulty. If the integral

$$\int_0^\infty e^{\tilde{s}_k t} f_0(t)\, \theta_n(t)\, dt$$

is needed, as in equation (7.7), it can be evaluated from the physical system shown in Figure 7.6 by noting that the impulse response $e^{\tilde{s}_k t} \theta_n(t)$ corresponds to a system function

$$\int_0^\infty e^{-st} e^{+\tilde{s}_k t}\, \theta_n(t)\, dt = \Theta_n(s - \tilde{s}_k) \qquad (7.16)$$
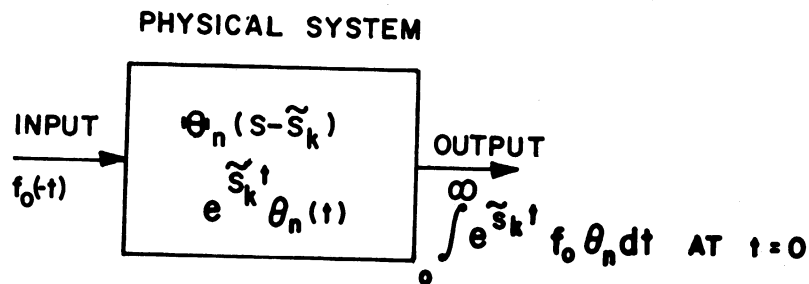
**PHYSICAL SYSTEM**



Figure 7.6   Physical System for Computation of Exponentially Weighted Integral

and that at $t = 0$ the convolution integral gives the indicated output.

When the physical system is an electrical network the change

from $\Theta_n(s)$ to $\Theta_n(s-\tilde{s}_k)$ is easily made if $\tilde{s}_k$ is real. All the dynamic

elements of the system are modified as shown in Figure 7.7. Note that

the change in impedance for the inductor or the capacitor always pro-

duces a negative $\tilde{s}_K$. Positive values of $\tilde{s}_k$ must be obtained with active

elements as the example in Figure 7.7d demonstrates.



Figure 7.7  Circuit Modifications for Changing s into $(s-\tilde{s}_k)$
 (a)  Capacitive Impedance    (b)  Inductive Impedance
 (c) and (d)  Integrator System Function

If the function $W(t)$ is not a sum of real exponentials (this

happens in equations (7.6) and (7.7) when the $\tilde{s}_k$ are complex), the com-

putation of the weighted integrals must assume the more direct approach

discussed earlier. The integral

$$\int_0^\infty W\, f_o\, \Theta_n\, dt$$

is calculated by forcing the system $\Theta_n(s)$ with $W(-t)f_o(-t)$ and measuring

the output at t = 0. The main complication, if it can be considered a

complication, is the generation of the function $W(-t)f_o(-t)$.

Error Analysis--Last of all, it is important to consider the effects of computer errors. As in most computational problems, digital or analog, complete error analysis is very difficult and must be replaced in part by a comparison of computer solutions with known analytical solutions. This has been done in several examples in the next section with results that appear quite satisfactory for average approximation problems. In a theoretical sense, it is also possible to say something about the relations between errors in the approximation coefficients d] and the WME of the resulting approximation. While such relations are important and are developed in the following paragraphs, it must be remembered that their final application depends on knowledge of the size of the errors in d] and the dependence of these errors on computer errors.

To discover the effect of errors in d] it is convenient to write the computed coefficients as d] + $\Delta$d] where the elements of $\Delta$d] represent the computer errors in $d_1$, - - - $d_N$. Because the coefficients are no longer the optimum coefficients, the WME will be increased from I by an amount $\Delta$I defined by

$$I + \Delta I = \int_{T_1}^{T_2} W[f_0 - \underline{\theta}, C^{-1} (d] + \Delta d])]^2 dt. \qquad [1] \qquad (7.17)$$

---

[1] Note that the expression $I = I_{max} - \underline{d}, C^{-1} d]$ has no meaning when the coefficients $d_1$, - - - $d_N$ are in error since it was derived on the basis of

$$d_n \equiv \int_{T_1}^{T_2} W f_0 \theta_n dt \qquad .$$

Expansion of equation (7.17) gives

$$I + \Delta I = \int_{T_1}^{T_2} W[(f_o - f_o^*) - \underline{\theta}\, C^{-1}\, \Delta d]\ ]^2\, dt$$

$$= I - 2\int_{T_1}^{T_2} W(f_o - f_o^*)\,\underline{\theta}\, C^{-1}\, \Delta d]\, dt + \int_{T_1}^{T_2} W(\,\underline{\theta}\, C^{-1}\, \Delta d])^2\, dt$$

$$= I - 2\int_{T_1}^{T_2} Wf_o\,\underline{\theta}\, dt C^{-1}\, \Delta d] + 2\int_{T_1}^{T_2} W\{\,\underline{d}\, C^{-1}\, \theta]\}\,\underline{\theta}\, dt\, C^{-1}\, \Delta d]$$

$$+ \int_{T_1}^{T_2} W\{\,\underline{\Delta d}\, C^{-1}\, \theta]\}\ \{\,\underline{\theta}\, C^{-1}\, \Delta d]\}\, dt$$

$$= I - 2\,\underline{d}\, C^{-1}\, \Delta d] + 2\,\underline{d}\, C^{-1}\, \Delta d] + \underline{\Delta d}\, C^{-1}\, \Delta d]$$

and

$$\Delta I = \underline{\Delta d}\, C^{-1}\, \Delta d] \qquad . \qquad\qquad (7.18)$$

Thus the increase in the WME is readily determined when $\Delta d]$ is known. If it is found that $\Delta I$ is small compared to $I$, then the computer errors in $d]$ can be considered negligible.

From the standpoint of ultimate approximation error, formula (7.18) is the logical one to apply. A less valid approach is to find the errors in the coefficients $a]$ by

$$\Delta a] = C^{-1}\, \Delta d] \quad . \qquad\qquad (7.19)$$

In many cases (especially when the signs of the terms in

$$\Delta a_n = \sum_{m=1}^{N} c_{nm}^{-1}\, \Delta d_m$$

alternate) $\Delta I$ is small even when $\Delta a_1,$ - - - $\Delta a_N$ are very large.  Such results are possible because the effect of a large error in one co-efficient may be compensated by corrective changes in the remaining coefficients.

While equation (7.18) is applicable when $\Delta d]$ is known it does not say anything about the sensitivity of $\Delta I$ to $\Delta d]$ when $\Delta d]$ is not known.  In this respect it is helpful to consider the relation-ship between

$$E = \frac{\sum_{n=1}^{N} \Delta d_n^2}{\sum_{n=1}^{N} d_n^2} \qquad (7.20)$$

the relative mean square error in $d_1,$ - - - $d_N,$ and

$$\Delta I_{rel} = \frac{\Delta I}{\int_{T_1}^{T_2} W f_o^2 \, dt} = \frac{\Delta I}{I_{max}} \qquad (7.21)$$

the relative mean square approximation error due to $\Delta d_1,$ - - - $\Delta d_N.$ If $k_1$ is the smallest eigenvalue of $C^{-1}$ and $k_N$ is the largest eigen-value of $C^{-1}$ then [8]

$$k_1 \sum_{n=1}^{N} \Delta d_n^2 \leq \underline{\Delta d} \, C^{-1} \, \Delta d] \leq k_N \sum_{n=1}^{N} \Delta d_n^2 \qquad (7.22)$$

and

$$k_1 \sum_{n=1}^{N} d_n^2 \leq \underline{d} \, C^{-1} \, d] \leq k_N \sum_{n=1}^{N} d_n^2 \qquad . \qquad (7.23)$$

But $I_{max} \stackrel{\sim}{=} \lfloor d, C^{-1} d \rfloor$ for $I << I_{max}$ so

$$\Delta I_{rel} = \frac{\lfloor \Delta d, C^{-1} \Delta d \rfloor}{I_{max}} \stackrel{\sim}{=} \frac{\lfloor \Delta d, C^{-1} \Delta d \rfloor}{\lfloor d, C^{-1} d \rfloor} \quad . \quad (7.24)$$

Hence,

$$\Delta I_{rel} \stackrel{\sim}{=} r \, E, \quad I << I_{max} \quad (7.25)$$

where

$$\frac{k_1}{k_N} \leq r \leq \frac{k_N}{k_1} \quad . \quad (7.26)$$

Equation (7.26) puts bounds on $r$, the sensitivity of the relative mean square approximation error to relative mean square co-efficient error. When $k_N/k_1$ is large the sensitivity may be very large or very small depending on the particular values of $d_1$, - - - $d_N$ and $\Delta d_1$, - - - $\Delta d_N$. From an uncertainty point of view this is undesirable and it would be better if $k_N/k_1 = 1$. This is the case when $\theta_1$, - - - $\theta_N$ are orthonormal. Thus, orthonormal approximating functions have a decided advantage over non-orthonormal functions for analog computation of approximation coefficients. If non-orthonormal functions must be employed, some check on final approximation accuracy should be made.[1] Usually, this involves a computer measurement of $e = f_0 - f_0^*$; the cal-culation of $k_N/k_1$ is more difficult and less conclusive.

To review, it has been shown that WME approximation by analog computation is practical and has many advantages. Among them are: the simple computation of weighted or unweighted approximation coefficient

---

[1] This is especially true for high order matrices where $k_N/k_1$ tends to be large.

integrals, convenient realization and simulation of H*(s), rapid error calculation, and a computer method of pole optimization. Although computer errors are an added problem, they are small in most practical problems and are minimized by the use of orthonormal approximating functions.

Examples

The examples of WME impulse response approximation given in this section are presented to illustrate the conclusions of the preceding sections and to show the practicality of the approximation procedures. They are summarized with pertinent approximation information in Table I and are described in detail in the following paragraphs. Tables II and III contain supporting numerical results while Figures 7.8 through 7.15 show plots of $f_o$ and $f_o$*.

The computing components employed in examples four, six, and eight through eleven were part of a general purpose, direct current electronic differential analyzer. All summing amplifiers, integrating amplifiers, and coefficient potentiometers were calibrated within 0.1 percent and were connected according to the diagrams of the previous section. A relay control device was incorporated to disconnect all integrator inputs at t = 0. This held the output voltages so they could be measured statically with an accuracy comparable to component accuracy.

The first example in Table I involves the analytic approximation of a fourth power pulse by six non-orthogonal Laguerre functions[1] whose poles are at s = -1. The pulse is defined by

---

[1] See the Appendix, section 1.

TABLE I.  EXAMPLE INFORMATION

| Example Number | Prescribed Response | Type of Approximation | Approximating* Functions | Pole Positions | N | W(t) | Figure |
|---|---|---|---|---|---|---|---|
| 1 | fourth power pulse, equation (7.27) | unconstrained analytic | non-orthogonal Laguerre, (1) | $\alpha_n = 1$ $\beta_n = 0$ | 6 | 1 | 7.8 |
| 2 | fourth power pulse, equation (7.27) | unconstrained analytic | non-orthogonal constrained Laguerre, (2) | $\alpha_n = 1$ $\beta_n = 0$ | 5 | 1 | 7.9 |
| 3 | fourth power pulse, equation (7.27) | constrained analytic, $f^{*\prime}(0) = 0$ | non-orthogonal constrained Laguerre, (2) | $\alpha_n = 1$ $\beta_n = 0$ | 5 | 1 | 7.10 |
| 4 | fourth power pulse, equation (7.27) | unconstrained computer | non-orthogonal Laguerre, (1) | $\alpha_n = 1$ $\beta_n = 0$ | 6 | 1 | 7.11 |
| 5 | fourth power pulse, equation (7.27) | unconstrained analytic | orthonormal Laguerre | $\alpha_n = 1$ $\beta_n = 0$ | 6 | 1 | -- |
| 6 | fourth power pulse, equation (7.27) | unconstrained computer | orthonormal Laguerre | $\alpha_n = 1$ $\beta_n = 0$ | 6 | 1 | -- |
| 7 | delayed pulse, equation (7.28) | unconstrained analytic | non-orthogonal Butterworth, (6) | $\alpha_1 = 1$, $\alpha_2 = .8$ $\alpha_4 = .3$, $\beta_2 = 1$ $\beta_4 = 1.044$ | 5 | 1 | 7.12 |

* Numbers refer to Appendix sections.

TABLE I.  (CONT'D)

| Example Number | Prescribed Response | Type of Approximation | Approximating* Functions | Pole Positions | N | W(t) | Figure |
|---|---|---|---|---|---|---|---|
| 8 | trapezoid pulse, equation (7.29), T = 7 | unconstrained computer | orthonormal | $\alpha_1 = 1$, $\alpha_2 = .75$, $\alpha_4 = .5$, $\alpha_6 = .25$, $\beta_n = 1$ | 7 | 1 | 7.13 |
| 9 | trapezoid pulse, equation (7.29), T = 4 | unconstrained computer | orthonormal | $\alpha_1 = 1$, $\alpha_2 = 2/3$, $\alpha_4 = 1/3$, $\beta_n = 1$ | 5 | 1 | 7.14 |
| 10 | trapezoid pulse, equation (7.29), T = 4 | unconstrained computer | orthonormal | $\alpha_1 = .75$, $\alpha_2 = 5/12$, $\alpha_4 = 1/12$, $\beta_2 = .856$, $\beta_4 = .946$ | 5 | $e^{t/2}$** | 7.14 |
| 11 | matched filter response | unconstrained computer | orthonormal | $\alpha_1 = 1$, $\alpha_2 = .809$, $\alpha_4 = .309$, $\beta_n = 1$ | 5 | 1 | 7.15 |

* Numbers refer to Appendix sections.

**Obtained by predistortion.  Poles of final approximation are those of example 9.

TABLE II.  APPROXIMATION DATA

| Example Number | $I_{max}$ | $\dfrac{I_{rel}^*}{I_{c_{rel}}}$ | $\Delta I_{rel}$ | $e_{max}$ | $e_{max}^*$ | $\overline{T}_2 - \overline{T}_1$ |
|---|---|---|---|---|---|---|
| 1 | 3.2507 | .00651 | -- | .114 | .065 | 10 |
| 2 | 3.2507 | .00684 | -- | .092 | .067 | 10 |
| 3 | 3.2507 | .00845 | -- | .095 | .074 | 10 |
| 4 | 3.2507 | .0102 | .000365 | -- | -- | -- |
| 5 | 3.2507 | .00651 | -- | -- | -- | -- |
| 6 | 3.2507 | .0032 | $5.92 \times 10^{-6}$ | -- | -- | -- |
| 7 | 2.9556 | .0999 | -- | .810 | .243 | 10 |
| 8 | 4.6667 | .0066 | -- | .054 | .083 | 9 |
| 9 | 2.6667 | .0098 | -- | .133 | .102 | 5 |
| 10 | 5.8100 | .0220 | -- | .178 | .179 | 8 |
| 11 | -- | -- | -- | -- | -- | -- |

\* In computer approximations $I_{rel}$ is obtained from computer coefficients.

TABLE III.   APPROXIMATION COEFFICIENTS

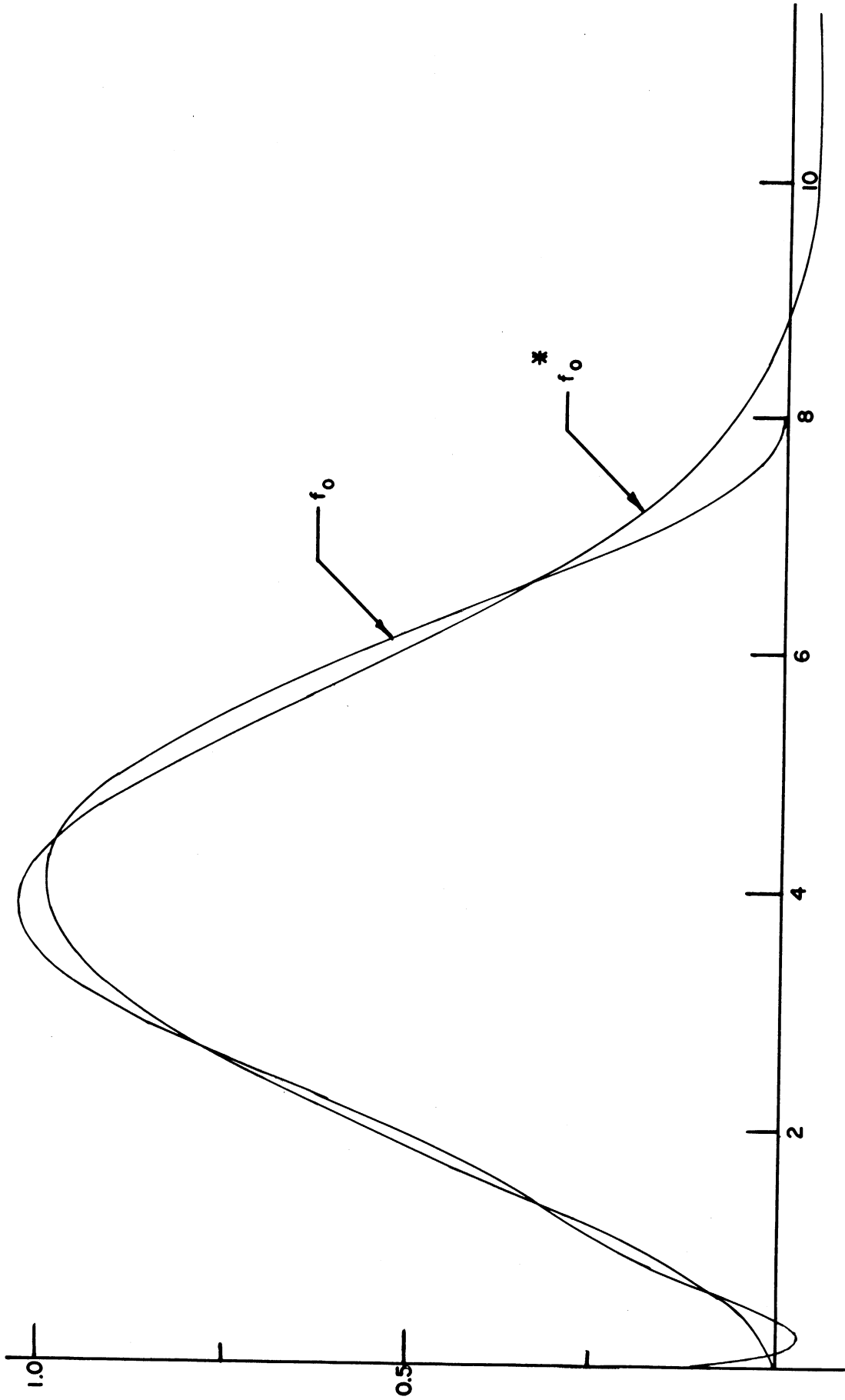| Example Number | Coefficient | Coefficient Number | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 1 | $d_n$ | 0.218425 | 0.464987 | 0.633805 | 0.689099 | 0.643712 | 0.534782 | -- |
| | $a_n$ | 0.1136 | -1.4244 | 7.7913 | -17.9408 | 26.9403 | -11.3127 | -- |
| 2 | $d_n$ | 0.464987 | 0.633805 | 0.689099 | 0.643712 | 0.534782 | -- | -- |
| | $a_n$ | -0.8565 | 6.2768 | -15.6691 | 25.1229 | -10.7069 | -- | -- |
| 3 | $d_n$ | 0.464987 | 0.633805 | 0.689099 | 0.643712 | 0.534782 | -- | -- |
| | $b_n$ | 0.0000 | 2.8509 | -9.5025 | 19.6415 | -8.7493 | -- | -- |
| 4 | $d_n$ | 0.220 | 0.467 | 0.634 | 0.690 | 0.644 | 0.535 | -- |
| | $a_n$ | -0.132 | 0.696 | 0.400 | -4.832 | 15.168 | -7.040 | -- |
| 5 | $a_n$ | 0.30890 | -1.00628 | 1.26388 | -0.67687 | -0.05930 | 0.24999 | -- |
| 6 | $a_n$ | 0.312 | -1.008 | 1.265 | -0.678 | -0.059 | 0.252 | -- |
| 7 | $d_n$ | 0.0183156 | 0.0352155 | 0.0351941 | 0.6062901 | 0.0170923 | -- | -- |
| | $a_n$ | 14.5785 | -20.2825 | 2.8474 | 4.9605 | -0.0038 | -- | -- |
| 8 | $a_n$ | 1.399 | -0.204 | 1.543 | -0.395 | -0.305 | -0.202 | -0.083 |
| 9 | $a_n$ | 1.329 | -0.599 | -0.706 | 0.130 | 0.004 | -- | -- |
| 10 | $a_n$ | 1.883 | -1.167 | -0.857 | 0.043 | 0.196 | -- | -- |
| 11 | $a_n$ | 0.338 | -0.4621 | -1.075 | 0.385 | 0.164 | -- | -- |

Figure 7.8 Unconstrained Analytic Approximation of Fourth Power Pulse
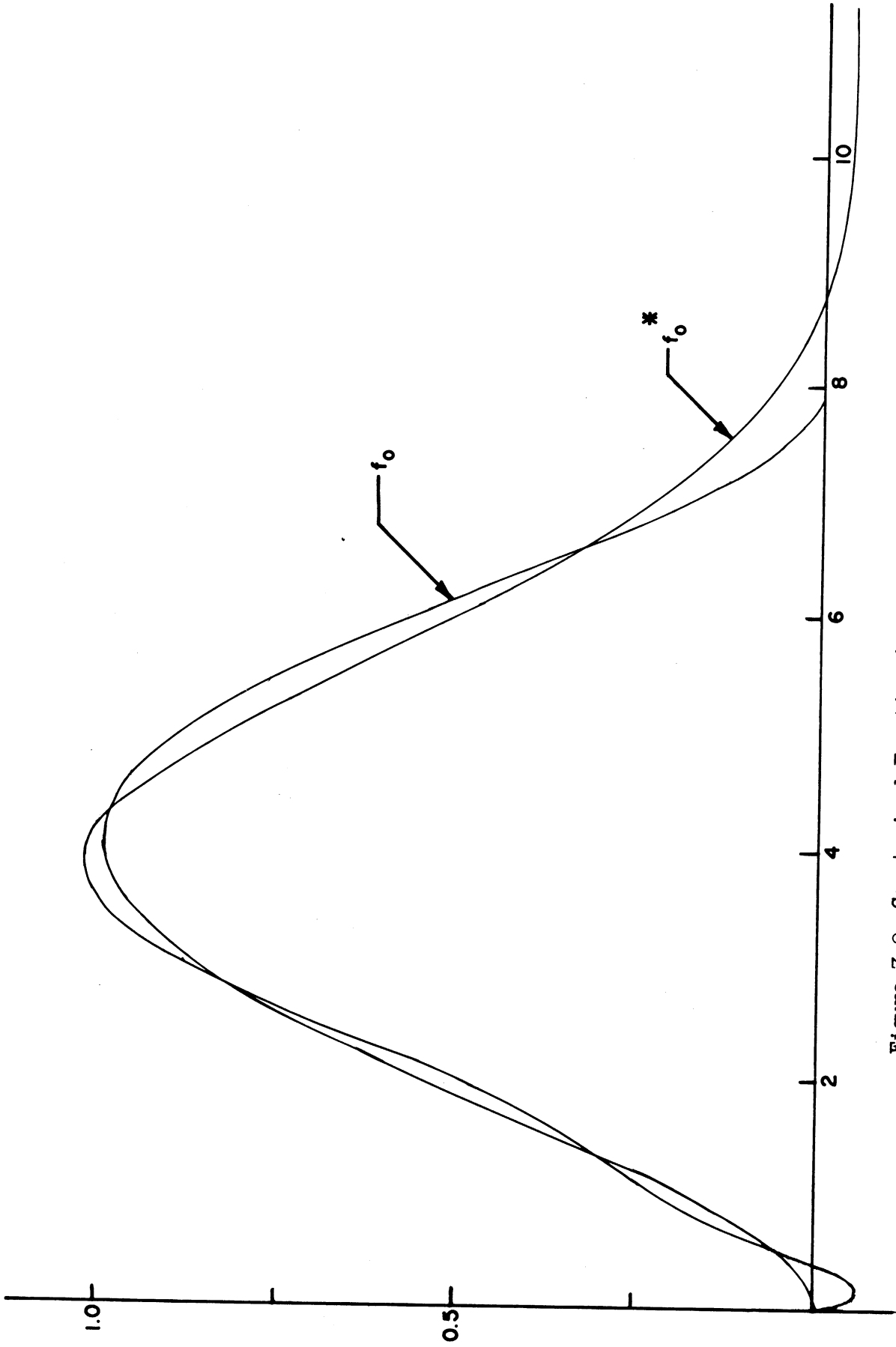
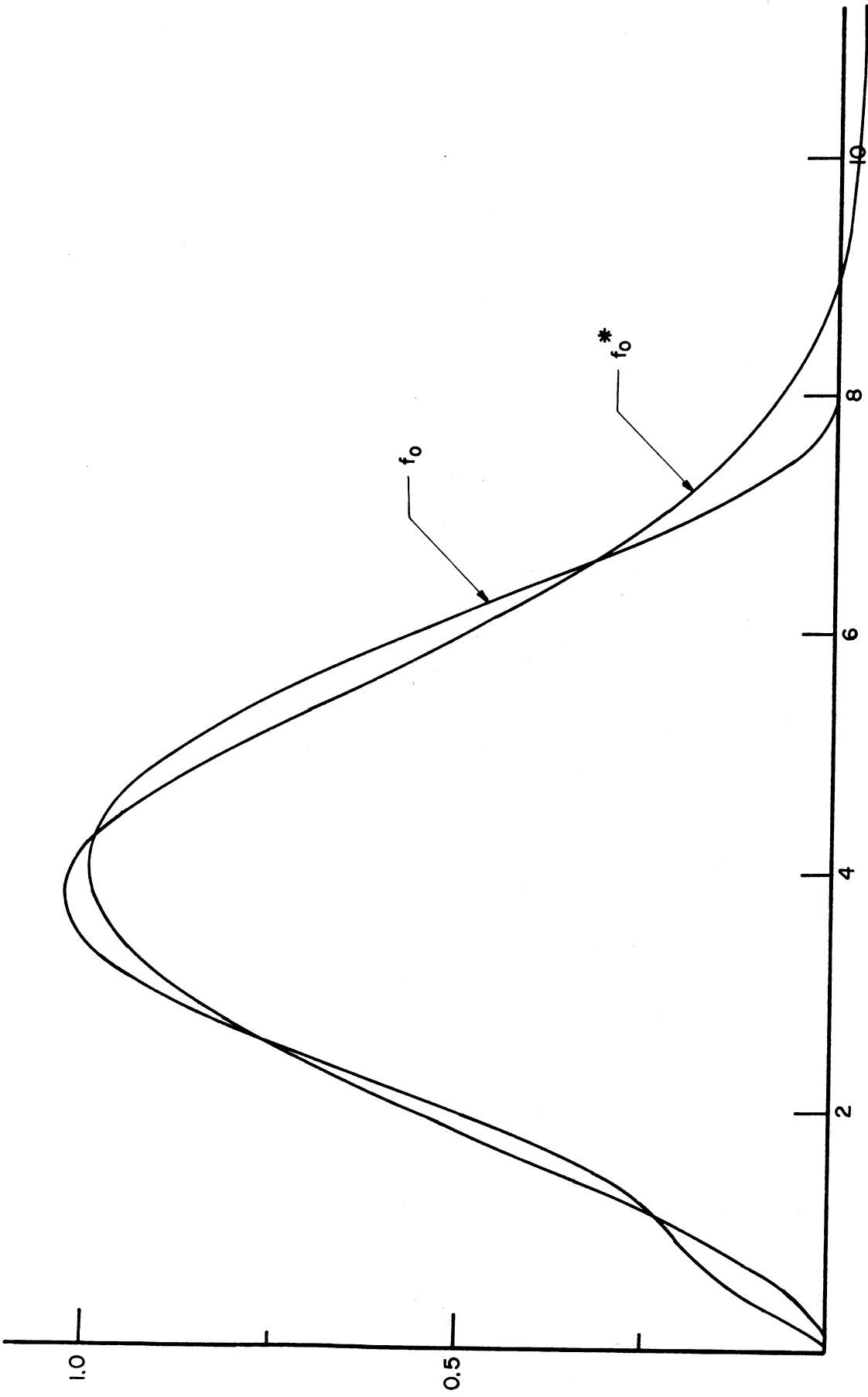Figure 7.9  Constrained Function Approximation of Fourth Power Pulse

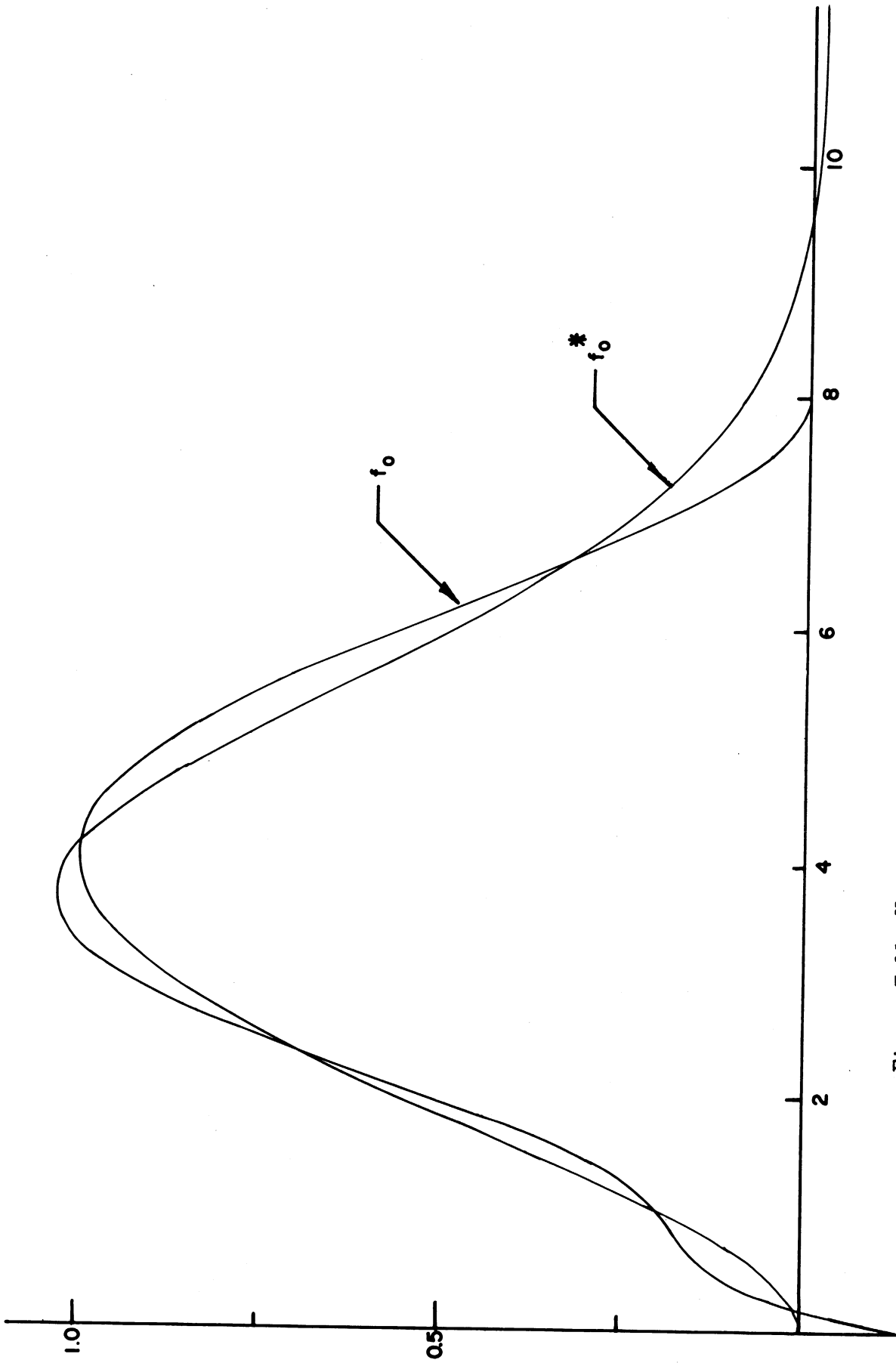Figure 7.10  Constrained Approximation of Fourth Power Pulse

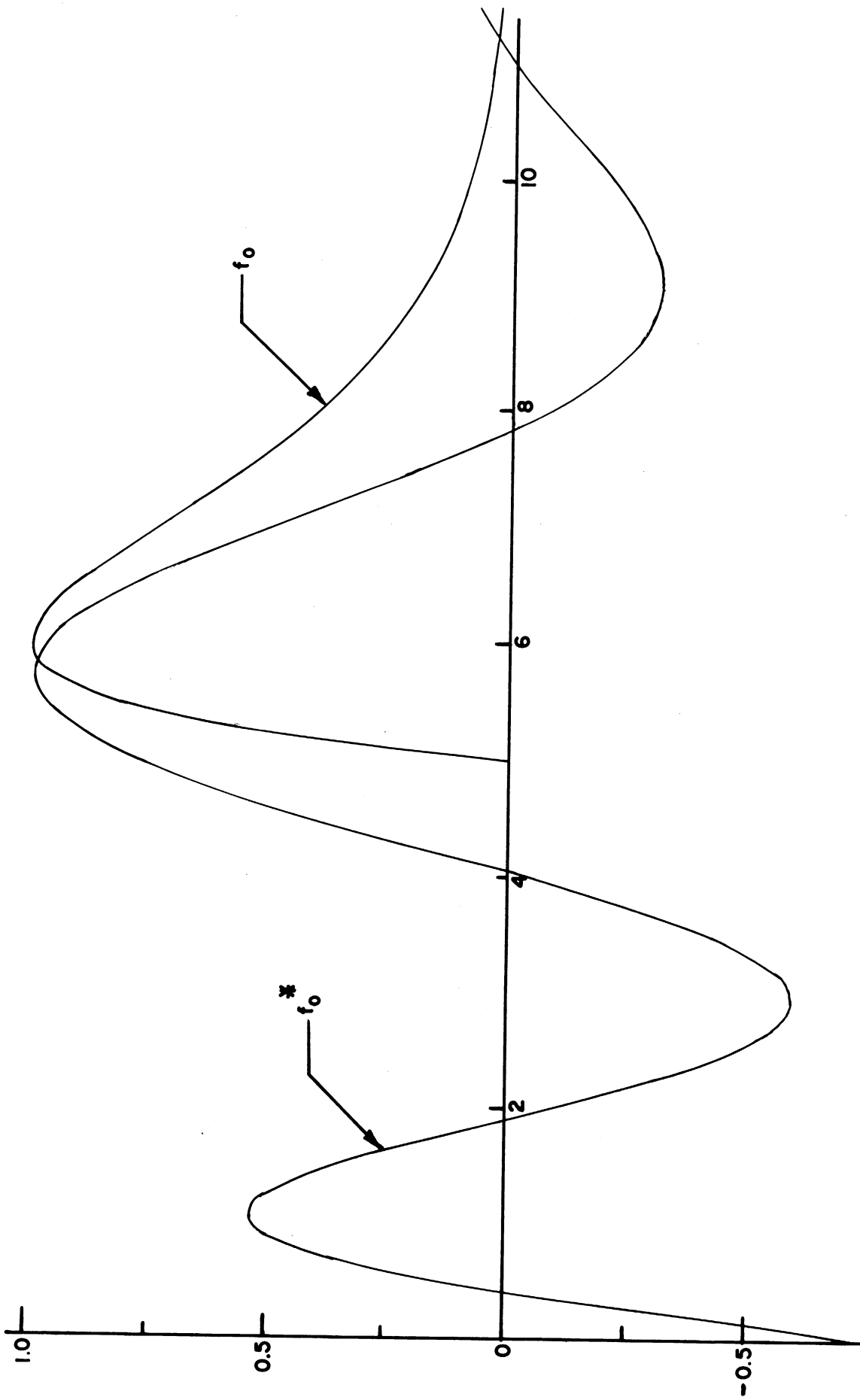Figure 7.11  Unconstrained Computer Approximation of Fourth Power Pulse

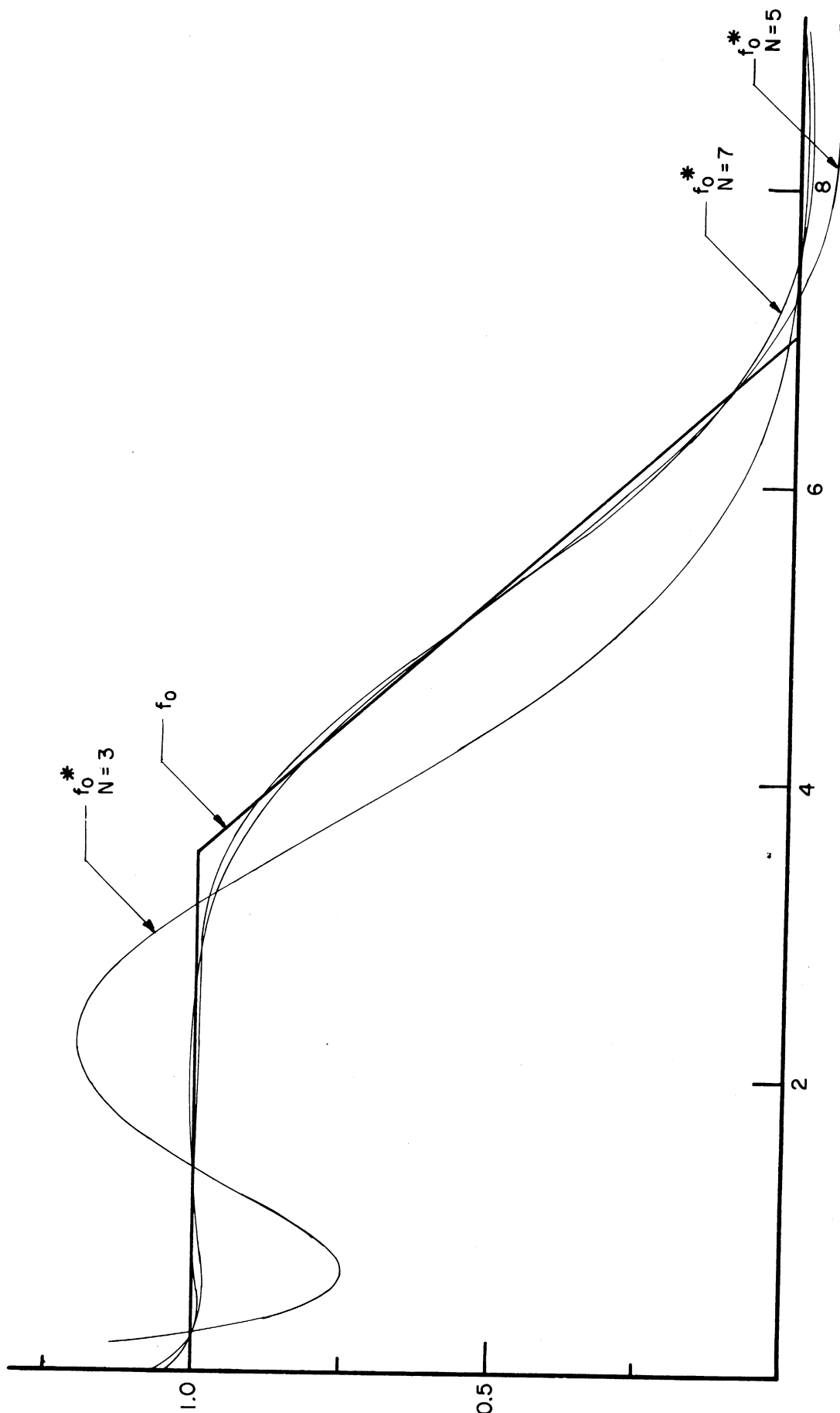Figure 7.12 Unconstrained Approximation of Delayed Pulse

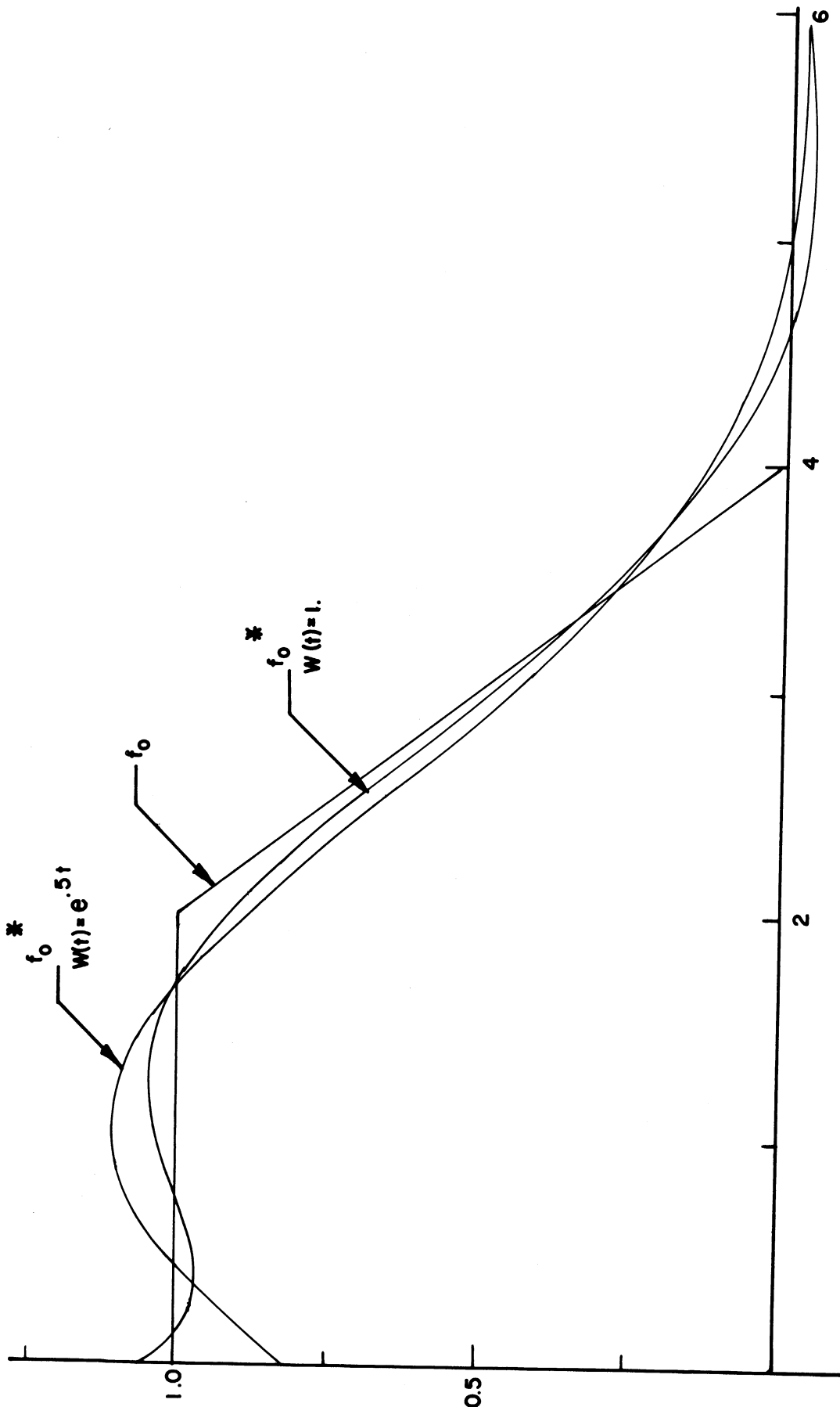Figure 7.13  Orthonormal Function Approximation of Trapezoid Pulse

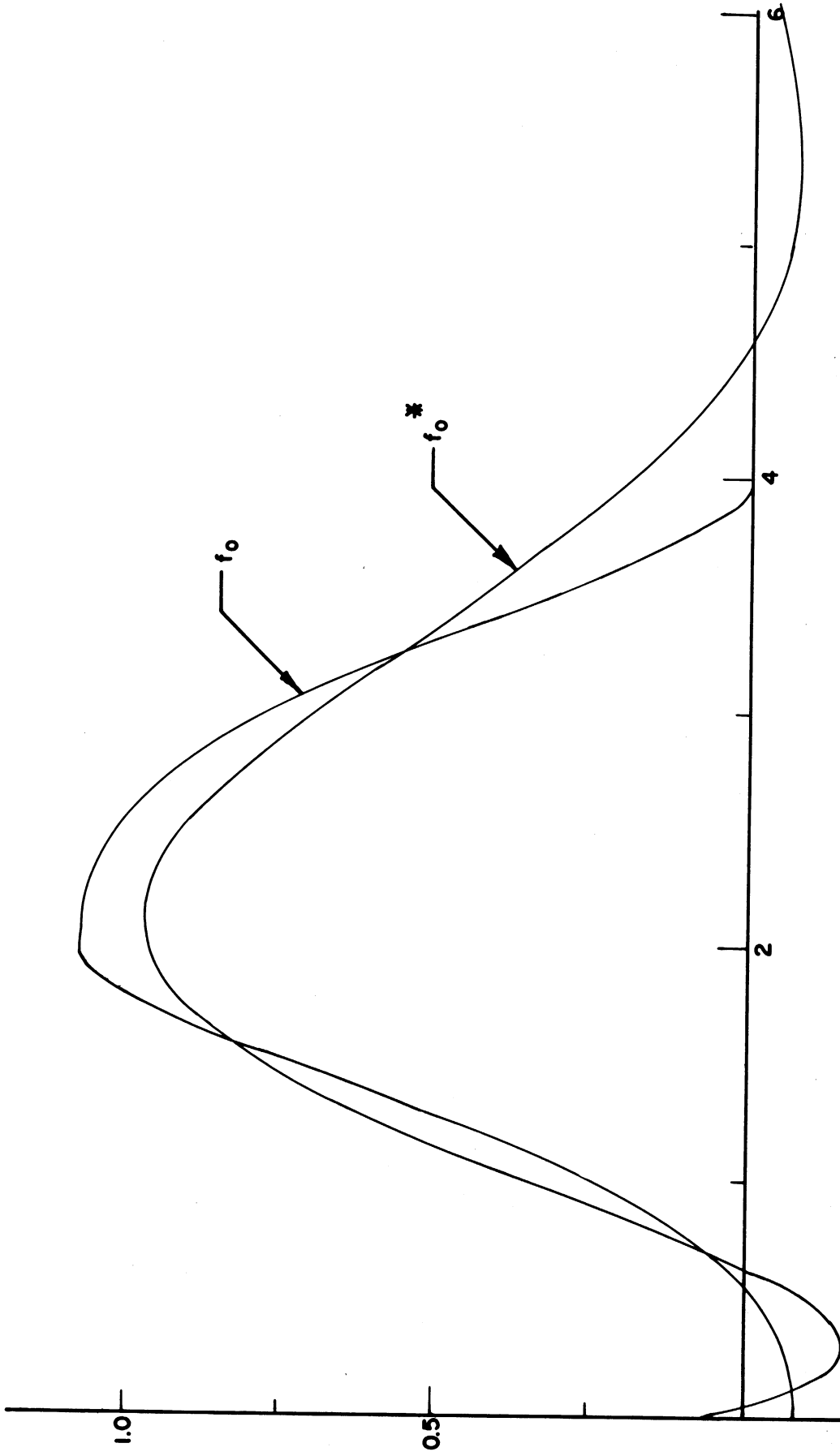Figure 7.14 Weighted and Unweighted Approximation of Trapezoid Pulse

Figure 7.15 Orthonormal Function Approximation of Matched Filter Response

$$f_0(t) = \frac{t^4}{256} - \frac{t^3}{16} + \frac{t^2}{4} \; , \quad 0 \le t < 8$$

$$= 0 \quad , \; t < 0, \quad t \ge 8 \; , \tag{7.27}$$

is symmetric about $t = 4$, has a maximum at $t = 4$ of $f_0(4) = 1$, and has zero value and slope at $t = 0$ and $t = 8$. Its length is taken to match approximately the length of the longest approximating function,

$$\theta_6 = \frac{t^5}{5!} e^{-t}$$

which has its maximum at $t = 5$. Coefficients $d_1$, - - - $d_6$ and $a_1$, - - $a_6$ are determined from $F_0(s)$ and its derivatives from the formulas of the Appendix and are listed in Table III. Figure 7.8 compares the approximation $f_0{}^*$ with the prescribed response $f_0$; the error is small with the greatest deviations occurring at $t = 0$ and $t = 8$, where $f_0$ has properties which are not characteristic of the approximating functions.

As the first step in improving the approximation at $t = 0$ the constrained Laguerre functions of Appendix section 2 are employed. These functions have the same poles as the unconstrained functions but satisfy in addition the condition $\theta_n(0) = 0$. For a six pole approximation five functions are required. Again, the coefficients are listed in Table III, in this case under example two. The plot of $f_0{}^*$ is shown in Figure 7.9 and is seen to differ little from the previous plot except at $t = 0$. The leniency of the imposed constraint is confirmed by the small increase in relative mean square error from 0.00651 to 0.00684.

As the final step in matching the approximation at $t = 0$, example three utilizes the constrained Laguerre functions in a

constrained approximation so that $f_o*(0) = f_o'*(0) = 0$. The constrained coefficients $b_1, - - - b_5$ are obtained from equations (4.30), and the approximation is plotted in Figure 7.10. A noticeable change in the error magnitude and distribution is observed. The severity of the constraint is indicated by the increase in $I_{rel}$ from 0.00684 to 0.00845.

In order to determine the accuracy of computer approximation, example four repeats the unconstrained pulse approximation of example one. Table III shows that the computer errors in $d_1, - - - d_6$ are very small, all less than 0.002. Even so, the errors in the calculated approximation coefficients $a_1, - - - a_6$ are extremely large. The apparent discrepancy is due to the sign alternations of the terms in the sum

$$a_n = \sum_{m=1}^{6} c_{nm}^{-1} d_m \quad .$$

At first inspection it would appear that the $a_n$ are worthless, but the approximation plotted in Figure 7.11 shows otherwise. Although e has increased, it is still tolerably small. This result is confirmed by the small increase in relative mean square error, $\Delta I_{rel} = 0.000365$.[1] These facts all agree with the remarks of the previous section.

While $\Delta I_{rel}$ is quite small the ratio $\Delta I_{rel}/E = r = 89.7$ is large. Hence it would seem advantageous to obtain the approximation by means of orthonormal functions. This has been done analytically in

[1] The value

$$I_{rel} = 1 - \frac{d]\, C^{-1}\, d]}{I_{max}}$$

given in Table II has questionable significance in computer approximations because d] is not precise. See footnote at bottom of page 85. A more satisfactory computer value could be obtained by direct integration of

$$\int_0^{} w\, [f_o - f_o*]^2\, dt \quad .$$

example five and by the computer in example six. Again the computer

errors in the coefficients are small, less than 0.003, but the increase

in relative mean square error $\Delta I = 5.92 \times 10^{-6}$ is negligibly small.

This graphically illustrates the superiority of orthonormal functions

for computer approximation.

Example seven is concerned with the analytic approximation

of the delayed exponential pulse

$$f_O(t) = 0 \quad , \quad t < 5$$
$$= e^6(t-5)e^{-t}, \quad t \geq 5 \quad . \tag{7.28}$$

The pulse reaches its maximum of one at $t = 6$ and decays slowly. It

is particularly difficult to approximate because it is zero for the

first five seconds.[1] The Butterworth functions of Appendix section 6

are used in the approximation because their poles assume a configuration

in the s-plane which is similar to the configuration obtained by a Padé

approximation of an ideal delay. . As in example one the length of the

prescribed response is chosen to approximate the length of the longest

approximating function. The approximation coefficients obtained from

$F_O(s)$ are shown in Table III; $f_O{}^*(t)$ is plotted in Figure 7.12. Ob-

viously, five approximating functions are insufficient in this problem.

This is not surprising since all finite representations of $e^{-Ts}$ with

good transient response have many terms.

In example eight a computer approximation of the trapezoid

pulse

$$f_O(t) = 0 \quad , \quad t < 0 \tag{7.29}$$
$$= 1 \quad , \quad 0 \leq t < \frac{T}{2}$$

_____

[1] Delay removal (see Chapter II) would permit $f_O$ to be approximated
with no error.

$$= 2(1 - \frac{t}{T}), \quad \frac{T}{2} \leq t < T$$

$$f_o(t) = 0 \qquad , \quad t \geq T$$

is made. The seven approximating functions are orthonormal with poles on the unit circle at $s = -1$, $-3/4 \pm j\sqrt{1-(3/4)^2}$, $-1/2 \pm j\sqrt{1-(1/2)^2}$, $-1/4 \pm j\sqrt{1-(1/4)^2}$. Different pulse lengths yield the following computer mean square errors: $T = 10$, $I_{rel} = 0.0188$; $T = 7$, $I_{rel} = 0.00658$; $T = 5$, $I_{rel} = 0.00654$; $T = 4$, $I_{rel} = 0.0121$.[1] Consequently, it appears that the length $T = 6$ is about optimum for the above approximating functions. This figure agrees closely with the length of the longest approximating functions. Figure 7.13 shows $f_o^*$ for $T = 7$ with seven, five, and three of the above poles. For $N = 5$, the last two poles are omitted; for $N = 3$ the last four poles are omitted. Omission of the last two poles, which contribute mainly to $f_o^*$ for $t > 7$, causes $I_{rel}$ to increase from 0.00658 to 0.0139.

Example nine repeats the problem of example eight with only five poles, $s = -1$, $-2/3 \pm j\sqrt{1 - (2/3)^2}$, $-1/3 \pm j\sqrt{1 - (1/3)^2}$. The pulse length is reduced to $T = 4$ to compensate for the change in approximating functions. As would be expected, Figure 7.14 indicates an increase in approximation error over the error in example eight, $N = 7$; the same conclusion is not valid in example eight, $N = 5$. Hence, the present poles are superior for a five pole approximation.

Example ten considers the same approximation with $W(t) = e^{0.5t}$. The method of predistortion is used (equations (7.8) through (7.11)) in conjunction with the pole shift procedures of Figures 7.6

---

[1] The computer mean square error $I_{rel}$, although differing somewhat from the true mean square error, should be a good comparative measure of the different approximation errors.

and 7.7. To fix the poles of the final approximation at s = -1,

$-2/3 \pm j \sqrt{1 - (2/3)^2}$, $-1/3 \pm j \sqrt{1 - (1/3)^2}$, the poles of the ortho-

normal functions are taken at $s = -3/4$, $-5/12 \pm j \sqrt{1 - (2/3)^2}$,

$-1/12 \pm j \sqrt{1 - (1/3)^2}$. The weighted approximation is plotted with

the unweighted approximation in Figure 7.14. As predicted by theory,

errors at large t are reduced at the expense of increased errors at

small t.

Example eleven concerns the approximation of a filter for

detecting a signal in white noise. If the detection condition speci-

fies that the ratio of peak signal to rms noise be a maximum, then

the matched-filter criterion of Van Vleck and Middleton [39] holds

and states that $h(t) = S(-t)$ where $S(t)$ is the signal to be detected.

In the present problem $S(t)$ is generated from a square pulse by a

transmitter which acts like a linear system of second order. Figure

7.16 shows the signal source and defines the signal parameters. Since

$h(t) = S(-t) = 0$ for $t > 0$, it cannot be approximated directly. It

must first be delayed. By letting $h(t) = S(-t-4)$ for $t > 0$ and

$h(t) = 0$ for $t < 0$ most of the signal is considered. Figure 7.15

shows $f_0(t) = h(t) = S(-t-4)$ and the corresponding computer approxi-

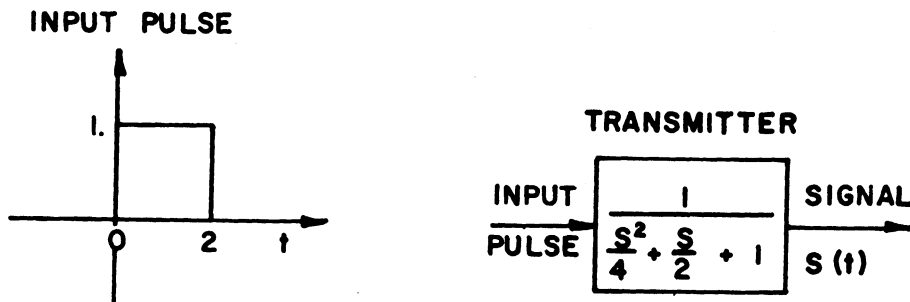mation $f_0^*(t)$. The five approximating functions have Butterworth poles



Figure 7.16 Signal Source for Detection Problem

at $s = e^{j\pi}$, $e^{j(\pi \pm \pi/5)}$, $e^{j(\pi \pm 2\pi/5)}$. Considering the nature of $f_o$, the five term approximation is fairly good. A time delay greater than four seconds in the function $S(-t)$ would reduce the error due to omission of $S(t)$ for $t > 0$, but it would also increase the mean square approximation error. In this respect the approximation of the matched filter impulse response is more difficult than the previous approximations.

The last three columns of Table II are included to demonstrate the application and accuracy of equation (3.5) for the estimation of the peak approximation error. $e_{max}$ is peak error obtained from $e = f_o - f_o*$; $e^*_{max}$ is the estimated peak error from equation (3.5); $\overline{T}_2 - \overline{T}_1$ is the time interval in which $W^{1/2} f_o$ and $W^{1/2} f_o*$ have appreciable value. Because the approximating functions have approximately the same length as $f_o$, $\overline{T}_2 - \overline{T}_1$ is taken to be slightly greater than the length of $f_o$. An exception is example ten where $W^{1/2} f_o$ is somewhat shorter than $\theta_4$ and $\theta_5$. Agreement between $e_{max}$ and $e^*_{max}$ is quite good. Examples seven through ten should be discounted to a certain degree because of the computer errors in $I_{rel}$.

# VIII.  THE ARBITRARY INPUT PROBLEM

Many practical applications of linear system approximation, such as the determination of linear system characteristics from experimentally obtained response data or the synthesis of amplifiers with pulse inputs, lead naturally to prescribed inputs which are arbitrary functions of time.  In these applications the approximation problem becomes quite complicated.  No longer are the response functions $\Theta_1(t)$, $-----\Theta_N(t)$ sums of exponential functions.  Consequently, most of the practical procedures of Chapter VII (for example, the simple orthogonalization of approximating functions and the computer evaluation of approximation coefficients) are not applicable.  Furthermore, error questions make doubtful the common practice of replacing the arbitrary input problem by the corresponding impulse response problem.  Not only may there be an error in the determination of $h(t)$ from $f_0(t)$ and $f_i(t)$, but as Chapter III indicates, direct approximation of $h(t)$ and $f_0(t)$ yield different error distributions in the frequency domain.  This chapter attempts to resolve some of the difficulties of the WME arbitrary input problem by presenting simplified analytic and computer techniques.

Since the theory of WME approximation, as it pertains to the arbitrary input problem, has been discussed fully in Chapter IV, the following sections will be concerned primarily with modifications of the theory and with detailed methods of coefficient evaluation.  The first section describes the reduction of the arbitrary input problem to the easily solved exponential input problem.  After several methods for analog computation of approximation coefficients are developed, various

procedures for handling experimental data are discussed. The chapter

concludes with a WME method for correlation function computation.

### The Equivalent Exponential Input Problem

Chapter II showed that the effect of an arbitrary input is to

weight the error $E_h = H - H^*$ by $F_i(j\omega)$. As a result minor changes in

$f_i(t)$ (more precisely, in $F_i(j\omega)$ ) do little to change $H^*$. Thus $f_i(t)$

may be replaced by an exponential series approximation, $\overline{f}_i(t)$, with

little change in the final approximation. Since the exponential input

yields exponential approximation functions, the new approximation pro-

blem is easier to solve.

Statement of the modified approximation problem is straight-

forward. $f_i(t)$ is replaced by $\overline{f}_i(t)$ and $f_0(t)$ is replaced by $\overline{f}_0(t) =$

$\int_0^\infty h(\tau)\overline{f}_i(t-\tau)d\tau$ ; i.e. $\overline{F}_0 = H\overline{F}_i = \dfrac{F_0}{F_i}\overline{F}_i$. The approximate system

function $\overline{H^*}$ which is obtained from the approximation to $\overline{F}_0$, $\overline{F_0^*} = \overline{H^*\overline{F}_i}$,

corresponds closely to $H^*$ of the original problem. It is important to

note that $\overline{f_0^*}(t)$ is made to approximate $\overline{f}_0(t)$ and not $f_0(t)$. Approxima-

tion of $f_0(t)$ would cause the difference between $\overline{H^*}$ and $H^*$ to be highly

dependent on the accuracy of the approximate input $\overline{f}_i(t)$.

The major difficulty in formulating the equivalent exponential

input problem is the determination of $\overline{f}_0(t)$. This is especially true

when $h(t)$ is given implicitly by specification of rather arbitrary $f_0(t)$

and $f_i(t)$. On the other hand, if $f_0(t)$ and $f_i(t)$ have known Laplace

transforms, no problem exists since $\overline{F}_0(s) = \dfrac{F_0(s)}{F_i(s)}\overline{F}_i(s)$. Many methods,

such as those of the previous chapter, are suitable for obtaining the

exponential series approximation $\overline{f}_i(t)$. When $W(t) = 1$ the freedom in

choice of $\overline{F}_i(t)$ is greater since only the magnitudes of $F_i(j\omega)$ and $\overline{F}_i(j\omega)$

must be matched.[1] This suggests various frequency domain approximation methods [41].

Once the preceding steps have been completed, actual solution of the approximation problem is simple. Since $\overline{F}_o^*(s) = \overline{H}^*(s) \, \overline{F}_i(s)$ is a rational function of $s$, all the terms of $\overline{f}_o^*(t)$ are exponential functions. Hence, the approximation procedures of Chapter VII are all valid. The system function is given by $\overline{H}^* = \dfrac{\overline{F}_o^*}{\overline{F}_i}$ . If $\overline{H}^*$ is not to contain as poles and zeros, the zeros and poles of $\overline{F}_i$ , the approximation $\overline{F}_o^*$ must be constrained to possess the factor $\overline{F}_i$. For example, if $\overline{F}_i = \dfrac{1}{s+1}$ and it is specified that $|s| \overset{\lim}{\to \infty} s \, \overline{H}^* = $ const. and that $\overline{H}^*$ does not have a zero at $s = -1$, then $\overline{F}_o^*$ must contain at least one pole at $s = -1$ and must be constrained so $\underset{|s| \to \infty}{\lim} s^2 \overline{F}_o^* = $ const.

While increasing the likelihood of more constraints, the above solution of the arbitrary input problem is little more complicated than the solution of the impulse response problem. Analytic solutions are perfectly feasible when $f_o(t)$ and $f_i(t)$ have known Laplace transforms. The steps are enumerated as follows:

1. Approximate $f_i(t)$ with as few exponential functions as considered necessary, to obtain $\overline{f}_i(t)$.

2. Calculate $\overline{F}_o(t)$ (or its transform $\overline{F}_o(s)$), the response of $H = \dfrac{F_o}{F_i}$ to $\overline{f}_i(t)$.

3. Approximate $\overline{F}_o(t)$ by the constrained exponential series $\overline{f}_o^*$ so that prescribed constraint conditions on $\overline{h}^*(t)$ (or $\overline{H}^* = \dfrac{\overline{F}_o^*}{\overline{F}_i}$) are met.

4. Obtain the approximate system from $\overline{H}^* = \dfrac{\overline{F}_o^*}{\overline{F}_i}$ .

---

[1] See Chapter III, page 23.

Some of the simpler arbitrary inputs do not require an exponential function representation. Suppose $f_i(t) = u_{-1}(t)$, the unit step function, and

$$f_o(t) = 0 \quad , \quad t \leq \frac{T}{2}$$
$$= (t-\frac{T}{2}) \frac{2}{T} \ , \ \frac{T}{2} < t \leq T$$
$$= 1 \quad , \quad T < t \ ,$$

(8.1)

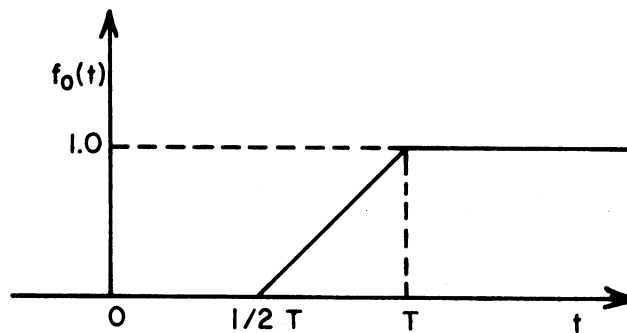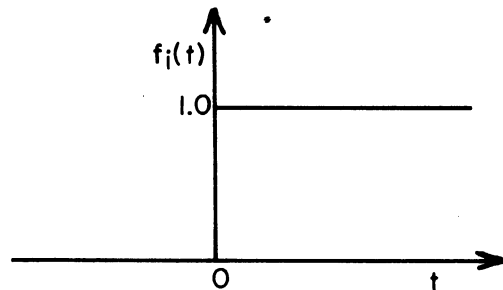as in Figure 8.1. Since $\int_0^\infty f_o^2 \, dt = \infty$, the problem cannot be solved



Figure 8.1 Prescribed Input and Response for Example

directly. By taking

$$\bar{f}_o(t) = u_{-1}(t) - f_o(t) = 0 \qquad , \qquad t \leq 0$$
$$= 1 \qquad , \qquad 0 < t \leq \frac{T}{2}$$
$$= 2 - \frac{2}{T}t \qquad , \qquad \frac{T}{2} < t \leq T$$
$$= 0 \qquad , \qquad T < t$$

$$(8.2)$$

the difficulty is removed.[1] Once $\bar{f}_o^*$ has been obtained (see examples

eight, nine, and ten of the previous chapter), $f_o^* = u_{-1} - \bar{f}_o^*$ and

$$H^* = \frac{F_o^*}{F_i} = s \left[\frac{1}{s} - \bar{F}_o^*\right] = 1 - s\bar{F}_o^* \qquad (8.3)$$

Unless $\bar{F}_o^*$ is a constrained approximation $|s| \overset{\lim}{\to} \infty H = const..$ Approxi-

mation for prescribed step function response has assured the condition

$H(0) = H^*(0) = 1..$

### Computer Evaluation of Coefficient Integrals

When it is impossible or impractical to reduce the arbitrary

input problem to a form where the approximating functions are exponential

functions, direct calculation of the integrals for $c_{nm}$ and $d_n$ is necessary.

Usually, the functions to be integrated are not elementary and some com-

puter technique must be employed. This section presents practical com-

puter diagrams for this purpose. Although the discussion is limited to

the integrals for $c_{nm}$ and $d_n$, similar procedures may also be applied to

other integrals, such as those which often appear in constraint equations.

For convenience, the equations defining $c_{nm}$ and $d_n$ are repeated

here.

---

1 This is really the preliminary simplification of extraction. See page
   15.

$$d_n = \int_{T_1}^{T_2} W(t)\ f_0(t)\ \theta_n(t)\ dt$$

$$= \int_{T_1}^{T_2} W(t)\ f_0(t) \int_0^\infty f_i(t-\tau_1)\ \varphi_n(\tau_1)\ d\tau_1\ dt \qquad (8.4)$$

$$c_{nm} = \int_{T_1}^{T_2} W(t)\ \theta_n(t)\ \theta_m(t)\ dt$$

$$= \int_{T_1}^{T_2} W(t) \int_0^\infty f_i(t-\tau_1)\ \varphi(\tau_1)\ d\tau_1 \int_0^\infty f_i(t-\tau_2)\ \varphi(\tau_2)\ d\tau_2\ dt$$

$$(8.5)$$

The physical systems $\Phi_n(s)$ which make up the final approximation $H^*$ are used, as shown in Figure 8.2, to obtain the required functions $\theta_n(t)$. As expected, the circuit for the computation of $c_{nm}$
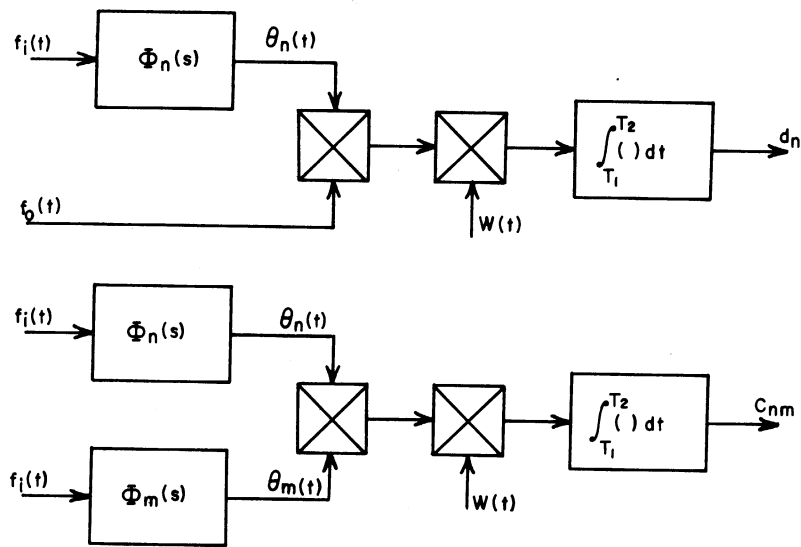


Figure 8.2  Diagrams for Direct Computation of $d_n$ and $c_{nm}$

requires only the prescribed input $f_i(t)$ and not $f_0(t)$. If all the integrals are computed simultaneously, $\dfrac{N(N+3)}{2}$ setups of the indicated type are required. By storing the functions $f_0$ and $f_i$ the integrations may be performed consecutively with only one setup, changing as required

the system functions $\Phi_n(s)$.

Sometimes the prescribed input and response functions are specified in terms of the auto-and cross-correlations functions

$$\psi_{ii}(\tau) = \frac{1}{T_2-T_1}\int_{T_1}^{T_1}f_i(t)\,f_i(t+\tau)\,dt \qquad \text{[1]}$$

(8.6)

$$\psi_{oi}(\tau) = \frac{1}{T_2-T_1}\int_{T_2}^{T_2}f_o(t)\,f_i(t+\tau)\,dt$$

(8.7)

This is especially true when $f_i(t)$ and $f_o(t)$ are the results of experimental measurements. The interval $(T_1,T_2)$ is assumed to be long enough so that changes in $T_1$ and $T_2$ which are as large as the maximum considered $|\tau|$ values, produce little change in $\psi_{ii}$ or $\psi_{oi}$. Under these conditions

$$\psi_{oi}(-\tau_1) = \frac{1}{T_2-T_1}\int_{T_1}^{T_2}f_o(t)\,f_i(t-\tau_1)\,dt$$

$$\psi_{ii}(\tau_2-\tau_1) = \frac{1}{T_2-T_1}\int_{T_1}^{T_2}f_i(t)\,f_i(t+\tau_2-\tau_1)\,dt$$

(8.8)

$$= \frac{1}{T_2-T_1}\int_{T_1+\tau_2}^{T_2+\tau_2}f_i(t-\tau_2)\,f_i(t-\tau_1)\,dt$$

$$\cong \frac{1}{T_2-T_1}\int_{T_1}^{T_2}f_i(t-\tau_2)\,f_i(t-\tau_1)\,dt$$

(8.9)

Upon changing the order of integration in equations (8.4) and (8.5) it is seen that

$$d_n = (T_2-T_1)\int_0^\infty \psi_{oi}(-\tau_1)\,\varphi_n(\tau_1)\,d\tau_1$$

(8.10)

---

1 See James, et al. [16]. Normally $T_1$ and $T_2$ go to $-\infty$ and $\infty$. Since practically they remain finite, the limit will not be taken here.

$$c_{nm} = (T_2 - T_1) \int_0^\infty \varphi_n(\tau_2) \int_0^\infty \psi_{ii}(\tau_2 - \tau_1) \varphi_m(\tau_1) \, d\tau_1 \, d\tau_2$$

$$(8.11)$$

when $W(t) = 1$.

Figure 8.3 gives the diagrams for evaluating $c_{nm}$ and $d_n$ when the known correlation functions are generated as time functions. As in Chapter VII, the linear systems $\Phi_n(s)$ perform the desired integrations by convolution. A different method for computing $c_{nm}$ is shown
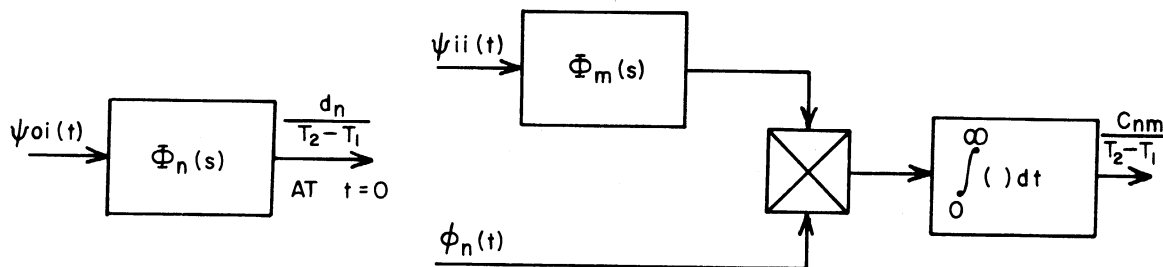


Figure 8.3  Diagrams for Computation of $d_n$ and $c_{nm}$ from Correlation Functions



Figure 8.4  Diagram for Computation of $c_{nm}$ from Correlation Function

in Figure 8.4. Here, the time function $g(-t)$ must be reproduced from $g(t)$ and used as a forcing function for the second system.[1]

The above integration schemes are, of course, subject to computer errors. Analysis of the errors as they influence $f_o^*$ is more complicated than in Chapter VII because both $c_{nm}$ and $d_n$ are inaccurate. Without going into detail, the least difficulty is again encountered when the eigenvalues of $C^{-1}$ are approximately equal. Unfortunately, the functions $\Theta_n(t)$ depend on both $f_i(t)$ and $\Phi_n(s)$ so orthogonalization of $\Theta_1, ----- \Theta_N$ is not attractive. The only thing that can be done is to choose $\Phi_1, ----- \Phi_N$ in such a manner as to make $\Theta_1, ----- \Theta_N$ approximately orthonormal. This approach is discussed further in the next section.

### Experimental Determination of System Characteristics

The problem of determining the system function (H) of an unknown physical system from input and response data ($f_i$ and $f_o$) is an important one which has been investigated by many authors. Although many methods have been proposed, most of them rely either on numerical analysis [24,32,33,42] or correlation function measurement[2] [12,23]. The purpose of this section is to show some of the relations between WME approximation and experimental system function determination.[3] No claim of originality is made for all of the material.

---

1 Laning and Battin [21] (page 218) describe a similar integral evaluation for a different application.

2 Truxal [38] gives an excellent review of the problems involved in determining correlation functions from experimental data.

3 Booton [4] describes a general mean square error method for measurement and representation of system characteristics.

Review of the Problem--Most practical measurement problems are complicated by the fact that the input to the physical system under test consists of two components: (1) the known test input $f_i(t)$, and (2) an undesired and usually unknown input $n_i(t)$. The function $n_i(t)$ may be an irremovable noise input, or it may be an input required for continued operation of the system under test. Correspondingly, the measured response consists of the two components, $f_o(t)$ and $n_o(t)$, resulting separately from each of the two inputs. From measurement of the operating records, $f_i(t)$ and $f_o(t) + n_o(t)$, the unknown system function must be approximated as closely as possible.

Auto- and cross-correlation functions offer one solution of the problem. The functions

$$\frac{1}{T_2-T_1} \int_{T_1}^{T_2} f_i(t)\, f_i(t+\tau)\, dt = \psi_{ii}(\tau)$$

$$\frac{1}{T_2-T_1} \int_{T_1}^{T_2} [f_o(t) + n_o(t)]\, f_i(t+\tau)\, d\tau = \psi_{oi}(\tau) + \frac{1}{T_2-T_1} \int_{T_1}^{T_2}$$

$$n_o(t)\, f_i(t+\tau)\, dt$$

are computed by one of the conventional methods. If the time interval $(T_2-T_1)$ is sufficiently long, and if $n_o(t)$ and $f_i(t)$ are uncorrelated, the second integral of the second equation vanishes so that $\psi_{ii}(\tau)$ and $\psi_{oi}(\tau)$ are known. Because $\psi_{ii}(\tau)$ and $\psi_{oi}(\tau)$ are related by

$$\psi_{oi}(-t) = \int_0^{\infty} h(\tau)\, \psi_{ii}(t-\tau)\, d\tau \qquad [1]$$

(8.12)

the unknown weighting function $h(t)$ is also known. Various methods for solving the integral equation have been employed. Most of them are standard linear system approximation procedures because $\psi_{ii}(t)$ and

---

1 See Truxal [38], page 437.

$\psi_{oi}(-t)$ are related, as inputs and outputs of the system $H(s)$ are related.

The WME Approximation--The most obvious method of WME approximation is to reduce the aforementioned correlation function problem to an equivalent arbitrary input-response problem, i.e., let $\overline{f}_i(t) = \psi_{ii}(t)$ and $\overline{f}_o(t) = \psi_{oi}(-t)$. This method has the deficiency that the new input $\overline{f}_i$ does not produce the same frequency weighting of approximation error that $f_i$ produces. A more direct and satisfactory procedure for approximating $h(t)$ is the coefficient evaluation scheme of equations (8.10) and (8.11). In this approach, the original test input $f_i(t)$ is retained even though $c_{nm}$ and $d_n$ are calculated from $\psi_{ii}$ and $\psi_{oi}$.

A still more direct approach does not require the added work and error of preliminary calculation of the correlation functions. The coefficients $c_{nm}$ and $d_n$ are obtained from $f_i(t)$, $f_o(t)$, and the physical systems $\Phi_n(s)$ as in Figure 8.2. The components for computation of $d_n$ ($W(t) = 1$) are shown in Figure 8.5 as they would appear in an actual test situation; the components for computation of $c_{nm}$ remain as in Figure 8.2. Presence of the undesired output $n_o(t)$ causes little error in $d_n$ if the
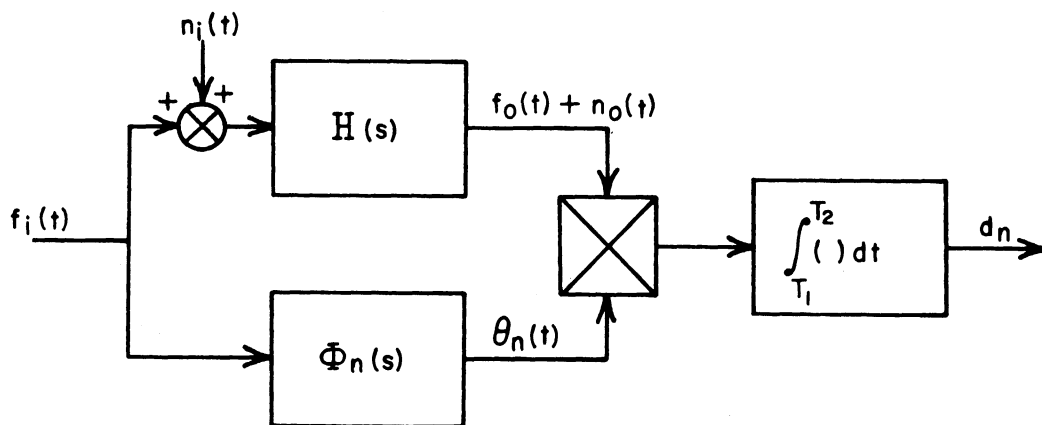


Figure 8.5 Diagram for Experimental Determination of $d_n$

period $(T_2-T_1)$ is large, and $n_o(t)$ and $f_i(t)$ are uncorrelated. When all coefficients are computed simultaneously, a great number of integrating

and multiplying components is necessary. Consecutive computation of

$c_{nm}$ and $d_n$ is possible if $f_i(t)$ and $f_0(t) + n_0(t)$ are stored or if $f_i$

is stationary so that any $(T_1, T_2)$ gives the same result for large $(T_2-T_1)$.

Long test periods are not mandatory when $n_i(t)$ is small or zero. A

period several times longer than the period for which $h(t)$ has appreciable

value should be adequate.

White Noise Testing and Function Orthogonality--All of the

above methods become greatly simplified when $W(t) = 1$, $(T_2-T_1)$ is large,

and $f_i(t)$ is a white noise. Since $\psi_{ii}(\tau) = Pu_0(\tau)$ [1] the coefficients

$c_{nm}$ defined by equation (8.11) become simply

$$c_{nm} = P (T_2-T_1) \int_0^\infty \varphi_n(\tau_1) \varphi_m(\tau_1) \, d\tau_1 \qquad (8.13)$$

But the functions $\varphi_n(t)$ are known sums of exponential functions so

the $c_{nm}$ can be calculated and do not need to be measured experimentally.

In fact, when the $\varphi_n$ are orthonormal

$$c_{nm}^{-1} = \frac{1}{P(T_2-T_1)} \qquad (8.14)$$

If the functions $\varphi_n$ are not orthonormal but are chosen from the appendix,

the tabulated inverse matrices yield $C^{-1}$ upon multiplication by $\frac{1}{P(T_2-T_1)}$.

The constants $d_n$ must still be obtained experimentally.

Arguments concerning the eigenvalues of $C^{-1}$ and sensitivity

of the approximation error to coefficient errors hold equally well when

applied to white noise testing or impulse response approximation. There-

fore, the use of orthonormal $\Phi_n$ functions is again indicated. If $f_i(t)$

---

1 P is the power per unit bandwidth. In practice $F_i(j\omega)$ is band limited and $\psi_{ii}$ is not an exact impulse function. If the spectral density is flat well beyond the response frequencies of the $\Phi_n$, then the representation is satisfactory.

is not white noise, approximate orthogonality of $\Theta_1, ----- \Theta_n$ can be had by letting $\Phi_n = G\overline{\Phi}_n$, where $\overline{\Phi}_1, ----- \overline{\Phi}_N$ are orthonormal and G is a rational system function which changes $f_i(t)$ into a function which approximates white noise.[1] The complexity of the functions $\Phi_1, ----- \Phi_N$ should not be too great since G can usually be kept rather simple. Because $\Theta_1, ----- \Theta_N$ are not exactly orthonormal, C must be inverted; however, the approximate orthogonality assures good accuracy.

Choice of Approximating Functions--The choice of the approximating functions $\Phi_1, ----- \Phi_N$ remains a difficulty in the experimental approximation problem. Generally, little is known about the exact pole positions of the unknown system function H(s). It is even possible that H(s) does not possess a finite number of poles. Sometimes, the physical system does give information about the form and approximate pole positions of H(s). This is often true of the linear equations which describe the small motion response of an airplane; known airplane data enable a reasonably good estimate of pole positions. In such systems, the estimated poles plus a few added poles should provide an excellent basis for choice of $\Phi_1, ----- \Phi_N$.

In other cases, the poles may often be obtained by examination of the transient response of the unforced system. Graphical determination of decay rates and oscillation frequencies is one such approach. Another is physical mechanization of the Prony method. If no success is had in pole determination, accuracy of the final approximation must rest primarily on a sufficiently large number of fairly arbitrary approximating functions.

---

1 The determination of G is straightforward. See Goldman [11], page 264.

Multiple Input and Output Systems --Multiple input and output

systems are approximated by simple extension of the above principles.

A separate approximation is found for each of the input-output pairs.

Although the calculations are no more complicated than those already

discussed, the number of calculations mounts rapidly with multiplicity

of inputs and outputs. When the same approximating functions are used

for all approximations in a single system, some duplication in approxi-

mation and realization is avoided. For example, single input, multiple

output systems lead to the realization of Figure 8.6. Only one set of

approximating systems is required. By introducing the coefficient gains

and summing operations at the inputs of the $\Phi_n(s)$ a similar reduction

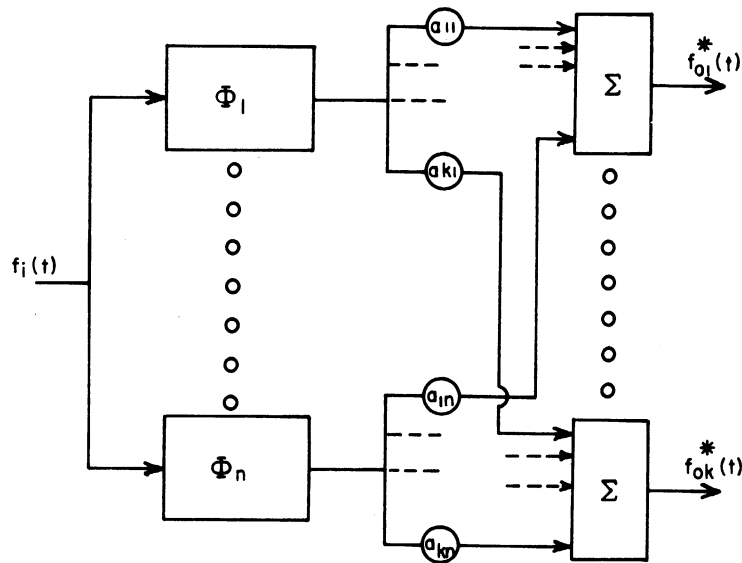in complexity is obtained for single output, multiple input systems.



Figure 8.6. Diagram for Realization of Single Input, Multiple Output System.

If multiple input, multiple output systems are to be realized, use of

the same schemes is not possible without duplication of each of the

approximating systems. However, the common poles in all approximations

do not permit the realization of a linear system with only M dynamic

elements (M is the number of poles in the approximating functions).

Optimum Filter Design[1]-Procedures which are valid for exper-imental determination of system characteristics are also valid for experimental design of optimum filters. The unknown system is simply replaced by an ideal system which performs exactly the prescribed filtering task. Even though such an ideal system is not mathematically attainable, the input and response functions which are required for mean square approximation of the system are available.

The problem is illustrated in Figure 8.7. The weighted mean square error I is minimized in the usual way through the approximation



IDEAL FILTER

$H(s)$

$f_0(t)$ DESIRED RESPONSE (DEPENDENT ON $f_i(t)$ )

$f_i(t) + n(t)$
INPUT PLUS NOISE

APPROXIMATE FILTER

$+$ $e(t)$
$-$

$H^*(s)$

$f_0^*(t)$ APPROXIMATION TO $f_0(t)$

Figure 8.7. Optimum Filter Problem.

coefficients $a_n$ defined by $c_{nm}$ and $d_n$. To obtain $f_0(t)$ it is presumed that $f_i(t)$ and $n(t)$ can be separated during the experimental design measurements. Typical optimum filter operations which might be spec-ified include: (1) smoothing, $f_0(t) = f_i(t)$ ; (2) delayed smoothing $f_0(t) = f_i(t-T)$ ; (3) prediction, $f_0(t) = f_i(t+T)^2$; (4) differentiation, $f_0(t) = f_i'(t)$. Since the approximate system is time invariant, the

---

1 Bose [5] has recently proposed methods similar to the ones of this section.

2 The time advance $f_0(t)=f_i(t+T)$ would be obtained by working with delayed input and output functions, $f_0(t-T)=f_i(t)$ and $f_i(t-T) + n(t-T)$.

statistical properties of $f_i(t)$ and $n(t)$ should also be time invariant, i.e. stationary. Moreover, the time averaging period $(T_2-T_1)$ must be sufficiently large to reduce the statistical deviations in $c_{nm}$ and $d_n$ to a reasonable point.

The optimum filter design of the preceding paragraph is only optimum in terms of the approximating functions $\Phi_1,----- \Phi_N$. With an infinite set of complete approximating functions the approximation error would be decreased even further. While this "finiteness" error is decreased by increasing N, there is no reason for going beyond the N where the finiteness error becomes small compared to the filtering error.

Experimental design of optimum filters has only been touched upon here. Certain classes of non-stationary inputs could be treated by constrained approximation [43]. Introduction of distributed element approximating functions would permit design of finite memory filters [43]. Many other questions, such as practicality, accuracy, and computational complexity can be settled only by completing some trial problems.

### Correlation Function Computation

Although the subject of correlation function computation cannot be classified under linear system approximation, the following method is so closely related to the material of the above sections that it merits mention.[1] Conventional methods of computing the cross-correlation function of $f(t)$ and $g(t)$

---

1 The same method of correlation function computation was discovered earlier by Lampard [20]. However, his work is restricted entirely to orthogonal Laguerre functions.

$$\psi_{fg}(\tau) = \frac{1}{T_2-T_1} \int_{T_1}^{T_2} f(t)\, g(t+\tau)\, dt$$

$$\cong \frac{1}{T_2-T_1} \int_{T_1}^{T_2} f(t-\tau)\, g(t)\, dt \tag{8.15}$$

require the operations of delay, multiplication, and integration. In the present section $\psi_{fg}$ is represented by a mean square approximation of exponential functions, either orthonormal or linearly independent.

For $\tau > 0$ ($\tau < 0$ will be considered later), let $\psi_{fg}$ be approximated by

$$\psi_{fg}^{*}(\tau) = \sum_{n=1}^{N} a_n \varphi_n(\tau) \qquad \tau > 0 \tag{8.16}$$

As in previous approximations, $a] = C^{-1} d]$ where $C$ is known from $\varphi_1, \text{-----}\ \varphi_N$ and $d]$ is defined by



Figure 8.8  Physical Systems for Computing Correlation Function Coefficients $d_n$

$$d_n = \int_0^\infty \psi_{fg}(\tau)\, \varphi_n(\tau)\, d\tau$$

$$= \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} g(t) \int_0^\infty f(t-\tau)\, \varphi_n(\tau)\, d\tau\, dt$$

$$= \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} g(t)\, \Theta_n(t)\, dt \tag{8.17}$$

The function $\Theta_n(t)$ is the response of the system $\Phi_n$ to the input $f(t)$.

Thus, the circuit of Figure 8.8 gives $d_1, \text{-----}\ d_N$. In order to

approximate $\psi_{fg}$ for $\tau < 0$ it is necessary to note that $\psi_{gf}(\tau) = \psi_{fg}(-\tau)$.

Hence the $\overline{d}_n$, which hold for the negative $\tau$ approximation, are computed

by reversing the positions of $f(t)$ and $g(t)$. The complete approximation

is

$$\psi_{fg}^* = \sum_{n=1}^N [a_n \varphi_n(\tau) + \overline{a}_n \varphi_n(-\tau)] \tag{8.18}$$

The cross-spectral density is the Fourier transform of equation (8.18).

$$\Psi_{fg}(\omega) = \sum_{n=1}^N [a_n \Phi_n(j\omega) + \overline{a}_n \text{conj.}\ \Phi_n(j\omega)] \tag{8.19}$$

When $g = f$ the auto-correlation function is obtained. Since

$\psi_{ff}(\tau) = \psi_{ff}(-\tau)$ equations (8.18) and (8.19) simplify to

$$\psi_{ff}^*(\tau) = \sum_{n=1}^N a_n [\varphi_n(\tau) + \varphi_n(-\tau)] \tag{8.20}$$

$$\Psi_{ff}(\omega) = 2 \sum_{n=1}^N a_n \text{Re}\ \Phi_n(j\omega) \tag{8.21}$$

Mean square error computation of correlation functions has the following advantages: (1) no delay device is required, (2) $f(t)$ and $g(t)$ need not be stored for repeated calculations since the $d_n$ may be computed simultaneously, (3) functional representations are obtained for the correlation functions and spectral densities. Disadvantages are:(1) lack of satisfactory method for choosing approximating functions, (2) inability to calculate $I_{max} = \int_{-\infty}^{\infty} \psi_{fg}^2 \, d\tau$ and thus the mean square approximation error. The last difficulty can be eased by checking the approximation error in $\psi_{fg}$ or $\psi_{ff}$ at $\tau = 0$ or the value of $\Psi_{ff}(\omega)$ at several frequencies.

# IX. CONCLUSION

The primary goal of the preceding chapters has been to extend the theory and practicality of linear system approximation by mean square error minimization in the time domain. Success in achieving this goal can be attributed to two considerations. First, emphasis has been placed on the general theory of mean square approximation rather than on a specialized aspect of it, such as approximation by restricted classes of orthogonal functions. Second, recognition has been made of the difficulty of practical computations by supplementing available analytic procedures by matrix tables and computer methods.

Since the effort to present a comprehensive exposition has necessarily required the inclusion of supporting material, it is proper to summarize the major results which are believed to be new. They are as follows:

1. A detailed investigation of the mean square error criterion, its meaning in the frequency domain, and its relation to the arbitrary input problem. A simple formula for estimating the peak approximation error from the mean square error.

2. A general statement and solution of the mean square error approximation problem for non-orthogonal approximating functions and an arbitrary input. The extension of constrained approximations to the general mean square approximation problem and the introduction of constrained approximating functions.

3. The generation of constrained orthonormal exponential functions and of more general unconstrained orthonormal exponential functions.

4. The approximation of prescribed responses having known Laplace transforms by means of non-orthogonal exponential functions and exponentially weighted non-orthogonal exponential functions.

5. Tables of inverse matrices to simplify approximation calculations for various families of exponential approximating functions.

6. The analog computation of approximation coefficients through linear system realization of approximating functions, including system simulation and approximation error evaluation, orthonormal function realization, pole optimization through orthonormal function computation, computation of exponentially weighted integrals, and error analysis.

7. A simplified solution of the arbitrary input problem through the use of an equivalent exponential input.

8. The computer evaluation of approximation coefficients for arbitrary or experimentally measured input and response functions. The experimental design of optimum linear filters.

9. The application of the general WME approximation theory to correlation function determination.

Although the above results are encouraging, the mean square approximation of linear system response has several important limitations. Most seriously, the pole determination problem has no completely satisfactory solution. The more exact methods of Chapter VI are lengthy and

offer no real assurance that the poles chosen are close to optimum.
When the prescribed system is lumped and of relatively low order (a
situation which is only of significance in experimental determination
of unknown system characteristics) and is excited by a simple input,
the Prony method or variations of it provide good pole determination.
For more complicated conditions the poles must be estimated as best as
possible and a sufficient number of them included to meet the prescribed
accuracy requirements.  In this respect the WME approximation procedure
is no more deficient than any other presently known approximation pro-
cedure.

Finally, the calculations tend to become excessively lengthy
when the number of approximating functions is large, when the weight
factor is not constant, when several constraints are applied, or when
the input is approximated by more than a few exponential functions.  This
difficulty should not be considered so much as a fault of the WME method,
as an indication of the complexity  of the problem being solved.  In
many cases the length of calculations could be reduced by a more complete
tabulation of inverse matrices for unconstrained, constrained and weighted
approximating functions.  A still more general approach would be the
programming of the approximation problem on a large scale digital com-
puter.

The following possibilities for further work in mean square
approximation of linear systems exist:  (1) approximation by distributed
element linear systems, (2) approximation of time-variant linear systems
by time-variant approximating systems, (3) investigation of the feasi-
bility of experimental approximation methods.  Preliminary study has
shown that (2) gives promise of useful results.

APPENDIX

Inverse Matrix Tables

The following families of non-orthogonal approximating functions are defined in the interval $(0,\infty)$ and are unweighted. In addition to the tabulation of the inverse C matrices, formulas are given for the calculation of $d_1,\text{-----}\ d_N$. To aid in the choice of time scaling, an indication of the length of the longest approximating function in each family is also included.

1. Unconstrained Laguerre Set

Poles: all at $s = -1$.

System functions: $\Theta_n = \dfrac{1}{(s+1)^n}$

Time functions: $\Theta_n = \dfrac{t^{n-1}}{(n-1)!} e^{-t}$

Time of maximum value for longest function: N-1 seconds

Coefficients: $c_{nm} = \dfrac{(n + m - 2)!}{(n-1)!(m-1)! \, 2^{n+m-1}}$

$$d_n = \int_0^\infty \Theta_n f_o \, dt = \int_0^\infty \frac{t^{n-1}}{(n-1)!} e^{-t} f_o(t) \, dt$$

$$= \frac{(-1)^{n-1}}{(n-1)!} \left. \frac{d^{n-1} F_o}{d s^{n-1}} \right|_{s=1}$$

Inverse Matrices:

$$\begin{vmatrix} 4 & -- \\ -4 & 8 \end{vmatrix} \qquad \begin{vmatrix} 6 & -- & -- \\ -12 & 40 & -- \\ 8 & -32 & 32 \end{vmatrix} \qquad \begin{vmatrix} 8 & -- & -- & -- \\ -24 & 112 & -- & -- \\ 32 & -176 & 320 & -- \\ -16 & 96 & -192 & 128 \end{vmatrix}$$

$$\begin{vmatrix} 10 & -- & -- & -- & -- \\ -40 & 240 & -- & -- & -- \\ 80 & -560 & 1{,}472 & -- & -- \\ -80 & 608 & -1{,}728 & 2{,}176 & -- \\ 32 & -256 & 768 & -1{,}024 & 512 \end{vmatrix}$$

$$\begin{vmatrix} 12 & -- & -- & -- & -- & -- \\ -60 & 440 & -- & -- & -- & -- \\ 160 & -1,360 & 4,672 & -- & -- & -- \\ -240 & 2,208 & -8,128 & 14,976 & -- & -- \\ 192 & -1,856 & 7,168 & -13,824 & 13,312 & -- \\ -64 & 640 & -2,560 & 5,120 & -5,120 & 2,048 \end{vmatrix}$$

2. Constrained Laguerre Set

Poles: all at $s = -1$.

System functions: $\Theta_n = \dfrac{1}{(s+1)^{n+1}}$

Time functions: $\Theta_n = \dfrac{t^n}{n!} e^{-t}$

Time of maximum value for longest function: N seconds

Coefficients: $c_{nm} = \dfrac{(n + m)!}{n! \; m! \; 2^{n+m+1}}$

$$d_n = \int_0^\infty \Theta_n f_0 \, dt = \int_0^\infty \frac{t^n}{n!} e^{-t} f_0(t) \, dt$$

$$= \frac{(-1)^n}{n!} \left. \frac{d^n F_0}{d s^n} \right|_{s=1}$$

Inverse Matrices:

$$\begin{vmatrix} 16 & -- \\ -16 & 21\frac{1}{3} \end{vmatrix}$$

$$\begin{vmatrix} 40 & -- & -- \\ -80 & 192 & -- \\ 48 & -128 & 96 \end{vmatrix}$$

$$\begin{vmatrix} 80 & -- & -- & -- \\ -240 & 832 & -- & -- \\ 288 & -1{,}088 & 1{,}536 & -- \\ -128 & 512 & -768 & 409\frac{6}{10} \end{vmatrix}$$

$$\begin{vmatrix} 140 & -- & -- & -- & -- \\ -560 & 2{,}538\frac{2}{3} & -- & -- & -- \\ 1{,}008 & -4{,}928 & 10{,}176 & -- & -- \\ -896 & 4{,}608 & -9{,}984 & 10{,}240 & -- \\ 320 & -1{,}706\frac{2}{3} & 3{,}840 & -4{,}096 & 1{,}706\frac{2}{3} \end{vmatrix}$$

## 3. Real Exponential Set

Poles: at $s = -n$.

System functions: $\Theta_n = \frac{1}{s+n}$

Time functions: $\theta_n = e^{-nt}$

Time constant of longest function: 1. second

Coefficients: $c_{nm} = \frac{1}{n+m}$

$$d_n = \int_0^\infty \theta_n f_o \, dt = \int_0^\infty e^{-nt} f_o(t) \, dt$$
$$= F_o(n)$$

Inverse Matrices:

$$\begin{vmatrix} 18 & -- \\ -24 & 36 \end{vmatrix}$$

$$\begin{vmatrix} 72 & -- & -- \\ -240 & 900 & -- \\ 180 & -720 & 600 \end{vmatrix}$$

$$\begin{vmatrix} 200 & -- & -- & -- \\ -1,200 & 8,100 & -- & -- \\ 2,100 & -15,120 & 30,000 & -- \\ -1,120 & 8,400 & -16,800 & 9,800 \end{vmatrix}$$

$$\begin{vmatrix} 450 & -- & -- & -- & -- \\ -4,200 & 44,100 & -- & -- & -- \\ 12,600 & -141,120 & 471,000 & -- & -- \\ -15,120 & 176,400 & -604,800 & 793,800 & -- \\ 6,300 & -75,600 & 264,600 & -352,800 & 158,760 \end{vmatrix}$$

$$\begin{vmatrix} 882 & -- & -- & -- & -- & -- \\ -11,760 & 176,400 & -- & -- & -- & -- \\ 52,920 & -846,720 & 4,234,200 & -- & -- & -- \\ -105,840 & 1,764,000 & -9,702,000 & 19,845,000 & -- & -- \\ 97,020 & -1,663,200 & 8,731,800 & -19,404,000 & 19,209,960 & -- \\ -33,264 & 582,120 & -3,104,640 & 6,985,440 & -6,985,440 & 2,561,328 \end{vmatrix}$$

4. <u>Fourier Set $(\alpha = 1)$</u>

Poles:  at  $s = -1,\ -1\pm j,\ -1\pm j2,\ -1\pm j3$.

System functions:  $\Theta_1 = \dfrac{1}{s+1}$

$$\Theta_n = \frac{s+1}{s^2+2s+1+.25n^2} \qquad ,\ n\ \text{even}$$

$$= \frac{.5(n-1)}{s^2+2s+1+.25(n-1)^2} \qquad ,\ n\ \text{odd}$$

Time functions:  $\Theta_1 = e^{-t}$

$$\Theta_n = e^{-t}\cos .5nt \qquad\qquad ,\ n\ \text{even}$$

$$= e^{-t}\sin .5(n-1)t \qquad ,\ n\ \text{odd}$$

Time constant of longest function:  1. second

Coefficients:
$$d_1 = \int_0^\infty \Theta_1 f_0\ dt = \int_0^\infty e^{-t}f_0(t)\ dt\ =\ F_0(1)$$

$$d_n = \int_0^\infty \Theta_n f_0\ dt = \int_0^\infty f_0(t)e^{-t}\cos.5nt\ dt = \mathrm{Re}F_0(1-j.5n),$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad n\ \text{even}$$

$$= \int_0^\infty f_0(t)e^{-t}\sin .5(n-1)t\ dt =$$

$$\qquad\qquad\qquad\qquad \mathrm{Im}F_0(1-j.5(n-1)),$$

$$\qquad\qquad\qquad\qquad\qquad\qquad n\ \text{odd}$$

Inverse Matrices:

$$\begin{vmatrix} 50 & -- & -- \\ -40 & 36 & -- \\ -40 & 28 & 44 \end{vmatrix}$$

$$\begin{vmatrix} 200 & -- & -- & -- & -- \\ -80 & 100 & -- & -- & -- \\ -293\tfrac{1}{3} & 77\tfrac{7}{9} & 477\tfrac{7}{9} & -- & -- \\ -100 & -13\tfrac{1}{3} & 182\tfrac{2}{9} & 96\tfrac{2}{3} & -- \\ 66\tfrac{2}{3} & -48\tfrac{8}{9} & -102\tfrac{2}{9} & -21\tfrac{1}{9} & 47\tfrac{7}{9} \end{vmatrix}$$

$$\begin{vmatrix} +417.277 & -- & -- & -- & -- & -- & -- \\ + 38.5090 & +227.768 & -- & -- & -- & -- & -- \\ -693.318 & -127.759 & +1216.63 & -- & -- & -- & -- \\ -414.060 & -164.428 & + 764.415 & +557.085 & -- & -- & -- \\ - 57.7696 & -217.769 & +106.650 & +126.430 & +280.658 & -- & -- \\ - 12.8379 & - 75.1836 & + 9.99718 & - 3.26160 & +116.442 & +74.3945 & -- \\ +154.067 & + 96.6591 & -281.097 & -218.654 & -108.437 & -22.7395 & +111.771 \end{vmatrix}$$

5. <u>Fourier Set ($\alpha = 2$)</u>

Poles: at $s = -2$, $-2\pm j$, $-2\pm j2$, $-2\pm j3$

System functions: $\Theta_1 = \dfrac{1}{s+2}$

$$\Theta_n = \frac{s+2}{s^2+4s+4+.25n^2} \quad , \quad n \text{ even}$$

$$= \frac{.5(n-1)}{s^2+4s+4+.25(n-1)^2} \quad , \quad n \text{ odd}$$

Time functions: $\theta_1 = e^{-2t}$

$$= e^{-2t}\cos .5nt \quad , \quad n \text{ even}$$

$$= e^{-2t}\sin .5(n-1)t \quad , \quad n \text{ odd}$$

Time constant of longest function: .5 second

Coefficients:
$$d_1 = \int_0^\infty \theta_1 f_o \, dt = \int_0^\infty e^{-2t}f_o(t)\,dt = F_o(2)$$

$$d_n = \int_0^\infty \theta_n f_o \, dt = \int_0^\infty f_o(t)e^{-2t}\cos .5nt \, dt = \text{ReF}(2-j.5n),$$

$$n \text{ even}$$

$$= \int_0^\infty f_o(t)e^{-2t}\sin .5(n-1)dt =$$

$$\text{ImF}(2-j.5(n-1)),$$

$$n \text{ odd}$$

Inverse Matrices:

$$\begin{vmatrix} 1155.93 & -- & -- \\ -1087.94 & 1031.94 & -- \\ -543.971 & 495.973 & 327.986 \end{vmatrix}$$

$$\begin{vmatrix} -114.216 & -- & -- & -- & -- \\ 63.4051 & -2.5357 & -- & -- & -- \\ 209.899 & -137.873 & 147.812 & -- & -- \\ 57.2765 & -56.6722 & -59.6550 & -0.1193 & -- \\ -91.6699 & 41.9000 & -167.979 & 25.4848 & 178.428 \end{vmatrix}$$

$$
\begin{vmatrix}
744.520 & -- & -- & -- & -- & -- & -- \\
42.927 & -\ .046 & -- & -- & -- & -- & -- \\
-\ 307.778 & -84.290 & 1309.53 & -- & -- & -- & -- \\
49.614 & -56.645 & -\ 58.250 & -\ .039 & -- & -- & -- \\
-1453.43 & 17.557 & -517.647 & 42.064 & 3950.60 & -- & -- \\
-\ 746.282 & 14.258 & 376.995 & 6.934 & 1283.85 & 654.807 & -- \\
647.493 & 12.743 & 190.339 & -7.974 & -1826.77 & -612.515 & 885.057
\end{vmatrix}
$$

## 6. Butterworth Set

Poles: approximately at $s = e^{j\frac{\pi}{2}\left(1+\frac{2k-1}{N}\right)}$, $k = 1, - - N$

System functions: $\Theta_n = \dfrac{s+\alpha_n}{s^2+2\alpha_n s+\alpha_n^2+\beta_n^2}$

$$\Theta_{n+1} = \dfrac{\beta_n}{s^2+2\alpha_n s+\alpha_n^2+\beta_n^2}$$

Time functions:

$$\Theta_n = e^{-\alpha_n t}\cos \beta_n t$$

$$\Theta_{n+1} = e^{-\alpha_n t}\sin \beta_n t$$

Time constant of longest function: $\dfrac{1}{\alpha_N}$ seconds

Coefficients

$$d_n = \int_0^\infty \Theta_n f_o \, dt = \int_0^\infty f_o e^{-\alpha_n t}\cos \beta_n t \, dt = \mathrm{Re}F_o(\alpha_n - j\beta_n)$$

$$d_{n+1} = \int_0^\infty \Theta_{n+1} f_o \, dt = \int_0^\infty f_o e^{-\alpha_n t}\sin \beta_n t \, dt = \mathrm{Im}F_o(\alpha_n - j\beta_n)$$

Inverse Matrices and Pole Positions:

$\alpha_1 = .7, \quad \beta_1 = .7$

$$\begin{vmatrix} 2.8 & -- \\ -2.8 & 8.4 \end{vmatrix}$$

$\alpha_1 = 1., \quad \beta_1 = 0, \quad \alpha_2 = 0.5, \quad \beta_2 = 0.9$

$$\begin{vmatrix} -4.19002 & -- & -- \\ 1.79293 & 1.23279 & -- \\ 4.45686 & -3.01823 & -1.50614 \end{vmatrix}$$

$\alpha_1 = 0.9, \quad \beta_1 = 0.4, \quad \alpha_3 = 0.4, \quad \beta_3 = 0.9$

$$\begin{vmatrix} 24.3357 & -- & -- & -- \\ -54.7549 & 270.734 & -- & -- \\ -14.9758 & 33.6953 & 10.8159 & -- \\ 6.65580 & -55.3273 & -4.80699 & 15.0888 \end{vmatrix}$$

$$\alpha_1 = 1, \beta_1 = 0, \alpha_2 = 0.8, \beta_2 = 0.6, \alpha_4 = 0.3, \beta_4 = 1.$$

| | | | | |
|---|---|---|---|---|
| 527.998 | -- | -- | -- | -- |
| -545.360 | 578.690 | -- | -- | -- |
| -198.029 | 184.086 | 138.573 | -- | -- |
| 49.8568 | -59.8808 | - 7.09480 | 10.4793 | -- |
| 50.0106 | -49.5157 | -34.1496 | 3.03990 | 10.6929 |

# BIBLIOGRAPHY

1. Aigrain, P. R., and Williams, E. M., "Synthesis n-Reactance Networks for Desired Transient Response," Journal of Applied Physics, Vol. 20, pp.597-600; June, 1949.

2. Aigrain, P. R., and Williams, E. M., "Design of Optimum Transient Response Amplifiers," Proceedings of the Institute of Radio Engineers, Vol 37, pp. 873-879; August, 1949.

3. Ba Hli, F., "A General Method for Time Domain Synthesis," Transactions of the Institute of Radio Engineers, Vol. CT-1, pp. 21-29; September, 1954.

4. Booton, R. C., Jr., "The Measurement and Representation of Nonlinear Systems," Transactions of the Institute of Radio Engineers, Vol. CT-1, pp. 32-35; December, 1954.

5. Bose, A. G., "A Theory for the Experimental Determination of Optimum Non-linear Systems," The Institute of Radio Engineers Convention Record, Part 4, pp. 21-31; 1956.

6. Carr, J. W., III, "An Analytic Investigation of Transient Synthesis by Exponentials," S. M. Thesis, Dept. of Elect. Engineering, Mass. Institute of Tech.; 1949.

7. Cerrillo, M. V., and Guillemin, E. A., "Rational Fraction Expansion for Network Functions," Proceedings of the Symposium on Modern Network Synthesis, pp. 84-127, New York; April, 1952.

8. Courant, R., and Hilbert, D., Methods of Mathematical Physics, Interscience Publishers, Inc. New York; 1953

9. Gabor, D., "Communication Theory and Cybernetics," Transactions of the Institute of Radio Engineers, Vol. CT-1, pp. 19-32; December, 1954.

10. Gardner, M. F., and Barnes, J. L., Transients in Linear Systems, John Wiley and Sons, New York; 1942.

11. Goldman, Stanford, Information Theory, Prentice-Hall, Inc., New York; 1953.

12. Goodman, T. P., and Reswick, J. B., "Determination of System Characteristics from Normal Operating Records," ASME-IRD Conference on Automatic Control, Ann Arbor, Michigan; April, 1955.

13. Graham, D., and Lathrop, R. C., "The Synthesis of Optimum Transient Response," Transactions of the American Institute of Electrical Engineers, Vol. 72, Part II, p. 273, November, 1953.

14. Guillemin, E., "What is Network Synthesis," Transactions of the Institute of Radio Engineers, Vol. PGCT-1, pp. 4-19; December, 1952.

15. Huggins, W. H., "Network Approximation in the Time Domain," Report E5048A, Air Force Research Laboratories, Cambridge, Mass.; October, 1949.

16. James, H. M., Nichols, N. B., and Phillips, R. S., Theory of Servo-mechanisms, McGraw-Hill Book Company; 1947.

17. Kautz, W. H., "Network Synthesis for Specified Transient Response," Technical Report No. 209, Research Laboratory of Electronics, Mass Institute of Tech.; April, 1952.

18. Kautz, W. H., "Transient Synthesis in the Time Domain," Transactions of the Institute of Radio Engineers, Vol. CT-1, pp. 29-39; September, 1954.

19. Korn, G. A., and Korn, T. M., Electronic Analog Computers, 2nd edition, McGraw-Hill Book Company, New York; 1956

20. Lampard, D. G. "A Method of Determining Correlation Functions of Stationary Time Series," Proceedings of the Institution of Electrical Engineers, Part c, No. 1, p.-35; March, 1955.

21. Laning, J. H., Jr., and Battin, R. H., Random Processes in Automatic Control, McGraw-Hill Book Company, New York; 1956.

22. Lee, Y. W., "Synthesis of Electrical Networks by Means of the Fourier Transforms of Laguerre's Functions," Journal of Mathamatics and Physics, Vol. 11, pp. 83-113 June, 1932.

23. Lee, Y. W., "Application of Statistical Methods to Communications Problems," Technical Report No. 181, Research Laboratory of Electronics, Mass. Institute of Tech.; September, 1950.

24. Lewis, N. W., "Waveform Computations by the Time Series Method," Proceedings of the Institution of Electrical Engineers, Vol 99, Part III, pp. 109-110; September, 1952.

25. Linvill, W. K., "Use of Sampled Functions for Time Domain Synthesis," Proceedings of the National Electronics Conference, Vol. 9, pp. 533-542; 1953

26. Mathers, G. W., "The Synthesis of Lumped-Element Circuits for Optimum Transient Response," Technical Report No. 28, Electronics Research Laboratories, Stanford University; November, 1951.

27. Mc Cool, W. A., Paper 14, Project Cyclone Symposium I; March, 1951

28. Moore, A. D., "An Application of Prony's Method to Time Domain Synthesis," Technical Report No. 42, Electronics Research Laboratories, Stanford University, February, 1952.

29. Nadler, M., "The Synthesis of Electrical Networks According to Prescribed Transient Conditions," Proceedings of the Institute of Radio Engineers, Vol. 37, pp. 627-629; June, 1949.

30. Otterman, J., "Time Domain Network Synthesis for an Analog Computer Setup," _Proceeding of National Simulation Conference_, pp. 24.1 - 24.5, Dallas, Texas; January, 1956.

31. Papoulis, A., "Network Response in Terms of Behavior at Imaginary Frequencies," _Proceedings of the Symposium on Modern Network Synthesis_, New York, N. Y.; April, 1955.

32. Shinbrot, Marvin, "A Least Squares Curve Fitting Method with Applications to the Calculation of Stability Coefficients from Transient Response Data," _Technical Note 2341_, National Advisory Committee for Aeronautics; April, 1951.

33. Shinbrot, Marvin, "On the Analysis of Linear and Nonlinear Dynamical Systems from Transient Response Data," _Technical Note 3288_, National Advisory Committee for Aeronautics; December, 1954.

34. Spencer, R. C., "Network Synthesis and the Moment Problem," _Transactions of the Institute of Radio Engineers_, Vol. CT-1 pp. 32-33; June, 1954.

35. Strieby, M., "A Fourier Method for Time Domain Synthesis," _Proceedings of the Symposium on Modern Network Synthesis_, pp. 197-211, New York; April, 1955.

36. Teasdale, R. D., "Time Domain Approximation by Use of Pade Approximates," _The Institute of Radio Engineers Convention Record_, Part 5, pp. 89-94; March, 1953

37. Thomson, W. E., "The Synthesis of a Network to Have a Sine-Squared Impulse Response," _Proceeding of the Institution of Electrical Engineers_, Vol. 99, Part III, pp. 373-376; November, 1952.

38. Truxal, J. R., _Automatic Feedback Control System Synthesis_, McGraw-Hill Book Company, Inc., New York; 1955.

39. Van Vleck, J. H., and Middleton, D., "A theoretical Comparison of the Visual, Aural, and Meter Reception of Pulsed Signals in the Presence of Noise," _Journal of Applied Physics_, Vol. 17, pp. 940-971; November, 1946.

40. Westcott, J. H., "The Introduction of Constraints into Feedback-System Designs," _Transactions of the Institute of Radio Engineers_, Vol. CT-1, pp. 39-49, September, 1954.

41. Winkler, S., "The Approximation Problem of Network Synthesis," _Transactions of the Institute of Radio Engineers_, Vol. CT-1, pp. 5-21; September, 1954.

42. Zabusky, N. J., "A Numerical Method for Determining a System Impulse Response from the Transient Response to Arbitrary Inputs," _Transactions of the Institute of Radio Engineers_, Vol. PGAC-1, pp. 40-56; May, 1956.

43. Zadeh, L. A., and Ragazzini, J. R., "An Extension of Wiener's Theory of Prediction," <u>Journal of Applied Physics</u>, Vol. 21, pp. 645-655; July, 1950.

43. Zadeh, L. A., and Ragazzini, J. R., "An Extension of Wiener's Theory of Prediction," <u>Journal of Applied Physics</u>, Vol. 21, pp. 645-655; July, 1950.