

①

AD A278 246

TECHNICAL REPORT

**GENERAL RESEARCH
IN FLIGHT SCIENCES**

JANUARY 1959 - JANUARY 1960

DTIC
ELECTE
APR 04 1994
S E D

**VOLUME IV
MATHEMATICS AND STATISTICS**

LMSD - 288139

JANUARY 1960

WORK CARRIED OUT UNDER
THE LOCKHEED GENERAL RESEARCH PROGRAM

Lockheed

Best Available Copy

MISSILES and SPACE DIVISION
LOCKHEED AIRCRAFT CORPORATION • SUNNYVALE, CALIF.

TECHNICAL REPORT

**GENERAL RESEARCH
IN FLIGHT SCIENCES**

JANUARY 1959 - JANUARY 1960

SEARCHED
SERIALIZED
APR 04 1994
S E D

**VOLUME IV
MATHEMATICS AND STATISTICS**

LMSD - 288139

JANUARY 1960

WORK CARRIED OUT UNDER
THE LOCKHEED GENERAL RESEARCH PROGRAM

Lockheed

MISSILES and SPACE DIVISION

LOCKHEED AIRCRAFT CORPORATION • SUNNYVALE, CALIF.

FOREWORD

The Lockheed Missiles and Space Division sponsors a comprehensive program of general research in connection with its defense contracts. A portion of the total research effort at LMSD is carried out in Flight Sciences in the fields of fluid mechanics, mechanics of deformable bodies, flight dynamics, space mechanics, and mathematics and statistics.

The results obtained from theoretical and experimental studies in Flight Sciences are felt to be of general interest and have been assembled into a collection of research reports. These reports comprise a series of papers which deal with various topics in the fields listed above.

This volume is concerned with mathematics and statistics.

Accession For	
NTIS CRA&I	<input checked="" type="checkbox"/>
DTIC TAB	
Unannounced Justification	
By	
Distribution /	
Availability Codes	
Dist	Avail and/or Special
A-1	

CONTENTS

FOREWORD	iii
	<u>Number</u>
INTERVAL ANALYSIS I	
R. E. Moore and C. T. Yang	1
CONVERGENCE OF APPROXIMATE EIGENVECTORS IN JACOBI METHODS	
R. L. Causey and P. Henrici	2
WORD CORRELATION STUDY	
R. P. Mitchell	3
EFFICIENT ESTIMATION OF REGRESSION PARAMETERS FOR SECOND ORDER STATIONARY PROCESSES	
C. T. Striebel	4
A MODEL OF TURBULENCE DISPLAYING AUTONOMOUS OSCILLATIONS	
R. J. Dickson	5

INTERVAL ANALYSIS 1

**R.E. Moore
C.T. Yang**



1

FOREWORD

This article is a reprint of the Technical Document
LMSD-285875, bearing the same title and dated September
1959. It discusses work carried out under the Lockheed
General Research Program.

CONTENTS

<u>Section</u>		<u>Page</u>
	Foreword	iii
	Introduction	1
1	Preliminaries	2
2	Addition	6
3	Multiplication	7
4	Subtraction	11
5	Division	12
6	Arithmetic Functions	14
7	Relation Between Arithmetic Functions and Rational Functions	19
8	First Approximation Theorem	21
9	Second Approximation Theorem	26
10	Approximation of a Continuous Function by Arithmetic Functions	29

INTERVAL ANALYSIS I

INTRODUCTION

Digital computations by computers consist of finite sequences of psuedo-arithmetic operations. On the other hand, the exact numerical solution of a mathematical problem, if computable at all, often requires an infinite sequence of exact arithmetic operations.

The study of approximation by digital computations is the underlying motivation for the present study. A digital computation and the analysis of its error as an approximation are usually carried out separately. However, in the present study an interval arithmetic is devised which forms a basis for a concomitant analysis of error in a digital computation. In this system, computations are performed with intervals and intervals are so produced to contain, by construction, the exact numerical solutions sought. Hence an approximation and its possible error will be obtained at the same time, choosing say the midpoint of an interval as the approximation.

This report is the first part of our study, in which we first examine some properties of exact or ideal interval arithmetic. After a preliminary discussion of the space of intervals (§ 1) we study addition, multiplication, subtraction and division of intervals (§§ 2-5). Then we construct arithmetic functions as compositions of these elementary operations (§ 6). As one may expect, arithmetic functions play exactly the same role in interval analysis as rational functions in real analysis, so that there is a relation between arithmetic functions and rational functions (§ 7).

In the present report we apply interval analysis to the study of the following problems. Let $f(x)$ be a continuous real-valued function defined on an interval $[a,b]$. (1) What is the image interval $f([a,b])$? (2) What is the definite integral $\int_a^b f(x) dx$? When $f(x)$ is a rational function, the approximations are given in theorems 1 and 2 (§§ 8,9). In general, if $f(x)$ is an arbitrary continuous function, we may still have approximation theorems 3 and 4 (§ 10), although they are not as precise as the first two theorems.

In forthcoming reports we shall apply interval analysis to differential equations and report on results of machine computations using a digital version of interval arithmetic modified to enable the computations to be carried out with pseudo-arithmetic operations. (See also, LMSD-48421, "Automatic Error Analysis in Digital Computation," by R. E. Moore.)

1. PRELIMINARIES

Throughout the whole study R denotes the real line. Whenever a and b are real numbers with $a \leq b$, $[a,b]$ denotes the subset of R consisting of all the real numbers x with $a \leq x \leq b$. In symbols,

$$[a,b] = \{x \in R \mid a \leq x \leq b\}.$$

Let \mathcal{I} be the set of all such sets $[a,b]$.

That means,

$$\mathcal{I} = \{[a,b] \mid a,b \in R \text{ and } a \leq b\}.$$

Then we have natural functions

$$p : R \rightarrow \mathcal{A} ,$$

$$\alpha : \mathcal{A} \rightarrow R ,$$

$$\beta : \mathcal{A} \rightarrow R ,$$

$$\gamma : \mathcal{A} \rightarrow R ,$$

$$\sigma : \mathcal{A} \rightarrow R ,$$

defined by

$$p(x) = [x, x] ,$$

$$\alpha([a, b]) = a ,$$

$$\beta([a, b]) = b ,$$

$$\gamma([a, b]) = \max \{ |a| , |b| \} ,$$

$$\sigma([a, b]) = b - a$$

respectively, where $x \in R$ and $[a, b] \in \mathcal{A}$.

Whenever $x, x' \in R$, we let

$$\rho(x, x') = |x - x'| .$$

Then ρ is a metric on R , that means, ρ has the following properties.

- (1) Whenever $x, x' \in R$, $\rho(x, x') = 0$ if and only if $x = x'$.
- (2) Whenever $x, x' \in R$, $\rho(x, x') = \rho(x', x)$.
- (3) Whenever $x, x', x'' \in R$, $\rho(x, x') + \rho(x', x'') \geq \rho(x, x'')$.

Whenever $A, A' \in \mathcal{A}$ we let

$$P(A, A') = \max \left\{ \rho(\alpha(A), \alpha(A')), \rho(\beta(A), \beta(A')) \right\}.$$

Then P has the same properties as those ρ has on R so that P is a metric on \mathcal{A} .

As direct consequences of the definitions of $\rho, P, \alpha, \beta, \gamma, \sigma$ we have

(1-1) The function $p: R \rightarrow \mathcal{A}$ is isometric; that means, for any $x, x' \in R$,

$$P(p(x), p(x')) = \rho(x, x').$$

Hence p maps R homeomorphically onto $p(R)$.

(1-2) Whenever $A, A' \in \mathcal{A}$,

$$\rho(\alpha(A), \alpha(A')) < P(A, A').$$

Hence the function $\alpha: \mathcal{A} \rightarrow R$ is uniformly continuous, that means, for any $\epsilon > 0$ there is a $\delta > 0$ such that whenever $A, A' \in \mathcal{A}$ with $P(A, A') < \delta$, we have $\rho(\alpha(A), \alpha(A')) < \epsilon$.

Since uniformly continuous functions are continuous, it follows that:

(1-3) The function $\alpha: \mathcal{A} \rightarrow R$ is continuous, that means, for any $A \in \mathcal{A}$ and any $\epsilon > 0$ there is a $\delta > 0$ such that whenever $A' \in \mathcal{A}$ with $P(A, A') < \delta$ we have $\rho(\alpha(A), \alpha(A')) < \epsilon$.

Just as (1-2) and (1-3), we have

(1-4) Whenever $A, A' \in \mathcal{A}$,

$$\rho(\beta(A), \beta(A')) \leq P(A, A').$$

Hence the function $\beta: \mathcal{A} \rightarrow \mathbb{R}$ is uniformly continuous and consequently it is continuous.

Using (1-2), (1-4) and well-known properties of real numbers, we have

(1-5) The functions $\gamma: \mathcal{A} \rightarrow \mathbb{R}$ and $\sigma: \mathcal{A} \rightarrow \mathbb{R}$ are uniformly continuous and hence they are continuous.

The following will be needed later

(1-6) Let $A, A' \in \mathcal{A}$ and let $\alpha > 0$. Then $P(A, A') < \alpha$ if and only if the following two conditions hold.

- (i) For every $x \in A$ there is some $x' \in A'$ with $\rho(x, x') < \alpha$.
- (ii) For every $y' \in A'$ there is some $y \in A$ with $\rho(y', y) < \alpha$.

Whenever $I \in \mathcal{A}$ we let

$$\mathcal{A}_I = \{A \in \mathcal{A} \mid A \subset I\} .$$

(1-7) Whenever $I \in \mathcal{A}$, \mathcal{A}_I is compact; that means, every sequence in \mathcal{A}_I contains a convergent subsequence.

Let \mathcal{J} be the subset of \mathcal{A} consisting of all the elements of \mathcal{A} not containing 0.

(1-8) The set \mathcal{J} is open in \mathcal{A} ; that means, for every $A \in \mathcal{J}$ there is a positive number γ_A such that every $A' \in \mathcal{A}$ with $P(A, A') < \gamma_A$ belongs to \mathcal{J} . In fact, we may choose $\gamma_A = \min \{|\alpha(A)|, |\beta(A)|\}$.

2. ADDITION

There is a function

$$\oplus : \mathcal{L} \times \mathcal{L} \rightarrow \mathcal{L}$$

defined by

$$\begin{aligned} \oplus (A,B) &= \{x + y \mid x \in A \text{ and } y \in B\} \\ &= [\alpha(A) + \alpha(B), \beta(A) + \beta(B)], \quad A, B \in \mathcal{L} . \end{aligned}$$

The set $\oplus(A,B)$ is also written $A \oplus B$. The function \oplus is called the addition on \mathcal{L} .

(2-1) The function $p : R \rightarrow \mathcal{L}$ preserves the addition; that means, for any $x, y \in R$,

$$p(x) \oplus p(y) = p(x + y) .$$

Because of (2-1), the addition \oplus on \mathcal{L} may be regarded as an extension of the addition $+$ on R . This is the reason for calling \oplus the "addition" on \mathcal{L} .

(2-2) The addition \oplus is commutative; that means, for any $A, B \in \mathcal{L}$,

$$A \oplus B = B \oplus A .$$

(2-3) The addition \oplus is associative; that means, for any $A, B, C \in \mathcal{L}$,

$$A \oplus (B \oplus C) = (A \oplus B) \oplus C .$$

(2-4) Whenever $(A,B), (A',B') \in \mathcal{L} \times \mathcal{L}$,

$$P(A \oplus B, A' \oplus B') \leq P(A, A') + P(B, B') .$$

Hence the addition \oplus is uniformly continuous; that means, for any $\epsilon > 0$ there is a $\delta > 0$ such that whenever $(A,B), (A',B') \in \mathcal{A}$ with $P(A,A') < \delta$ and $P(B,B') < \delta$, we have $P(A \oplus B, A' \oplus B') < \epsilon$.

Proof. Let

$$A = [a,b], A' = [a',b'], B = [c,d], B' = [c',d'] .$$

Then

$$\begin{aligned} P(A \oplus B, A' \oplus B') &= P\left([a+c, b+d], [a'+c', b'+d']\right) \\ &= \max\{\rho(a+c, a'+c'), \rho(b+d, b'+d')\} \\ &\leq \max\{\rho(a,a') + \rho(c,c'), \rho(b,b') + \rho(d,d')\} \\ &\leq \max\{\rho(a,a'), \rho(b,b')\} + \max\{\rho(c,c'), \rho(d,d')\} \\ &= P(A,A') + P(B,B') . \end{aligned}$$

To prove the uniform continuity we have only to pick $\delta = \epsilon/2$. q.e.d.

Since uniformly continuous functions are continuous, it follows that

(2-5) The addition \oplus is continuous; that means, for any $(A,B) \in \mathcal{A}$ and any $\epsilon > 0$ there is a $\delta > 0$ such that whenever $(A',B') \in \mathcal{A}$ with $P(A,A') < \delta$ and $P(B,B') < \delta$, we have $P(A \oplus B, A' \oplus B') < \epsilon$.

3. MULTIPLICATION

Whenever $A, B \in \mathcal{A}$ we let

$$\otimes(A,B) = \{xy \mid x \in A \text{ and } y \in B\} .$$

We claim that $\otimes(A,B) \in \mathcal{A}$.

Let $A = [a,b]$ and $B = [c,d]$. Then we have the following cases

(1) If $a \geq 0$ and $c > 0$, then

$$\otimes (A,B) = [ac,bd] .$$

(2) If $a \geq 0$ and $c < 0 < d$, then

$$\otimes (A,B) = [bc,bd] .$$

(3) If $a > 0$ and $d \leq 0$, then

$$\otimes (A,B) = [bc,ad] .$$

(4) If $a < 0 < b$ and $c \geq 0$, then

$$\otimes (A,B) = [ad,bd] .$$

(5) If $a < 0 < b$ and $c < 0 < d$, then

$$\otimes (A,B) = [\min \{bc,ad\} , \max \{ac,bd\}] .$$

(6) If $a < 0 < b$ and $d \leq 0$, then

$$\otimes (A,B) = [bc,ac] .$$

(7) If $b \leq 0$ and $c \geq 0$, then

$$\otimes (A,B) = [ad,bc] .$$

(8) If $b \leq 0$ and $c < 0 < d$, then

$$\otimes (A,B) = [ad,ac] .$$

(9) If $b \leq 0$ and $d \leq 0$, then

$$\otimes (A,B) = [bd,ac] .$$

Because of this result, we have a function

$$\otimes : \mathcal{A} \times \mathcal{A} \rightarrow \mathcal{A}$$

called the multiplication on \mathcal{A} . Whenever $A, B \in \mathcal{A}$, the set $\otimes(A, B)$ is also written $A \otimes B$.

(3-1) The function $p : R \rightarrow \mathcal{A}$ preserves the multiplication.

Because of (3-1) the multiplication \otimes on \mathcal{A} may be regarded as an extension of the multiplication on R .

(3-2) The multiplication \otimes is commutative; that means, for any $A, B \in \mathcal{A}$,

$$A \otimes B = B \otimes A.$$

(3-3) The multiplication \otimes is associative; that means, for any $A, B, C \in \mathcal{A}$,

$$A \otimes (B \otimes C) = (A \otimes B) \otimes C$$

(3-4) For any $A, B, C \in \mathcal{A}$,

$$A \otimes (B \oplus C) \subset (A \otimes B) \oplus (A \otimes C)$$

but both sides may not be equal. Hence the distributive law does not hold.

(3-5) Whenever $(A, B), (A', B') \in \mathcal{A} \times \mathcal{A}$

$$P(A \otimes B, A' \otimes B') \leq \gamma(B) P(A, A') + \gamma(A') P(B, B').$$

Hence the multiplication is continuous.

Proof. Let $xy \in A \otimes B$, where $x \in A$ and $y \in B$. By (1-6), there is some $x' \in A'$ and some $y' \in B'$ such that $\rho(x, x') \leq P(A, A')$ and $\rho(y, y') \leq P(B, B')$.

Therefore $x'y' \in A' \otimes B'$ and

$$\begin{aligned} \rho(xy, x'y') &= |xy - x'y'| \\ &= |y(x - x') + x'(y - y')| \\ &\leq |y| \cdot |x - x'| + |x'| \cdot |y - y'| \\ &\leq \gamma(B) P(A, A') + \gamma(A') P(B, B') . \end{aligned}$$

Let $x'y' \in A' \otimes B'$, where $x' \in A'$ and $y' \in B'$. Similarly there is some $xy \in A \otimes B$ such that $x \in A$, $y \in B$ and

$$\rho(x'y', xy) < \gamma(B) P(A, A) + \gamma(A') P(B, B') .$$

Making use of (1-6) again, these results imply the first part of (3-5).

To prove the continuity of \otimes , we let $(A, B) \in \mathcal{A} \times \mathcal{A}$ and let $\epsilon > 0$.
Take

$$\delta = \min \left\{ \epsilon / (\gamma(A) + \gamma(B) + 1), 1 \right\} .$$

Then for any $(A', B') \in \mathcal{A} \times \mathcal{A}$ with $P(A, A') < \delta$ and $P(B, B') < \delta$ we have

$$\begin{aligned} P(A \otimes B, A' \otimes B') &< \gamma(B)\delta + \gamma(A')\delta \leq \gamma(B)\delta + (\gamma(A) + \delta)\delta \\ &\leq (\gamma(A) + \gamma(B) + 1)\delta < \epsilon . \quad \text{q.e.d.} \end{aligned}$$

It is not hard to see that the multiplication is not uniformly continuous. However, since the restriction of a continuous function on a compact set is uniformly continuous, it follows from (1-7) and (3-5) that

(3-6) Let I and J be fixed elements of \mathcal{A} . Then the multiplication

$$\otimes : \mathcal{A}_I \times \mathcal{A}_J \rightarrow \mathcal{A}$$

is uniformly continuous. In fact, for any $A, A' \in I$ and $B, B' \in J$ we have

$$P(A \oplus B, A' \oplus B') < \gamma(J) P(A, A') + \gamma(I) P(B, B') .$$

4. SUBTRACTION

By (3-5) we have

(4-1) Let E be a fixed element of \mathcal{A} . Then the function \otimes_E of \mathcal{A} into \mathcal{A} , defined by

$$\otimes_E (A) = E \otimes A, \quad A \in \mathcal{A},$$

is uniformly continuous.

In particular, if $E = [-1, -1]$, we have

(4-2) The function $\otimes_{[-1, -1]} : \mathcal{A} \rightarrow \mathcal{A}$, defined by

$$\otimes_{[-1, -1]} (A) = [-1, -1] \otimes A, \quad A \in \mathcal{A},$$

is uniformly continuous.

Whenever $A \in \mathcal{A}$, we shall abbreviate $\otimes_{[-1, -1]} (A)$ by $-A$. Since $-(-A) = A$, $\otimes_{[-1, -1]}$ is a homeomorphism.

Combining the addition \oplus and the function $\otimes_{[-1, -1]}$ we define the subtraction

$$\ominus : \mathcal{A} \times \mathcal{A} \rightarrow \mathcal{A}$$

by

$$\ominus (A, B) = A \oplus (-B) = \{x - y \mid x \in A \text{ and } y \in B\}, \quad A, B \in \mathcal{A} .$$

The set $\ominus(A, B)$ is also written $A \ominus B$.

From (2-1) and (3-1), it follows

(4-3) The function $p : R \rightarrow \mathcal{A}$ preserves the subtraction.

(4-4) Whenever $A, B \in \mathcal{A}$,

$$(-A) \otimes B = A \otimes (-B) = -(A \otimes B),$$

$$(-A) \otimes (-B) = A \otimes B.$$

From (2-4) and (4-2) it follows

(4-5) The subtraction \ominus is uniformly continuous and hence is continuous. In fact, for any $(A, B), (A', B') \in \mathcal{A} \times \mathcal{A}$,

$$P(A \ominus B, A' \ominus B') \leq P(A, A') + P(B, B').$$

5. DIVISION

For every $A = [a, b] \in \mathcal{J}$ (see § 1),

$$A^{-1} = \{x^{-1} \mid x \in A\} = [b^{-1}, a^{-1}]$$

is in \mathcal{J} . Hence we have a function $\tau : \mathcal{J} \rightarrow \mathcal{J}$ defined by

$$\tau(A) = A^{-1}.$$

(5-1) Whenever $J \in \mathcal{J}$, $\mathcal{A}_J \subset \mathcal{J}$ and for any $A, A' \in \mathcal{A}_J$,

$$P(A^{-1}, A'^{-1}) \leq \gamma (J^{-1})^2 P(A, A').$$

Hence τ is uniformly continuous on \mathcal{A}_J , $J \in \mathcal{J}$, and consequently it is continuous on \mathcal{J} .

Proof. Let $A, A' \in \mathcal{J}$. For any $x \in A$ there is, by (1-6), some $x' \in A'$ such that $\rho(x, x') \leq P(A, A')$ so that

$$\begin{aligned} \rho(x^{-1}, x'^{-1}) &= |x^{-1} - x'^{-1}| = |x^{-1}| |x'^{-1}| |x - x'| \\ &\leq \gamma(J)^2 P(A, A'). \end{aligned}$$

Similarly, for any $y' \in A'$ there is some $y \in A$ such that

$$\rho(y'^{-1}, y^{-1}) \leq \gamma(J)^2 P(A, A'). \text{ Hence the first part is proved.}$$

To prove the uniform continuity of τ on \mathcal{J} , $J \in \mathcal{J}$, we let $\epsilon > 0$ and take

$$\delta = \epsilon / \gamma(J^{-1})^2.$$

It is clear that for any $A, A' \in \mathcal{J}$ with $P(A, A') \leq \delta$, we have

$$P(A^{-1}, A'^{-1}) < \gamma(J^{-1})^2 \delta = \epsilon.$$

To prove the continuity of τ we let $A \in \mathcal{J}$ and let $\epsilon > 0$. Take a $J \in \mathcal{J}$ such that

$$\alpha(J) < \alpha(A) \leq \beta(A) < \beta(J).$$

Then for any $A' \in \mathcal{J}$ with $P(A, A') < \min\{\alpha(A) - \alpha(J), \beta(J) - \beta(A)\}$, $A' \in \mathcal{J}$. Hence the continuity of τ on \mathcal{J} implies the continuity of τ at A . q.e.d.

Since $(A^{-1})^{-1} = A, A \in \mathcal{J}$, τ is a homeomorphism.

Combining the multiplication \otimes and the function τ we define the division

$$\oplus : \mathcal{J} \times \mathcal{J} \rightarrow \mathcal{J}$$

by

$$\oplus (A,B) = A \otimes B^{-1} = \{xy^{-1} \mid x \in A \text{ and } y \in B\}, A \in \mathcal{I} \text{ and } B \in \mathcal{J} .$$

The set $\oplus (A,B)$ is also written A/B .

Since for any real number $x \neq 0$,

$$p(x)^{-1} = p(x^{-1}) ,$$

it follows from (3-1) that

(5-2) The function p preserves the division \oplus .

From (3-6) and (5-1) it follows

(5-3) Whenever $I \in \mathcal{I}$ and $J \in \mathcal{J}$,

$$\oplus : \mathcal{I} \times \mathcal{J} \rightarrow \mathcal{I}$$

is uniformly continuous. In fact, for any (A,B) , $(A',B') \in \mathcal{I} \times \mathcal{J}$ we have

$$P(A/B, A'/B') \leq \gamma(J^{-1}) P(A, A') + \gamma(J^{-1})^2 \gamma(I) P(B, B') .$$

Hence $\oplus : \mathcal{I} \times \mathcal{J} \rightarrow \mathcal{I}$ is continuous.

6. ARITHMETIC FUNCTIONS

Before giving the definition of an arithmetic function we remark that every arithmetic function has a domain contained in \mathcal{I} , an order which is a non-negative integer, and a finite number of parameters which are elements of \mathcal{I} . An arithmetic function of order n with parameters A_1, A_2, \dots, A_m is written $F_{A_1 A_2 \dots A_m}^{(n)}$ or simply F . It is a function from its domain $\mathcal{D}(F)$ to \mathcal{I} .

An arithmetic function of order n is defined by induction on n . When $n = 0$, the number m of parameters is either 0 or 1. In the former case it is the identity function $F^{(0)} : \mathcal{A} \rightarrow \mathcal{A}$ given by

$$F^{(0)}(X) = X, \quad X \in \mathcal{A}.$$

In the latter case it is the constant function $F_A^{(0)} : \mathcal{A} \rightarrow \mathcal{A}$ given by

$$F_A^{(0)}(X) = A, \quad X \in \mathcal{A},$$

where A is an arbitrary element of \mathcal{A} . Notice that every arithmetic function of order 0 has \mathcal{A} as its domain.

Let n be a positive integer and suppose that arithmetic functions of order $< n$ have been defined. Then every arithmetic function $F_{A_1 A_2 \dots A_m}^{(n)}$ of order n is defined as follows. There is an arithmetic

function $F_{A_1 A_2 \dots A_\ell}^{(s)}$ of order s , $0 \leq s \leq n-1$, with parameters

A_1, A_2, \dots, A_ℓ , $0 \leq \ell \leq m$, and an arithmetic function $F_{A_{\ell+1} A_{\ell+2} \dots A_m}^{(n-1-s)}$

of order $n-1-s$ with parameters $A_{\ell+1}, A_{\ell+2}, \dots, A_m$ such that

$F_{A_1 A_2 \dots A_m}^{(n)} = F_{A_1 A_2 \dots A_\ell}^{(s)} \circ F_{A_{\ell+1} A_{\ell+2} \dots A_m}^{(n-1-s)}$ is given by

$$F_{A_1 A_2 \dots A_m}^{(n)}(X) = F_{A_1 A_2 \dots A_\ell}^{(s)}(X) \circ F_{A_{\ell+1} A_{\ell+2} \dots A_m}^{(n-1-s)}(X),$$

where \circ is one of $\oplus, \ominus, \otimes, \oplus$ and X is such an element of \mathcal{A} that the right side is well-defined. It is clear that if \mathcal{D}_1 and \mathcal{D}_2 are the

respective domains of $F_{A_1 A_2 \dots A_\ell}^{(s)}$ and $F_{A_{\ell+1} A_{\ell+2} \dots A_m}^{(n-1-s)}$, then the

domain \mathcal{D} of $F_{A_1 A_2 \dots A_m}^{(n)}$ is given by

$$\mathcal{D} = \begin{cases} \mathcal{D}_1 \cap \mathcal{D}_2 \cap \left(F_{A_{\ell+1} A_{\ell+2} \dots A_m}^{(n-1-s)} \right)^{-1}(\mathcal{D}) & \text{if } \circ \text{ is } \oplus, \\ \mathcal{D}_1 \cap \mathcal{D}_2 & \text{if otherwise.} \end{cases}$$

Let F be an arithmetic function of parameters A_1, A_2, \dots, A_m and domain $\mathcal{D}(F)$. We may write

$$F(A_1, A_2, \dots, A_m; X)$$

instead of $F(X)$ and consider F as a function of $\mathcal{I}^m \times \mathcal{D}(F)$ into \mathcal{I} . Notice that $\mathcal{D}(F)$ depends on A_1, A_2, \dots, A_m .

(6-1) If $F(B_1, B_2, \dots, B_m; Y)$ is defined, then for any $C_1 \in \mathcal{D}_{B_1}, C_2 \in \mathcal{D}_{B_2}, \dots, C_m \in \mathcal{D}_{B_m}, Z \in \mathcal{D}_Y$, $F(C_1, C_2, \dots, C_m; Z)$ is defined and is contained in $F(B_1, B_2, \dots, B_m; Y)$.

Proof. If F is of order 0, our assertion is trivial. Hence our assertion holds for arithmetic functions of order 0.

Now we proceed by induction. Let n be a positive integer and assume our assertion for all arithmetic functions of order $< n$. By definition every arithmetic function F of order n with parameters A_1, A_2, \dots, A_m is given by

$$F(A_1, A_2, \dots, A_m; X) = F_1(A_1, \dots, A_l; X) \circ F_2(A_{l+1}, \dots, A_m; X),$$

where F_1 and F_2 are arithmetic functions of order $< n$ and \circ is one of $\oplus, \ominus, \otimes, \oslash$.

If $F(B_1, B_2, \dots, B_m; Y)$ is defined, then $F_1(B_1, \dots, B_l; Y)$ and $F_2(B_{l+1}, \dots, B_m; Y)$ is in \mathcal{I} when \circ is \oplus . By the induction hypothesis, for any $C_1 \in \mathcal{D}_{B_1}, C_2 \in \mathcal{D}_{B_2}, \dots, C_m \in \mathcal{D}_{B_m}, Z \in \mathcal{D}_Y$, $F_1(C_1, \dots, C_l; Z)$ and $F_2(C_{l+1}, \dots, C_m; Z)$ are defined and

$$F_1(C_1, \dots, C_l; Z) \subset F_1(B_1, \dots, B_l; Y),$$

$$F_2(C_{l+1}, \dots, C_m; Z) \subset F_2(B_{l+1}, \dots, B_m; Y).$$

So $F_2(C_{l+1}, \dots, C_m; Z) \in \mathcal{J}$ when \circ is \oplus . Hence

$$F(C_1, C_2, \dots, C_m; Z) = F_1(C_1, \dots, C_l; Z) \circ F_2(C_{l+1}, \dots, C_m; Z)$$

is defined and is contained in $F(B_1, B_2, \dots, B_m; Y)$. q.e.d.

(6-2) If $F(B_1, \dots, B_m; Y)$ is defined, then there is a number $k > 0$ such that whenever $A_1, A'_1 \in \mathcal{A}_{B_1}, \dots, A_m, A'_m \in \mathcal{A}_{B_m}, X, X' \in \mathcal{A}_Y$, we have

$$\begin{aligned} & P \left(F(A_1, \dots, A_m; X), F(A'_1, \dots, A'_m; X') \right) \\ & \leq k \left(P(A_1, A'_1) + \dots + P(A_m, A'_m) + P(X, X') \right). \end{aligned}$$

Hence F is uniformly continuous on $\mathcal{A}_{B_1} \times \dots \times \mathcal{A}_{B_m} \times \mathcal{A}_Y$.

Proof. When F is of order \circ it is evident that the inequality holds with $k = 1$. Therefore we may proceed by induction and assume the inequality for arithmetic functions of order $< n$.

Every arithmetic function F of order n is given by

$$F(A_1, \dots, A_m; X) = F_1(A_1, \dots, A_l; X) \circ F_2(A_{l+1}, \dots, A_m; X),$$

where F_1 and F_2 are arithmetic functions of order $< n$ and \circ is one of $\oplus, \ominus, \otimes, \odot$. By the induction hypothesis, there exist positive numbers k_1 and k_2 such that whenever $A_1, A'_1 \in \mathcal{A}_{B_1}, \dots, A_m, A'_m \in \mathcal{A}_{B_m}, X, X' \in \mathcal{A}_Y$, we have

$$\begin{aligned} & P \left(F_1(A_1, \dots, A_l; X), F_1(A'_1, \dots, A'_l; X') \right) \\ & \leq k_1 \left(P(A_1, A'_1) + \dots + P(A_l, A'_l) + P(X, X') \right) \\ & P \left(F_2(A_{l+1}, \dots, A_m; X), F_2(A'_{l+1}, \dots, A'_m; X') \right) \\ & \leq k_2 \left(P(A_{l+1}, A'_{l+1}) + \dots + P(A_m, A'_m) + P(X, X') \right). \end{aligned}$$

By (2-4), (3-6), (4-5), (5-3), there exists a positive number k for which our desired inequality holds. In fact, we can let

$$k = \begin{cases} 2(k_1 + k_2) \left(\gamma(J^{-1}) + \gamma(J^{-1})^2 \gamma(I) \right) & \text{if } \circ \text{ is } \oplus, \\ 2(k_1 + k_2) \left(\gamma(I) + \gamma(J) + 1 \right) & \text{if otherwise,} \end{cases}$$

where

$$I = F_1(B_1, \dots, B_\ell; Y), \quad J = F_2(B_{\ell+1}, \dots, B_m; Y).$$

q.e.d.

(6-3) If $F(B_1, \dots, B_m; Y)$ is defined, then there is a $\delta > 0$ such that whenever $A_1, \dots, A_m, X \in \mathcal{A}$ with

$$P(A_1, B_1) < \delta, \dots, P(A_m, B_m) < \delta, P(X, Y) < \delta,$$

$F(A_1, \dots, A_m; X)$ is defined and is continuous.

Proof. If F is of order 0, then our assertion is trivial. Therefore, we may proceed by induction and assume our assertion for arithmetic functions of order $< n$, $n > 0$.

Every arithmetic function F of order n is given by

$$F(A_1, \dots, A_m; X) = F_1(A_1, \dots, A_\ell; X) \circ F_2(A_{\ell+1}, \dots, A_m; X),$$

where F_1 and F_2 are arithmetic functions of order $< n$ and \circ is one of $\oplus, \ominus, \otimes, \oplus$.

By the induction hypothesis there is a positive number δ such that whenever $A_1, \dots, A_m, X \in \mathcal{A}$ with $P(A_1, B_1) < \delta, \dots, P(A_m, B_m) < \delta, P(X, Y) < \delta$, both $F_1(A_1, \dots, A_\ell; X)$ and $F_2(A_{\ell+1}, \dots, A_m; X)$ are defined and continuous. If $F_2(B_{\ell+1}, \dots, B_m; Y) \in \mathcal{J}$, we may choose δ so small that $F_2(A_{\ell+1}, \dots, A_m; X) \in \mathcal{J}$. Hence $F(A_1, \dots, A_m; X)$ is defined and is continuous.

q.e.d.

It follows from (6-1), (6-2) and (6-3) that

(6-4) Let F be an arithmetic function with m parameters and let

$$U = \{(A_1, \dots, A_m; X) \in \mathcal{A}^{m+1} \mid F(A_1, \dots, A_m; X) \text{ defined}\} .$$

Then U is open in \mathcal{A}^{m+1} and F is a continuous function of U into \mathcal{A} . Moreover, for every $(B_1, \dots, B_m; Y) \in U$, $\mathcal{A}_{B_1} \times \dots \times \mathcal{A}_{B_m} \times \mathcal{A}_Y \subset U$ and F is uniformly continuous on $\mathcal{A}_{B_1} \times \dots \times \mathcal{A}_{B_m} \times \mathcal{A}_Y$. Furthermore, $(A_1, \dots, A_m; X) \in \mathcal{A}_{B_1} \times \dots \times \mathcal{A}_{B_m} \times \mathcal{A}_Y$ implies

$$F(A_1, \dots, A_m; X) \subset F(B_1, \dots, B_m; Y) .$$

When we write $F(X)$ in place of $F(A_1, \dots, A_m; X)$ it is understood that parameters A_1, \dots, A_m are fixed. Since

$$\mathcal{D}(F) = \{X \in \mathcal{A} \mid F(X) \text{ defined}\} ,$$

it follows from (6-4) and (6-2) that

(6-5) For every arithmetic function F , $\mathcal{D}(F)$ is open in \mathcal{A} and $F: \mathcal{D}(F) \rightarrow \mathcal{A}$ is continuous. Let $I \in \mathcal{D}(F)$. Then $\mathcal{A}_I \subset \mathcal{D}(F)$ and $X \in \mathcal{A}_I$ implies $F(X) \subset F(I)$. Moreover, there is a positive number k such that whenever $X, X' \in \mathcal{A}_I$,

$$P(F(X), F(X')) \leq kP(X, X') .$$

7. RELATION BETWEEN ARITHMETIC FUNCTIONS AND RATIONAL FUNCTIONS.

In the construction of arithmetic functions, if we replace \mathcal{A} by R and replace $\oplus, \ominus, \otimes, \odot$ by corresponding operations on R , then we obtain rational functions in place of arithmetic functions. Therefore, we can establish a relation between arithmetic functions and rational function.

An arithmetic function is called special if all of its parameters belong to $p(R)$.

(7-1) If F is a special arithmetic function of domain $\mathcal{D}(F)$, then $p^{-1}(\mathcal{D}(F))$ is open in R and for every $x \in p^{-1}(\mathcal{D}(F))$, $Fp(x) \in p(R)$. Moreover, there is a unique rational function f whose domain contains $p^{-1}(\mathcal{D}(F))$ such that $pf = Fp$ or $f = p^{-1}Fp$.

$$\begin{array}{ccc} R & \xrightarrow{p} & \mathcal{A} \\ \uparrow f & & \uparrow F \\ p^{-1}(\mathcal{D}(F)) & \xrightarrow{p} & \mathcal{D}(F) \end{array}$$

As the converse of (7-1), we have

(7-2) Given any rational function f of domain $D(f)$ there is a special arithmetic function F of domain $\mathcal{D}(F) \supset p(D(f))$ such that $f = p^{-1}Fp$.

Remark. It is possible to have two distinct special arithmetic functions F_1 and F_2 such that $p^{-1}F_1p = p^{-1}F_2p$. For example,

$$\begin{aligned} F_1(X) &= (X \otimes X) \oplus X, \\ F_2(X) &= X \otimes (X \oplus p(1)), \end{aligned} \quad X \in \mathcal{A},$$

give two arithmetic functions F_1 and F_2 of domain \mathcal{A} . It is clear that

$$(p^{-1}F_1p)(x) = x^2 + x = x(x+1) = (p^{-1}F_2p)(x).$$

Since

$$F_1([-1, 0]) = [-1, 1], \quad F_2([-1, 0]) = [-1, 0],$$

F_1 and F_2 are distinct.

Let F be an arithmetic function with parameters A_1, \dots, A_m . Let G be the arithmetic function with B_1, \dots, B_m in place of A_1, \dots, A_m respectively; that means, it is given by

$$G(X) = F(B_1, \dots, B_m; X) .$$

If $B_1 \subset A_1, \dots, B_m \subset A_m$, then, by (6-4), $\mathcal{D}(F) \subset \mathcal{D}(G)$ and for every $X \in \mathcal{D}(F)$, $G(X) \subset F(X)$. Hence the relation that $B_1 \subset A_1, \dots, B_m \subset A_m$ will be written

$$G \subset F .$$

Let f be a rational function and let F be an arithmetic function. If there is a special arithmetic function $G \subset F$ with $f = p^{-1}Gp$, we say that f is an associated rational function of F or that F is an associated arithmetic function of f . In particular, if $G = F$ and then $f = p^{-1}Fp$, we say that F is an associated special arithmetic function of f and f is the associated rational function of F .

8. FIRST APPROXIMATION THEOREM

Let $X \in \mathcal{A}$. By a subdivision of X we mean

$$\xi = \{ \xi_1, \xi_2, \dots, \xi_r \}$$

such that

$$\begin{aligned} \alpha(X) = \alpha(\xi_1) < \beta(\xi_1) = \alpha(\xi_2) < \beta(\xi_2) = \alpha(\xi_3) < \\ \dots < \beta(\xi_{r-1}) = \alpha(\xi_r) < \beta(\xi_r) = \beta(X) . \end{aligned}$$

For every subdivision $\xi = \{ \xi_1, \xi_2, \dots, \xi_r \}$ we let

$$\sigma(\xi) = \max \{ \sigma(\xi_1), \sigma(\xi_2), \dots, \sigma(\xi_r) \} .$$

Let f be a rational function of domain $D(f)$; let G be an associated special arithmetic function of f and let F be an associated arithmetic function of f with $F \supset G$ (see § 7).

Let I be an element of \mathcal{A} contained in $D(f)$. Then for every $x \in I$, $Gp(x) = pf(x)$ is well-defined. Since, by (6-5), the domain $\mathcal{D}(G)$ of G is open, there is, for every $x \in I$, a positive number r_x such that whenever $Y \in \mathcal{A}$ with $P(p(x), Y) < r_x$, $G(Y)$ is well-defined. Let

$$I_x = [x - r_x/2, x + r_x/2] .$$

Then $G(I_x)$ is well-defined.

Since I is compact, there exist a finite number of points of I , say x_1, x_2, \dots, x_t , such that I is contained in the union of the interior

$$Q_i = (x_i - r_{x_i}/2, x_i + r_{x_i}/2)$$

of I_{x_i} , $i = 1, \dots, t$. We abbreviate I_{x_i} by I_i .

Let B_1, \dots, B_m be the parameters of G and let A_1, \dots, A_m be the parameters of F . By definition,

$$B_1, \dots, B_m \in p(R) ;$$

$$B_1 \subset A_1, \dots, B_m \subset A_m ;$$

$$G(Y) = F(B_1, \dots, B_m; Y) , \quad Y \in \mathcal{D}(G) ;$$

$$F(Y) = F(A_1, \dots, A_m; Y) , \quad Y \in \mathcal{D}(F) .$$

By (6-3) there is, for every $i = 1, \dots, t$, a $\delta_i > 0$ such that whenever $\sigma(A_1) < \delta_i, \dots, \sigma(A_m) < \delta_i$, $F(I_i)$ is defined. Let

$$3\delta = \min \{ \delta_1, \dots, \delta_t \} .$$

Then whenever $\sigma(A_1) < 3\delta$, ..., $\sigma(A_m) < 3\delta$, $F(I_1)$ is defined for all i .

By a well-known theorem on compact metric spaces there exists a $\delta' > 0$ such that whenever $Y \in \mathcal{I}_I$ with $\sigma(Y) < \delta'$, Y is contained in one of I_1, \dots, I_t so that $F(Y)$ is defined by (6-5).

Let $X \in \mathcal{I}_I$ and let $\xi = \{\xi_1, \dots, \xi_r\}$ be a subdivision of X with $\sigma(\xi) < \delta'$. Then

$$F(X, \xi) = F(\xi_1) \cup \dots \cup F(\xi_r),$$

$$\Sigma(F, X, \xi) = (F(\xi_1) \otimes p\sigma(\xi_1)) \oplus \dots \oplus (F(\xi_r) \otimes p\sigma(\xi_r))$$

are well-defined when the parameters A_1, \dots, A_m satisfy

$$\sigma(A_1) < 3\delta, \dots, \sigma(A_m) < 3\delta.$$

Clearly $\Sigma(F, X, \xi) \in \mathcal{I}$. Since for every $j = 2, \dots, r$,

$$F(\xi_{j-1}) \cap F(\xi_j) \subset F(p\alpha(\xi_j)) \neq \emptyset,$$

it follows that $F(X, \xi) \in \mathcal{I}$.

By (6-2), there is, for every $i = 1, \dots, t$, a positive number k_1 such that for any $Y, Y' \in \mathcal{I}_{I_i}$,

$$P(G(Y), G(Y')) \leq k_1 P(Y, Y'),$$

$$P(F(Y), G(Y)) \leq k_1 (\sigma(A_1) + \dots + \sigma(A_m))$$

hold where the parameters satisfy

$$\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta,$$

and k_1 is independent of the parameters.

Let $k = \max \{k_1, \dots, k_t\}$.

Then for any $Y, Y' \in \mathcal{I}_I$ with $\sigma(Y) < \delta$ and $\sigma(Y') < \delta'$

$$P(F(Y), G(Y')) \leq k(P(Y, Y') + \sigma(A_1) + \dots + \sigma(A_m)) .$$

In fact, there is a finite sequence

$$Y_1 = Y, Y_2, \dots, Y_s = Y'$$

in \mathcal{I}_I such that

$$\sigma(Y_2) = \dots = \sigma(Y_s) = \sigma(Y') ,$$

$$\alpha(Y_2) < \dots < \alpha(Y_s) = \alpha(Y')$$

and Y_{j-1} and Y_j are contained in the same \mathcal{I}_i for some i . Hence

$$\begin{aligned} kP(Y, Y') &= k(P(Y_1, Y_2) + \dots + P(Y_{s-1}, Y_s)) \\ &\geq P(G(Y_1), G(Y_2)) + \dots + P(G(Y_{s-1}), G(Y_s)) \\ &\geq P(G(Y), G(Y')) . \end{aligned}$$

Moreover

$$k(\sigma(A_1) + \dots + \sigma(A_m)) \geq P(F(Y), G(Y)) .$$

Our assertion thus follows.

Now we are ready to prove

Theorem 1. Let f be a rational function of domain $D(f)$, let G be an associated special arithmetic function of f and let \mathcal{I} be an element of \mathcal{A} contained in $D(f)$. Then there are positive numbers δ, δ', k such that whenever F is an associated arithmetic function of f with

parameters A_1, \dots, A_m such that $F \supset G$ and $\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta,$
 X is an element of \mathcal{A}_I and $\xi = \{\xi_1, \dots, \xi_r\}$ is a subdivision of X
with $\sigma(\xi) < \delta',$

$$F(X, \xi) = F(\xi_1) \cup \dots \cup F(\xi_r)$$

is defined and satisfies

$$f(X) \subset F(X, \xi) \subset f(X) \oplus [-\kappa, \kappa],$$

where

$$f(X) = \{f(x) \mid x \in X\},$$

$$\kappa = k(\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m)).$$

Proof. Let δ, δ', k be chosen as above. Let F be an associated arithmetic function of f with parameters A_1, \dots, A_m such that $F \subset G$ and $\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta$. Let $X \in \mathcal{A}_I$ and let ξ be a subdivision of X with $\sigma(\xi) < \delta'$. We have shown that $F(X, \xi)$ is defined.

For every $y \in f(X)$ there is an $x \in X$ with $f(x) = y$. Let $x \in \xi_j$. Then, by (6-5),

$$y = f(x) \in Fp(x) \subset F(\xi_j) \subset F(X, \xi).$$

Hence $f(X) \subset F(X, \xi)$.

For every $y \in F(X, \xi)$, there is a ξ_j with $y \in F(\xi_j)$. Let $x \in \xi_j$. Then $P(\xi_j, p(x)) \leq \sigma(\xi_j) < \sigma(\xi)$. It follows that

$$P(F(\xi_j), Gp(x)) \leq k(\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m))$$

or

$$P(F(\xi_j), pf(x)) \leq \kappa.$$

Therefore

$$\begin{aligned} y \in F(\xi_j) &\subset pf(x) \oplus [-\kappa, \kappa] \\ &\subset f(X) \oplus [-\kappa, \kappa] \end{aligned}$$

Hence $F(X, \xi) \subset f(X) \oplus [-\kappa, \kappa]$. q.e.d.

Corollary 1. Let f be a rational function of domain $D(f)$, let G be an associated special arithmetic function of f and let X be an element of \mathcal{I} contained in $D(f)$. Let F be an associated arithmetic function of f with parameters A_1, \dots, A_m such that $F \supset G$. Then whenever ξ is a subdivision of X with small $\sigma(\xi)$ and $\sigma(A_1), \dots, \sigma(A_m)$ are small, $F(X, \xi)$ is defined. Moreover, as $\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m) \rightarrow 0$,

$$\lim F(X, \xi) = f(X) ,$$

that means,

$$\lim P(F(X, \xi), f(X)) = 0 .$$

9. SECOND APPROXIMATION THEOREM

Theorem 2. Let f be a rational function of domain $D(f)$, let G be an associated special arithmetic function of f and let I be an element of \mathcal{I} contained in $D(f)$. Then there are positive numbers δ, δ', k such that whenever F is an associated arithmetic function of F with parameters A_1, \dots, A_m such that $F \supset G$ and $\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta, X = [a, b]$ is an element of \mathcal{I} and ξ is a subdivision of X ,

$$\Sigma(F, X, \xi) = (F(\xi_1) \otimes p\sigma(\xi_1)) \oplus \dots \oplus (F(\xi_r) \otimes p\sigma(\xi_r))$$

is defined and satisfies

$$p \left(\int_a^b f(x) dx \right) \subset \Sigma(F, X, \xi) \subset p \left(\int_a^b f(x) dx \right) + [-\kappa\sigma(X), \kappa\sigma(X)] ,$$

where

$$\kappa = k \left(\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m) \right) .$$

Proof. Let δ, δ', k be as before. Let F be an associated arithmetic function of F with parameters A_1, \dots, A_m such that $F \supset G$ and $\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta$. Let $X = [a, b]$ be an element of \mathcal{A}_I and let $\xi = \{ \xi_1, \dots, \xi_r \}$ be a subdivision of X with $\sigma(\xi) < \delta'$. We have shown that $\Sigma(F, X, \xi)$ is defined.

Let

$$m_j = \inf_{x \in \xi_j} f(x) \quad , \quad M_j = \sup_{x \in \xi_j} f(x) .$$

Then

$$\sum_{j=1}^r m_j \sigma(\xi_j) \leq \int_a^b f(x) dx \leq \sum_{j=1}^r M_j \sigma(\xi_j) .$$

Since ξ_j is compact, there is a point x_j of ξ_j with $f(x_j) = m_j$.

Then

$$m_j = f(x_j) \in F p(x_j) \subset F(\xi_j)$$

so that

$$m_j \sigma(\xi_j) \in F(\xi_j) \otimes p\sigma(\xi_j) .$$

Hence

$$\sum_{j=1}^r m_j \sigma(\xi_j) \in \Sigma(F, X, \xi) .$$

Similarly we can show that

$$\sum_{j=1}^r M_j \sigma(\xi_j) \in \Sigma(F, X, \xi) .$$

Since $\Sigma(F, X, \xi) \in \mathcal{A}$, it follows that

$$p\left(\int_a^b f(x)dx\right) \subset \Sigma(F, X, \xi) .$$

Let x_j be as above. Since

$$P\left(\xi_j, p(x_j)\right) \leq \sigma(\xi_j) \leq \sigma(\xi) ,$$

it follows that

$$P\left(F(\xi_j), Gp(x_j)\right) \leq k\left(\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m)\right) = \kappa$$

so that

$$F(\xi_j) \subset p(m_j) \oplus [-\kappa, \kappa]$$

Therefore

$$F(\xi_j) \oplus p\sigma(\xi_j) \subset p\left(m_j \sigma(\xi_j)\right) \oplus [-\kappa \sigma(\xi_j), \kappa \sigma(\xi_j)] .$$

Hence

$$\Sigma(F, X, \xi) \subset p\left(\sum_{j=1}^r m_j \sigma(\xi_j)\right) \oplus [-\kappa \sigma(X), \kappa \sigma(X)] .$$

Similarly we can prove that

$$\Sigma(F, X, \xi) \subset p\left(\sum_{j=1}^r m_j \sigma(\xi_j)\right) \oplus [-\kappa \sigma(X), \kappa \sigma(X)] .$$

Since $\Sigma(F, X, \xi) \in \mathcal{A}$, it follows that

$$\Sigma(F, X, \xi) \subset p\left(\int_a^b f(x)dx\right) \oplus [-\kappa \sigma(X), \kappa \sigma(X)] .$$

q.e.d.

Corollary. Let f be a rational function of domain $D(f)$, let G be an associated special arithmetic function of f and let $X = [a, b]$ be an

element of \mathcal{A} contained in $D(f)$. Let F be an associated arithmetic function of f with parameters A_1, \dots, A_m such that $F \supset G$. Then whenever ξ is a subdivision of X with small $\sigma(\xi)$ and $\sigma(A_1), \dots, \sigma(A_m)$ are small, $\Sigma(F, X, \xi)$ are defined. Moreover, as $\sigma(\xi) + \sigma(A_1) + \dots + \sigma(A_m) \rightarrow 0$,

$$\lim \Sigma(F, X, \xi) = \int_a^b f(x) dx ,$$

that means

$$\lim P\left(\Sigma(F, X, \xi), p\left(\int_a^b f(x) dx\right)\right) = 0 .$$

10. APPROXIMATION OF A CONTINUOUS FUNCTION BY ARITHMETIC FUNCTIONS

Let I be an element of \mathcal{A} and let $f = I \rightarrow R$ be a continuous function. Let

$$\{F_n\} = \{F_1, F_2, \dots\}$$

be a sequence of arithmetic functions such that $p(I)$ is contained in the domain $D(F_n)$ for all n . If for every $x \in I$,

$$F_1 p(x) \supset F_2 p(x) \supset \dots$$

and

$$\bigcap_{n=1}^{\infty} F_n p(x) = p f(x) ,$$

we say that $\{F_n\}$ converges to f . In symbols,

$$\lim_{n \rightarrow \infty} F_n = f .$$

(10-1) Let f be a rational function of domain $D(f)$ and let I be an element of \mathcal{A} contained in $D(f)$. Let F be an associated arithmetic

function of f with parameters A_1, \dots, A_m . Then for every $\epsilon > 0$ there is a $\delta > 0$ such that if $\sigma(A_1) < \delta, \dots, \sigma(A_m) < \delta$, then for every $x \in I$, $Fp(x)$ is defined and $\sigma(Fp(x)) < \epsilon$.

Proof. Let G be the associated special arithmetic function of f with $F \supset G$ and let B_1, \dots, B_m be the parameters of G . By (6-3) there is, for every $x \in I$, a $\delta_x > 0$ such that if $P(A_1, B_1) < \delta_x, \dots, P(A_m, B_m) < \delta_x$ and X is an element of \mathcal{I} with $P(X, p(x)) < \delta_x$, then $F(X)$ is defined and $P(F(X), Gp(x)) < \epsilon/2$. Since I is compact there exist a finite number of points of I , say x_1, \dots, x_t , such that the union of

$$\left(x_j - \delta_{x_j}, x_j + \delta_{x_j} \right), \quad j = 1, \dots, t,$$

contains I . Let

$$\delta = \min \left\{ \delta_{x_1}, \dots, \delta_{x_t} \right\}.$$

It follows that if $P(A_1, B_1) < \delta, \dots, P(A_m, B_m) < \delta$, then for every $x \in I$,

$$\sigma(Fp(x)) < \epsilon.$$

In fact, there is an x_j with $P(p(x), p(x_j)) < \delta_j$ so that $P(Fp(x), pf(x)) < \epsilon/2$. Hence

$$\sigma(Fp(x)) < \epsilon. \quad \text{q.e.d.}$$

Let f be a rational function of domain $D(f)$ and let I be an element of \mathcal{I} contained in $D(f)$. By (10-1) we can easily construct a sequence

of arithmetic functions

$$F_1 \supset F_2 \supset \dots$$

such that for every $x \in I$,

$$\sigma(F_n p(x)) < 1/n, \quad n = 1, 2, \dots$$

Hence

$$\lim_{n \rightarrow \infty} F_n = f.$$

This result can be extended as follows.

(10-2) For every continuous function $f : I \rightarrow R$, $I \in \mathcal{I}$, there is a sequence of arithmetic functions

$$F_n, \quad n = 1, 2, \dots$$

with

$$\lim_{n \rightarrow \infty} F_n = f.$$

Proof. It is well-known that every continuous function can be approximated by polynomials. For every integer $n > 1$ we let f_n be a polynomial such that for every $x \in I$, $\rho(f_n(x), f(x)) < 1/(6 \cdot 3^n)$. By (10-1), there is an arithmetic function G_n such that for every $x \in I$, $G_n p(x)$ is defined and

$$P(G_n p(x), p f_n(x)) \leq 1/(6 \cdot 3^n)$$

Let F_n be the arithmetic function such that

$$F_n(x) = G_n(x) \oplus \left[-1/(6 \cdot 3^n), 1/(6 \cdot 3^n) \right], \quad x \in \mathcal{D}(G_n).$$

Then for every $x \in I$, $F_n p(x)$ is defined. Since

$$\begin{aligned} P(G_n p(x), pf(x)) &\leq P(G_n p(x), pf_n(x)) + P(pf_n(x), pf(x)) \\ &\leq 1/(6 \cdot 3^n) + 1/(6 \cdot 3^n) = 2/(6 \cdot 3^n), \end{aligned}$$

it follows that

$$\begin{aligned} pf(x) &\subset G_n p(x) \oplus \left[-2/(6 \cdot 3^n), 2/(6 \cdot 3^n) \right], \\ G_n p(x) &\subset pf(x) \oplus \left[-2/(6 \cdot 3^n), 2/(6 \cdot 3^n) \right]. \end{aligned}$$

Hence

$$\begin{aligned} pf(x) \oplus \left[-2/(6 \cdot 3^n), 2/(6 \cdot 3^n) \right] &\subset F_n p(x) \\ &\subset pf(x) \oplus \left[-1/3^n, 1/3^n \right] \end{aligned}$$

Consequently

$$F_{n+1} p(x) \subset pf(x) \oplus \left[-1/3^{n+1}, 1/3^{n+1} \right] \subset F_n p(x).$$

Since $\sigma(F_n p(x)) < 2/3^n$ and $\lim_{n \rightarrow \infty} 2/3^n = 0$, it follows that $\lim F_n = f$.

q.e.d.

Let I be an element of \mathcal{I} and let

$$\xi = \{ \xi_1, \dots, \xi_r \}, \quad \eta = \{ \eta_1, \dots, \eta_s \}$$

be subdivisions of I . If there exist integers $1 \leq j(1) < j(2) < \dots < j(r) = s$ such that for every $i = 1, \dots, r$, $\{ \eta_{j(i-1)+1}, \dots, \eta_{j(i)} \}$ is a subdivision of ξ_i , we write $\xi \prec \eta$ and call η a refinement of ξ .

Let $I \in \mathcal{I}$ and let F be an arithmetic function with $\mathcal{D}(F) \supset p(I)$. As before, there is a positive number δ such that whenever $X \in \mathcal{I}$ with

$\sigma(X) < \delta$, $F(X)$ is defined. Hence if ξ is a subdivision of I with $\sigma(\xi) < \delta$, both $F(I, \xi)$ and $\Sigma(F, I, \xi)$ are defined.

(10-3) Let F be an arithmetic function and let $\xi = \{\xi_1, \dots, \xi_r\}$ be a subdivision of $I \in \mathcal{A}$ such that $F(\xi_1), \dots, F(\xi_r)$ are defined. Then for every refinement $\eta = \{\eta_1, \dots, \eta_s\}$ of ξ , $F(\eta_1), \dots, F(\eta_s)$ are defined so that $F(I, \eta)$ and $\Sigma(F, I, \eta)$ are defined. Moreover,

$$F(I, \eta) \subset F(I, \xi), \quad \Sigma(F, I, \eta) \subset \Sigma(F, I, \xi) .$$

(10-4) Let $I \in \mathcal{A}$ and let F be an arithmetic function. Let $\xi = \{\xi_1, \dots, \xi_r\}$ be a subdivision of I such that $F(\xi_1), \dots, F(\xi_r)$ are defined, and let

$$\xi = \xi^{(1)} \prec \xi^{(2)} \prec \dots$$

be a sequence of subdivisions of I . Then $F(I, \xi^{(n)})$ and $\Sigma(F, I, \xi^{(n)})$ are defined for all n and

$$\begin{aligned} F(I, \xi^{(1)}) \supset F(I, \xi^{(2)}) \supset \dots , \\ \Sigma(F, I, \xi^{(1)}) \supset \Sigma(F, I, \xi^{(2)}) \supset \dots . \end{aligned}$$

If $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$, then $\bigcap_{n=1}^{\infty} F(I, \xi^{(n)})$ and $\bigcap_{n=1}^{\infty} \Sigma(F, I, \xi^{(n)})$ are independent of the choice of ξ and $\{\xi^{(n)}\}$.

Proof. Let $\eta = \{\eta_1, \dots, \eta_s\}$ be a subdivision of I such that $F(\eta_1), \dots, F(\eta_s)$ are defined and let $\zeta = \{\zeta_1, \dots, \zeta_t\}$ be a refinement of η . By definition, there are integers $1 \leq j(1) < \dots < j(s) = t$

such that for every $i = 1, \dots, s$, $\{\xi_{j(i-1)+1}, \dots, \xi_{j(i)}\}$ is a sub-division of η_i . It follows from (6-5) that

$$F(\xi_{j(i-1)+1}) \subset F(\eta_i), \dots, F(\xi_{j(i)}) \subset F(\eta_i) .$$

Hence

$$\begin{aligned} F(I, \xi) &= F(\xi_1) \cup \dots \cup F(\xi_t) \\ &\subset F(\eta_1) \cup \dots \cup F(\eta_s) = F(I, \eta) \end{aligned}$$

As a consequence of this result we have

$$F(I, \xi^{(1)}) \supset F(I, \xi^{(2)}) \supset \dots$$

Since

$$\begin{aligned} &\alpha \left(\left(F(\xi_{j(i-1)+1}) \otimes p\sigma(\xi_{j(i-1)+1}) \right) \oplus \dots \oplus \left(F(\xi_{j(i)}) \otimes p\sigma(\xi_{j(i)}) \right) \right) \\ &= \alpha \left(F(\xi_{j(i-1)+1}) \right) \sigma(\xi_{j(i-1)+1}) + \dots + \alpha \left(F(\xi_{j(i)}) \right) \sigma(\xi_{j(i)}) \\ &\geq \alpha \left(F(\eta_i) \right) \sigma(\xi_{j(i-1)+1}) + \dots + \alpha \left(F(\eta_i) \right) \sigma(\xi_{j(i)}) \\ &= \alpha \left(F(\eta_i) \right) \sigma(\eta_i) = \alpha \left(F(\eta_i) \otimes p\sigma(\eta_i) \right) \end{aligned}$$

and similarly

$$\begin{aligned} &\beta \left(\left(F(\xi_{j(i-1)+1}) \otimes p\sigma(\xi_{j(i-1)+1}) \right) \oplus \dots \oplus \left(F(\xi_{j(i)}) \otimes p\sigma(\xi_{j(i)}) \right) \right) \\ &\leq \beta \left(F(\eta_i) \otimes p\sigma(\eta_i) \right) , \end{aligned}$$

it follows that

$$\begin{aligned} & \left(F(\xi_{j(i-1)+1}) \otimes p\sigma(\xi_{j(i-1)+1}) \right) \oplus \dots \oplus \left(F(\xi_{j(i)}) \otimes p\sigma(\xi_{j(i)}) \right) \\ & \subset F(\eta_i) \otimes p\sigma(\eta_i) . \end{aligned}$$

Hence

$$\begin{aligned} \Sigma(F, I, \xi) &= \left(F(\xi_1) \otimes p\sigma(\xi_1) \right) \oplus \dots \oplus \left(F(\xi_t) \otimes p\sigma(\xi_t) \right) \\ &\subset \left(F(\eta_1) \otimes p\sigma(\eta_1) \right) \oplus \dots \oplus \left(F(\eta_s) \otimes p\sigma(\eta_s) \right) = \Sigma(F, I, \eta) . \end{aligned}$$

From this result it follows that

$$\Sigma(F, I, \xi^{(1)}) \supset \Sigma(F, I, \xi^{(2)}) \supset \dots$$

Now we assume $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$. Let $\eta = \{\eta_1, \dots, \eta_s\}$ be a subdivision of I such that $F(\eta_1), \dots, F(\eta_s)$ are defined, and let

$$\eta = \eta^{(1)} \prec \eta^{(2)} \prec \dots$$

be a sequence of subdivisions with $\lim_{n \rightarrow \infty} \sigma(\eta^{(n)}) = 0$.

Let $\epsilon > 0$. For any $y \in \bigcap_{n=1}^{\infty} F(I, \xi^{(n)})$ there is a sequence

$$\xi_{i(1)}^{(1)} \supset \xi_{i(2)}^{(2)} \supset \dots$$

in \mathcal{J} such that for every $n = 1, 2, \dots$, $\xi_{i(n)}^{(n)} \in \xi^{(n)}$ and $F(\xi_{i(n)}^{(n)}) \ni y$.

Since $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$, $\bigcap_{n=1}^{\infty} \left(\xi_{i(n)}^{(n)} \right)$ contains a single point x . By

(6-5) there is a $\delta > 0$ such that whenever $X \in \mathcal{J}_I$ with $P(X, p(x)) < \delta$,

$F(X)$ is defined and $P(F(X), Fp(x)) < \varepsilon$. Since $\lim_{n \rightarrow \infty} \sigma(\eta^{(n)}) = 0$, $\sigma(\eta^{(n)}) < \delta$ holds for all large n . Let $\eta_i^{(n)}$ be the element of $\eta^{(n)}$ containing x . Then $P(\eta_i^{(n)}, p(x)) < \delta$ so that $P(F(\eta_i^{(n)}) Fp(x)) < \varepsilon$. Hence for all large n

$$y \in Fp(x) \subset F(\eta_i^{(n)}) \oplus [-\varepsilon, \varepsilon] \subset F(I, \eta^{(n)}) \oplus [-\varepsilon, \varepsilon]$$

This proves that $y \in \bigcap_{n=1}^{\infty} F(I, \eta^{(n)}) \oplus [-\varepsilon, \varepsilon]$. Since y is an arbitrary point of $\bigcap_{n=1}^{\infty} F(I, \xi^{(n)})$, it follows that

$$\bigcap_{n=1}^{\infty} F(I, \xi^{(n)}) \subset \bigcap_{n=1}^{\infty} F(I, \eta^{(n)}) \oplus [-\varepsilon, \varepsilon].$$

Similarly,

$$\bigcap_{n=1}^{\infty} F(I, \eta^{(n)}) \subset \bigcap_{n=1}^{\infty} F(I, \xi^{(n)}) \oplus [-\varepsilon, \varepsilon].$$

Applying (1-6), we have

$$\bigcap_{n=1}^{\infty} F(I, \xi^{(n)}) = \bigcap_{n=1}^{\infty} F(I, \eta^{(n)}).$$

In order to prove that $\bigcap_{n=1}^{\infty} \Sigma(F, I, \xi^{(n)}) = \bigcap_{n=1}^{\infty} \Sigma(F, I, \eta^{(n)})$, we may assume $\sigma(I) > 0$. It is sufficient to prove that for every integer $m \geq 1$ and every $\varepsilon > 0$,

$$\Sigma(F, I, \xi^{(m)}) \oplus [-\varepsilon, \varepsilon] \supset \Sigma(F, I, \eta^{(n)})$$

holds for large n . In fact, if this is proved, then

$$\Sigma(F, I, \xi^{(m)}) \oplus [-\varepsilon, \varepsilon] \supset \bigcap_{n=1}^{\infty} \Sigma(F, I, \eta^{(n)}).$$

Since ε is arbitrary, it follows that $\Sigma(F, I, \xi^{(m)}) \supset \bigcap_{n=1}^{\infty} \Sigma(F, I, \eta^{(n)})$.

Since m is arbitrary, it follows that

$$\bigcap_{m=1}^{\infty} \Sigma(F, I, \xi^{(m)}) \supset \bigcap_{n=1}^{\infty} \Sigma(F, I, \eta^{(n)}) .$$

Similarly,

$$\bigcap_{n=1}^{\infty} \Sigma(F, I, \eta^{(n)}) \supset \bigcap_{m=1}^{\infty} \Sigma(F, I, \xi^{(m)}) .$$

Hence our assertion follows.

Now we let m be an arbitrary integer ≥ 1 and let $\xi^{(m)} = \{A_1, \dots, A_u\}$.

Let ϵ be an arbitrary positive number < 1 . By (6-3), there is, for

every $x \in I$, a positive number r_x such that whenever $X \in \mathcal{A}$ with $P(X, p(x)) < r_x$. $F(X)$ is defined and $P(F(X), Fp(x)) < \epsilon/u \left(\gamma \left(F(I, \xi^{(m)}) \right) + 1 \right)$.

Since I is compact, there is a $\delta > 0$ such that whenever $X \in \mathcal{A}_I$ with

$\sigma(X) < \delta$, $P(X, p(x)) < r_x$ for some $x \in I$.

Since $\lim_{n \rightarrow \infty} \sigma(\eta^{(n)}) = 0$, there is an integer n_0 such that

$$\sigma(\eta^{(n)}) < \min \left(\delta, \epsilon/u \left(\gamma \left(F(I, \xi^{(m)}) \right) + 1 \right) \right)$$

for all integers $n > n_0$. Let $n > n_0$ and let

$$\eta^{(n)} = \{B_1, \dots, B_v\} .$$

For every $i = 1, \dots, u$, there is a largest integer $j(i)$ with $\beta(A_i) \in B_{j(i)}$. Clearly

$$1 \leq j(1) \leq j(2) \leq \dots \leq j(u) = v .$$

Since $P(B_{j(i)}, p\beta(A_i)) < \delta$, it follows that

$$P(F(B_{j(i)}), Fp\beta(A_i)) < \epsilon/u \left(\gamma \left(F(I, \xi^{(m)}) \right) + 1 \right) < 1 .$$

Therefore

$$F(B_{j(i)}) \subset F_{p\sigma}(A_i) \oplus [-1, 1] \subset F(A_i) \oplus [-1, 1]$$

and then

$$\begin{aligned} \gamma(F(B_{j(i)}) \otimes p\sigma(B_{j(i)})) &\leq (\gamma(F(A_i)) + 1) \sigma(B_{j(i)}) \\ &\leq (\gamma(F(I, \xi^{(m)})) + 1) \sigma(B_{j(i)}) . \end{aligned}$$

Hence

$$\begin{aligned} \sum_{i=1}^u \gamma(F(B_{j(i)}) \otimes p\sigma(B_{j(i)})) &\leq (\gamma(F(I, \xi^{(m)})) + 1) \sum_{i=1}^m \sigma(B_{j(i)}) \\ &\leq (\gamma(F(I, \xi^{(m)})) + 1) u \sigma(\eta^{(n)}) < \varepsilon \end{aligned}$$

and consequently

$$(F(B_{j(1)}) \otimes p\sigma(B_{j(1)})) \oplus \dots \oplus (F(B_{j(u)}) \otimes p\sigma(B_{j(u)})) \subset [-\varepsilon, \varepsilon]$$

Since, for every $k = j(i-1) + 1, \dots, j(i) - 1$, $B_k \subset A_i$, it follows that

$$\begin{aligned} (F(B_{j(i-1)+1}) \otimes p\sigma(B_{j(i-1)+1})) \oplus \dots \oplus (F(B_{j(i)-1}) \otimes p\sigma(B_{j(i)-1})) \\ \subset F(A_i) \otimes p\sigma(A_i) , \quad i = 1, \dots, u . \end{aligned}$$

Hence, by adding these equations, we have

$$\Sigma(F, I, \eta^{(n)}) \subset \Sigma(F, I, \xi^{(m)}) \oplus [-\varepsilon, \varepsilon] .$$

This completes our proof.

q.e.d.

Let F be an arithmetic function with domain $\mathcal{D}(F)$. Let \mathcal{A}^F be the subset of \mathcal{A} consisting of all the elements X of \mathcal{A} with $p(X) \subset \mathcal{D}(F)$. Then, by (10-4), we may define a function

$$\bar{F} : \mathcal{A}^F \longrightarrow \mathcal{A}$$

by

$$\bar{F}(X) = \bigcap_{n=1}^{\infty} F(X, \xi^{(n)}),$$

where $X \in \mathcal{A}^F$ and $\xi^{(1)} \prec \xi^{(2)} \prec \dots$ is a sequence of subdivisions of X with $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$.

(10-5) Let X and Y be elements of \mathcal{A}^F with $\beta(X) = \alpha(Y)$. Then $X \cup Y \in \mathcal{A}^F$ and

$$\bar{F}(X \cup Y) = \bar{F}(X) \cup \bar{F}(Y).$$

Proof. Since $\beta(X) = \alpha(Y)$, $X \cup Y \in \mathcal{A}$. Since $X \cup Y \in \mathcal{A}^F$, $p(X) \subset \mathcal{D}(F)$ and $p(Y) \subset \mathcal{D}(F)$. Hence $p(X \cup Y) \subset \mathcal{D}(F)$ and consequently $X \cup Y \in \mathcal{A}^F$.

If $Y \in p(R)$, then $\bar{F}(Y) = F(Y) \subset \bar{F}(X)$ so that our assertion is obvious. If $Y \notin p(R)$, then we may have a sequence of subdivisions

$$\xi^{(n)} = \left\{ \xi_1^{(n)}, \dots, \xi_{r(n)}^{(n)} \right\}$$

of $X \cup Y$ such that $\xi^{(1)} \prec \xi^{(2)} \prec \dots$, $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$ and for every integer n there is an integer $s(n)$ with $\beta\left(\xi_{s(n)}^{(n)}\right) = \beta(X)$. Therefore

$$\eta^{(n)} = \left\{ \xi_1^{(n)}, \dots, \xi_{s(n)}^{(n)} \right\}$$

is a sequence of subdivisions of X such that $\eta^{(1)} \prec \eta^{(2)} \prec \dots$ and

$\lim_{n \rightarrow \infty} \sigma(\eta^{(n)}) = 0$; and $\zeta^{(n)} = \left\{ \xi_{s(n)+1}^{(n)}, \dots, \xi_{r(n)}^{(n)} \right\}$ is a sequence of

subdivisions of Y such that $\zeta^{(1)} \prec \zeta^{(2)} \prec \dots$ and $\lim_{n \rightarrow \infty} \sigma(\zeta^{(n)}) = 0$.
 Since

$$F(X \cup Y, \xi^{(n)}) = F(X, \eta^{(n)}) \cup F(Y, \zeta^{(n)}),$$

it follows that

$$\overline{F}(X \cup Y) = \overline{F}(X) \cup \overline{F}(Y) \quad . \quad \text{q.e.d.}$$

(10-6) $\overline{F} : \mathcal{A}^F \rightarrow \mathcal{A}$ is continuous.

Proof. Let $Y \in \mathcal{A}^F$ and let $\varepsilon > 0$. Since $F(\alpha(Y))$ is defined, there is a $\delta > 0$ such that whenever $Z \in \mathcal{A}$ with $\sigma(Z) < \delta$, $P(Z, p\alpha(Y)) < \delta$, $F(Z)$ is defined and $P(F(Z), Fp\alpha(Y)) < \varepsilon/4$. Since

$$F(Z) \supset \overline{F}(Z) \supset Fp\alpha(Y) = \overline{F}p\alpha(Y) \quad ,$$

we have

$$P(\overline{F}(Z), \overline{F}p\alpha(Y)) < \varepsilon/4 \quad .$$

Similarly there is a $\delta' > 0$ such that whenever $Z' \in \mathcal{A}$ with $P(Z', p\beta(Z)) < \delta'$, $F(Z')$ is defined and

$$P(\overline{F}(Z'), \overline{F}p\beta(Z)) < \varepsilon/4 \quad .$$

For every $X \in \mathcal{A}_F$ with $P(X, Y) < \min(\delta, \delta')$ we have $Z, Z' \in \mathcal{A}$ such that

$$P(Z, p\alpha(Y)) < \delta \quad , \quad P(Z', p\beta(Y)) < \delta'$$

and one of the following holds:

- (1) $\alpha(X) = \alpha(Z)$, $\beta(Z) = \alpha(Y)$, $\beta(Y) = \alpha(Z')$, $\beta(Z') = \beta(X)$;
- (2) $\alpha(X) = \alpha(Z)$, $\beta(X) = \alpha(Z')$, $\beta(Z) = \alpha(Y)$, $\beta(Z') = \beta(Y)$;
- (3) $\alpha(Z) = \alpha(Y)$, $\beta(Z) = \alpha(X)$, $\beta(Y) = \alpha(Z')$, $\beta(X) = \beta(Z')$;
- (4) $\alpha(Z) = \alpha(Y)$, $\beta(Z) = \alpha(X)$, $\beta(X) = \alpha(Z')$, $\beta(Z') = \beta(Y)$.

In case (1) we have, by (10-5),

$$\bar{F}(X) = \bar{F}(Z) \cup \bar{F}(Y) \cup \bar{F}(Z')$$

Then $\bar{F}(X) \supset \bar{F}(Y)$. Since

$$\bar{F}(Z) \subset \bar{F}_p \alpha(Y) \oplus [-\epsilon/2, \epsilon/2] \subset \bar{F}(Y) \oplus [-\epsilon/2, \epsilon/2],$$

and similarly,

$$\bar{F}(Z') \subset \bar{F}(Y) \oplus [-\epsilon/2, \epsilon/2],$$

it follows that

$$\bar{F}(X) \subset \bar{F}(Y) \oplus [-\epsilon, \epsilon].$$

Hence $P(\bar{F}(X), \bar{F}(Y)) < \epsilon$.

In case (2), we have, by (10-5),

$$\bar{F}(X) \cup \bar{F}(Z') = \bar{F}(Z) \cup \bar{F}(Y).$$

Since $\bar{F}(Z) \subset \bar{F}(Y) \oplus [-\epsilon/2, \epsilon/2]$ and $\bar{F}(Z') \subset \bar{F}(Y) \oplus [-\epsilon/2, \epsilon/2]$, it follows that $\bar{F}(X) \subset \bar{F}(Y) \oplus [-\epsilon, \epsilon]$. On the other hand,

$$\begin{aligned} P(\bar{F}(Z'), \bar{F}_p \beta(X)) &\leq P(\bar{F}(Z'), \bar{F}_p \beta(Y)) \oplus P(\bar{F}_p \beta(X), \bar{F}_p \beta(Y)) \\ &< \epsilon/4 + \epsilon/4 = \epsilon/2 \end{aligned}$$

so that

$$\bar{F}(Z') \subset \bar{F}(X) \oplus [-\epsilon/2, \epsilon/2].$$

Hence $\bar{F}(Y) \subset \bar{F}(X) \oplus [-\epsilon, \epsilon]$. This again proves that $P(\bar{F}(X), \bar{F}(Y)) < \epsilon$.

Similar argument shows that our assertion also holds for the other two cases. q.e.d.

(10-7) Let I be an element of \mathcal{A} and let $f : I \rightarrow \mathbb{R}$ be a continuous function. Let $\{F_n\}$ be a sequence of arithmetic functions with $\lim_{n \rightarrow \infty} F_n = f$. Then

$$\overline{F}_1(I) \supset \overline{F}_2(I) \supset \dots$$

$$\text{and } \bigcap_{n=1}^{\infty} \overline{F}_n(I) = f(I).$$

Proof. Let n be an integer ≥ 1 . Let $\varepsilon > 0$. For every $x \in \mathcal{A}$ there is an $r_x > 0$ such that whenever $X \in \mathcal{A}$ with $P(X, p(x)) < r_x$, $F_n(X)$ and $F_{n+1}(X)$ are defined and $P(F_{n+1}(X), F_{n+1}p(x)) < \varepsilon$. Let δ be a positive number such that whenever $X \in \mathcal{A}_I$ with $\sigma(X) < \delta$, we have $P(X, p(x)) < r_x$ for some $x \in X$.

Let $\xi = \{\xi_1, \dots, \xi_r\}$ be any subdivision of I with $\sigma(\xi) < \delta$. For every $i = 1, \dots, r$, there is an $x_i \in I$ such that $P(\xi_i, p(x_i)) < r_{x_i}$.

Then

$$F_{n+1}(\xi_i) \subset F_{n+1} p(x_i) \oplus [-\varepsilon, \varepsilon] \subset F_n p(x_i) \oplus [-\varepsilon, \varepsilon]$$

$$F_n(I, \xi) \oplus [-\varepsilon, \varepsilon]$$

so that

$$F_{n+1}(I, \xi) \subset F_n(I, \xi) \oplus [-\varepsilon, \varepsilon].$$

Hence

$$\overline{F}_{n+1}(I) \subset \overline{F}_n(I) \oplus [-\varepsilon, \varepsilon].$$

Since ϵ is arbitrary, it follows that

$$\overline{F}_{n+1}(I) \subset \overline{F}_n(I) .$$

This proves that

$$\overline{F}_1(I) \supset \overline{F}_2(I) \supset \dots .$$

Clearly $\overline{F}_n(I) \supset f(I)$ for all n . Let $\epsilon > 0$.

For every $n = 1, 2, \dots$ we let

$$I_n = \left\{ x \in I \mid P(F_n p(x), p f(x)) \geq \epsilon \right\} .$$

It follows from the continuity of F_n that I_n is closed in I . Since $I_1 \supset I_2 \supset \dots$ and $\bigcap_{n=1}^{\infty} I_n = \emptyset$, we infer that there is an integer n_0 such that whenever $n > n_0$, $P(F_n p(x), p f(x)) < \epsilon$ holds for all $x \in I$. Let $n > n_0$. Then for any subdivision $\xi = \{\xi_1, \dots, \xi_r\}$ of I we have

$$F_n(\xi_i) \subset f(\xi_i) \oplus [-\epsilon, \epsilon] , \quad i = 1, \dots, r .$$

Hence

$$\overline{F}_n(I) \subset F_n(I, \xi) \subset f(I) \oplus [-\epsilon, \epsilon] .$$

Since ϵ is arbitrary, we infer that

$$\bigcap_{n=1}^{\infty} \overline{F}_n(I) = f(I) \quad \text{q.e.d.}$$

Theorem 3. Let $I \in \mathcal{A}$ and let $f : I \rightarrow R$ be a continuous function.
Let $\{F_n\}$ be a sequence of arithmetic functions with $\lim_{n \rightarrow \infty} F_n = f$.

Then there is a sequence $\xi^{(1)} \prec \xi^{(2)} \prec \dots$ of subdivisions of I such that $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$ and

$$\bigcap_{n=1}^{\infty} F_n(I, \xi^{(n)}) = f(I) .$$

If, moreover, for every $x \in I$ and every integer $n \geq 1$, $F_{n+1} p(x)$ is contained in the interior of $F_n p(x)$, then there is a sequence

$\xi^{(1)} \prec \xi^{(2)} \prec \dots$ of subdivisions of I such that $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$, $F_1(I, \xi^{(1)}) \supset F_2(I, \xi^{(2)}) \supset \dots$ and $\bigcap_{n=1}^{\infty} F_n(I, \xi^{(n)}) = f(I)$.

Proof. By (10-7), we have $\bar{F}_1(I) \supset \bar{F}_2(I) \supset \dots$ and $\bigcap_{n=1}^{\infty} \bar{F}_n(I) = f(I)$.

Let $\xi^{(1)}$ be a subdivision of I such that $\sigma(\xi^{(1)}) < 1$ and

$P(F(I, \xi^{(1)}), \bar{F}(I)) < 1$. Suppose that we have subdivisions

$\xi^{(1)} \prec \xi^{(2)} \prec \dots \prec \xi^{(k)}$ of I such that $\sigma(\xi^{(n)}) < 1/n$ and

$P(F(I, \xi^{(n)}), \bar{F}(I)) < 1/n$, $n = 1, \dots, k$. We let $\xi^{(k+1)}$ be a refine-

ment of $\xi^{(k)}$ with $P(F_{k+1}(I, \xi^{(k+1)}), \bar{F}_{k+1}(I)) < 1/(k+1)$. By induction,

we have a sequence $\xi^{(1)} \prec \xi^{(2)} \prec \dots$ of subdivisions of I with

$P(F_n(I, \xi^{(n)}), \bar{F}_n(I)) < 1/n$. Hence

$$\bigcap_{n=1}^{\infty} F_n(I, \xi^{(n)}) = f(I) .$$

If for every $x \in I$ and every integer $n \geq 1$, $F_{n+1} p(x)$ is contained in the interior of $F_n p(x)$, as in the proof of (10-7), there is a $\delta_n > 0$

such that whenever ξ is a subdivision of I with $\sigma(\xi) < \delta_n$,

$F_n(I, \xi) \supset F_{n+1}(I, \xi)$, $n = 1, 2, \dots$. Now we construct a sequence

$\xi^{(1)} < \xi^{(2)} < \dots$ of subdivisions of I , just as above, satisfying the additional condition that

$$\sigma(\xi^{(n)}) < \delta_n, \quad n = 1, 2, \dots$$

Then our conclusion follows.

q.e.d.

Let F be an arithmetic function with domain $D(F)$ and \mathcal{A}^F be the subset of \mathcal{A} consisting of all the elements I of \mathcal{A} with $p(I) \subset D(F)$. Then, by (10-4), we may define a function

$$\Sigma_F : \mathcal{A}^F \rightarrow \mathcal{A}$$

by

$$\Sigma_F(I) = \bigcap_{n=1}^{\infty} \Sigma(F, I, \xi^{(n)}),$$

where $I \in \mathcal{A}^F$ and $\xi^{(1)} < \xi^{(2)} < \dots$ is a sequence of subdivisions of I with $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$.

As (10-5), (10-6) and (10-7), we have

(10-8) Let I and J be elements of \mathcal{A}^F with $\beta(I) = \alpha(J)$. Then $I \cup J \in \mathcal{A}^F$ and

$$\Sigma_F(I \cup J) = \Sigma_F(I) \oplus \Sigma_F(J).$$

(10-9) $\Sigma_F : \mathcal{A}^F \rightarrow \mathcal{A}$ is continuous.

(10-10) Let $I = [a, b] \in \mathcal{A}$ and let $f : I \rightarrow \mathbb{R}$ be a continuous function. Let $\{F_n\}$ be a sequence of arithmetic functions with $\lim_{n \rightarrow \infty} F_n = f$.

Then

$$\Sigma_{F_1}(I) \supset \Sigma_{F_2}(I) \supset \dots$$

and $\bigcap_{n=1}^{\infty} \Sigma_{F_n}(I) = \int_a^b f(x) dx$

Making use of (10-10) and the definition of Σ_F we can prove

Theorem 4. Let $I = [a, b] \in \mathcal{I}$ and let $f : I \rightarrow R$ be a continuous function. Let $\{F_n\}$ be a sequence of arithmetic functions with
 $\lim_{n \rightarrow \infty} F_n = f$. Then there is a sequence $\xi^{(1)} \prec \xi^{(2)} \prec \dots$ of subdivisions of I such that $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$ and

$$\bigcap_{n=1}^{\infty} \Sigma(F_n, I, \xi^{(n)}) = \int_a^b f(x) dx .$$

If, moreover, for every $x \in I$ and every integer $n \geq 1$ $F_{n+1} p(x)$ is contained in the interior of $F_n p(x)$, then there is a sequence $\xi^{(1)} \prec \xi^{(2)} \prec \dots$ of subdivisions of I such that $\lim_{n \rightarrow \infty} \sigma(\xi^{(n)}) = 0$,
 $\Sigma(F_1, I, \xi^{(1)}) \supset \Sigma(F_2, I, \xi^{(2)}) \supset \dots$ and $\bigcap_{n=1}^{\infty} \Sigma(F_n, I, \xi^{(n)}) = \int_a^b f(x) dx$.

**CONVERGENCE OF APPROXIMATE
EIGENVECTORS IN JACOBI METHODS**

R.L. Causey
P. Henrici



2

FOREWORD

The work covered in this paper was begun by the authors at the Space Technology Laboratories, Los Angeles, under the sponsorship of the U.S. Air Force. It was concluded at Lockheed Missiles and Space Division, Sunnyvale, California, under the Lockheed General Research Program.

The paper is being published in the Numerische Mathematik, Vol. 2, Heidelberg, Springer-Verlag, 1960.

The authors wish to express their appreciation to Prof. George E. Forsythe for calling their attention to some of the problems discussed in this paper.

ABSTRACT

The paper constitutes a study of the convergence of infinite products of unitary matrices connected with Jacobi methods for computing eigenvalues of Hermitian matrices. Bounds are obtained for the error in approximate eigenvectors resulting from a finite number of steps of three Jacobi processes, namely, the classical Jacobi method, quasi-cyclic restricted Jacobi methods and threshold cyclic Jacobi methods. The results apply only to Hermitian matrices which do not have repeated eigenvalues. Also, certain questions are answered concerning the representation of an arbitrary unitary matrix as an infinite product.

CONTENTS

<u>Section</u>		<u>Page</u>
	Foreword	iii
	Abstract	v
1	Introduction	1
2	Norms of Matrices	5
3	Preliminary Theorem on Infinite Products	7
4	Convergence of Eigenvectors and Error Estimates	10
5	Remarks	18
6	References	21

CONVERGENCE OF APPROXIMATE EIGENVECTORS
IN JACOBI METHODS

1. INTRODUCTION. In this paper we are interested in studying the convergence of infinite products of unitary matrices connected with Jacobi methods [4] for computing eigenvalues of Hermitian matrices. We obtain bounds for the error in approximate eigenvectors resulting from a finite number of steps of a Jacobi process. Also, we are able to answer certain questions concerning the representation of an arbitrary unitary matrix as an infinite product. Use is made of the classical results of Jacobi [5] and of the recent results of Henrici [4] and Pope and Tompkins [7]. Basically, we adhere to the notation and terminology used in [4].

Consider the following computational algorithm. Let $A = A_0 = (a_{rc})$ be a Hermitian matrix of order n with eigenvalues $\lambda_r (r = 1, 2, \dots, n)$. One calculates a sequence of matrices $A_1, A_2, \dots, A_k = (a_{rc}^{(k)})$, ... which are unitarily similar to A by the recurrence relation

$$A_{k+1} = U_k^* A_k U_k \quad (k = 0, 1, 2, \dots). \quad (1.1)$$

The $U_k = (u_{rc}^{(k)})$ are special unitary matrices of order n . For every value of k there is specified a pair $\pi_k = (i_k, j_k) = (i, j)$ of indices (we omit the subscript k in the sequel for notational simplicity) satisfying $1 \leq i < j \leq n$, such that the 2×2 matrix

* Bracketed numbers indicate References, page 21.

$$V_k = \begin{pmatrix} u_{ii}^{(k)} & u_{ij}^{(k)} \\ u_{ji}^{(k)} & u_{jj}^{(k)} \end{pmatrix}, \quad (1.2)$$

which is a principal submatrix of U_k , is unitary. All other elements of U_k satisfy

$$u_{rc}^{(k)} = \delta_{rc} = \begin{cases} 1, & r = c \\ 0, & r \neq c \end{cases}. \quad (1.3)$$

The matrices U_k are completely determined by the pairs π_k and the 2×2 unitary matrices V_k .

Any set of rules for choosing the U_k is called a Jacobi method. A Jacobi method is said to be convergent if

$$A_k \rightarrow \Lambda \quad (k \rightarrow \infty), \quad (1.4)$$

for all A under the adopted set of rules, where Λ is a diagonal matrix whose diagonal elements are the λ_r in some order. Either of the quantities

$$S_k = \sum_{r \neq c} |a_{rc}^{(k)}|^2 \quad (1.5)$$

or

$$M_k = \max_{r \neq c} |a_{rc}^{(k)}| \quad (1.6)$$

may be used as a measure of the closeness of A_k to Λ . For the Jacobi methods to be considered in this paper, a necessary and sufficient condition for convergence is that $S_k \rightarrow 0$ ($k \rightarrow \infty$) or,

equivalently, $\mu_k \rightarrow 0$ ($k \rightarrow \infty$). See either [2] or [4] for the proof. If the infinite product of matrices

$$U = \prod_{k=0}^{\infty} U_k = U_0 U_1 U_2 \dots \quad (1.7)$$

converges, then

$$U^*AU = \Lambda, \quad (1.8)$$

and the columns of U are a complete set of normalized eigenvectors of A . We shall investigate the convergence of (1.7) for three convergent Jacobi methods: (a) the classical Jacobi method, (b) the quasicyclic restricted Jacobi methods [4] and (c) the threshold cyclic Jacobi method [7]. Throughout the paper we shall restrict V_k to be of the form

$$V_k = \begin{pmatrix} \cos \theta_k & -e^{i\phi_k} \sin \theta_k \\ e^{-i\phi_k} \sin \theta_k & \cos \theta_k \end{pmatrix} \quad (1.9)$$

By the classical Jacobi method we shall mean a Jacobi method with the following set of rules for determining the U_k of (1.1). Choose π_k such that

$$|a_{ij}^{(k)}| = \mu_k, \quad (1.10)$$

and choose the V_k of (1.9) such that

$$a_{ij}^{(k+1)} = 0. \quad (1.11)$$

The equation (1.11) will always be satisfied if the relations

$$\phi_k = \arg a_{ij}^{(k)}, \quad (1.12)$$

$$\tan 2\theta_k^* = \frac{2 |a_{ij}^{(k)}|}{a_{ii}^{(k)} - a_{jj}^{(k)}} \quad (\theta_k = \theta_k^*) \quad (1.13)$$

hold simultaneously. Since Jacobi dealt only with real symmetric matrices A and orthogonal matrices U_k , the set of rules (1.10) - (1.13) constitute a generalization of his original method.

A Jacobi method is called cyclic if in every segment of $N = n(n-1)/2$ consecutive elements of the sequence $\{\pi_k\}$ every pair (p, q) ($1 \leq p < q \leq n$) occurs exactly once. This is a special case of a quasicyclic Jacobi method with period K where the following condition is imposed. In every segment of $K \geq N$ consecutive elements of the sequence $\{\pi_k\}$ every pair (p, q) ($1 \leq p < q \leq n$) occurs at least once. We next define a restricted Jacobi method. Let θ_k^* be an angle satisfying (1.13) and belonging to the closed interval $[-\pi/4, \pi/4]$, and let $\psi > 0$ be a constant angle not depending on k . Then any method in which ϕ_k is chosen to satisfy (1.12) and θ_k is chosen such that

$$\text{sign } \theta_k = \text{sign } \theta_k^* , \quad |\theta_k| = \text{Min} (|\theta_k^*|, \psi) \quad (1.14)$$

is called a restricted Jacobi method with bound ψ . In [4] it was proved that any quasicyclic, restricted Jacobi method with suitable bound converges. Furthermore it is shown in [4] that, if A has n distinct eigenvalues, any quasicyclic Jacobi method in which θ_k , ϕ_k satisfy (1.12), (1.13) and $\theta_k^* \in [-\pi/4, \pi/4]$ converges quadratically provided the offdiagonal elements are already sufficiently small. More precisely, if $4n\mu_k \leq d$, then*

* Mr. E. R. Hansen has shown (oral communication) that the exponent in the exponential function may be replaced by $\sqrt{2nKd}^{-1}\mu_k$.

$$\mu_{k+K} \leq 2^{\frac{1}{2}} n(K-1)e^{1.334nK^{\frac{3}{2}} d^{-1}\mu_k} d^{-1}\mu_k^2, \quad (1.15)$$

where

$$d = \text{Min}_{r \neq c} |\lambda_r - \lambda_c|.$$

The threshold cyclic Jacobi method is a modification of the cyclic Jacobi method due to Pope and Tompkins [7]. It depends on the choice of a sequence of positive threshold values t_ν ($\nu = 0, 1, 2, \dots$) such that $t_{\nu+1} < t_\nu$ holds for all ν . A fixed cyclic ordering is adopted for the π_k ; ϕ_k always satisfies (1.12), and for each ν , one chooses

$$\theta_k = \begin{cases} \theta_k^* & , \text{ for } |a_{ij}^{(k)}| \geq t_\nu \\ 0 & , \text{ for } |a_{ij}^{(k)}| < t_\nu \end{cases}. \quad (1.16)$$

When all $|a_{rc}^{(k)}| < t_\nu$ ($r \neq c$), t_ν is replaced by $t_{\nu+1}$. Pope and Tompkins prove that the threshold method converges if $\lim_{\nu \rightarrow \infty} t_\nu = 0$. In the present paper we shall take the sequence $\{t_\nu\}$ to be a geometric progression; i.e., we shall assume

$$t_\nu = \alpha q^\nu \quad (0 < q < 1). \quad (1.17)$$

2. NORMS OF MATRICES. In this section we summarize for convenience of reference certain known results concerning norms of matrices. Let $H = (h_{rc})$ and G denote square matrices of order n . The norm (cf. [1]) of a square matrix H is defined to be a non-negative number $N(H)$ satisfying the conditions

$$N(H) > 0 \text{ for } H \neq 0, \quad (2.1)$$

$$N(cH) = |c| \cdot N(H) \text{ for all scalars } c, \quad (2.2)$$

$$N(H + G) \leq N(H) + N(G) , \quad (2.3)$$

$$N(HG) \leq N(H) \cdot N(G) . \quad (2.4)$$

This definition is equivalent to that given by Ostrowski [6] for his "multiplicative norm." By (2.4)

$$N(H) = N(HI) \leq N(H) \cdot N(I) ,$$

hence

$$N(I) \geq 1 \quad (2.5)$$

for any norm. There exist norms for which $N(I) = 1$, an example being the familiar spectral norm

$$N_1(H) = \text{Max}_{x \neq 0} \left(\frac{x^* H x}{x^* x} \right)^{\frac{1}{2}} . \quad (2.6)$$

Note that for any unitary matrix U , $N_1(U) = 1$. The Euclidean norm is defined by the relation

$$N_2(H) = \left(\sum_{r,c=1}^n |h_{rc}|^2 \right)^{\frac{1}{2}} . \quad (2.7)$$

It is well known that, for all H ,

$$n^{-\frac{1}{2}} \leq \frac{N_1(H)}{N_2(H)} \leq 1 . \quad (2.8)$$

Inequalities such as (2.8) hold quite generally. In fact, it was shown in [6] that if $N(H)$ and $N'(H)$ are any two matrix norms, then

$$a(N, N') < \frac{N(H)}{N'(H)} < b(N, N') \quad (2.9)$$

where $a(N, N')$ and $b(N, N')$ are positive constants which do not depend on H . It is easy to show that

$$\text{Max } |h_{rc}| \leq N_1(H), \quad (2.10)$$

and by (2.9) there exists a positive constant $b(N_1, N)$, independent of H , so that

$$|h_{rc}| \leq b(N_1, N) \cdot N(H) \quad (2.11)$$

holds for all r and c and for any norm.

3. PRELIMINARY THEOREM ON INFINITE PRODUCTS. A sequence of $n \times n$

matrices $B_1, B_2, \dots, B_m, \dots$ with $B_k = (b_{rc}^{(k)})$ is said to be

convergent if a limit exists for each element of B_m . That is, if

$\lim_{m \rightarrow \infty} b_{rc}^{(m)} = b_{rc}$, then the matrix $B = (b_{rc})$ is said to be the limit of

the sequence $\{B_k\}$. If the sequence $\{B_k\}$ is not convergent, it is said

to be divergent. An infinite product

$$(I + E_1)(I + E_2) \dots (I + E_m) \dots \quad (3.1)$$

of $n \times n$ matrices is here* called convergent if the partial product

$$P_m = \prod_{k=1}^m (I + E_k) \quad (3.2)$$

converges to a limit. We now prove the following

THEOREM 1. Let $N(H)$ denote some norm of the matrix H . Let the factors of the infinite product (3.1) be $n \times n$ unitary matrices, and let the series $\sum N(E_k)$ converge. Then (3.1) converges to a unitary matrix U . Furthermore, for all m

$$N(U - P_m) \leq c(N) \sum_{k=m+1}^{\infty} N(E_k), \quad (3.3)$$

where $c(N) > 0$ depends only on N and n .

* In defining convergence of infinite products of matrices in general, it is necessary to exclude nonsingular factors in the product and also to distinguish between singular and nonsingular limits. However, since we are dealing only with unitary matrices in the present paper, it is unnecessary to complicate the discussion with such distinctions.

Proof: We first establish that if the sequence $\{P_m\}$ has a limit U , then U must be a unitary matrix. Let $P_m = U + Q_m$. We have that $Q_m \rightarrow 0$ ($m \rightarrow \infty$). Since each factor $I + E_k$ is unitary, P_m is unitary for all m . Thus

$$P_m^* P_m = I = (U + Q_m)^*(U + Q_m),$$

or

$$I - U^*U = U^*Q_m + Q_m^*U + Q_m^*Q_m.$$

The right hand side of the last expression tends to the null matrix as $m \rightarrow \infty$, and this shows that U is unitary. Let $r > m$ and note that

$$P_r - P_m = P_m [(I + E_{m+1}) \dots (I + E_r) - I]. \quad (3.4)$$

We shall prove that P_m has a limit with the help of the following

LEMMA. If H_1, H_2, \dots, H_p and G_1, G_2, \dots, G_p are any two sets of $n \times n$ matrices, then the following identity holds

$$\begin{aligned} H_1 H_2 \dots H_p - G_1 G_2 \dots G_p &= (H_1 - G_1) G_2 \dots G_p \\ &+ H_1 (H_2 - G_2) G_3 \dots G_p + \dots + H_1 \dots H_{p-1} (H_p - G_p). \end{aligned} \quad (3.5)$$

The truth of (3.5) is easily established by induction on p . The details are left to the reader. Letting $H_i = I + E_{m+i}$, $G_i = I$ and combining (3.4) and (3.5) we have

$$\begin{aligned} P_r - P_m &= P_m [E_{m+1} + (I + E_{m+1})E_{m+2} \\ &+ (I + E_{m+1})(I + E_{m+2})E_{m+3} + \dots + (I + E_{m+1}) \dots (I + E_{r-1})E_r]. \end{aligned} \quad (3.6)$$

Products of the matrices $(I + E_k)$ are unitary, hence by (2.4), (3.6)

$$N_1(P_r - P_m) \leq N_1(E_{m+1}) + \dots + N_1(E_r), \quad (3.7)$$

since the spectral norm of a unitary matrix is one. Now for an arbitrary norm N , there exist* $a(N) > 0$ and $b(N) > 0$ such that

$$a(N) \cdot N(H) \leq N_1(H) \leq b(N) \cdot N(H)$$

holds for all H . We have therefore that

$$a(N) \cdot N(P_r - P_m) \leq N_1(P_r - P_m) \quad (3.8)$$

and

$$\sum_{k=m+1}^r N_1(E_k) \leq b(N) \sum_{k=m+1}^r N(E_k) \quad (3.9)$$

Letting $c(N) = b(N)/a(N)$ and combining (3.7) - (3.9), we get

$$N(P_r - P_m) \leq c(N) \sum_{k=m+1}^r N(E_k) \quad (3.10)$$

By hypothesis the series $\sum N(E_k)$ is convergent. Therefore, given any $\epsilon > 0$, there exists an integer $\mu > 0$ such that

$$\sum_{k=\mu+1}^{\infty} N(E_k) < \frac{\epsilon}{b(N) \cdot c(N)} \quad (3.11)$$

We see from (3.10) and (3.11) that

$$N(P_{m+\rho} - P_m) < \frac{\epsilon}{b(N)} \quad (3.12)$$

holds for all $m \geq \mu$ and for $\rho = 1, 2, \dots$; but (3.12) is precisely the condition that

$$|(P_{m+\rho} - P_m)_{rc}| < \epsilon \quad (3.13)$$

holds for all r and c by (2.11). Hence P_m tends to a limit U as $m \rightarrow \infty$. Finally letting $r \rightarrow \infty$ in (3.10) we obtain (3.3) and this proves Theorem 1.

* The spectral norm N_1 is fixed; hence, the $a(N_1, N)$ and $b(N_1, N)$ of (2.9) do not depend on it. Therefore, for simplicity, we have here written $a(N)$ for $a(N_1, N)$ and $b(N)$ for $b(N_1, N)$.

4. CONVERGENCE OF EIGENVECTORS AND ERROR ESTIMATES. We first give a rather general result about convergence and estimation of the error in all the eigenvectors. Following this we give estimates for the three Jacobi methods considered in this paper.

THEOREM 2. Let the $n \times n$ Hermitian matrix A have n distinct eigenvalues, and let the sequence of matrices (1.1) be generated by a Jacobi method such that the angles θ_k satisfy (1.13), every $\theta_k \in [-\pi/4, \pi/4]$, and such that the series

$$\sum_{k=0}^{\infty} S_k^2$$

converges. Then the infinite product (1.7) converges. Furthermore, for all $m > m_0$ we have

$$N_2(U - P_m) \leq \frac{2n^2}{d} \sum_{k=m}^{\infty} S_k^2, \quad (4.1)$$

where

$$d = \min_{r \neq c} |\lambda_r - \lambda_c|, \quad (4.2)$$

and where m_0 is the smallest integer such that for all $m > m_0$

$$S_m^2 < \frac{d}{4}. \quad (4.3)$$

Proof: We shall make use of the fact that, for some ordering of the eigenvalues λ_r of A ,

$$|\lambda_r - a_{rr}^{(k)}| \leq S_k^2 \quad (r = 1, 2, \dots, n) \quad (4.4)$$

hold for all k . The inequalities (4.4) follow from Theorem 1 of [4].

Note that, while (4.4) is valid for any value of k , the ordering of the

eigenvalues λ_r is generally not the same for all k . However, our use of (4.4) does not depend on the ordering of the λ_r .

Let the U_k of (1.1) be expressed in the form

$$U_k = I + E_k, \quad (4.5)$$

and let $E_k = (e_{rc}^{(k)})$. By (1.2), (1.3) and (1.9) we see that the 2×2 principal submatrix of E_k corresponding to V_k is given by

$$\begin{pmatrix} e_{ii}^{(k)} & e_{ij}^{(k)} \\ e_{ji}^{(k)} & e_{jj}^{(k)} \end{pmatrix} = \begin{pmatrix} \cos \theta_k - 1 & -e^{i\phi_k} \sin \theta_k \\ e^{-i\phi_k} \sin \theta_k & \cos \theta_k - 1 \end{pmatrix}, \quad (4.6)$$

and all other elements of E_k are zero. From the definition (2.7) we get

$$N_2(E_k) = [2(\cos \theta_k - 1)^2 + 2 \sin^2 \theta_k]^{1/2},$$

and with the help of some standard trigonometric identities, we obtain

$$N_2(E_k) = 8^{1/2} \left| \sin \frac{\theta_k}{2} \right|. \quad (4.7)$$

By the definitions (1.5) and (1.6)

$$\mu_k \leq \frac{1}{\sqrt{2}} S_k^{1/2} \quad (4.8)$$

Hence, by (4.3) we have

$$\mu_m < \frac{1}{\sqrt{2}} \frac{d}{4} < \frac{d}{4} \quad (m > m_0). \quad (4.9)$$

Using (1.6), (1.3) and (4.2)–(4.4) it follows that

$$\left| \tan 2 \theta_k \right| = \frac{2 |a_{ij}^{(k)}|}{|a_{ii}^{(k)} - a_{jj}^{(k)}|} < \frac{2\mu_k}{\frac{d}{2}} = \frac{\mu_k}{\frac{d}{4}} \quad (4.10)$$

holds for $k > m_0$. Hence, by (4.9), $|\tan 2\theta_k| < 1$ and $|2\theta_k| < \pi/4$ ($k > m_0$), since $\theta_k \in [-\pi/4, \pi/4]$. For these small angles it is easily verified that $4|\tan(\theta/2)| \leq |\tan 2\theta|$. Hence for $k > m_0$, we have from (4.7), (4.8) and (4.10)

$$N_2(E_k) = 8^{1/2} |\sin \frac{\theta_k}{2}| \leq \frac{8^{1/2}}{4} |\tan 2\theta_k| < \frac{8^{1/2}}{d} \mu_k \leq \frac{2}{d} S_k^{1/2}. \quad (4.11)$$

By hypothesis the series $\sum S_k^{1/2}$ converges. Therefore $\sum N_2(E_k)$ converges and by Theorem 1 the infinite product (1.7) converges. Now applying (3.3) to the case where $N = N_2$, we see from (2.8) and (2.9) that

$$c(N_2) = \frac{b(N_1, N_2)}{a(N_1, N_2)} = n^{\frac{1}{2}}. \quad (4.12)$$

Combining (3.3), (4.11) and (4.12) we obtain (4.1) and this completes the proof of Theorem 2.

As before we let N denote the number of elements on one side of the diagonal of an $n \times n$ matrix, viz.

$$N = \frac{1}{2} n(n-1), \quad (4.13)$$

and for convenience we define

$$h = (1 - N^{-1})^{\frac{1}{2}}. \quad (4.14)$$

We let $[x]$ denote the integral part of x , i.e., the largest integer not exceeding x . Finally we let $\sigma = [N/q^2]$ and we state our main result as follows.

THEOREM 3. Let the $n \times n$ Hermitian matrix A have n distinct eigenvalues, and let the sequence of matrices (1.1) be generated by

either Method (a): the classical Jacobi method, Method (b): a quasicyclic restricted Jacobi method, or Method (c): a threshold cyclic Jacobi method in which the thresholds t_v satisfy (1.17). Then the infinite product of matrices (1.7) converges. Furthermore, for Method (a)

$$N_2(U - P_m) \leq \frac{2}{d} (nS_0)^{\frac{1}{2}} \frac{h^{m-1}}{1-h} \quad (m > m_0) \quad , \quad (4.15)$$

for Method (b)

$$N_2(U - P_m) \leq \frac{4Kn^{\frac{3}{2}}}{d} \mu_m \quad (m > \rho) \quad , \quad (4.16)$$

where ρ is given by (4.23) below, and for Method (c)

$$N_2(U - P_m) \leq \frac{2\sigma}{d} q^{[m/\sigma]} (2nN\mu_0)^{\frac{1}{2}} \frac{1}{1-q} \quad (m > m_0) \quad , \quad (4.17)$$

where in each case the integer m_0 has the same meaning as it had in Theorem 2.

Proof: It is easy to show that, for any step in a Jacobi method in which the parameters θ_k, ϕ_k satisfy (1.12) and (1.13), the following equation holds

$$S_{k+1} = S_k - 2|a_{ij}^{(k)}|^2 \quad . \quad (4.18)$$

In the classical Jacobi method we have

$$|a_{ij}^{(k)}|^2 \geq \frac{1}{2N-1} S_k \quad (4.19)$$

always satisfied. Hence

$$S_{k+1} \leq (1 - N^{-1})S_k = h^2 S_k \quad (k = 0, 1, 2 \dots) \quad (4.20)$$

holds for Method (a). It follows therefore that

$$\sum_{k=0}^{\infty} S_k^* = S_0^* (1 + h + h^2 + \dots) = S_0^* \frac{1}{1-h} , \quad (4.21)$$

hence by Theorem 2 the product (1.7) converges. By (4.20)

$$\sum_{k=m}^{\infty} S_k^* = S_0^* \frac{h^{m-1}}{1-h}$$

holds for all m . For $m > m_0$ we may combine the last expression with (4.1) and obtain the inequality (4.15).

Now consider Method (b). We note that $\mu_k \rightarrow 0$ as $k \rightarrow \infty$, since a quasicyclic restricted Jacobi method is convergent. Since A has n distinct eigenvalues, there exists a smallest integer p_1 such that

$$\text{Min} (|\theta_k^*|, \psi) = |\theta_k^*| \quad (k \geq p_1) ,$$

and by (1.14) the restricted method behaves like an ordinary quasicyclic method from the step p_1 on. There also exists a smallest integer

$p_2 \geq p_1$ such that $4n\mu_k \leq d$ ($k \geq p_2$) and hence (1.15) holds for $k \geq p_2$.

Since the sequence $\{S_k^*\}$ is monotone decreasing and $\mu_k \rightarrow 0$, there is a smallest integer p_3 such that (4.3) is satisfied, viz.

$$S_k^* \leq \frac{d}{4} \quad (k \geq p_3) ,$$

and therefore (4.11) holds for $k \geq p_3$. Finally, a smallest integer p_4 exists such that

$$\mu_k < \frac{1}{2g}^{-1} \quad (k \geq p_4) , \quad (4.22)$$

where g is the coefficient of μ_k^2 in (1.15). Now let

$$\rho = \text{Max} (p_2, p_3, p_4) . \quad (4.23)$$

For $k \geq \rho$ we have that (1.15) and (4.11) hold, hence

$$\begin{aligned}
\sum_{k=p}^{\infty} N_2(E_k) &< \frac{g^{1/2}}{d} \sum_{k=p}^{\infty} \mu_k, \\
&= \frac{g^{1/2}}{d} \sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \mu_{p+1+jK}, \\
&\leq \frac{g^{1/2}}{d} \{ (\mu_p + g\mu_p^2 + g^3\mu_p^4 + \dots) \\
&\quad + (\mu_{p+1} + g\mu_{p+1}^2 + g^3\mu_{p+1}^4 + \dots) + \dots \\
&\quad + (\mu_{p+K-1} + g\mu_{p+K-1}^2 + g^3\mu_{p+K-1}^4 + \dots) \}, \\
&\leq \frac{g^{1/2}}{dg} \left(\frac{g\mu_p}{1-g\mu_p} + \frac{g\mu_{p+1}}{1-g\mu_{p+1}} + \dots + \frac{g\mu_{p+K-1}}{1-g\mu_{p+K-1}} \right),
\end{aligned}$$

and using (4.22) we get

$$\sum_{k=p}^{\infty} N_2(E_k) \leq \frac{4\sqrt{2}}{d} (\mu_p + \mu_{p+1} + \dots + \mu_{p+K-1}). \quad (4.24)$$

Since S_k is monotone decreasing,

$$\mu_{p+v} \leq \frac{1}{\sqrt{2}} S_{p+v} \leq \frac{1}{\sqrt{2}} S_p \leq \sqrt{n} \mu_p \leq \frac{n}{\sqrt{2}} \mu_p \quad (v = 0, 1, 2, \dots),$$

we combine this with (4.24) to obtain

$$\sum_{k=p}^{\infty} N_2(E_k) \leq \frac{4nK}{d} \mu_p. \quad (4.25)$$

Now combining (3.3), (4.12) and (4.25) we get (4.16), and since

$\mu_m \rightarrow 0$ ($m \rightarrow \infty$) this proves the convergence of (1.7) for Method (b).

Turning our attention to Method (c), we see by (1.16) that $U_k = I$, $E_k = 0$ whenever $|a_{1j}^{(k)}| < t_v$ so that a step in the process (1.1) counts only when a non-trivial transformation is made, i.e., whenever

$|a_{ij}^{(k)}| \geq t_v$. We let the scalar α in (1.17) be equal to μ_0 so that

$$|a_{rc}| \leq \mu_0 = \mu_0 q^0 = t_0 \quad (r \neq c).$$

There may be no non-trivial transformations made at any given threshold level, but if there are, we can easily bound their number. Assume that after the r -th transformation $\mu_r < t_{v-1}$. Then

$$S_r < 2Nt_{v-1}^2, \quad (4.26)$$

and by (4.18)

$$S_{k+1} \leq S_k - 2t_v^2 \quad (4.27)$$

holds for each transformation made at the v -th threshold level, since

$|a_{ij}^{(k)}| \geq t_v$. Combining (4.26) and (4.27) we see that

$$S_{r+p} \leq 2Nt_{v-1}^2 - 2pt_v^2$$

holds after p transformations at the v -th level. Letting $[x]$ denote the integral part of x , we have that there are at most

$$\sigma_v = \left[\frac{Nt_{v-1}^2}{t_v^2} \right]$$

steps at the v -th level, for

$$\begin{aligned} S_{r+\sigma_v} &< 2Nt_{v-1}^2 - 2 \left[\frac{Nt_{v-1}^2}{t_v^2} \right] t_v^2, \\ &< 2Nt_{v-1}^2 - 2 \left(\frac{Nt_{v-1}^2}{t_v^2} - 1 \right) t_v^2, \\ &< 2t_v^2, \end{aligned}$$

so that $\mu_{r+\sigma_v} < t_v$, i.e., the next threshold level has been reached.

From (1.17)

$$\sigma_v = \left[\frac{N}{q^2} \right] = \sigma \quad (v = 1, 2, \dots) .$$

We have proved that

$$S_{v\sigma} \leq 2Nt_v^2 \quad (v = 1, 2, \dots) ,$$

so that, using (1.17)

$$S_{v\sigma}^{1/2} \leq (2N\mu_0)^{1/2} q^v . \quad (4.28)$$

Now, since $s[t/s] \leq t$ for all positive integers s and t , we may write

$$\sum_{k=m}^{\infty} S_k^{1/2} \leq \sum_{k=m}^{\infty} S_{\sigma}^{1/2} \left[\frac{k}{\sigma} \right] , \quad (4.29)$$

$$\leq \sum_{k=\sigma}^{\infty} S_{\sigma}^{1/2} \left[\frac{k}{\sigma} \right] .$$

Let i be one of the integers $0, 1, \dots, \sigma-1$ and let

$$k = \sigma \left[\frac{m}{\sigma} \right] + i + j\sigma .$$

We have

$$\left[\frac{k}{\sigma} \right] = \left[j + \left[\frac{m}{\sigma} \right] + \frac{i}{\sigma} \right] ,$$

and since $i < \sigma$, we have

$$\left[\frac{k}{\sigma} \right] = j + \left[\frac{m}{\sigma} \right] \quad (i = 0, 1, \dots, \sigma-1) .$$

Thus we may write (4.29) as follows

$$\begin{aligned} \sum_{k=m}^{\infty} S_k^{1/2} &\leq \sum_{i=0}^{\sigma-1} \sum_{j=0}^{\infty} S_{\sigma}^{1/2} \left(j + \left[\frac{m}{\sigma} \right] \right) , \\ &= \sigma \sum_{j=0}^{\infty} S_{\sigma}^{1/2} \left(j + \left[\frac{m}{\sigma} \right] \right) , \end{aligned}$$

and using (4.28) we have

$$\begin{aligned} \sum_{k=m}^{\infty} S_k^{1/2} &\leq \sigma(2N\mu_0)^{1/2} \sum_{j=0}^{\infty} q^{j + \left[\frac{m}{\sigma} \right]}, \\ &= \sigma q^{\left[\frac{m}{\sigma} \right]} (2N\mu_0)^{1/2} \frac{1}{1-q}. \end{aligned} \quad (4.30)$$

Finally, combining (4.1) with (4.30) we get immediately the expression (4.17). Obviously,

$$q^{\left[\frac{m}{\sigma} \right]} \rightarrow 0 \quad (m \rightarrow \infty),$$

and thus the product (1.7) converges also for Method (c). This completes the proof of Theorem 3.

5. REMARKS. If we refer to unitary matrices of the general type used in the process (1.1) as elementary unitary matrices, then it is a rather easy matter to show that an arbitrary unitary matrix can be decomposed into a product of a finite number of elementary unitary matrices. In [3] (footnote (5), page 8) Givens gave a proof that every orthogonal matrix with positive determinant can be expressed as a product of a finite number of elementary orthogonal matrices. We shall not give the details here, but a proof of the corresponding result on the decomposition of an arbitrary unitary matrix can be constructed by a method very much the same as that used in [3].

In spite of the result just mentioned, it is perhaps not without some interest that we can use Theorem 3 to prove that an arbitrary unitary matrix can be expressed as an infinite product of elementary

unitary matrices. The decomposition is, of course, not unique. We sketch a proof as follows.

Suppose W is any $n \times n$ unitary matrix. Let D be a fixed $n \times n$ diagonal matrix with n distinct and real diagonal elements. Let the Hermitian matrix

$$A = WDW^* \quad (5.1)$$

be diagonalized by (say) Method (a) of Theorem 3, so that

$$U^*AU = D',$$

where U is the limit of the infinite product (1.7). The matrix D' has the same diagonal elements as D but possibly with a different arrangement. By a finite number of elementary unitary transformations with V_k of the form

$$\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix},$$

we can transform D' into D , so that

$$V^*AV = D, \quad (5.2)$$

where V is an infinite product of elementary unitary matrices. From (5.1) and (5.2) we see that

$$VDV^* = WDW^*. \quad (5.3)$$

It is easy to show from (5.3) that W can be expressed as VY where Y is a diagonal unitary matrix, i.e.

$$y_{rr} = e^{i\alpha_r} \quad (r = 1, 2, \dots, n). \quad (5.4)$$

Since any diagonal matrix Y satisfying (5.4) can be decomposed as a

product of at most $n-1$ elementary unitary matrices with V_k of the form

$$\begin{pmatrix} e^{i\beta} & 0 \\ 0 & e^{i\beta} \end{pmatrix},$$

the desired result is established.

Section 6
REFERENCES

1. V. N. Faddeeva, Vyčislitel'nye Metody Lineinoi Algebry, Gostekhizdat, Moscow-Leningrad, 1950. (Translated by Curtis D. Benster under title Computational Methods of Linear Algebra, Dover, New York, 1959)
2. G. E. Forsythe and P. Henrici, "The Cyclic Jacobi Method for Computing the Principal Values of a Complex Matrix," Trans. Amer. Math. Soc., (to be published)
3. W. Givens, Numerical Computation of the Characteristic Values of a Real Symmetric Matrix, Oak Ridge National Laboratory, Rep. ORNL 1574, 1954. (Purchasable from Office of Technical Services, U. S. Dept. of Commerce, Washington 25, D. C.)
4. P. Henrici, "On the Speed of Convergence of Cyclic and Quasicyclic Jacobi Methods for Computing Eigenvalues of Hermitian Matrices," J. Soc. Indust. Appl. Math., Vol. 6, 1958, pp. 144-162
5. C. G. J. Jacobi, "Über ein leichtes Verfahren, die in der Theorie der Säkularstörungen vorkommenden Gleichungen numerisch aufzulösen," J. Reine Angew. Math., Vol. 30, 1846, pp. 51-95
6. A. Ostrowski, "Über Normen von Matrizen," Math. Z., Vol. 63, 1955, pp. 2-18
7. D. A. Pope and C. B. Tompkins, "Maximizing Functions of Rotations-- Experiments Concerning Speed of Diagonalization of Symmetric Matrices Using Jacobi's Method," J. Assoc. Comput. Mach., Vol. 4, 1957, pp. 459-466

WORD-CORRELATION STUDY

R.P. Mitchell



3

FOREWORD

This document, originally published as LMSD 311696, September 1959, was the final report for Word Correlation Study, Contract No. AF 30(602)-1889, Rome Air Development Center.

The author would like to express his appreciation to members of Applied Mathematics for their suggestions and encouragement in this research.

ABSTRACT

A method of syntactic analysis of declarative English sentences, proposed by Bar-Hillel and Lambek, was operationally extended in such a manner that the product of the analysis is a machine retrieval language. Syntactic analysis thus becomes a recursive translation from the natural language to the proposed machine language. This method is called "microsyntactic analysis," and the smallest meaningful linguistic units, or morphemes, of the retrieval language are termed "microsyntactic indices."

A three-value storage logic was devised for storage and retrieval operations. Various relations between syntactic strings and algebraic strings were defined. A characterization of semantic correspondences in terms of algebraic relations was proposed and is currently under study.

CONTENTS

<u>Section</u>		<u>Page</u>
	Foreword	iii
	Abstract	v
1	Introduction	1
2	A Storage-Dependent Algebra of Indices	6
4	Syntactic Theory of Definitions	11
5	Technical Summary and Conclusions	15
6	References	18
 <u>Appendix</u>		
A	Word Correlation Study Program	19

Section 1

INTRODUCTION

The first quarter of research on this contract resulted in a proposal for studying word correlation in a sublanguage of English and presented a method by which correlations could be computed for the given sublanguage. The method of syntactic analysis conceived first by Y. Bar-Hillel [5] and studied by J. Lambek [6] was arbitrarily extended to a larger part of the language than intended by these authors. This extension, termed by us "macrosyntax" because it provides a basis for classification of lexical units into function-sets, may not be justified on logical grounds; it was taken only as a working basis for further research.

The second and third quarters of research resulted in a refinement of syntactic analysis called "microsyntax." The chief advance consisted in the abstraction of syntactic forms from semantic and/or functional interpretations given them in the macrosyntax, and opening the way to mathematical characterizations of syntactic forms and their construction. Microsyntactic indexing was also introduced, providing a means for very fine semantic distinctions between words whenever such interpretations are desired by the analyst. Concurrently with these studies, we made serious efforts to find number-theoretic analogues of syntactic structure in the hope that syntactic analysis could be accomplished on standard digital computers. In this we were only partially successful, and these efforts had been abandoned at the close of the third quarter.

A list of basic English words was abstracted from West [7], their syntactic forms listed, and frequencies of occurrence per million

computed from data given by West [7]. This list was computer-programmed for future work; the program is described in the appendix of this report.

The fourth quarter of research has resulted in the conception and outline of a storage algebra of microsyntactic indices. The most important result of this work is the characterization of semantic correspondences in terms of relations defined in this algebra. We feel that this characterization gives us a firm basis for the derivation of a correlation function.

Section 2

A GENERALIZED SYNTAX FOR L

In [3] we presented a method for indexing strings of words which are analyzable in terms of a microsyntax of English. This method has been generalized and is presented here in a more systematic form. In the reference cited we defined heuristically the order of a syntactic form as the highest number of strokes separating elements of the form. It will be observed that forms of order greater than zero are constructed by repeated application of left or right stroke functions. Thus, we need only two mappings and the concept of order of syntactic forms to characterize all possible constructions. In [2], when we defined the microsyntax of L, we defined two sets of elements, $A = \{n_1, s_1, e_1\}$ and $B = \{n_2, s_2, e_2\}$, and constructed all elements of the microsyntax from these sets by means of appropriate mappings. The indicial feature of this method was extended in [3], and the method of microsyntactic indexing was introduced. In the process of systematizing the method of syntactic analysis originally proposed by Bar-Hillel in [4], we have abstracted the constructive properties of the method and discarded the interpretations of strings of these elements as "noun phrases," "verb phrases," "sentences," etc. These interpretations may serve a convenient classification purpose in descriptive semantics, but they are meaningless syntactically. In microsyntactic indexing, n and s serve an indicial function only, and may be replaced by any convenient symbol. Let Z be the set of zero-order indices. Indices of higher order may, in general, be either entirely associative or non-associative between certain elements. We denote by X_n the set constructed by n consecutive applications of the left stroke function on elements of Z ; thus, X_1 is the set of elements of the form $a \backslash b$, $X_2 = \{a \backslash b \backslash c\}$, etc. Similarly, we denote by Y_m the

set of elements formed by m successive applications of the right stroke function. All members of X_n and Y_m are first-order and entirely associative. Forms containing elements which are non-associative relative to a left or a right stroke are constructed from the sets X_n and Y_m . These forms are the second-order forms, and are in turn used in the construction of third and higher order forms. We denote the many possible sets of second-order forms by $X_{i,j,k,\dots}^{m,n,p,\dots}$ and $Y_{i,j,k,\dots}^{m,n,p,\dots}$ where X or Y indicates that the left or right stroke, respectively, is the stroke function relative to which the elements are non-associative; and the elements of the superscript sequence indicate the number of right or left strokes, respectively, in the forms following the forms indicated in the paired subscript sequence. For example, $X_{1,0}^2 = \{a \backslash b \backslash \backslash c / d / e \backslash \backslash f\}$, $Y^{1,1} = \{a \backslash b / / c \backslash d\}$, $Y_1^{1,1} = \{a / b // c \backslash d // e \backslash f\}$, $X_{,1}^{1,0} = \{a / b \backslash \backslash c / d \backslash \backslash e\}$, etc. Most of the possible second-order forms do not occur in practice; among these are the forms which are non-associative relative to both stroke functions. Where such forms should be necessary, and for forms of higher order, it is notationally simpler to write the required sets in terms of the left and right strokes, as was done in the above examples. A slightly more cumbersome notation is to define the stroke functions as functions of ordered pairs of sets (as we have implicitly done above), then write the actual function as required. For example, an element belonging to $Y^{1,1}$ would be written $\rho[\lambda(a,b), \lambda(c,d)]$, where λ and ρ are the left and right stroke functions respectively, and a, b, c, d are all elements of Z . Some convention must be used to characterize association between elements; capital Greek letters, Λ and P , may be used for this purpose.

We require an identity element $e \in Z$, such that

$$\lambda(e, a) = \rho(a, e) = a$$

where a is an element of any set. Sequences of indices are said to "reduce" to an index when any of the following sequences occur:

$$\lambda(a,b) \lambda(b,c) = \lambda(a,c)$$

$$\Lambda(a,b) \Lambda(b,c) = \Lambda(a,c)$$

$$\rho(a,b) \rho(b,c) = \rho(a,c)$$

$$P(a,b) P(b,c) = P(a,c)$$

where a, b, c are elements of any set. We define a "syntactic string" as any order sequence of indices. A string is said to reduce to an index when sequential pairs of indices reduce according to the transitive property of the stroke functions given above. Conversely, an expansion of any index is any string obtained by the converse application of the transitive property, providing the indices used in the expansion are members of sets required for analysis of the language. In particular, if q is any index, we denote by $\rightarrow q$ the set of allowable expansions of q .

Section 3

A STORAGE-DEPENDENT ALGEBRA OF INDICES

The central idea in developing a storage-dependent algebra of indices is the assignment of values to indices and algebraic strings of indices according to the results of a storage-verification procedure. We call the values thus assigned "s-values" for the indices, and denote them by "0," "1," and " δ ." It is assumed that the storage of indices is a continually changing storage, but that storage-verification occurs for any fixed interval while the storage is static. The s-values 0, 1 are interpreted to mean negative and positive results of the verification procedures, while δ is interpreted as the (intentional or unintentional) non-application of the verification-procedure.

The stipulation of s-values for algebraic strings of indices is to some extent arbitrary, depending upon the purpose of the algebraic analysis. Once this purpose has been clarified, then the algebraic relations among indices can be defined in terms of assignments of s-values for the primitive relations. For the heuristic reason that, for a sufficiently large collection of sentences in storage, the negations of a large number of words and phrases are likely also to be in storage, we define, for each index p in L , an index $-p$ in L such that the following conditions of storage verification hold; if $p = 1$, $-p = 1$; if $p = 0$, $-p = \delta$; and if $p = \delta$, $-p = 0$. The element $-p$ is the antonym of p in a descriptive semantics.

We define a binary relation D between indices of L in terms of the following stipulation of s-values.

p	q	pDq
0	0	0
0	1	1
0	δ	δ
1	0	1
1	1	1
1	δ	1
δ	0	δ
δ	1	1
δ	δ	δ

We also define a binary relation C among indices of L such that C satisfies the identity $pCq = -(-pD-q)$. From the definitions it follows that pCq satisfies the following table of s-values.

p	q	pCq
0	0	0
0	1	1
0	δ	0
1	0	1
1	1	1
1	δ	1
δ	0	0
δ	1	1
δ	δ	δ

Examination of the s-table for pDq shows that

$$pDq = qDp . \tag{1}$$

Further, from the definition for -p , we have

$$-(-p) = p . \tag{2}$$

Using the definition for C and (1), we have

$$pCq = qCp . \quad (3)$$

Through the s-tables, it can be shown that

$$pD(qDr) = (pDq)Dr , \quad (4)$$

and it is easily proved that

$$pC(qCr) = (pCq)Cr \quad (5)$$

Again through the s-tables, the following property may be verified:

$$pC(qDr) = (pCq)D(pCr) \quad (6)$$

We have shown that, for the given stipulation of s-values the relations C and D are commutative, associative, and distributive with respect to C .

Algebraic strings are linearly related to the stroke functions λ and ρ . Let x, y, z be any indices such that, if α denotes either λ or ρ , $\alpha(x,y), \alpha(x,z), \alpha(y,z)$ are all indices in L . Then we define:

$$\alpha(xCy, z) = \alpha(x, z)\alpha(y, z) \quad (7)$$

$$\alpha(x, yCz) = \alpha(x, y)\alpha(x, z) \quad (8)$$

$$\alpha(xDy, z) = \alpha(x, z)D\alpha(y, z) \quad (9)$$

$$\alpha(x, yDz) = \alpha(x, y)D\alpha(x, z) \quad (10)$$

Since C and D are commutative, we have

$$\alpha(xCy, z) = \alpha(yCx, z) \quad (11)$$

$$\alpha(x, yCz) = \alpha(x, zCy) \quad (12)$$

$$\alpha(xDy, z) = \alpha(yDx, z) \quad (13)$$

$$\alpha(x, yDz) = \alpha(x, zDy) \quad (14)$$

We define

$$\alpha(-x, z) = \alpha(x, -z) = -\alpha(x, z) \quad (15)$$

and show that this definition is consistent with (7) - (10).

$$\begin{aligned}
\alpha(xCy, z) &= \alpha[-(-xD-y), z] \\
&= -\alpha(-xD-y, z) = -[\alpha(-x, z)D\alpha(-y, z)] \\
&= -[-\alpha(x, z)D-\alpha(y, z)] = \alpha(x, z)C\alpha(y, z),
\end{aligned}$$

and similarly for $\alpha(x, yCz)$. From (15), we have the identities:

$$\alpha(xCy, z) = \alpha(-xD-y, -z) \quad (16)$$

$$\alpha(xDy, z) = \alpha(-xC-y, -z) \quad (17)$$

$$\alpha(x, yCz) = \alpha(-x, -yD-z) \quad (18)$$

$$\alpha(x, yDz) = \alpha(-x, -yC-z) \quad (19)$$

Further, from (7) - (10) and associativity, we have

$$\begin{aligned}
\alpha(uCv, xCy) &= \\
&\alpha(u, x)C\alpha(u, y)C\alpha(v, x)C\alpha(v, y) \quad (20)
\end{aligned}$$

$$\alpha(uDv, xDy) = \alpha(u, x)D\alpha(u, y)D\alpha(v, x)D\alpha(v, y) \quad (21)$$

Indices containing both C and D are allowable if and only if all indices in the resulting decomposition have the same s-value. To prove this take, for example, the index $\alpha(uCv, xDy)$. We have

$$\begin{aligned}
\alpha(uCv, xDy) &= \alpha[-(uCv), -xC-y] = -\alpha(uCv, -xC-y) = \\
&-[-\alpha(u, x)C-\alpha(u, y)C-\alpha(v, x)C-\alpha(v, y)]
\end{aligned}$$

On the other hand, we also have

$$\begin{aligned}
\alpha(uCv, xDy) &= \alpha[-uD-v, -(xDy)] = -\alpha(-uD-v, xDy) = \\
&-[-\alpha(u, x)D-\alpha(u, y)D-\alpha(v, x)D-\alpha(v, y)]
\end{aligned}$$

The two decompositions can be equal only if all indices in the strings are either 0, δ , or 1.

From the transitive property of the stroke functions, we have

$$\alpha(xCy, z)\alpha(z, u) = \alpha(xCy, u) \quad (22)$$

Since $\alpha(xCy, z) = \alpha(x, z)C\alpha(y, z)$, the sequence of indices must be distributive with respect to C and D; thus

$$\begin{aligned} \alpha(xCy, z)\alpha(z, u) &= [\alpha(x, z)C\alpha(y, z)]\alpha(z, u) = \\ &\alpha(x, z)\alpha(z, u)C\alpha(y, z)\alpha(z, u) = \alpha(x, u)C\alpha(y, u) = \\ &\alpha(xCy, u) \end{aligned}$$

However, $\alpha(u, xCy)\alpha(xCy, v) = \alpha(u, v)$ if and only if $x = y$, and similarly if C is replaced by D in one or both indices.

Using (2) and (15), we have

$$\begin{aligned} \alpha(-x, y) \alpha(y, -z) &= \alpha(x, -y) \alpha(-y, z) = \\ [-\alpha(x, y)] [-\alpha(y, z)] &= \alpha(x, z) \end{aligned} \tag{23}$$

So far, we have not used any property of δ in our results, and hence the definitions and consequent properties of both algebraic and syntactic strings hold for a Boolean algebra as well. An interesting application of these results would be the analysis of the predicate calculus in terms of this theory; however, we have only fragmentary results on this application, and shall not include them here.

Section 4

SYNTACTIC THEORY OF DEFINITIONS

It was proposed, in the first report on this study [1], to compile a list of words to be taken as undefined and define all other words in the general lexicon in terms of the primitive words. The novel feature of our proposal was to define words in such a manner that the string of the definition is an allowable expansion of the syntactic form of the word defined. This feature would permit free substitution of definitions in the context of the words defined, and hence would permit correlation of sentences. The theory of definition is the heart of descriptive semantics, and it is in this area of research that any method or proposal for word correlation must be conceived. As was indicated in [4], we have encountered considerable difficulty in writing satisfactory definitions within the proposed syntactic restrictions. A "satisfactory" definition is simply a matter of the conscience and experience of the person writing the definition, and part of the difficulty lies herein--we have no criteria for an acceptable definition. The other part of the difficulty is that the theory of definition at this level is not properly a scientific problem, even for those who choose to recognize it as a problem. It is a philosophical problem in the present state of our knowledge of human behaviour and communication, and subject to endless debate.

The purpose of a definition in a natural language is to establish a correspondence between a word and whatever the word names or whatever function it serves in the language. The manner in which this purpose is accomplished is by stating a correspondence between the word and a collection of other words in the language, placing reliance on the human interpreter to perform the actual process of naming objects of experience,

and whatever other functions are required. Since the construction of a dictionary or thesaurus involves these complex human activities to a sophisticated degree, and since in our present state of knowledge these activities are so little understood, the proper objects of our research are simply those correspondences between words and collections of words which are already established in conventional dictionaries, rather than the construction of a new dictionary. We have limitations imposed by the method of syntactic analysis, of course, and these restrictions will induce changes in the precise wording of definitions. But we can study the problems in abstracto, and the results obtained ought to tell us whether any useful concept of word correlation is feasible.

Generally, a definition is a correspondence between a syntactic string and a sequence of syntactic strings. We consider here the correspondence between a single word and a single string, and write $p \rightarrow \sigma$ for the correspondence, where p is the word defined and σ the string constituting the definition. p denotes a set of indices each of which stands in the given relation to the string σ . Since we have no criterion for substitutability of σ in any string containing an index belonging to p , the relation \rightarrow is not an equivalence relation. We define the equivalence relation \equiv between indices as follows:

$$p \equiv \sigma \text{ if and only if } p \rightarrow \sigma \text{ and } \sigma \rightarrow p,$$

the relation holding for all indices belonging to p . We define the relation \rightarrow in terms of the storage-dependent algebra:

$$p \rightarrow \sigma = - pD\sigma \tag{24}$$

The relation \rightarrow may hold for any two indices in L . In the thesaurus-type of definition, where a word is "defined" in terms of a collection of words, the relation holds between the word defined and each word in the collection. From the definitions, we have the following table of s -values for \rightarrow :

p	q	$p \rightarrow q$
0	0	δ
0	1	1
0	δ	δ
1	0	1
1	1	1
1	δ	1
δ	0	0
δ	1	1
δ	δ	δ

The first set of values reflects the fact that the semantic correspondence between any two indices does not depend upon the storage of the indices. Further, if $q = \sigma$, a string, then the correspondence may implicitly be in storage through the words in δ , though σ itself need not be in storage.

The semantic correspondence relation has the following elementary properties:

$$p \rightarrow q = - (pC-q) \quad (25)$$

$$\alpha(x, y \rightarrow z) = \alpha(x, y) \rightarrow \alpha(x, z) \quad (26)$$

$$\alpha(x \rightarrow y, z) = \alpha(x, z) \rightarrow \alpha(y, z) \quad (27)$$

The form $\alpha(x \rightarrow y, u \rightarrow v)$, if it should occur, implies $x = y$ and $u = v$; this is easily seen by application of (26) and (27).

The purpose of the three-valued algebra we have outlined is to provide a basis for semantic correlation relative to a given set of correspondences between words and strings of words. It is our intention to define the s-value δ in such a manner that it provides a numerical measure of correlation relative to the given correspondences. This function will also provide a basis for semantic inference in L relative

to that sub-language of L which is in storage. Our speculations regarding the form of this numerical function center around statistical properties of indices subject to the syntactic restraints of acceptable strings and to the imposed structure of the storage algebra. Results obtained at this time are too indefinite for inclusion in this report.

Section 5

TECHNICAL SUMMARY AND CONCLUSIONS

We present in this section a summary of the more important technical results obtained for this contract.

The subject of analysis is the English language which occurs in meaningful discourse; in particular, the analysis of meaningful declarative sentences. The syntax of this language was analyzed by the method proposed by Bar-Hillel [5] and Lambek [6]. The result was a sublanguage of the original language consisting of sentences for which the method applies and the words which occur in those sentences. This method of analysis was extended operationally to include a larger set of sentences for experimental purposes. This extension is termed "macrosyntax" in our reports.

One of the weaknesses of the macrosyntax is the ambiguity of association in the more complex forms. For example, the form $n \setminus s / n$ can be taken as either of the forms $(n \setminus s) / n$ or $n \setminus (s / n)$. A resolution of this ambiguity was found by defining two semigroups $(A, /)$ and (B, \setminus) , with $A = \{n_1, s_1, e_1\}$ and $B = \{n_2, s_2, e_2\}$, then defining mappings which carry cross-products of the two sets into sets whose members are the non-associative forms. This syntax was termed "microsyntax."

It was then recognized that the indicial feature of the microsyntax could be used to index words in context. Thus, $n_1, n_2, n_1 / n_1, n_1 \setminus \setminus s_2, n_2 \setminus s_2, s_1, s_2$, etc. are all sets whose elements are words or strings of words in the natural language used in particular sentences with meanings intended in those sentences. This method of indexing was called "microsyntactic

indexing," and represents a recursive translation of the natural language into a machine language suitable for storage and retrieval purposes. The analysis of the syntax of the natural language, with this recursive translation into the microsyntactic indices, was felt to be an essential part of the groundwork for word correlation. This is because the meanings of words depend upon their usage and functions in context, and a method for preserving context while simultaneously permitting algebraic operations with individual words was essential. We feel that microsyntactic indexing satisfies this technical requirement.

An algebra suitable for operations among indices in storage is outlined in this report. A three-value algebra was chosen so that one of the values could be used as a correlation function between any pair of indices. Relations between the stroke functions of the microsyntax and the algebraic relations were defined and simple consequences of the definitions were derived. A characterization of the correspondence between a word and a string which defines the word was proposed and defined in terms of the algebra. This is felt to be the most important technical result of our research, as it is through the analysis of these correspondences that the theory of word correlation must be constructed.

To each word in the natural language there corresponds a set of microsyntactic indices, at most one index for each occurrence of the word in the material analyzed. In addition, each sentence which is analyzable in the microsyntax is indexed. The storage can be manipulated at will using a Boolean algebra; for example, various hypotheses could be constructed and consequences of the hypotheses derived automatically by machine. However, the technique of using a Boolean algebra presupposes a precise knowledge of the contents of the storage, and when storage is large, the Boolean algebra becomes impractical. An automatic program for word correlation is essential as a guide for the analyst. If a correlation function were

available its application would tell the analyst in which direction to proceed in storage manipulation without reading out the entire storage. In other words, the correlation function serves as a tool for plausible inference; from the value of the correlation function in each particular instance the analyst is able to infer the contents of storage and construct the Boolean search function accordingly. These are the reasons we have chosen a three-value storage algebra.

It is unclear at present what form the correlation function must have to be useful in the intended way. It must be independent of the cardinality of index sets corresponding to words, yet it seems reasonable that the function be related to the probability that a given word is in storage. On the other hand, since we are seeking a semantic correlation function, we have no a priori reason to assume that probability enter the picture at all. A scheme for computing word correlation similar to the simple ideas originally proposed in [1] may serve the intended purpose, with suitable criteria for substitution.

Section 6

REFERENCES

1. Lockheed Missiles and Space Division, Word Correlation Study, I, by R. P. Mitchell and E. Greer, Contract No. AF 30(602)-1889, Sunnyvale, Calif.
2. Lockheed Missiles and Space Division, Word Correlation Study, II, by R. P. Mitchell, Contract No. AF 30(602)-1889, Sunnyvale, Calif.
3. Lockheed Missiles and Space Division, Word Correlation Study, III, by R. P. Mitchell, Contract No. AF 30(602)-1889, Sunnyvale, Calif.
4. Lockheed Missiles and Space Division, Proposal for Continuation of Word Correlation Study, Contract No. AF 30(602)-1889, LMSD-49901, Sunnyvale, Calif., 27 Jul 1959.
5. Bar-Hillel, Y. A., "Quasi-Arithmetical Notation for Syntactic Description," Language, pp. 47-58, 1953.
6. Lambek, J., "The Mathematics of Sentence Structure," American Mathematical Monthly, pp. 154-170, Mar 1958.
7. West, M. P., A General Service List of English Words, New York, Longmans, Green and Co., 1957.

Appendix A
WORD CORRELATION STUDY PROGRAM

General Description:

Given a list of words with their frequencies of occurrence in five million words and the percentages of occurrence as different syntactic types, this program will produce a readable list consisting of the words, their frequencies of occurrence in one million words, and the frequencies of occurrence as the different syntactic types. Where the percentage is unknown, the program will provide the number -- llll llll - ll as an indicative symbol that the word is known to occur as a certain syntactic type but the programmer must put in the symbol in the desired location in the output buffer. See Note 2.

The program requires all of first core and up to location 10240 of second core. The remainder of second core may be used to enlarge the program or make comparisons between the words and/or numbers associated with them.

Output Format:

Each 'page' is headed by the page number and five words. The syntactic types are represented by numbers from 30 to 60, and appear to the left. Each word heads a column of numbers, the first of which is the frequency per million words, and the others, opposite the syntactic types, represent the occurrence of the word as the various types.

Input:

Input to the program is by tape. For every page, the five words to be printed on that page must be punched on one card in the following format:

column 1	begin first word
column 19	second word
column 37	third word
column 55	fourth word
column 73	fifth word

(Note: If the fifth word has more than eight letters, start in the number of column to the left of column 73 that equals the number of letters greater than eight.) e.g., a nine-letter word must start in column 72; a twelve-letter word should start in column 69. If possible, words with more than twelve letters should not be the fifth word on a page.

The numbers associated with every set of five words are punched on cards separate from the word card. Using the standard VARAB or VARCAR format (the same), punch the five frequencies, followed by the five percentages for the first syntactic type, next the five percentages for the second syntactic type, and so on. Zeros must be punched when a word doesn't occur as a particular type. There should be five cards of numbers of output desired, for each one card of words for each page. This is for 15 syntactic types. If more types are given, the argument word in location 00131)B must be changed from 03 00005 00010 to 03 000nn 00010, n = number if octal, of cards to be read for each set of five words.

Operation:

Card-to-tape program deck I; card-to-tape 'data deck II' (the word cards); card-to-tape 'data deck III' (the number cards). Place II on uniservo No. 4, III on uniservo No. 3; load octal program tape number. Output is on uniservo No. 5. The program starts at cell 00250)B and should stop at 05043)B. The original program provides for 100 pages of output, or 500 words. For fewer pages an MJ to 05043B) should be provided at the end of

the number of pages required. This may be accomplished by an octal corrector with the program deck, in which case the program must be card-to-taped; or by a four-field octal loader corrector, which must be loaded after the program tape (before program is started).

Notes:

1. The output buffer consists of an 84 word BUFF and 324 word HOLD. The syntactic type symbols go to HOLD + 60, HOLD + 66, HOLD + 72, HOLD + 78, HOLD + 84; HOLD + 96, HOLD + 102, HOLD + 108, HOLD + 114, HOLD + 120, HOLD + 132, etc. The numbers associated with these types are in locations called HOLD + 61 through HOLD + 65; HOLD + 67 through HOLD + 71; and so on.

2. If the indicative number - 1111 1111 - 11 is desired, for example, opposite type 31, under the second word on a page, then TP STAR HOLD + 68 must be inserted by programmer in 20 word buffer provided for that purpose for each page before PREDIT is called to output the page.

**EFFICIENT ESTIMATION OF
REGRESSION PARAMETERS FOR CERTAIN
SECOND-ORDER STATIONARY PROCESSES**

C.T. Striebel



4

FOREWORD

This reprint was published originally
as a technical report, LMSD-288079,
under the same title, November 1959.

ABSTRACT

The problem considered is that of estimating the single regression parameter k in the presence of a second-order stationary disturbing process $X(t)$ with rational spectral density. If the process

$$y(t) = k\phi(t) + X(t)$$

is observed in the continuous interval $0 \leq t \leq T$, a linear unbiased estimate \bar{k}^T of k is said to be efficient if its variance is asymptotically minimum among all linear unbiased estimates.

For known mean value functions of the form

$$\phi(t) = t^2 e^{at} \sum_{\alpha=1}^n \phi_{\alpha} e^{i\lambda_{\alpha} t}$$

efficient estimates are given for the two cases $a = 0$ and $a > 0$. They are shown to be economical of information concerning the spectral density in that no estimate exists which is efficient for a wider class of spectra.

CONTENTS

<u>Section</u>		<u>Page</u>
	Foreword	iii
	Abstract	v
1	The Problem and Its Background	1
2	Summary	7
3	Efficient Estimates	12
4	Efficiency Classes	30
5	References	40

Section 1
THE PROBLEM AND ITS BACKGROUND

Let $x(t)$ be a second order stationary stochastic process with mean value function zero and covariance

$$E[x(t) \overline{x(s)}] = R(t-s) , \quad (1.1)$$

and suppose that the process

$$y(t) = k\phi(t) + x(t) \quad (1.2)$$

is observed for the continuous time parameter t in the interval $0 \leq t \leq T$. The function $\phi(t)$ is known, and the unknown parameter k is to be estimated. This process can be thought to consist of a systematic component $k\phi(t)$, which is completely predictable except for the unknown scale factor k , plus a random disturbing component $x(t)$. The problem is to estimate this scale factor from the observed process.

Only linear unbiased estimates will be considered, and they will be represented as linear functionals

$$\overline{k}^T = \overline{k}^T [y(t), 0 \leq t \leq T] \quad (1.3)$$

which are limits in quadratic mean of finite linear combinations of the process $y(t_i)$, $i = 1, \dots, N$ where the arguments t_i are in the interval $0 \leq t_i \leq T$. Such an estimate \overline{k}^T is said to be

asymptotically efficient, or simply efficient, for the problem

$[\varphi(t), R(t)]$ if it satisfies the condition

$$\frac{\text{variance } \bar{k}^T}{\text{variance } \hat{k}^T} \rightarrow 1 \quad \text{as } T \rightarrow \infty, \quad (1.4)$$

where \hat{k}^T is the minimum variance linear unbiased estimate of k computed for the problem determined by the mean value function $\varphi(t)$ and the covariance $R(t)$.

Interest in efficient estimates arises from the fact that the "best" estimate \hat{k}^T is usually very inconvenient. If the estimate \hat{k}^T is represented by the linear functional $\hat{k}^T[y(t), 0 \leq t \leq T]$ then it must satisfy the linear equation

$$\hat{k}^T[R(t-s), 0 \leq t \leq T] = \hat{M}^T \overline{\varphi(s)} \quad 0 \leq s \leq T, \quad (1.5)$$

where \hat{M}^T is a constant to be determined by the condition of unbiasedness,

$$\hat{k}^T[\varphi(t), 0 \leq t \leq T] = 1. \quad (1.6)$$

For most combinations of functions $\varphi(t)$ and $R(t)$ it is difficult to exhibit this solution explicitly; and provided it can be exhibited at all, the solution usually depends on complete knowledge of $R(t)$ and is cumbersome to compute. Efficient estimates are provided by linear functional \bar{k}^T which are in a sense asymptotic solutions to the linear equation (1.5).

The principal alternative estimate, which has been proposed, is the least square estimate,

$$\bar{k}_L^T = \int_0^T \overline{\varphi(t)y(t)} dt / \int_0^T |\varphi(t)|^2 dt . \quad (1.7)$$

This estimate has the advantages that it is easy to compute and requires no knowledge whatever of the covariance $R(t)$. Previous work on the problem of efficient estimates has been primarily devoted to determining those combinations of functions $\varphi(t)$ and $R(t)$ for which the least square estimate is efficient.

For the Ornstein Uhlenbeck process, that is,

$$R(t) = e^{-\rho|t|} , \quad (1.8)$$

and for mean value functions

$$\varphi(t) = t^r \text{ or } e^{-i\lambda_0 t} , \quad (1.9)$$

where r is a non-negative integer and λ_0 is a real frequency, Mann and Moranda in reference [10] proved that the least square estimate is efficient. The author in reference [13] extended this result to include mean value functions of the form

$$\varphi(t) = t^r e^{-i\lambda_0 t} \quad (1.10)$$

and showed further that for the more general function

$$\varphi(t) = \sum_{\alpha=1}^n \varphi_{\alpha} t^r e^{-i\lambda_{\alpha} t} \quad (1.11)$$

where the φ_{α} are non-zero constants and $n > 1$, the least square estimate is not efficient.

For a much broader class of covariance functions $R(t)$ and essentially the same mean value functions $\varphi(t)$, this problem was first discussed by Grenander in reference [3]. Further work was carried out by Grenander and Rosenblatt in references [4] and [5]. Rosenblatt considered some of the same problems in the case of vector valued time series in reference [11] and extended his results in [12]. Most of these results together with some examples appear in Chapter 7 of reference [6]. In this work only the discrete parameter case is considered, and the mean value function $\varphi(t)$ is assumed to be "slowly increasing." This assumption requires that the function $\varphi(t)$ be essentially of form (1.11). All restrictions on the class of covariances are imposed on the equivalent class of spectral densities $f(\lambda)$, which by assumption exist and satisfy the relation

$$R(t) = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\lambda t} f(\lambda) d\lambda \quad (1.12)$$

for a discrete parameter process and

$$R(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} f(\lambda) d\lambda \quad (1.13)$$

for a continuous parameter process. Thus a problem $[\varphi, R]$ can be referred to by the equivalent $[\varphi, f]$. For positive continuous spectral density and "slowly increasing" mean value function, a necessary and

sufficient condition is given in reference [6] for the least square estimate to be efficient. The condition, which is the same as that found in reference [13] for the continuous parameter Ornstein Uhlenbeck process and mean value function of the form (1.11), is that $f(-\lambda_{\alpha})$ be constant for $\alpha = 1, \dots, n$.

In Chapter 1.3 of reference [7] Grenander reproduces a few of these results using the methods of Toeplitz forms. In Chapter 1.4, under certain regularity conditions on $f(\lambda)$, he extends his results to the continuous parameter case for the single example

$$\varphi(t) = 1. \quad (1.14)$$

With the exception of those in reference [7], all the above mentioned results are derived for the more general problem

$$E[y(t)] = \sum_{i=1}^p k_i \varphi_i(t) \quad (1.15)$$

where the k_i are unknown parameters and the $\varphi_i(t)$ are known functions. For $p > 1$, the definition of efficiency used by Mann, Moranda, and Striebel is different from that used by Rosenblatt and Grenander. For the case $p = 1$, both agree with definition (1.4) made above. In the present paper only the case $p = 1$ will be considered though it is believed that the results obtained could be generalized to larger p using the methods and definitions of efficiency given by Rosenblatt and Grenander.

In the present paper $f(\lambda)$ is assumed to be a rational function with no real zeros. Again it is believed that this requirement could be slightly relaxed so as to include the case considered in Chapter 1.4 of reference [7] though this has not been attempted.

Under these two limitations, for the continuous parameter process the results in references [10], [13], and [7], mentioned above, are special cases of the theory to be presented here. The continuous parameter equivalents of the important examples discussed in reference [6] are also covered.

Section 2

SUMMARY

In the results summarized above it has been shown that in certain circumstances the least square estimate is efficient for a large class of spectral densities. It shall be the purpose here to discuss efficiency classes of spectral densities defined as follows: For a fixed mean value function $\varphi(t)$, a class of spectral densities \mathcal{F} and a member $f(\lambda)$ of \mathcal{F} , the efficiency class $\mathcal{F}(f, \varphi)$ of the density $f(\lambda)$ in the class \mathcal{F} with respect to the mean value function $\varphi(t)$ is defined to be all members $g(\lambda)$ of \mathcal{F} for which there exists an estimate (not necessarily the least square) that is efficient for the two problems $[\varphi, f]$ and $[\varphi, g]$. Thus for $\varphi(t) = 1$ and a certain class of densities \mathcal{F} , Grenander in Chapter 1.3 of reference [7] shows that $\mathcal{F}(f, \varphi) = \mathcal{F}$ for all $f \in \mathcal{F}$, since in this case he establishes that the least square estimate is efficient for all members of his class \mathcal{F} .

Section 3 is devoted to establishing the efficiency of certain proposed estimates (2.10), (2.11), and (2.13). It is also shown that for questions of efficiency the problem based on the process $y(t)$ for t in the finite interval $0 \leq t \leq T$ can be extended to an equivalent problem with t in the semi-infinite interval $-\infty < t \leq T$. This property is used in

section 4 to apply known results for the latter problem in obtaining a necessary condition that a density $g(\lambda)$ belong to an efficiency class $\mathcal{F}(f, \varphi)$.

The mean value functions $\varphi(t)$ which will be considered are of the form

$$\varphi(t) = \sum_{\gamma=1}^m \sum_{j=0}^{r_{\gamma}} \varphi_{\gamma j} t^j e^{-i\lambda_{\gamma} t} \quad (2.1)$$

where r_{γ} is a non-negative integer, φ_{γ} and λ_{γ} are complex,

$$\max_{1 \leq \gamma \leq m} \operatorname{Re}(\lambda_{\gamma}) = a \geq 0, \quad (2.2)$$

and $\varphi_{\gamma r_{\gamma}} \neq 0$. The spectral densities considered are assumed to be rational functions with no real zeros. For spectral densities of this type, the solution of the linear equation (1.5) for the finite interval $0 \leq t \leq T$ is outlined by Laning and Battin in Chapter 8.4 of reference [8]. However, this estimate $\hat{k}^T[\varphi(t), 0 \leq t \leq T]$ will not be used here, since it will be shown that efficient estimates can instead be compared with the much more convenient estimate $\hat{k}^T[y(t), -\infty < t \leq T]$.

In order to exhibit the estimates, whose efficiency is considered in section 3, some notation must first be developed. The one-sided Laplace transform,

$$\phi(\lambda) = \int_0^{\infty} e^{-i\lambda t} \varphi(t) dt, \quad (2.3)$$

of $\varphi(t)$ given by (2.1) can be written

$$\phi(\lambda) = \sum_{\gamma=1}^m \sum_{j=1}^{r_{\gamma}+1} \frac{\phi_{r_{\gamma}j}}{(\lambda+\lambda_{\gamma})^j}, \quad (2.4)$$

where $\phi_{r_{\gamma}j} = (-1)^{j+1} j! \phi_{\gamma j+1}$. Let r be the maximum of r_{γ} among subscripts γ for which $\Re \lambda_{\gamma} = a$, and indicate by $\alpha = 1, \dots, n$ those subscripts for which $\Re(\lambda_{\alpha}) = a$ and $r_{\alpha} = r$. It will be seen that for questions of efficiency only those terms of the mean value function whose absolute value increases as $t^r e^{at}$ will be pertinent. Thus the mean value function (2.1) can in effect be replaced by

$$\varphi(t) = t^r e^{at} \sum_{\alpha=1}^n \phi_{\alpha r} e^{-it\Re(\lambda_{\alpha})}. \quad (2.5)$$

Using the method described by Doob on page 542 of reference [2], the spectrum can be factored and written in the form

$$f(\lambda) = |F(\lambda)|^2, \quad (2.6)$$

where $F(\lambda)$ is a rational function with zeros and poles in the upper half-plane, ($\Re \lambda > 0$). By the division algorithm the quotient of polynomials $1/F(\lambda)$ can be expanded as

$$\frac{1}{F(\lambda)} = E(\lambda) + M(\lambda) \quad (2.7)$$

where $E(\lambda)$ is a polynomial of degree e and $M(\lambda)$ is a proper rational function with poles in the upper half-plane. By proper it is meant that the degree of the numerator is less than the degree of the denominator. Similarly, the following expansions can be made:

$$\frac{\bar{\phi}_{\alpha r+1}}{F(-\bar{\lambda}_{\alpha})(\lambda+\bar{\lambda}_{\alpha})F(\lambda)} = E_{\alpha}(\lambda) + M_{\alpha}(\lambda) \quad \alpha = 1, \dots, n \quad (2.8)$$

where $E_{\alpha}(\lambda)$ is a polynomial of degree $e-1$; and, provided $a > 0$, $M_{\alpha}(\lambda)$ is a proper rational function with poles in the upper half-plane.

Let

$$E(\lambda) = \sum_{j=0}^e e_j \lambda^j \quad (2.9)$$

$$E_{\alpha}(\lambda) = \sum_{j=0}^{e-1} e_{\alpha j} \lambda^j \quad \alpha = 1, \dots, n$$

$$m(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} M(\lambda) d\lambda$$

$$m_{\alpha}(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} M_{\alpha}(\lambda) d\lambda \quad \alpha = 1, \dots, n.$$

In section 3 the following estimates are proved efficient:

$$\begin{aligned} & \tilde{M}^T k_1^T [y(t), 0 \leq t \leq T; a > 0] \\ &= \sum_{\alpha=1}^n e^{i\bar{\lambda}_{\alpha} T} \left[\sum_{j=0}^{e-1} (-i)^j e_{\alpha j} y^{(j)}(T) + \int_0^T m_{\alpha}(T-t) y(t) dt \right]; \end{aligned} \quad (2.10)$$

$$\begin{aligned} & \tilde{M}^T k_2^T [y(t), 0 \leq t \leq T; a > 0] \\ &= \sum_{j=0}^e (-i)^j e_j Y_T^{(j)}(T) + \int_0^T m(T-t) Y_T(t) dt, \end{aligned} \quad (2.11)$$

where

$$Y_T(t) = \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{F(-\lambda_{\alpha})} e^{i\bar{\lambda}_{\alpha}(T-t)} \int_0^t e^{i\bar{\lambda}_{\alpha}u} y(u) du ; \quad (2.12)$$

and

$$\begin{aligned} \tilde{M}^T k^T [y(t), 0 \leq t \leq T; a = 0] & \quad (2.13) \\ &= \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{f(-\lambda_{\alpha})} \int_0^T t^r e^{i\lambda_{\alpha}t} y(t) dt . \end{aligned}$$

In each case \tilde{M}^T denotes a constant to be determined so that the estimate is unbiased. It is evaluated by replacing $y(t)$ with $\varphi(t)$ in the right side of the expressions (2.10), (2.11), and (2.13). As the notation indicates, estimates (2.10) and (2.11) are efficient in the case $a > 0$, and estimate (2.13) is efficient when all the frequencies λ_{α} are real, that is, $a = 0$. In order to compute the estimates (2.10) and (2.11), $\tilde{k}_1^T[a > 0]$ and $\tilde{k}_2^T[a > 0]$, it is necessary that the entire spectral density $f(\lambda)$ be known. Estimate $\tilde{k}^T[a = 0]$, (2.13), requires knowledge of the spectrum only at the points $-\lambda_{\alpha}$, that is, $f(-\lambda_{\alpha})$, $\alpha = 1, \dots, n$.

Results concerning the least square estimate can be obtained by letting $F(\lambda) = 1$, for which expressions (2.10), (2.11) and (2.13) reduce to the least square estimate for the equivalent mean value function (2.5).

Section 3
EFFICIENT ESTIMATES

Since $f(\lambda)$ is a positive rational function it is clear that the condition

$$\int_{-\infty}^{\infty} \frac{\log f(\lambda)}{1 + \lambda^2} d\lambda > -\infty \quad (3.1)$$

is satisfied, and hence by Theorem 10.12b of reference [7] the process $x(t)$ can be represented as a moving average

$$x(t) = \int_{-\infty}^t \gamma(t-s) d\xi(s) \quad -\infty < t < \infty \quad (3.2)$$

where $\xi(s)$ is a fundamental random process with zero mean value and $\gamma(t)$ is the Fourier inverse of $F(\lambda)$,

$$\gamma(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} F(\lambda) d\lambda . \quad (3.3)$$

Suppose there exists a function $\psi(t)$ which satisfies the conditions

$$\int_{-\infty}^t \gamma(t-s) \dot{\psi}(s) ds = \phi(t) \quad 0 \leq t < \infty \quad (3.4)$$

and

$$\int_{-\infty}^T |\dot{\psi}(s)|^2 ds < \infty \quad \text{all } T < \infty , \quad (3.5)$$

then the process $y(t)$ can be represented as the moving average

$$y(t) = \int_{-\infty}^t \gamma(t-s) d\eta(s) \quad 0 \leq t < \infty \quad (3.6)$$

where $\eta(s)$ is the process defined by

$$\eta(s) = k\psi(t) + \xi(t) \quad -\infty < t < \infty . \quad (3.7)$$

Thus there are three processes which will be considered

I	$y(t) = k\phi(t) + x(t)$	$0 \leq t < \infty$
II	$y(t) = k\xi(t) + x(t)$	$-\infty < t < \infty$
III	$\eta(t) = k\psi(t) + \xi(t)$	$-\infty < t < \infty .$

The $x(t)$ processes are second order stationary with zero mean value and spectral density $f(\lambda)$. The $\xi(t)$ process is a fundamental random process. The mean value functions satisfy the conditions

$$\int_{-\infty}^t \gamma(t-s) \dot{\psi}(s) ds = \zeta(t) \quad -\infty < t < \infty \quad (3.8)$$

$$\zeta(t) = \phi(t) \quad 0 \leq t < \infty$$

and condition (3.5). Their transforms $\Phi(\lambda)$, given by (2.3), (3.9)

$$Z(\lambda) = \int_{-\infty}^{\infty} e^{-i\lambda t} \zeta(t) dt , \quad (3.10)$$

and

$$\Psi(\lambda) = \int_{-\infty}^{\infty} e^{-i\lambda t} \dot{\psi}(t) dt , \quad (3.11)$$

are assumed to converge to proper rational functions in some non-degenerate strip. They satisfy conditions (2.4) and from (3.8)

$$Z(\lambda) = \Psi(\lambda)F(\lambda) . \quad (3.12)$$

The integral equation (3.4) can be solved using Laplace transforms and the Wiener-Hopf technique. (See, for example, page 313 of [14].) Taking two-sided Laplace transforms of equation (3.4) yields the equation

$$\Phi(\lambda) + H(\lambda) = \Psi(\lambda)F(\lambda) \quad (3.13)$$

where

$$H(\lambda) = \int_{-\infty}^0 e^{-i\lambda t} \int_{-\infty}^t \gamma(t-s)\dot{\psi}(s)dsdt . \quad (3.14)$$

Solving (3.13), the transform $\Psi(\lambda)$ must satisfy

$$\Psi(\lambda) = \frac{\Phi(\lambda) + H(\lambda)}{F(\lambda)} \quad (3.15)$$

The integral $\Phi(\lambda)$, (2.3), converges in the region $\Re(\lambda) < -a$. In order that the Laplace inverse of $\Psi(\lambda)$ in (3.15) exist and satisfy conditions (3.4) and (3.5), it is sufficient that $H(\lambda)$ be a proper rational function with poles in the region $\Re(\lambda) < -b$ for some $-b < -a$ and that $\Psi(\lambda)$ also be a proper rational function. Such a function $H(\lambda)$ can be found, for example, as follows: Let $D_H(\lambda)$ be any polynomial with zeros whose imaginary parts are less than $-b$. By the division algorithm the rational function $\Phi(\lambda)D_H(\lambda)$ can be written

$$\Phi(\lambda)D_H(\lambda) = Q(\lambda) + R(\lambda) , \quad (3.16)$$

where $R(\lambda)$ is a proper rational function and $Q(\lambda)$ is a polynomial of degree equal to that of $\Phi(\lambda)D_H(\lambda)$ —that is, the degree of the numerator of $\Phi(\lambda)D_H(\lambda)$ minus the degree of its denominator. Since $\Phi(\lambda)$

is proper, the degree of $Q(\lambda)$ is less than that of $D_H(\lambda)$. Now define $H(\lambda)$ by

$$H(\lambda) = -\frac{Q(\lambda)}{D_H(\lambda)}, \quad (3.17)$$

so that

$$H(\lambda) + \phi(\lambda) = \frac{R(\lambda)}{D_H(\lambda)}. \quad (3.18)$$

By taking the degree of $D_H(\lambda)$ larger than e , the degree of $1/F(\lambda)$, a proper function $\Psi(\lambda)$ is obtained. From (3.12) $Z(\lambda)$ can be defined by

$$Z(\lambda) = \phi(\lambda) + H(\lambda). \quad (3.19)$$

Thus, given a process I, the existence of processes II and III has been established.

In reference [9], Mann has solved the problem of estimating the parameter k for the process III. The minimum variance linear unbiased estimate is given by

$$\hat{M}^T \hat{k}^T[\eta(t), -\infty < t \leq T] = \int_{-\infty}^T \overline{\psi(t)} d\eta(t), \quad (3.20)$$

and

$$\text{variance } \hat{k}^T[\eta(t), -\infty < t \leq T] = 1 / \int_{-\infty}^T |\overline{\psi(t)}|^2 dt. \quad (3.21)$$

The remainder of the section will be devoted to obtaining approximate expressions for the estimate (3.20), which are more convenient but are still efficient for the problem. The efficiency of proposed estimates $\tilde{k}_1^T[a > 0]$, $\tilde{k}_2^T[a > 0]$, and $\tilde{k}^T[a = 0]$ given by (2.10), (2.11), and (2.13) is proved in Theorem 1. However, before this is done, a non-rigorous derivation of the estimates will be given by way of motivation.

First equation (3.6) will be solved for $d\eta(s)$ by transform methods. Since the poles of $F(\lambda)$ are in the upper half-plane, its inverse transform $\gamma(t)$ vanishes for $t < 0$. Thus the process $y(t)$ can be written

$$y(t) = \int_{-\infty}^{\infty} \gamma(t-s)d\eta(s), \quad (3.22)$$

and

$$\begin{aligned} \int_{-\infty}^T e^{-i\lambda t} y(t) dt &= \int_{-\infty}^T e^{-i\lambda t} \int_{-\infty}^{\infty} \gamma(t-s) d\eta(s) dt \\ &= \int_{-\infty}^T \left[\int_{-\infty}^{\infty} e^{-i\lambda(\tau+s)} \gamma(\tau) d\tau \right] d\eta(s) \\ &= F(\lambda) \int_{-\infty}^T e^{-i\lambda s} d\eta(s). \end{aligned} \quad (3.23)$$

Substituting the Laplace inversion formula for both $\dot{\psi}(t)$ and $d\eta(t)$ in formula (3.20) gives the alternative form for the estimate

$$\begin{aligned}
\hat{M}^T \hat{k}^T [\eta(t), -\infty < t \leq T] &= \int_{-\infty}^T \overline{\psi(t)} d\eta(t) \\
&= \frac{1}{(2\pi)^2} \int_{-\infty}^T \left\{ \left[\int_{-\infty-ic}^{+\infty-ic} e^{i\lambda t} \frac{\overline{\Phi(-\lambda)} + \overline{H(-\lambda)}}{\overline{F(-\lambda)}} d\lambda \right] \right. \\
&\quad \left. \left[\int_{-\infty+id}^{+\infty+id} \frac{e^{i\omega t}}{\overline{F(\omega)}} \left(\int_{-\infty}^T e^{-i\omega u} y(u) du \right) d\omega \right] \right\} dt, \tag{3.24}
\end{aligned}$$

where c is in the region of convergence of the integral $\overline{\psi(\lambda)}$ in (3.11), $b > c > a$; d is taken in the interval $0 < d < c$ where the mean and variance of the process $\int_{-\infty}^T e^{-i\omega t} y(t) dt$ and the integral $\int_{-\infty}^T e^{(i\omega+i\lambda)t} dt$ are all finite. The notation $\overline{\Phi(\lambda)}$ indicates that complex conjugates of the coefficients in the rational function $\Phi(\lambda)$ are taken.

The most important change that will be made in (3.24) is truncation of $y(t)$ for $t < 0$. Thus (3.24) becomes

$$\hat{M}^T \hat{k}^T [y(t), 0 \leq t \leq T] = \frac{1}{2\pi} \int_{-\infty+id}^{+\infty+id} \frac{\overline{\overline{\psi^T(\omega)}}}{\overline{F(\omega)}} \left[\int_0^T e^{-i\omega t} y(t) dt \right] d\omega \tag{3.25}$$

where

$$\begin{aligned}
\overline{\overline{\psi^T(\omega)}} &= \frac{1}{2\pi} \int_{-\infty-ic}^{+\infty-ic} \frac{e^{T(i\omega+i\lambda)}}{i(\omega+\lambda)} \left[\frac{\overline{\Phi(-\lambda)} + \overline{H(-\lambda)}}{\overline{F(-\lambda)}} \right] d\lambda \\
&= \int_{-\infty}^T e^{-i\omega t} \overline{\psi(t)} dt. \tag{3.26}
\end{aligned}$$

Additional changes in the estimate (3.25) are made by expanding $\overline{\Psi^T(\omega)}$ as a sum of residues at the poles of $\Phi(\lambda)$, $H(\lambda)$ and $1/F(\lambda)$ and then picking out those terms which dominate at $T \rightarrow \infty$. When terms of $\overline{\Psi^T(\omega)}$ are omitted, the complex integral (3.25) will in general cease to exist in the common sense. It must be replaced by a Cesàro limit of order $e + 1$. This limit is defined and its important properties are stated in Theorem 2. The final estimates $\tilde{k}_1^T[a > 0]$, $\tilde{k}_2^T[a > 0]$ and $\tilde{k}^T[a = 0]$ given by (2.10), (2.11), and (2.13) are then obtained by interpreting the resulting Cesàro limit as the Laplace inverse of appropriate functions of $y(t)$ for $0 \leq t \leq T$.

The simplification of $\overline{\Psi^T(\omega)}$ will be different for the two cases $a > 0$ and $a = 0$. If $a > 0$, the dominant terms of $\overline{\Psi^T(\omega)}$ have magnitude $T^r e^{(a+\omega)T}$. They are contributed by the poles of $\overline{\Phi(-\lambda)}$, $\lambda = \lambda_\alpha$, $\alpha = 1, \dots, n$ and are of the form

$$\overline{\Psi^T(\omega)} \sim \frac{T^r}{r!} \sum_{\alpha=1}^n \frac{\overline{\Phi_{\alpha r+1}} e^{iT(\omega + \bar{\lambda}_\alpha)}}{F(-\bar{\lambda}_\alpha) (\omega + \bar{\lambda}_\alpha)} \quad (3.27)$$

Making this substitution in (3.25) gives the estimate

$$\begin{aligned} \tilde{M}^T \tilde{k}^T[y(t), 0 \leq t \leq T; a > 0] \\ = \frac{1}{2\pi} \int_{-\infty + id}^{+\infty + id} \left[\sum_{\alpha=1}^n \frac{\overline{\Phi_{\alpha r+1}} e^{i(\bar{\lambda}_\alpha + \omega)T}}{F(-\bar{\lambda}_\alpha) (\omega + \bar{\lambda}_\alpha) F(\omega)} \right] \left[\int_0^T e^{-i\omega t} y(t) dt \right] d\omega \quad (3.28) \end{aligned}$$

Estimates (2.10) and (2.11) are two slightly different interpretations of this Cesàro limit as a Laplace inverse. The existence in quadratic mean of the derivatives of $y(t)$ involved in estimates (2.10) and (2.11) can easily be established. In order to show that (3.28) itself provides a bona fide estimate, it would first have to be shown that the Cesàro limit involved exists in quadratic mean.

In the case $a = 0$, the magnitude of the dominant terms in the variance of \tilde{k}^T is T^{2r+1} . They are contributed by terms in $\tilde{\psi}^T(\omega)$ of the form

$$\overline{\tilde{\psi}^T(\omega)} \sim \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1} e^{iT(\lambda_{\alpha}+\omega)}}{F(-\lambda_{\alpha})(\omega+\lambda_{\alpha})^{r+1}} \quad (3.29)$$

In this case the dominant terms in (3.25) are the residues at the poles $\omega = -\lambda_{\alpha}$, which are evaluated simply in the estimate $\tilde{k}^T[a = 0]$ given by (2.13).

Theorem 1. Let $y(t)$, $0 \leq t \leq T$, be a process of type I where $x(t)$ is a second order stationary process whose spectral density $f(\lambda)$ is a rational function with no real zeros, and the mean value function $\varphi(t)$ is of the form (2.1). Then for $a > 0$, the estimates $\tilde{k}_1^T[y(t), 0 \leq t \leq T; a > 0]$ and $\tilde{k}_2^T[y(t), 0 \leq t \leq T; a > 0]$ given by (2.10) and (2.11) are efficient for estimating the parameter k . If $a = 0$, then estimate $\tilde{k}^T[y(t), 0 \leq t \leq T; a = 0]$ given by (2.13) is efficient.

Proof: Since any linear unbiased estimate based on process I

$(y(t), 0 \leq t \leq T)$ is also a linear unbiased estimate for process III

$(\eta(t), -\infty < t \leq T)$, it follows from (3.21) that the inequality

$$\text{variance } \hat{k}^T[y(t), 0 \leq t \leq T] \geq 1 / \int_{-\infty}^T |\dot{\psi}(t)|^2 dt \quad (3.30)$$

must hold. Thus it is sufficient to show that the limit

$$[\text{variance } \tilde{k}^T][\int_{-\infty}^T |\dot{\psi}(t)|^2 dt] \rightarrow 1 \quad \text{as } T \rightarrow \infty \quad (3.31)$$

holds in order to establish efficiency of an estimate \tilde{k}^T . The remainder of the proof is devoted to obtaining asymptotic expressions for the quantities $\int_{-\infty}^T |\dot{\psi}(t)|^2 dt$, \tilde{M}^T , and variance $\tilde{M}^T \tilde{k}^T$ computed for the three estimates considered. The notation

$$A^T \sim B^T \quad (3.32)$$

will be used to mean

$$\frac{A^T}{B^T} \rightarrow 1 \quad \text{as } T \rightarrow \infty. \quad (3.33)$$

In addition to the standard Laplace transform inversion formula, the following result due to Amerio [1] will be used:

Theorem 2. If the function $\varphi(t)$ and its first n derivatives are absolutely continuous in the interval $0 \leq t \leq b$, then the derivative $\varphi^{(n)}(t)$ satisfies the inversion formula

$$\varphi^{(n)}(t) = \lim_{R \rightarrow \infty} \frac{1}{2\pi} \int_{-R+ic}^{R+ic} (i\lambda)^n e^{i\lambda t} \varphi(\lambda) \left(1 - \frac{\lambda^2}{R^2}\right)^{n+1} d\lambda \quad a \leq t \leq b \quad (3.34)$$

where $\varphi(\lambda)$ is the transform of $\varphi(t)$

$$\varphi(\lambda) = \int_{-\infty}^{\infty} e^{-i\lambda t} \varphi(t) dt \quad (3.35)$$

and the line $\lambda = ic$ is in the strip of convergence of this integral for $\varphi(\lambda)$.

In addition, if $\varphi(\lambda)$ is a proper rational function, the Cesàro limit (3.34) can be evaluated by the contour integral on a path C around the poles of $\varphi(\lambda)$ above ic if $t > 0$ and below ic if $t < 0$,

$$\lim_{R \rightarrow \infty} \frac{1}{2\pi} \int_{-R+ic}^{+R+ic} (i\lambda)^n e^{i\lambda t} \varphi(\lambda) \left(1 - \frac{\lambda^2}{R^2}\right)^{n+1} d\lambda = \frac{1}{2\pi} \oint (i\lambda)^n e^{i\lambda t} \varphi(\lambda) d\lambda \quad (3.36)$$

From the Laplace inversion formula

$$\begin{aligned} \int_{-\infty}^T |\dot{\psi}(t)|^2 dt &= \frac{1}{(2\pi)^2} \int_{-\infty}^T \int_{-\infty-ic}^{+\infty-ic} \int_{-\infty-ic}^{+\infty-ic} e^{(i\lambda+i\omega)t} \bar{\psi}(\lambda) \bar{\psi}(-\omega) d\lambda d\omega dt \\ &= \frac{1}{(2\pi)^2} \int_{-\infty-ic}^{+\infty-ic} \int_{-\infty-ic}^{+\infty-ic} \frac{e^{T(i\lambda+i\omega)}}{i(\lambda+\omega)} \left(\frac{\varphi(\lambda)+H(\lambda)}{F(\lambda)} \right) \left(\frac{\bar{\varphi}(-\omega)+\bar{H}(-\omega)}{F(-\omega)} \right) d\lambda d\omega \end{aligned} \quad (3.37)$$

The complex integrals are taken in the region of analyticity which separates the poles of $\varphi(\lambda)/F(\lambda)$ from those of $H(\lambda)$. These integrals can be taken as contour integrals closing them upward around the poles of $\varphi(\lambda)/F(\lambda)$. Examination of the residues at the poles of $\varphi(\lambda)$ and

of $1/F(\lambda)$ shows that for large T , the dominant terms are contributed by the poles with minimum imaginary part and among these, by the poles with maximum order—that is, by $\lambda = -\lambda_\alpha$, $\alpha = 1, \dots, n$. The residue of $H(\lambda)/F(\lambda)$ at poles of $\Phi(\lambda)$ is zero, so the asymptotic expression

$$\int_{-\infty}^T |\dot{\psi}(t)|^2 dt \sim \frac{1}{(2\pi)^2} \oint \oint \frac{e^{T(i\lambda+i\omega)}}{i(\lambda+i\omega)} \frac{\Phi(\lambda)}{F(\lambda)} \frac{\bar{\Phi}(-\omega)}{\bar{F}(-\omega)} d\lambda d\omega \quad (3.38)$$

is obtained, where the contours include only the poles $\lambda = -\lambda_\alpha$ and $\omega = \bar{\lambda}_\alpha$, $\alpha = 1, \dots, n$. Evaluating this integral for $a > 0$ and picking out the terms of magnitude $T^{2r} e^{2a}$, gives the asymptotic expression

$$\int_{-\infty}^T |\dot{\psi}(t)|^2 dt \sim \frac{T^{2r} e^{2a}}{(r!)^2} \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\Phi_{\alpha r+1} \bar{\Phi}_{\beta r+1} e^{iT\phi(\lambda_\beta - \lambda_\alpha)}}{i(\bar{\lambda}_\beta - \lambda_\alpha) F(-\lambda_\alpha) \bar{F}(-\lambda_\beta)} \quad (3.39)$$

This can be summarized by the form

$$\int_{-\infty}^T |\dot{\psi}(t)|^2 dt \sim T^{2r} e^{2a} C(T) \quad (3.40)$$

where $C(T)$ is bounded from above and away from zero. This is clear because $C(T)$ is a positive definite form in the vector

$$\frac{\Phi_{\alpha r+1} e^{-iT\phi\lambda_\alpha}}{F(-\lambda_\alpha)} \quad \alpha = 1, \dots, n \quad (3.41)$$

whose terms have finite positive absolute values which do not depend on T . For $a = 0$, terms of order T^{2r+1} dominate and asymptotically the integral (3.38) becomes

$$\int_{-\infty}^T |\dot{\psi}(t)|^2 dt \sim \frac{T^{2r+1}}{r!(2r+1)} \sum_{\alpha=1}^n \frac{|\phi_{\alpha r+1}|^2}{F(-\lambda_{\alpha})} \quad (3.42)$$

For $a > 0$, $M_1^T[a > 0]$ can be evaluated from (2.10),

$$\tilde{M}_1^T[a > 0] = \sum_{\alpha=1}^n e^{i\bar{\lambda}_{\alpha}T} \left[\sum_{j=0}^{e-1} (-1)^j e_{\alpha j} \phi^{(j)}(T) + \int_0^T m_{\alpha}(T-t)\phi(t)dt \right] \quad (3.43)$$

From Theorem 2 this can be written

$$\begin{aligned} \tilde{M}_1^T[a > 0] &= \frac{1}{2\pi} \sum_{\alpha=1}^n e^{i\bar{\lambda}_{\alpha}T} \int_{-\infty-ic}^{+\infty+ic} e^{i\lambda T} \left[\sum_{j=0}^{e-1} (-1)^j e_{\alpha j} (i\lambda)^j \phi(\lambda) \right. \\ &\quad \left. + M_{\alpha}(\lambda)\phi(\lambda) \right] d\lambda \quad (3.44) \\ &= \frac{1}{2\pi} \int_{-\infty-ic}^{+\infty+ic} \phi(\lambda) \sum_{\alpha=1}^n e^{iT(\bar{\lambda}_{\alpha}+\lambda)} (E_{\alpha}(\lambda) + M_{\alpha}(\lambda)) d\lambda \end{aligned}$$

From (2.8) this becomes

$$\tilde{M}_1^T[a > 0] = \frac{1}{2\pi} \int_{-\infty-ic}^{+\infty+ic} \phi(\lambda) \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1} e^{iT(\bar{\lambda}_{\alpha}+\lambda)}}{F(-\lambda_{\alpha})(\lambda+\bar{\lambda}_{\alpha})F(\lambda)} d\lambda \quad (3.45)$$

This can be evaluated by a contour integral. Only terms of maximum order $T^r e^{2aT}$ are retained,

$$\tilde{M}_1^T[a > 0] \sim \frac{(i)^{r+1} T^r e^{2aT}}{r!} \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\bar{\phi}_{\alpha r+1} \bar{\phi}_{\beta r+1}}{(\bar{\lambda}_{\beta}-\lambda_{\alpha})F(-\lambda_{\alpha})} \frac{e^{iT(\lambda_{\beta}-\lambda_{\alpha})}}{F(-\lambda_{\beta})} \quad (3.46)$$

Applying Theorem 2 to (2.11), the same expression is obtained for

$$\tilde{M}_2^T[a > 0] \quad .$$

Next, variance $\tilde{M}_{k_1}^{\tau \sim \tau}[a > 0]$ will be computed for estimate (2.10).

variance $\tilde{M}_{k_1}^{\tau \sim \tau}[a > 0]$

$$= \sum_{\alpha=1}^n \sum_{\beta=1}^n e^{iT(\bar{\lambda}_\alpha - \lambda_\beta)} E \left[\sum_{j=0}^{e-1} (-1)^j e_{\alpha j} x^{(j)}(T) + \int_0^T m_\alpha(T-u)x(u)du \right] \quad (3.47)$$

$$\left[\sum_{k=0}^{e-1} (1)^k e_{\beta k} \overline{x^{(k)}(T)} + \int_0^T \overline{m_\beta(T-v) x(v)} dv \right]$$

$$= \sum_{\alpha=1}^n \sum_{\beta=1}^n e^{iT(\bar{\lambda}_\alpha - \lambda_\beta)} \left\{ \sum_{j=0}^{e-1} \sum_{k=0}^{e-1} (-1)^j (1)^k e_{\alpha j} \bar{e}_{\beta k} \frac{\partial^j}{\partial u^j} \frac{\partial^k}{\partial v^k} R(u-v) \Big|_{u=v=T} \right.$$

$$+ \sum_{j=0}^{e-1} (-1)^j e_{\alpha j} \int_0^T \overline{m_\beta(T-v)} \frac{\partial^j}{\partial u^j} R(u-v) \Big|_{u=T} dv$$

$$+ \sum_{k=0}^{e-1} (1)^k \bar{e}_{\beta k} \int_0^T m_\alpha(T-u) \frac{\partial^k}{\partial v^k} R(u-v) \Big|_{v=T} du$$

$$\left. + \int_0^T \int_0^T m_\alpha(T-u) \overline{m_\beta(T-v)} R(u-v) dudv \right\}$$

$$= \sum_{\alpha=1}^n \sum_{\beta=1}^n e^{iT(\bar{\lambda}_\alpha - \lambda_\beta)} \frac{1}{2\pi} \int_{-\infty}^{\infty} f(\lambda) \quad (C, e+1)$$

$$\left\{ \sum_{j=0}^{e-1} \sum_{k=0}^{e-1} (-1)^j (1)^k e_{\alpha j} \bar{e}_{\beta k} \frac{\partial^j}{\partial u^j} \frac{\partial^k}{\partial v^k} e^{i\lambda(u-v)} \Big|_{u=v=T} \right.$$

$$+ \sum_{j=0}^{e-1} (-1)^j e_{\alpha j} \int_0^T \overline{m_\beta(T-v)} \frac{\partial^j}{\partial u^j} e^{i\lambda(u-v)} \Big|_{u=T} dv$$

$$\left. + \sum_{k=0}^{e-1} (1)^k \bar{e}_{\beta k} \int_0^T m_\alpha(T-u) \frac{\partial^k}{\partial v^k} e^{i\lambda(u-v)} \Big|_{v=T} du + \int_0^T \int_0^T m_\alpha(T-u) \overline{m_\beta(T-v)} e^{i\lambda(u-v)} dudv \right\} d\lambda$$

$$= \sum_{\alpha=1}^n \sum_{\beta=1}^n e^{iT(\bar{\lambda}_\alpha - \lambda_\beta)} \frac{1}{2\pi} \int_{-\infty}^{\infty} f(\lambda) [E_\alpha(\lambda) e^{i\lambda T} + \int_0^T m_\alpha(T-u) e^{i\lambda u} du] [E_\beta(\lambda) e^{i\lambda T} + \int_0^T m_\beta(T-v) e^{i\lambda v} dv] d\lambda .$$

The factor in the square bracket will be examined in detail. Using the substitutions

$$e^{i\lambda T} E_\alpha(\lambda) = \frac{1}{2\pi} \int_{-\infty - id}^{+\infty - id} \frac{E_\alpha(\lambda) e^{ipT}}{(ip - i\lambda)} dp , \quad (3.48)$$

(C, e+1)

$$\int_0^T m_\alpha(T-u) e^{i\lambda u} du = \frac{1}{2\pi} \int_{-\infty - id}^{+\infty - id} \frac{e^{ipT} M_\alpha(p)}{(ip - i\lambda)} dp ,$$

where $d > 0$, and (2.8); the bracket can be written

$$E_\alpha(\lambda) e^{i\lambda T} + \int_0^T m_\alpha(T-u) e^{i\lambda u} du = \frac{1}{2\pi} \int_{-\infty - id}^{+\infty - id} \frac{e^{ipT} (E_\alpha(p) + M_\alpha(p))}{(ip - i\lambda)} dp \quad (3.49)$$

(C, e+1)

$$= \frac{1}{2\pi} \int_{-\infty - id}^{+\infty - id} \frac{\bar{\phi}_{\alpha+1} e^{ipT}}{F(-\bar{\lambda}_\alpha)(p + \bar{\lambda}_\alpha) F(p) (ip - i\lambda)} dp .$$

(C, e+1)

Let β_1, \dots, β_k be the zeros of $F(p)$. By assumption $\Re \beta_\gamma > 0$. The contour in (3.49) will be closed upward around the poles $p = \lambda, -\bar{\lambda}_\alpha, \beta_1, \dots, \beta_k$ and evaluated by residues,

$$E_\alpha(\lambda) e^{i\lambda T} + \int_0^T m_\alpha(T-u) e^{i\lambda u} du = \frac{\bar{\phi}_{\alpha+1}}{F(-\bar{\lambda}_\alpha)} \left[\frac{e^{i\lambda T}}{(\lambda + \bar{\lambda}_\alpha) F(\lambda)} - \frac{e^{-i\bar{\lambda}_\alpha T}}{F(-\bar{\lambda}_\alpha)(\bar{\lambda}_\alpha + \lambda)} + \frac{C_\gamma \delta \nu \alpha^T e^{i\beta_\gamma T}}{(\beta_\gamma - \lambda)^v} \right] . \quad (3.50)$$

Summation on γ , δ , and ν is understood. The constants $C_{\gamma\delta\nu\alpha}$ will not be explicitly evaluated. Making the substitutions (2.6) and (3.50), the integral (3.47) can now be evaluated,

$$\begin{aligned}
 & \text{variance } \tilde{M}^{\tilde{T}\tilde{T}} k_1 [a > 0] \\
 &= \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\bar{\phi}_{\alpha r+1} \phi_{\beta r+1} e^{i\mathbb{T}(\bar{\lambda}_\alpha - \lambda_\beta)}}{F(-\lambda_\alpha) F(-\lambda_\beta)} \frac{1}{(C, e+1)} \int_{-\infty}^{\infty} F(\lambda) \overline{F(\lambda)} \left[\frac{1}{(\lambda + \bar{\lambda}_\alpha)(\lambda + \lambda_\beta) F(\lambda) \overline{F(\lambda)}} \right. \\
 &+ \frac{e^{i\mathbb{T}(\lambda_\beta - \bar{\lambda}_\alpha)}}{F(-\bar{\lambda}_\alpha) \overline{F(-\lambda_\beta)} (\bar{\lambda}_\alpha + \lambda)(\lambda_\beta + \lambda)} \frac{C_{\gamma\delta\nu\alpha} C_{\gamma'\delta'\nu'} \beta^{\mathbb{T}\delta + \delta'} e^{i\mathbb{T}(\beta_\gamma - \bar{\beta}_{\gamma'})}}{(\beta_\gamma - \lambda)^\nu (\bar{\beta}_{\gamma'} - \lambda)^{\nu'}} \quad (3.51) \\
 &- \frac{e^{i\mathbb{T}(\lambda + \lambda_\alpha)}}{(\lambda + \bar{\lambda}_\alpha) F(\lambda) \overline{F(-\lambda_\beta)} (\lambda_\beta + \lambda)} - \frac{e^{-i\mathbb{T}(\bar{\lambda}_\alpha + \lambda)}}{\overline{F(\lambda)} (\lambda + \lambda_\beta) F(-\bar{\lambda}_\alpha) (\bar{\lambda}_\alpha + \lambda)} \\
 &+ \frac{C_{\gamma'\delta'\nu'} \beta^{\mathbb{T}\delta'} e^{i\mathbb{T}(\lambda - \bar{\beta}_{\gamma'})}}{(\lambda + \bar{\lambda}_\alpha) F(\lambda) (\bar{\beta}_{\gamma'} - \lambda)^{\nu'}} + \frac{C_{\gamma\delta\nu\alpha} \mathbb{T}^\delta e^{i\mathbb{T}(\beta_\gamma - \lambda)}}{\overline{F(\lambda)} (\lambda + \lambda_\beta) (\beta_\gamma - \lambda)^\nu} \\
 &\left. - \frac{C_{\gamma'\delta\nu\beta} \mathbb{T}^\delta e^{-i\mathbb{T}(\bar{\lambda}_\alpha + \bar{\beta}_{\gamma'})}}{F(-\bar{\lambda}_\alpha) (\bar{\lambda}_\alpha + \lambda) (\bar{\beta}_{\gamma'} - \lambda)^{\nu'}} - \frac{C_{\gamma\delta\nu\alpha} \mathbb{T}^\delta e^{i\mathbb{T}(\lambda_\beta + \beta_\gamma)}}{\overline{F(-\lambda_\beta)} (\lambda_\beta + \lambda) (\beta_\gamma - \lambda)^\nu} \right] d\lambda .
 \end{aligned}$$

Examination of the residues of each term shows that, with the exception of the first, all terms contain one of the factors $e^{i\mathbb{T}(\lambda_\beta - \bar{\lambda}_\alpha)}$, $e^{i\mathbb{T}(\beta_\gamma - \bar{\beta}_{\gamma'})}$, $e^{-i\mathbb{T}(\bar{\lambda}_\alpha + \bar{\beta}_{\gamma'})}$, or $e^{i\mathbb{T}(\lambda_\beta + \beta_\gamma)}$. Thus they go to zero as $\mathbb{T} \rightarrow \infty$. Hence, the asymptotic expression,

$$\text{variance } \tilde{M}^{\tilde{T}\tilde{T}} k_1 [a > 0] \sim e^{2a\mathbb{T}} \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\bar{\phi}_{\alpha r+1} \phi_{\beta r+1} e^{i\mathbb{T}(\lambda_\alpha - \lambda_\beta)}}{F(-\lambda_\alpha) F(-\lambda_\beta) i(\bar{\lambda}_\alpha - \lambda_\beta)} \quad (3.52)$$

is provided by the first term of (3.51). Computation of variance $\tilde{M}^T k_2^T [a > 0]$ for estimate (2.11) is entirely analogous and yields the same result.

Asymptotic expressions (3.39), (3.46), and (3.52) can now be combined to compute the efficiency of estimates (2.10) and (2.11),

$$\begin{aligned} \text{variance } \tilde{k}^T [a > 0] &= \int_{-\infty}^T |\dot{\psi}(t)|^2 dt & (3.53) \\ &= \frac{(\text{variance } \tilde{M}^T k^T [a > 0]) \left(\int_{-\infty}^T |\psi(t)|^2 dt \right)}{|\tilde{M}^T [a > 0]|^2} \sim 1. \end{aligned}$$

Next, $\tilde{M}^T [a = 0]$ will be computed from (2.13),

$$\begin{aligned} \tilde{M}^T [a = 0] &= \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{f(-\lambda_{\alpha})} \int_0^T t^r e^{i\bar{\lambda}_{\alpha} t} \phi(t) dt & (3.54) \\ &= \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{f(-\lambda_{\alpha})} \sum_{\gamma=1}^m \frac{\gamma^{r+1} \phi_{\beta j}}{2\pi} \int_{-\infty-ic}^{+\infty-ic} \frac{\left[\int_0^T t^r e^{(i\bar{\lambda}_{\alpha} + i\lambda) t} dt \right]}{(\lambda + \lambda_{\beta})^j} d\lambda \\ &= \sum_{\alpha=1}^n \sum_{\gamma=1}^m \frac{\gamma^{r+1} \bar{\phi}_{\alpha r+1}}{f(-\lambda_{\alpha})} \phi_{\beta j} \frac{i}{(j-1)!} \int_0^T t^r (it)^{j-1} e^{i(\lambda_{\alpha} - \lambda_{\beta}) t} dt \end{aligned}$$

The dominant terms are for $\gamma = \alpha$ and $j = r + 1$, so that

$$\tilde{M}^T [a = 0] \sim \frac{i^{r+1} T^{2r+1}}{r!(2r+1)} \sum_{\alpha} \frac{|\bar{\phi}_{\alpha r+1}|^2}{f(-\lambda_{\alpha})} \quad (3.55)$$

For estimate (2.13) the variance can be computed as follows:

$$\begin{aligned}
 & \text{variance } \hat{M}^T \hat{k}^T [a = 0] \\
 &= \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\bar{\phi}_{\alpha r+1} \phi_{\beta r+1}}{f(-\lambda_{\alpha}) f(-\lambda_{\beta})} \int_0^T \int_0^T u^r v^r e^{i\lambda_{\alpha} u - i\lambda_{\beta} v} R(u-v) du dv \\
 &= \sum_{\alpha=1}^n \sum_{\beta=1}^n \frac{\bar{\phi}_{\alpha r+1} \phi_{\beta r+1}}{f(-\lambda_{\alpha}) f(-\lambda_{\beta})} \frac{1}{2\pi} \int_{-\infty}^{\infty} f(\lambda) \\
 & \quad \left[\int_0^T \int_0^T u^r v^r e^{i\lambda_{\alpha} u - i\lambda_{\beta} v + iu\lambda - iv\lambda} du dv \right] d\lambda \\
 &= \sum_{\alpha=1}^n \sum_{\beta=1}^n \bar{\phi}_{\alpha r+1} \phi_{\beta r+1} \frac{1}{2\pi} \int_{-\infty}^{\infty} f(\lambda) \\
 & \quad \left[\sum_{j=0}^r \frac{e^{iT(\lambda+\lambda_{\alpha})} T^{r-j} r! (-1)^j}{(i\lambda_{\alpha} + i\lambda)^{j+1} (r-j)!} + \frac{r! (-1)^{r+1}}{(i\lambda_{\alpha} + i\lambda)^{r+1}} \right] \\
 & \quad \left[\sum_{k=0}^r \frac{e^{-iT(\lambda+\lambda_{\beta})} T^{r-k} r!}{(i\lambda_{\beta} + i\lambda)^{k+1} (r-k)!} + \frac{r!}{(i\lambda_{\beta} + i\lambda)^{r+1}} \right] d\lambda .
 \end{aligned} \tag{3.56}$$

Evaluating each term by residues, the only terms of order T^{2r+1} are the cross product terms for $\alpha = \beta$. The contour must be deformed around the poles and then closed upward for one cross product term and downward for the other. Thus only one of them contains the poles of interest. As before, the poles of $F(\lambda)$ and $\bar{F}(\lambda)$ will contribute terms containing the factors $e^{iT(\eta_{\gamma} + \lambda_{\alpha})}$ where η_{γ} are the poles of $F(\lambda)$ and hence $\Im(\eta_{\gamma}) > 0$, so that terms of this type are asymptotically negligible. Thus the asymptotic expression becomes

$$\begin{aligned}
& \text{variance } \tilde{M}^T \tilde{k}^T [a = 0] \\
& \sim \sum_{\alpha=1}^n \frac{|\phi_{\alpha r+1}|^2}{[f(-\lambda_{\alpha})]^2} f(-\lambda_{\alpha}) \sum_{j=0}^n \frac{T^{r-j} (iT)^{j+r-1} (r!)^2 (-1)^j}{(j+r+1)! i^{r+j+2} (r-j)!} \quad (3.57) \\
& = \frac{T^{2r+1}}{2r+1} \sum_{\alpha=1}^n \frac{|\phi_{\alpha r+1}|^2}{f(-\lambda_{\alpha})}
\end{aligned}$$

Combining expressions (3.42), (3.55), and (3.57), it is seen that estimate (2.13) is also efficient. This concludes the proof of Theorem 1.

Corollary 1. Under the assumptions of Theorem 1, it follows that

$$\frac{\text{variance } \hat{k}^T [y(t), 0 \leq t \leq T]}{\text{variance } \hat{k}^T [y(t), -\infty < t \leq T]} \rightarrow 1 \text{ as } T \rightarrow \infty \quad (3.58)$$

Proof: In the proof of Theorem 1 it was shown that

$$(\text{variance } \tilde{k}^T [y(t), 0 \leq t \leq T]) \left\{ \int_0^T |\dot{\psi}(t)|^2 dt \right\} \rightarrow 1 \quad (3.59)$$

and hence

$$\begin{aligned}
1 & \leq \frac{\text{variance } \hat{k}^T [y(t), 0 \leq t \leq T]}{\text{variance } \hat{k}^T [y(t), -\infty < t \leq T]} \leq \frac{\text{variance } \tilde{k}^T [y(t), 0 \leq t \leq T]}{\text{variance } \hat{k}^T [\eta(t), -\infty < t \leq T]} \\
& = (\text{variance } \tilde{k}^T [y(t), 0 \leq t \leq T]) \left\{ \int_{-\infty}^T |\dot{\psi}(t)|^2 dt \right\} \rightarrow 1. \quad (3.60)
\end{aligned}$$

Corollary 2. If a linear unbiased estimate is efficient for process I, then it is efficient for processes II and III.

Proof: This follows immediately from Corollary 1.

Section 4
EFFICIENCY CLASSES

In this section a necessary condition that a density g be in the efficiency class $\mathcal{F}(f, \varphi)$ will be established.

Lemma 1. If there exists a linear unbiased estimate \bar{k}^T that is efficient for either problem I or II, then it follows that

$$\int_{-\infty}^{\infty} \left| \frac{\bar{\Psi}^T(\lambda)}{b^T} - \frac{F(\lambda)K^T(\lambda)}{a^T} \right|^2 d\lambda \rightarrow 0 \quad \text{as } T \rightarrow \infty \quad (4.1)$$

where

$$K^T(\lambda) = \bar{k}^T [e^{i\lambda t}, -\infty < t \leq T], \quad (4.2)$$

$$\bar{\Psi}^T(\lambda) = \int_{-\infty}^T e^{-i\lambda t} \dot{\psi}(t) dt, \quad (4.3)$$

$$b^T = \sqrt{\int_{-\infty}^{\infty} |\bar{\Psi}^T(\lambda)|^2 d\lambda}, \quad (4.4)$$

and

$$a^T = \sqrt{\int_{-\infty}^{\infty} |F(\lambda)K^T(\lambda)|^2 d\lambda}. \quad (4.5)$$

Proof: From reference [2] page 534 it is clear that the linear estimate \bar{k}^T can be represented as

$$\begin{aligned} \bar{k}^T[y(t), -\infty < t \leq T] &= \bar{k}^T \left[\frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} Z^T(\lambda) d\lambda \right. \\ &\quad \left. + \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\lambda t} dz(\lambda), -\infty < t \leq T \right] \quad (4.6) \\ &= k \frac{1}{2\pi} \int_{-\infty}^{\infty} K^T(\lambda) Z^T(\lambda) d\lambda + \frac{1}{2\pi} \int_{-\infty}^{\infty} K^T(\lambda) dz(\lambda) \end{aligned}$$

where $dz(\lambda)$ is the spectral process of $x(t)$ with $E[|dz(\lambda)|^2] = f(\lambda)d\lambda$, $K^T(\lambda)$ is given by (4.2), and $Z^T(\lambda)$ by

$$Z^T(\lambda) = \int_{-\infty}^T e^{-i\lambda t} \psi(t) dt = \Phi^T(\lambda) + H(\lambda). \quad (4.7)$$

Unbiasedness of the estimate implies that

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} K^T(\lambda) Z^T(\lambda) d\lambda = 1. \quad (4.8)$$

The variances of the estimates \bar{k}^T and $\hat{k}^T[\eta(t)]$ are given by

$$\begin{aligned} \text{variance } \bar{k}^T &= \frac{1}{(2\pi)^2} \int_{-\infty}^{\infty} |K^T(\lambda)|^2 f(\lambda) d\lambda \quad (4.9) \\ &= \frac{1}{(2\pi)^2} \int_{-\infty}^{\infty} |K^T(\lambda) F(\lambda)| d\lambda \end{aligned}$$

and

$$1/\text{variance } \hat{k}^T[\eta(t)] = \int_{-\infty}^T |\dot{\psi}(t)|^2 dt = \frac{1}{(2\pi)^2} \int_{-\infty}^{\infty} |\dot{\Psi}^T(\lambda)|^2 d\lambda \quad (4.10)$$

Applying the argument used to derive equation (3.23) to equation (3.8), it is seen that

$$Z^T(\lambda) = F(\lambda)\Psi^T(\lambda) . \quad (4.11)$$

From Corollary 2 the estimate \hat{k}^T is also efficient for process III, so that

$$a^T b^T = \frac{\text{variance } \hat{k}^T[y(t), -\infty < t \leq T]}{\text{variance } \hat{k}^T[n(t), -\infty < t \leq T]} \rightarrow 1 \quad \text{as } T \rightarrow \infty . \quad (4.12)$$

The result can now be derived as follows:

$$\begin{aligned} \int_{-\infty}^{\infty} \left| \frac{\overline{\Psi^T(\lambda)}}{b^T} - \frac{F(\lambda)K^T(\lambda)}{a^T} \right|^2 d\lambda &= \int_{-\infty}^{\infty} \left[\frac{\overline{\Psi^T(\lambda)}}{b^T} - \frac{F(\lambda)K^T(\lambda)}{a^T} \right] \left[\frac{\overline{\Psi^T(\lambda)}}{b^T} - \frac{F(\lambda)K^T(\lambda)}{a^T} \right] d\lambda . \\ &= \left(\frac{b^T}{b^T} \right)^2 + \left(\frac{a^T}{a^T} \right)^2 - \frac{2}{a^T b^T} \Re \int_{-\infty}^{\infty} F(\lambda)K^T(\lambda)\overline{\Psi^T(\lambda)} d\lambda \\ &= 2 - \frac{2}{a^T b^T} \rightarrow 1 . \end{aligned} \quad (4.13)$$

Lemma 2. If there exists a linear unbiased estimate that is efficient for either process I or II and for the two problems $[\varphi(t), f(\lambda)]$ and $[\varphi(t), g(\lambda)]$; that is, if $g \in \mathcal{F}(\varphi, f)$, then

$$\int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda)G(\lambda)}}{b_F^T} - \frac{a_G^T \overline{\Psi_G^T(\lambda)} F(\lambda)}{a_F^T b_G^T} \right|^2 d\lambda \rightarrow 0 \quad \text{as } T \rightarrow \infty . \quad (4.14)$$

Proof: From an elementary inequality, it is clear that

$$\int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda)G(\lambda)}}{b_F^T} - \frac{a_G^T \overline{\Psi_G^T(\lambda)F(\lambda)}}{a_F^T b_F^T} \right|^2 d\lambda$$

$$\cong 2 \int_{-\infty}^{\infty} \left\{ \left| G(\lambda) \left[\frac{\overline{\Psi_F^T(\lambda)}}{b_F^T} - \frac{K^T(\lambda)F(\lambda)}{a_F^T} \right] \right|^2 + \left| \frac{a_G^T F(\lambda)}{a_F^T} \left[\frac{\overline{\Psi_G^T(\lambda)}}{b_G^T} - \frac{K^T(\lambda)G(\lambda)}{a_G^T} \right] \right|^2 \right\} d\lambda \quad (4.15)$$

Since $G(\lambda)$ and $F(\lambda)$ are rational functions with no real poles, they are bounded

$$|G(\lambda)| \leq M$$

$$|F(\lambda)| \leq M \quad (4.16)$$

for λ real and hence

$$\int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda)G(\lambda)}}{b_F^T} - \frac{a_G^T \overline{\Psi_G^T(\lambda)F(\lambda)}}{a_F^T b_F^T} \right|^2 d\lambda$$

$$\cong 2 \left\{ M^2 \int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda)}}{b_F^T} - \frac{K^T(\lambda)F(\lambda)}{a_F^T} \right|^2 d\lambda + \left(\frac{a_G^T}{a_F^T} \right)^2 M^2 \int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_G^T(\lambda)}}{b_G^T} - \frac{K^T(\lambda)G(\lambda)}{a_G^T} \right|^2 d\lambda \right\} \quad (4.17)$$

From (3.40) and (3.42)

$$\left(\frac{a_F^T}{b_F^T} \right)^2 \sim T^{2r} e^{2a} C_F(T) \quad \text{or} \quad T^{2r+1} C \quad (4.18)$$

where $C(T)$ is bounded from above and away from zero. From efficiency

$\frac{a_F^T}{b_F^T} > 1$ so that

$$\frac{a_F^T}{b_F^T} \sim \frac{T^{-2r} e^{-2a}}{C_F(T)} \quad \text{or} \quad \frac{T^{-(2r+1)}}{C} \quad (4.19)$$

Thus the expression

$$\begin{pmatrix} a_G^T \\ \frac{a_G}{T} \\ a_F^T \end{pmatrix} \sim \frac{C_F(T)}{C_G(T)} \quad \text{or} \quad \frac{C_F}{C_G} \quad (4.20)$$

is bounded and the result follows from Lemma 1.

Theorem 3. If there exists a linear unbiased estimate that is efficient for either process I or II and for the two problems $[\varphi(\lambda), f(\lambda)]$ and $[\varphi(\lambda), g(\lambda)]$; that is, if $g \in \mathcal{F}[\varphi, f]$, then (i) if $a > 0$, it follows that

$$\left| \frac{C_1}{G(\lambda)} \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{G(-\lambda_{\alpha})} \frac{e^{iT\varphi\lambda_{\alpha}}}{(\bar{\lambda}_{\alpha} + \lambda)} - \frac{1}{F(\lambda)} \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha r+1}}{F(-\lambda_{\alpha})} \frac{e^{iT\varphi\lambda_{\alpha}}}{(\bar{\lambda}_{\alpha} + \lambda)} \right| \rightarrow 0 \quad (4.21)$$

as $T \rightarrow \infty$,

and (ii) if $a = 0$,

$$F(-\lambda_{\alpha}) = C_2 G(-\lambda_{\alpha}) \quad \alpha = 1, \dots, n, \quad (4.22)$$

where C_1 and C_2 are constants.

Proof: (i) For any sequence T_n there exists a subsequence T_j , and complex numbers l_{α} such that

$$e^{iT_j \varphi \lambda_{\alpha}} \rightarrow l_{\alpha} \quad \text{as } j \rightarrow \infty, \quad \alpha = 1, \dots, n. \quad (4.23)$$

For this subsequence, it follows from (3.27), (4.18) and (4.19) that

$$\begin{aligned}
& \left| \frac{\overline{\Psi_F^T(\lambda)G(\lambda)}}{b_F^{T_j}} - \frac{a_G^{T_j} \overline{\Psi_G^T(\lambda)F(\lambda)}}{a_F^{T_j} b_G^{T_j}} \right| \\
& \sim \left| \frac{T_j^r e^{iT_j \lambda} e^{aT_j}}{r!} \left[\frac{G(\lambda)}{\sqrt{C_F(T_j)}} \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}} e^{iT_j \lambda \alpha}}{(\lambda + \overline{\lambda_\alpha}) F(-\lambda_\alpha)} \right. \right. \\
& \qquad \qquad \qquad \left. \left. - \frac{\sqrt{C_F(T_j)}}{C_G(T_j)} F(\lambda) \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}} e^{iT_j \lambda \alpha}}{(\lambda + \overline{\lambda_\alpha}) G(-\lambda_\alpha)} \right] \right| \\
& \rightarrow \left| \frac{G(\lambda)}{\sqrt{C_F}} \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}}}{(\lambda + \overline{\lambda_\alpha})} \frac{l_\alpha}{F(-\lambda_\alpha)} - \frac{\sqrt{C_F}}{C_G} F(\lambda) \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}} l_\alpha}{G(-\lambda_\alpha) (\overline{\lambda_\alpha} + \lambda)} \right|
\end{aligned} \tag{4.24}$$

From Lemma 2 this limit must be identically zero.

(ii) For $a = 0$, the integral (4.14) will be evaluated. Making the substitution $\overline{\Psi^T(\lambda)} = Z^T(\lambda) / F(\lambda)$ for $\overline{\Psi_F^T(\lambda)}$ and $\overline{\Psi_G^T(\lambda)}$, this becomes

$$\begin{aligned}
& \int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda)G(\lambda)}}{b_F^T} - \frac{a_G^T \overline{\Psi_G^T(\lambda)F(\lambda)}}{a_F^T b_G^T} \right|^2 d\lambda \\
& = \int_{-\infty}^{\infty} \left| \frac{G(\lambda)}{F(\lambda)} - \frac{a_G^T b_F^T}{a_F^T b_G^T} \frac{F(\lambda)}{G(\lambda)} \right|^2 \left| \frac{Z^T(\lambda)}{b_F^T} \right|^2 d\lambda .
\end{aligned} \tag{4.25}$$

It is clear from the discussion in section 3 that the extension of $\phi(t)$ to $\mathcal{J}(t)$ depends only on the degree of $F(\lambda)$ and hence the same extension $Z^T(\lambda)$ will do for both $F(\lambda)$ and $G(\lambda)$. The function $Z(\lambda)$ will be evaluated by closing the contour in the integral

$$Z^T(\lambda) = \frac{1}{2\pi} \int_{-\infty-ic}^{+\infty-ic} \frac{e^{(i\omega-1\lambda)T}}{(i\omega-1\lambda)} [\Phi(\omega) + H(\omega)] d\omega, \quad (4.26)$$

upward around the poles of $\Phi(\omega)$ and $\omega = \lambda$ and computing the residues

$$\begin{aligned} Z^T(\lambda) &= \Phi(\lambda) + H(\lambda) + \sum_{\gamma=1}^m \sum_{j=1}^{r_\gamma+1} \frac{\phi_{\gamma j} e^{-i\lambda T}}{(j-1)!} \left. \frac{\partial^{j-1}}{\partial \omega^{j-1}} \frac{e^{i\omega T}}{(\omega-\lambda)} \right|_{\omega = -\lambda_\gamma} \\ &= H(\lambda) + \sum_{\gamma=1}^m \sum_{j=1}^{r_\gamma+1} \phi_{\gamma j} \left[\frac{1}{(\lambda+\lambda_\gamma)^j} \right. \\ &\quad \left. + \sum_{k=0}^{j-1} \frac{e^{-i\lambda T}}{(j-1)!} \binom{j-1}{k} \frac{(iT)^k e^{-i\lambda_\gamma T}}{(\lambda+\lambda_\gamma)^{j-k}} (j-1+k)! (-1)^k \right]. \end{aligned} \quad (4.27)$$

Since $b_F^T \sim T^{2r+1} C_F$, terms of magnitude T^{2r+1} must be found in (4.25). Terms containing the factor $H(\lambda)$ will not be of this magnitude. Thus making the substitution (4.27) in (4.25)

$$\begin{aligned} &\int_{-\infty}^{\infty} \left| \frac{\overline{\Psi_F^T(\lambda) G(\lambda)}}{b_F^T} - \frac{a_G^T}{a_F^T b_G^T} \frac{\overline{\Psi_G^T(\lambda) F(\lambda)}}{\Psi_G^T(\lambda) F(\lambda)} \right|^2 d\lambda \\ &\sim \frac{1}{(b_F^T)^2} \sum_{\gamma=1}^m \sum_{\delta=1}^m \sum_{i=1}^{r_\gamma+1} \sum_{j=1}^{r_\delta+1} \phi_{\gamma i} \overline{\phi_{\delta j}} \int_{-\infty}^{\infty} \left| \frac{G(\lambda)}{F(\lambda)} \right. \\ &\quad \left. - \frac{a_G^T b_F^T F(\lambda)}{a_F^T b_G^T G(\lambda)} \left[\frac{1}{(\lambda+\lambda_\gamma)^i} + e^{-i\lambda T} \sum_{k=0}^{i-1} \frac{(iT)^k e^{-i\lambda_\gamma T}}{(\lambda+\lambda_\gamma)^{i-k}} \frac{(-1)^k}{k!} \right] \right. \\ &\quad \left. \left[\frac{1}{(\lambda+\lambda_\delta)^{j+1}} + e^{i\lambda T} \sum_{\ell=0}^{j-1} \frac{(-iT)^\ell e^{i\lambda_\delta T}}{(\lambda+\lambda_\delta)^{j-\ell}} \frac{(-1)^\ell}{\ell!} \right] \right|^2 d\lambda \end{aligned} \quad (4.28)$$

Terms of maximum order are obtained when $\gamma = \delta = \alpha$ and $i = j = r + 1$ for the cross product terms. The contour may be deformed downward around the poles at $\lambda = -\lambda_\alpha$, $\alpha = 1, \dots, n$. Again the poles of

$$\left| \frac{G(\lambda)}{F(\lambda)} - \frac{a_G^T b_F^T}{a_F^T b_G^T} \frac{F(\lambda)}{G(\lambda)} \right|^2 \text{ will not contribute terms of maximum order.}$$

Evaluating by residues,

$$\begin{aligned} & \int_{-\infty}^{\infty} \left| \frac{\overline{\psi_F^T(\lambda)} G(\lambda)}{b_F^T} - \frac{a_G^T}{a_F^T b_G^T} \overline{\psi_G^T(\lambda)} F(\lambda) \right|^2 d\lambda \\ & \sim \frac{2\pi i}{C_F^T} \frac{1}{2r+1} \sum_{\alpha=1}^n |\phi_{\alpha r+1}|^2 \left| \frac{G(-\lambda_\alpha)}{F(-\lambda_\alpha)} - \frac{C_F}{C_G} \frac{F(-\lambda_\alpha)}{G(-\lambda_\alpha)} \right|^2 \sum_{\ell=0}^r \frac{(-i\pi)^\ell (-1)^\ell (i\pi)^{2r+1-\ell}}{\ell! (2r+1-\ell)!} \\ & = \frac{2\pi}{(2r+1)(r!)^2} \sum_{\alpha=1}^n |\phi_{\alpha r+1}|^2 \left| \frac{G(-\lambda_\alpha)}{F(-\lambda_\alpha)} - \frac{C_F}{C_G} \frac{F(-\lambda_\alpha)}{G(-\lambda_\alpha)} \right|^2. \end{aligned} \tag{4.29}$$

From Lemma 2 this limit must be zero, and since $\phi_{\alpha r+1} \neq 0$, it follows that (4.22) is satisfied.

In the case $a = 0$, it is clear that the estimate $\tilde{k}^T[y(t), 0 \leq t \leq T; a = 0]$ given by (2.13) is efficient for all densities $g(\lambda)$ which satisfy condition (ii) of the theorem. Thus condition (ii) is a necessary and sufficient condition that $g(\lambda)$ belong to $\mathcal{F}(f, \varphi)$, and the estimate (2.13) is economical of information concerning the spectrum $f(\lambda)$ in the sense that there exists no linear unbiased estimate which is efficient for a wider class of spectra.

In the case $a > 0$, the theorem states that if there exists an estimate that is efficient for $[\varphi(t), f(\lambda)]$ and $[\varphi(t), g(\lambda)]$ on a sequence T_j for which condition (4.23) holds, then the spectra $f(\lambda)$ and $g(\lambda)$ are related by

$$\frac{C}{G(\lambda)} \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha+1} l_{\alpha}}{G(-\lambda_{\alpha})(\lambda + \bar{\lambda}_{\alpha})} = \frac{1}{F(\lambda)} \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha+1} l_{\alpha}}{F(-\lambda_{\alpha})(\lambda + \bar{\lambda}_{\alpha})} \quad (4.30)$$

For $n = 1$ or 2 it can be shown that this implies that $g(\lambda) = Cf(\lambda)$ and hence that estimates $\tilde{k}_1^T[y(t), 0 \leq t \leq T; a > 0]$ and $\tilde{k}_2^T[y(t), 0 \leq t \leq T; a > 0]$ given by (2.10) and (2.11) are economical in the same sense stated above for $\tilde{k}^T[y(t), 0 \leq t \leq T, a = 0]$. However, the economy of estimates (2.10) and (2.11) is of little practical interest since it is very unlikely that the function

$$\frac{1}{F(\lambda)} \sum_{\alpha=1}^n \frac{\bar{\phi}_{\alpha+1} l_{\alpha}}{F(-\lambda_{\alpha})(\lambda + \bar{\lambda}_{\alpha})} \quad (4.31)$$

would in practice be known if the function $F(\lambda)$ were unknown. This is not the case for the estimate $\tilde{k}^T[y(t), 0 \leq t \leq T; a = 0]$, since here a real reduction is made in the information required concerning $f(\lambda)$. For instance, if $n = 1$, $\tilde{k}^T[a = 0]$ becomes essentially the least square estimate and no knowledge of $f(\lambda)$ is required. For larger n , if the spectrum is to be estimated from an independent experiment, then estimates need only be made at the frequencies of interest $-\lambda_{\alpha}$, $\alpha = 1, \dots, n$, considerably reducing the work.

For the least square estimate \bar{k}_L^T given by (1.4), the function $K^T(\lambda)$ becomes

$$K^T(\lambda) = \frac{\int_0^T \overline{\varphi(t) e^{i\lambda t}} dt}{\int_0^T |\varphi(t)|^2 dt} = \frac{\overline{\phi^T(\lambda)}}{\int_0^T |\varphi(t)|^2 dt} \quad (4.32)$$

For the pseudo-spectral density $g(\lambda) = 1$,

$$\Psi_G^T(\lambda) = Z^T(\lambda) = \phi^T(\lambda) + H(\lambda), \quad (4.33)$$

and

$$K^T(\lambda) = \frac{\overline{\Psi_G^T(\lambda)} - \overline{H(\lambda)}}{\int_0^T |\varphi(t)|^2 dt}. \quad (4.34)$$

If this expression is substituted in (4.1) and terms involving $H(\lambda)$ are neglected, the expression (4.25) is obtained and the following results concerning the least square estimate can be deduced: If the least square estimate is efficient for $f(\lambda)$; then if $a = 0$

$$f(-\lambda_\alpha) = c_2 \quad \alpha = 1, \dots, n, \quad (4.35)$$

and if $a > 0$

$$\left| \frac{1}{F(\lambda)} \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}} e^{iT\varphi(\lambda_\alpha)}}{F(-\lambda_\alpha)(\lambda + \overline{\lambda_\alpha})} - c_2 \sum_{\alpha=1}^n \frac{\overline{\phi_{\alpha r+1}} e^{iT\varphi(\lambda_\alpha)}}{(\lambda + \overline{\lambda_\alpha})} \right| \rightarrow 0 \quad (4.36)$$

as $T \rightarrow \infty$.

Section 5

REFERENCES

- [1] Luigi Amerio, "Sull' inversione della trasformata di Laplace," Rendiconti delle Acc. della Scienze Fisiche e Matematiche, Napoli, Series IV, Vol. X (1940), p. 232
- [2] J. L. Doob, Stochastic Processes, John Wiley and Sons, Inc., New York, 1953.
- [3] U. Grenander, "On the estimation of regression coefficients in the case of an autocorrelated disturbance," Annals of Math. Stat. Vol. 25 (1954), p. 252.
- [4] U. Grenander and M. Rosenblatt, "An extension of a theorem of G. Szegő and its application to the study of stochastic processes," Trans. Amer. Math. Soc. Vol. 76 (1954), p. 112.
- [5] U. Grenander and M. Rosenblatt, "Regression analysis of time series with stationary residuals," Proc. Nat. Acad. Sci. Vol 40 (1954), p. 812.
- [6] U. Grenander and M. Rosenblatt, Statistical Analysis of Stationary Time Series, John Wiley and Sons, Inc., New York, 1957 .
- [7] U. Grenander and Gabor Szegő, Toeplitz Forms and Their Application, University of California Press, Berkeley and Los Angeles, 1958
- [8] J. H. Laning, Jr. and R. H. Battin, Random Processes in Automatic Control, McGraw-Hill Book Company, New York, 1956
- [9] H. B. Mann, "A theory of estimation for the fundamental random process and the Ornstein Uhlenbeck process," Sankhyā, Vol. 13 (1954), p. 325.
- [10] H. B. Mann and P. B. Moranda, "On the efficiency of the least square estimates of parameters in the Ornstein Uhlenbeck process," Sankhyā, Vol. 13 (1954), p. 351.

- [11] M. Rosenblatt, "On the estimation of regression coefficients of a vector-valued time series with stationary residuals," Annals of Math. Stat., Vol. 27 (1956), p. 99.
- [12] M. Rosenblatt, "Some regression problems in time series analysis," Third Berkeley Symposium, Vol. 1, p. 165, University of California Press, Berkeley and Los Angeles, 1956.
- [13] C. T. Striebel, "On the efficiency of estimates of trend in the Ornstein Uhlenbeck process," Annals of Math. Stat., Vol. 29 (1958), p. 192.
- [14] B. van der Pol and H. Bremmer, Operational Calculus, Cambridge University Press, 1955.

**A MODEL OF TURBULENCE
DISPLAYING AUTONOMOUS OSCILLATIONS**

R.J. Dickson



5

FOREWORD

This paper was published originally under the same title, as LMSD-48446, 25 March 1959.

The principal result of Section 3 was presented at the annual meeting of the American Mathematical Society, Philadelphia, Penn., 20-22 January 1959. A technical abstract (#553-84) appears in Notices of the A. M. S., Vol. 5, No. 7, Dec 1958, pp. 822-823.

ABSTRACT

This report presents the results of the author's attempt to find and treat simplified examples of problems governed by nonlinear space-time systems. The work was carried out under the General Research Program of Lockheed Missiles and Space Division.

The introduction sketches the general setting of a "Model of Turbulence" in the sense of J. M. Burgers, and directs attention to the character of the particular problem considered here.

Section 2 introduces a special class of nonlinear ordinary differential equations and explains its relationship to a partial differential system studied by Burgers.

Section 3 takes the first step in a detailed study of this class of equations. Specifically, it is proved that a certain fourth-order autonomous system of ordinary equations possesses an unstable equilibrium solution and a stable periodic solution; and that with increasing time all other solutions tend asymptotically to one of these two possibilities.

CONTENTS

<u>Section</u>		<u>Page</u>
	FOREWORD	iii
	ABSTRACT	v
1	INTRODUCTION	1-1
2	A MODIFIED MODEL	2-1
	Introduction	2-1
	Theorem 2.1	2-2
3	THE CASE $n = 3$	3-1
	Introduction	3-1
	Theorem 3.1	3-2
	Theorem 3.2	3-2
	Theorem 3.3	3-5
	Theorem 3.4	3-6
4	REFERENCES	4-1

Section 1
INTRODUCTION

Hydrodynamic turbulence is a variety of parasitic oscillation. Its qualitative features are easily observed in nature and in the laboratory, yet their mathematical description has remained for many years at a primitive level. Almost no progress has been made toward proving - or disproving - that the dynamical (Navier-Stokes) equations of fluid mechanics possess solutions consistent with observed turbulent flows.

To gain insight into the mathematical structure of turbulence, J. M. Burgers initiated in 1939 the study of "models of turbulence" (Refs. 1 and 2). Typically, one of Burgers' models is a mathematical problem which preserves selected features of the hydrodynamical problem yet remains simple enough to admit analytical study; it does not correspond to any real physical problem, but is studied formally as a mathematical problem for its suggestiveness and for possible insight into the nature of its more complex parent.

The relation of Burgers' work to the existing literature on the subject of "isotropic turbulence" may be clarified by the following quotation (Ref. 1, p. 47):

It should be observed that the theories developed by Taylor, von Kármán and others are concerned mainly with a different problem, viz., the character of the turbulence found, e.g., in an air current in the wake of a grid, and its gradual decay. These theories do not consider the side of the problem to which attention has been given here, viz., the development of a dissipative secondary phenomenon, which grows by detracting energy from a given primary phenomenon and in this process gains such an intensity, that finally a balance is obtained between the energy detracted and the energy dissipated. This latter subject is the one which has been traced here through stages of successive complexity . . . and it will be clear that it has a

generality much wider than the field of hydrodynamics. The problem can be compared to that of the development of relaxational oscillations, investigated by van der Pol. The classical example treated by van der Pol, however, refers to a system with a single variable, whereas the present case is characterized by the appearance of a dissipative secondary phenomenon embracing an infinite number of degrees of freedom . . .

Referring to phase flows generated by the Navier-Stokes equations, E. Hopf (Ref. 3) described the motivation for the study of models when he wrote:

How do the solutions which represent the observed turbulent motions fit into the phase picture? The great mathematical difficulties of these important problems are well known and at present the way to a successful attack on them seems hopelessly barred. There is no doubt, however, that many characteristic features of the hydrodynamical phase flow occur in a much larger class of similar problems governed by nonlinear space-time systems. In order to gain insight into the nature of hydrodynamical phase flows we are, at present, forced to find and to treat simplified examples within that class. The study of such models has been originated by J. M. Burgers in a well known memoir . . .

This paper presents results of the author's attempt "to find and to treat simplified examples within that class." More specifically, the attempt was to find a model which compromised the shortcomings of two of Burgers' models. The first of these, which will be called here the "preliminary model," is the following differential boundary value problem for a (real valued) function $u(x,t)$ defined for positive "time," $t > 0$, on the one-dimensional "space," $0 \leq x \leq 1$:

$$u_t = \nu u_{xx} + u - 2uu_x \quad (1.1)$$

$$u(0,t) = u(1,t) = 0, \quad t > 0 \quad (1.2)$$

The subscripts here denote partial derivatives, while ν is a positive constant which Burgers calls the viscosity. In Burgers' "interpretation," the function u is thought of as the turbulent perturbation of an Eulerian velocity field; thus, the obvious solution $u(x,t) = 0$ is referred to as the "laminar" solution. If Eq. (1.1) is multiplied

by u and integrated over $0 \leq x \leq 1$, the contribution of the nonlinear term vanishes because of Eq.(1.2) The contribution of the viscosity term can be integrated by parts, yielding formally the relation

$$\frac{1}{2} \frac{d}{dt} \int_0^1 u^2 dx = -\nu \int_0^1 u_x^2 dx + \int_0^1 u^2 dx \quad (1.3)$$

which Burgers calls the "energy equation" because it reveals that the "energy of turbulence," $\int_0^1 u^2 dx$, tends to be diminished (or dissipated) by the viscosity term and increased by the "driving" term u of Eq. (1.1), which plays a role analogous to the pressure gradient term of the Navier-Stokes equations. (The question of which of these influences predominates is discussed in Ref. 4, where it is shown, in particular, that Eq. (1.1) belongs to a class of equations for which a modified maximum principle holds.)

If ν is sufficiently large, it can be proved that the viscous dissipation is predominant; i. e., all solutions tend to the stationary, laminar solution as $t \rightarrow \infty$. As ν is decreased, however, a positive critical value is reached at which this total stability of the laminar solution is lost. On the basis of hydrodynamical evidence, it would be reasonable to conjecture that this loss of stability marks the entry into a "turbulent regime" in which a time-dependent periodic solution appears and that further decrease of ν will lead to successive branchings at which more and more complicated oscillatory solutions become possible. (For a detailed conjecture on the nature of the branching process and illustrative examples, see Ref. 3.) So far as is known, however, the model under consideration is not sufficiently sophisticated to display such behavior. A branching process does indeed occur, but the "turbulent" solutions which appear are stationary, i. e., time independent. (For details see Refs. 1 and 2.)

In constructing a more complicated turbulence model, Burgers ingeniously excluded the possibility of stationary solutions other than the laminar one. Specifically, he

proposed the system

$$v_t = \nu v_{xx} + v - w + (w^2 - v^2)_x \quad (1.4)$$

$$w_t = \nu w_{xx} + v + w + (2vw)_x$$

for two functions $v(x,t)$, $w(x,t)$ on $0 \leq x \leq 1$, $t \geq 0$ subject to

$$v(0,t) = v(1,t) = w(0,t) = w(1,t) = 0, \quad t \geq 0 \quad (1.5)$$

It is convenient to combine Eqs. (1.4) into one equation for the complex valued function

$$u(x,t) = v(x,t) + iw(x,t)$$

Burgers' model is then

$$u_t = \nu u_{xx} + (1+i)u - 2\bar{u}u_x, \quad \bar{u} = v - iw \quad (1.6)$$

$$u(0,t) = u(1,t) = 0 \quad (1.7)$$

Multiplication of Eq. (1.6) by \bar{u} , integration over $0 \leq x \leq 1$, and integration by parts yield

$$\int_0^1 \bar{u} u_t dx = -\nu \int_0^1 |u_x|^2 dx + (1+i) \int_0^1 |u|^2 dx \quad (1.8)$$

The real and imaginary parts of Eq. (1.8) are

$$\frac{1}{2} \frac{d}{dt} \int_0^1 |u|^2 dx = -\nu \int_0^1 |u_x|^2 dx + \int_0^1 |u|^2 dx \quad (1.9)$$

and

$$\frac{1}{2i} \int_0^1 (\bar{u} u_t - u \bar{u}_t) dx = \int_0^1 |u|^2 dx \quad (1.10)$$

Equation (1.9) is the energy equation, exactly analogous to Eq. (1.3); Eq. (1.10) shows that $u = 0$ is the only solution independent of the time. If the model has any turbulent solutions at all, they are not stationary.

But are there any turbulent solutions? As with the preliminary model, it is found that the laminar solution loses stability below a positive critical value of viscosity. For subcritical ν , there is presumably a turbulent regime. Although Burgers presented heuristic arguments favoring the existence and boundedness of solutions for all time, essentially nothing is known with certainty about the behavior of solutions in the turbulent range of ν .

Whereas the preliminary model was too simple to display time-dependent turbulent solutions, this model is not simple enough to permit analysis. It was an attempt to formulate models of intermediate complexity which generated the results of this paper.

Section 2
A MODIFIED MODEL

INTRODUCTION

Numerical approximations to the solution of a partial differential boundary value problem are ordinarily obtained with digital or electrical analog computers. Such machines will accept the problem only in a modified or analogous form. For a digital computer, the problem must be replaced by one of its approximate "discretized" versions requiring only algebraic operations. For an analog machine, the discretization need not be complete; in a space-time system, for example, it is sufficient to "lump" space alone, leaving time as a continuous variable.

In either case, metamorphosis of the problem requires that functions of a continuous variable be replaced by functions defined on a finite set of points, and that derivatives be replaced by algebraic differences. There is freedom of choice in the number "n" of points, i. e., the "finess" of the discretization, and in the particular difference scheme employed. No criteria for making these choices are dictated by the original differential problem. Mathematical problems of considerable depth arise from the questions of how to make these choices in an optimal way and how to assess the goodness of the resulting approximation.

The more or less arbitrary choice of a particular difference scheme is often justified by the observation that the distinction between formally admissible schemes vanishes as $n \rightarrow \infty$, i. e., in the limit of refinement of the discretization. It may occur, however, that for small values of n certain special schemes are most effective in preserving some kind of analogy between the exact and approximate problems. In discretizing the "space" $0 \leq x \leq 1$ involved in Burgers' model Eq. (1.6), the most obvious and straightforward choice is poor since it does not yield for each n an exact analog of Eq. (1.8).

Specifically, let $u_k(t)$ denote an approximation to $u(k/n, t)$, $k = 1, \dots, n-1$ and set $u_0(t) = u(0, t) \equiv 0$, $u_n(t) = u(1, t) \equiv 0$. Suppose the u_k to be solutions of the following system of ordinary differential equations obtained from Eq. (1.6) by replacing spatial derivatives by central divided differences:

$$\frac{du_k}{dt} \equiv \dot{u}_k = \nu \frac{u_{k+1} - 2u_k + u_{k-1}}{1/n^2} + (1+i)u_k - 2\bar{u}_k \frac{\bar{u}_{k+1} - \bar{u}_{k-1}}{2/n}$$

$$k = 1, \dots, n-1$$

Multiplying by \bar{u}_k and summing on k , it is easily found that the exact analog of Eq. (1.8) does not hold because the contribution from the nonlinear terms does not vanish identically, i. e.,

$$\sum_{k=1}^{n-1} u_k^2 (u_{k+1} - u_{k-1}) \neq 0.$$

In the attempt to find an approximating system without this imperfection, it has been found sufficient to approximate the function \bar{u} in the nonlinear term $\bar{u}\bar{u}_x$ by the average $(\bar{u}_{k+1} + \bar{u}_k + \bar{u}_{k-1})/3$. Properties of the resulting system will be summarized as Theorem 2.1.

THEOREM 2.1:

If

$$\nu > 0, n \geq 2, u_0 = u_n = 0, \quad (2.1)$$

the system

$$\dot{u}_k = \nu n^2 (u_{k+1} - 2u_k + u_{k-1}) + (1+i)u_k - \frac{n}{3} (\bar{u}_{k-1} + \bar{u}_k + \bar{u}_{k+1}) (\bar{u}_{k+1} - \bar{u}_{k-1}) \quad (2.2)$$

for the (complex valued) functions $u_k(t)$, $k = 1, 2, \dots, n-1$ has the following properties:

- Every solution satisfies the relations

$$\begin{aligned} \text{Re} \sum_{k=1}^{n-1} \bar{u}_k \dot{u}_k &\equiv \frac{1}{2} \frac{d}{dt} \sum_{k=1}^{n-1} |u_k|^2 \\ &= -\nu n^2 \sum_{k=1}^n |u_k - u_{k-1}|^2 + \sum_{k=1}^{n-1} |u_k|^2 \end{aligned} \quad (2.3)$$

$$\begin{aligned} \text{Im} \sum_{k=1}^{n-1} \bar{u}_k \dot{u}_k &\equiv \frac{1}{2i} \sum_{k=1}^{n-1} (\bar{u}_k \dot{u}_k - u_k \dot{\bar{u}}_k) \equiv \sum_{k=1}^{n-1} |u_k|^2 \dot{\theta}_k \\ &= \sum_{k=1}^{n-1} |u_k|^2 \end{aligned} \quad (2.4)$$

where

$$u_k = |u_k| \exp(i\theta_k)$$

- The solution taking arbitrarily assigned initial values exists uniquely and is bounded over any finite time interval $0 \leq t \leq T < \infty$.
- $u_k(t) = 0, k = 1, 2, \dots, n-1$ is the only time-independent, i. e., equilibrium, solution. Define $\nu_{mn} = [2n \sin(m\pi/2n)]^{-2}, m = 1, \dots, n-1$. (As $n \rightarrow \infty \nu_{1n} \rightarrow \pi^{-2}, \nu_{n-1n} \rightarrow 0$). If $\nu > \nu_{1n}$, every solution approaches the zero solution as $t \rightarrow \infty$. If $0 < \nu < \nu_{n-1n}$, every nonzero solution approaches infinity, i. e.,

$$\sum_{k=1}^{n-1} |u_k|^2 \rightarrow \infty \quad \text{as } t \rightarrow \infty.$$

Proof.

1. After multiplying Eq. (2.1) by u_k and summing on k , Eqs. (2.3) and (2.4) follow from the

Lemma: If u_k , $k = 0, 1, \dots, n$ are complex numbers with $u_0 = u_n = 0$,

then

$$\sum_{k=1}^{n-1} u_k (u_{k+1} - 2u_k + u_{k-1}) = - \sum_{k=1}^n |u_k - u_{k-1}|^2 \quad (2.5)$$

and

$$\sum_{k=1}^{n-1} u_k (u_{k-1} + u_k + u_{k+1})(u_{k+1} - u_{k-1}) = 0. \quad (2.6)$$

Equation (2.5) follows at once from the Abel partial summation formula, quite analogously to the integration by parts used in obtaining Eq. (1.8). Equation (2.6) follows from the observation that the k^{th} summand can be written as $\alpha_k - \alpha_{k-1} + \beta_k - \beta_{k-1}$, where $\alpha_k = u_k^2 u_{k+1}$, $\beta_k = u_{k+1}^2 u_k$, so that the sum telescopes to $\alpha_{n-1} - \alpha_0 + \beta_{n-1} - \beta_0 = 0$.

2. It is apparent from classical theorems that a unique solution to the initial value problem exists for some interval $0 \leq t < T$. That T may be assumed arbitrarily large follows from the classical extension theorem once it is shown that

$$c(t) = \sum_{k=1}^{n-1} |u_k|^2$$

is bounded over any finite interval. To prove this, observe that the right member of Eq. (2.3) is a hermitian form (call it H) in the variables u_1, \dots, u_{n-1} .

In the sequel, it will appear that the eigenvalues of this form are given explicitly, in decreasing order, by $\lambda_k = 1 - 4\nu n^2 \sin^2(k\pi/2n)$, $k=1, 2, \dots, n-1$. By a well-known characterization of the least and greatest eigenvalues,

$$\lambda_{n-1}c \leq H \leq \lambda_1 c$$

yielding, from Eq. (2.3), the differential inequality

$$\lambda_{n-1} c(t) \leq \frac{1}{2} \dot{c}(t) \leq \lambda_1 c(t)$$

Integration gives

$$c(0) e^{2\lambda_{n-1}t} \leq c(t) \leq c(0) e^{2\lambda_1 t} \quad (2.7)$$

which exhibits the desired bound.

3. The uniqueness of the equilibrium solution is obvious from Eq. (2.4). The remaining assertions are implied by Eq. (2.7) upon noting that $\nu > \nu_{1n}$ implies $\lambda_1 < 0$ and $\nu < \nu_{n-1n}$ implies $\lambda_{n-1} > 0$.

As in the models considered by Burgers, Eq. (2.2) exhibits a "laminar" regime, $\nu > \nu_{1n}$, in which all solutions tend to a totally stable unique laminar solution. For fixed n and small ν , however, the model encounters a deficiency inherent in finite discretization schemes. Of necessity, a function defined on a finite mesh cannot adequately approximate spatial fluctuations on a scale smaller than the mesh spacing; in a model of turbulence like Burgers', however, fluctuations on a finer and finer spatial scale are to be expected with diminishing viscosity. It is, therefore, not surprising that solutions of a model like Eq. (2.2) are found to "blow up" if ν is too small. As a consequence, it is apparent that Eq. (2.2) can display "interesting" behavior only

in some range $\nu' < \nu < \nu_{1n}$, and then Eq. (2.1) gives ν_{n-1n} as a lower bound on ν' .

In the case $n = 2$, the model is essentially trivial, since Eq. (2.2) reduces to a single linear equation and the interval of interest degenerates to a point, i. e., $\nu_{n-1n} = \nu_{1n}$. It is suggestive, however, that for $\nu = \nu_{1n}$, all solutions are periodic.

The simplest nontrivial case, $n = 3$, will be analyzed in Section 3. In the remainder of this section, the effect upon Eq. (2.2) of a particular change of dependent variables will be described.

One of the tools used by Burgers in his study of Eq. (1.8) and similar equations is the expansion of the solution function in Fourier series; the partial differential equation yields formally an infinite system of ordinary differential equations for the determination of the Fourier coefficients as functions of time. At a fixed time, the coefficients may be called a "spectrum" which describes the harmonic content of the spatial distribution of "turbulent" fluctuations; the significance of this spectrum and various questions concerned with its change with time are discussed in Burgers' papers (Refs. 1 and 2).

To obtain an analogous (but finite) spectrum for the discretized model Eq. (2.2), the solution $u_k(t)$ may be expanded (with respect to the spatial variable k) in a finite sine series; in other words, new complex coordinates ξ_k may be introduced in place of the u_k in the (complex, $n - 1$ dimensional) phase space of the system (2.2) by means of the transformation

$$u_k = \sum_{j=1}^{n-1} \xi_j \sin j \frac{k\pi}{n}, \quad k = 1, 2, \dots, n - 1 \quad (2.8)$$

Abbreviating $\sin j \frac{k\pi}{n} = \sigma_{kj}$ and noting the orthogonality relation

$$\sum_{k=1}^{n-1} \sigma_{ki} \sigma_{kj} = \frac{n}{2} \delta_{ij} = \begin{cases} n/2 & i = j \\ 0 & i \neq j \end{cases}$$

the inverse transformation is easily found to be

$$\zeta_i = \frac{2}{n} \sum_{k=1}^{n-1} u_k \sigma_{ki} \quad (2.9)$$

Inserting Eq. (2.8) in Eq. (2.2) and performing algebraic reductions lead to the following systems of equations for the ζ_k

$$\dot{\zeta}_k = (\lambda_k + i)\zeta_k - \frac{n}{3} \sum_{1 \leq i < i+k \leq n-1} \alpha_{i, i+k} \ddot{\zeta}_i \ddot{\zeta}_{i+k} - \frac{n}{6} \sum_{1 \leq i, k-i \leq n-1} \beta_{i, k-i} \ddot{\zeta}_i \ddot{\zeta}_{k-i} \quad (2.10)$$

$$k = 1, 2, \dots, n-1$$

where

$$\alpha_{jk} = \sin j \frac{\pi}{n} (1 + 2 \cos k \frac{\pi}{n}) - \sin k \frac{\pi}{n} (1 + 2 \cos j \frac{\pi}{n})$$

$$\beta_{jk} = \sin j \frac{\pi}{n} (1 + 2 \cos k \frac{\pi}{n}) + \sin k \frac{\pi}{n} (1 + 2 \cos j \frac{\pi}{n})$$

$$\lambda_k = 1 - 4\nu n^2 \sin^2 \frac{k\pi}{2n}$$

i. e., the linear part of the system has been diagonalized by the change of coordinates.

Section 3
THE CASE $n = 3$

INTRODUCTION

If $n = 3$, Eqs. (2.2), (2.8), (2.9), and (2.10) have the following forms:

$$\begin{aligned} u_1 &= (1 - 18\nu + i) u_1 + 9\nu u_2 - (\bar{u}_1 + \bar{u}_2) \bar{u}_2 \\ \dot{u}_2 &= 9\nu u_1 + (1 - 18\nu + i) u_2 + (\bar{u}_1 + \bar{u}_2) \bar{u}_1 \end{aligned} \quad (3.1)$$

$$\begin{aligned} u_1 &= \frac{\sqrt{3}}{2} (\xi_1 + \xi_2) \\ u_2 &= \frac{\sqrt{3}}{2} (\xi_1 - \xi_2) \end{aligned} \quad (3.2)$$

$$\begin{aligned} \xi_1 &= \frac{1}{\sqrt{3}} (u_1 + u_2) \\ \xi_2 &= \frac{1}{\sqrt{3}} (u_1 - u_2) \end{aligned} \quad (3.3)$$

$$\begin{aligned} \dot{\xi}_1 &= (1 - 9\nu + i) \xi_1 + \sqrt{3} \bar{\xi}_1 \bar{\xi}_2 \\ \dot{\xi}_2 &= (1 - 27\nu + i) \xi_2 - \sqrt{3} \bar{\xi}_1^2 \end{aligned} \quad (3.4)$$

The most evident property of the system (3.4) will be stated as Theorem 3.1.

THEOREM 3.1

If $\xi_1(t)$, $\xi_2(t)$ is a solution of Eq.(3.4) in an interval containing a point t_0 such that $\xi_1(t_0) = 0$, then

$$\xi_1(t) \equiv 0 \quad , \quad \xi_2(t) = \xi_2(t_0) \exp [(1 - 27\nu + i)(t - t_0)] \quad (3.5)$$

is the unique extension of the given solution to $-\infty < t < +\infty$.

This theorem implies that the phase point sets $\xi_1 = 0$ and $\xi_1 \neq 0$ are invariant manifolds; solution trajectories never pass from one to the other. Solutions of the form Eq. (3.5) will be subsequently referred to as trivial. For $\nu > 1/27$, it is apparent that all trivial solutions approach the zero solution.

Because of the complex notation, Eq.(3.5) is a convenient condensation for a system of four first-order differential equations for four real-valued functions. For discussion of the nontrivial solutions, real coordinates will be introduced which permit the reduction of the initial-value problem for Eq. (3.4) to the initial-value problem for a system of three real differential equations plus a quadrature.

THEOREM 3.2

The solution of Eq. (3.4) taking initial values $\xi_1(0) \neq 0$, $\xi_2(0)$ is given by

$$\xi_1(t) = \frac{\xi_1(0)}{\sqrt{3} |\xi_1(0)|} \sqrt{z(t)} \exp \left\{ i \int_0^t [1 - y(s)] ds \right\} \quad (3.6)$$

$$\xi_2(t) = \frac{\overline{\xi_1(0)}}{\sqrt{3} \xi_1(0)} [x(t) + i y(t)] \exp \left\{ -2i \int_0^t [1 - y(s)] ds \right\}$$

in which $x(t)$, $y(t)$, $z(t)$ is the solution of the system

$$\dot{x} = (1 - 27\nu)x - 3y - z + 2y^2 \quad (3.7a)$$

$$\dot{y} = 3x + (1 - 27\nu)y - 2xy \quad (3.7b)$$

$$\dot{z} = 2(1 - 9\nu)z + 2xz \quad (3.7c)$$

satisfying the initial conditions dictated by the special case $t = 0$ of the transformation relations

$$x(t) + iy(t) = \sqrt{3} \frac{\xi_1(t) \xi_2(t)}{\xi_1(t)}, \quad z(t) = 3 |\xi_1(t)|^2 \quad (3.8)$$

Proof. The proof consists of straightforward verification and will be omitted. In its place the formal processes leading to Eqs. (3.7) will be sketched. ("Formal" here refers specifically to the fact that possible discontinuities of the function $\varphi_2(t)$ defined below are ignored.)

Into Eqs. (3.4) introduce real polar coordinates by

$$\sqrt{3} \xi_k = \rho_k \exp(i \varphi_k), \quad k = 1, 2 \quad (3.9)$$

The following differential equations are found for ρ_k , φ_k :

$$\begin{aligned} \dot{\rho}_1 &= (1 - 9\nu) \rho_1 + \rho_1 \rho_2 \cos(2\varphi_1 + \varphi_2) \\ \dot{\rho}_2 &= (1 - 27\nu) \rho_2 - \rho_1^2 \cos(2\varphi_1 + \varphi_2) \end{aligned} \quad (3.10)$$

$$\begin{aligned}\rho_1 \dot{\varphi}_1 &= \rho_1 - \rho_1 \rho_2 \sin(2\varphi_1 + \varphi_2) \\ \rho_2 \dot{\varphi}_2 &= \rho_2 + \rho_1^2 \sin(2\varphi_1 + \varphi_2)\end{aligned}\tag{3.11}$$

By defining a new variable $\psi = 2\varphi_1 + \varphi_2$, Eqs. (3.11) yield formally

$$\dot{\psi} = 3 - \left(2\rho_2 - \frac{\rho_1^2}{\rho_2}\right) \sin \psi\tag{3.12}$$

Now Eqs. (3.10) and (3.12) constitute a third-order system for ρ_1, ρ_2, ψ as functions of t . Given a solution of this system the functions φ_1, φ_2 can be found by quadrature from Eqs. (3.11).

The system of Eqs. (3.7) now results from transformation of the system (3.10) and (3.12) according to

$$\begin{aligned}x &= \rho_2 \cos \psi \\ y &= \rho_2 \sin \psi \\ z &= \rho_1^2\end{aligned}$$

The combined transformation from Eqs. (3.4) to (3.7) is given by Eq. (3.8). The first of Eqs. (3.11) yields

$$\varphi_1(t) = \varphi_1(0) + \int_0^t [1 - y(s)] ds$$

or

$$\frac{\xi_1(t)}{|\xi_1(t)|} = \frac{\xi_1(0)}{|\xi_1(0)|} \exp \int_0^t [1 - y(s)] ds$$

The transformation employed can be interpreted so as to yield also a correspondence between the trivial solutions of Eq. (3.4) and solutions of Eqs. (3.7) with $z(t) = 0$; simplicity is lost in the process, however, since the restriction of Eq. (3.4) to $\xi_1 = 0$ leaves a linear system while the restriction of Eqs. (3.7) to $z = 0$ is non-linear.

It is easily seen that in the (invariant) half-space $z > 0$ the system (3.7) can have at most one critical point. In fact, setting the right-hand members equal to zero, Eq. (3.7) implies

$$x = x_c = 9\nu - 1$$

Eq. (3.7b) then requires

$$y = y_c = \frac{1 - 9\nu}{1 - 15\nu}$$

and Eq. (3.7a) requires

$$z = z_c = (1 - 9\nu)(27\nu - 1) [1 + (1 - 15\nu)^{-2}]$$

This point (x_c, y_c, z_c) belongs to the half-space $z > 0$ if ν lies in either of the intervals $1/27 < \nu < 1/15$ or $1/15 < \nu < 1/9$. That the first interval represents values of ν too small to yield "interesting" results in the case $n = 3$ is suggested by the following negative result.

THEOREM 3.3

If $1/27 < \nu < 1/15$, every solution of Eqs. (3.7) which intersects the region $z > 0$, $y < 3/2$ is unbounded as $t \rightarrow +\infty$.

Proof. Using Eqs.(3.7), it is readily verified that for any solution

$$\frac{d}{dt} \left[z \left(y - \frac{3}{2} \right) e^{-3(1-15\nu)t} \right] = \frac{3}{2} (1-27\nu) z e^{-3(1-15\nu)t}$$

For a nontrivial solution ($z(t) > 0$) and $\nu > 1/27$, this implies

$$\frac{d}{dt} \left[z \left(y - \frac{3}{2} \right) e^{-3(1-15\nu)t} \right] < 0$$

Hence for $t_1 < t_2$

$$z(t_2) \left(y(t_2) - \frac{3}{2} \right) < z(t_1) \left(y(t_1) - \frac{3}{2} \right) e^{-3(1-15\nu)(t_1-t_2)}$$

If $y(t_1) - 3/2 < 0$, $\nu < 1/15$, and $z(t)$ is bounded, this inequality implies $y(t_2) \rightarrow -\infty$ as $t_2 \rightarrow \infty$. Under the hypothesis of the theorem, therefore, at least one of the functions $z(t)$, $y(t)$ is not bounded.

Against this negative result for $1/27 < \nu < 1/15$, the theorem following reveals interesting behavior for solutions of Eqs.(3.7) for at least a part of the interval $1/15 < \nu < 1/9$.

THEOREM 3.4

There is a number ν' , $1/15 \leq \nu' < 0.0836 < 1/9$ such that for $\nu' < \nu < 1/9$ every nontrivial solution of the system (3.7) approaches the equilibrium solution (x_c, y_c, z_c) as $t \rightarrow \infty$: Eqs. (3.6) then imply that every nontrivial solution of Eqs. (3.4) approaches a periodic solution.

Proof. The proof depends on the following properties of the function

$$V(x, y, z) = (x - x_c)^2 + (y - y_c)^2 + (z - z_c) - z_c \log \frac{z}{z_c} \quad (3.13)$$

defined in the half-space $z > 0$, $-\infty < x, y < +\infty$.

- $V \geq 0$ $V = 0$ if and only if $x = x_c$, $y = y_c$, $z = z_c$.
- $V \rightarrow \infty$ if either $z \rightarrow 0$ or $x^2 + y^2 \rightarrow \infty$ in $z > 0$.
- If $0 \leq \alpha < \infty$, the set $E(\alpha) = \{ (x, y, z) \mid V(x, y, z) \leq \alpha \}$ is a closed and bounded subset of $z > 0$ and $E(\alpha) \subset E(\beta)$ if $\alpha < \beta$.
- Along a solution curve of the system (3.7).

$$-\frac{1}{2} \frac{dV}{dt} = (2 + 3x_c)(x - x_c)^2 - 2y_c(x - x_c)(y - y_c) + (2 + 5x_c)(y - y_c)^2 \quad (3.14)$$

The right-hand member of this relation is a positive definite quadratic form in the variables $(x - x_c)$, $(y - y_c)$ provided that $\nu' < \nu < 1/9$, where $\nu' = 0.0836$ is the largest real root of the discriminant (which is a rational function of ν). The only implication of Eq.(3.14) needed in the proof is the inequality

$$-\frac{1}{2} \frac{dV}{dt} \leq B \left[(x - x_c)^2 + (y - y_c)^2 \right] \quad (3.15)$$

where

$$B = 2(1 - 2x_c) - \sqrt{x_c^2 + y_c^2}$$

is the positive smallest eigenvalue of the form.

These properties are either obvious or can be proved by direct calculation.

If $x(t)$, $y(t)$, $z(t)$ is the solution of Eqs. (3.7) with initial values $x_0, y_0, z_0 > 0$ and if $\alpha_0 = V(x_0, y_0, z_0)$, the phase point $(x(t), y(t), z(t))$ belongs to $E(\alpha_0)$ for $t \geq 0$ since $V = V[x(t), y(t), z(t)]$ is a nonincreasing function of t by Eq. (3.15). Thus, the phase point is bounded; from Eqs. (3.7), it follows that the phase-point velocity $(\dot{x}^2 + \dot{y}^2 + \dot{z}^2)^{1/2}$ is also uniformly bounded for $t \geq 0$.

It will be shown first that $x(t) \rightarrow x_c$ and $y(t) \rightarrow y_c$ as $t \rightarrow \infty$, or, equivalently, that no point of the limit set of a solution lies off the half-line $z > 0, (x - x_c)^2 + (y - y_c)^2 = 0$. Assuming the contrary, there is point P of the limit set, a neighborhood, $N(P)$, of P , and an $\epsilon > 0$ such that $N(P)$ does not intersect the cylinder $(x - x_c)^2 + (y - y_c)^2 \leq \epsilon$. During the time intervals that the phase point lies in $N(P)$, Eq. (3.15) implies $dV/dt \leq -2B\epsilon$. Since V is nonnegative and nonincreasing, this inequality implies that the time which a phase point spends in $N(P)$ is bounded. The boundedness of the phase point velocity and the assumption that P belongs to the limit set imply, however, that the time spent in $N(P)$ is unbounded. The contradiction proves that the limit set lies on the half-line and hence that $x(t) \rightarrow x_c$ and $y(t) \rightarrow y_c$.

It remains to prove that $\lim z(t)$ exists and equals z_c . Along a trajectory, V decreases and thus has a limit. From Eq. (3.13) and the existence of limits for V , x , and y , there follows the existence of $\lim (z - z_0 \log z)$ and, hence, that of $\lim z$. From Eq. (3.7a), the existence of $\lim \dot{x}$ is then immediate. But if $\lim x$ and $\lim \dot{x}$ exist, it follows at once that $\lim \dot{x} = 0$. Hence, the right-hand member of Eq. (3.7a) converges to zero as $t \rightarrow \infty$, which implies that $\lim z = z_c$ and concludes the proof.

Section 4
REFERENCES

1. J. M. Burgers, "Mathematical Examples Illustrating Relations Occurring in the Theory of Turbulent Fluid Motion," Verhandel. Kon. Nederl. Akad. Wetenschappen Amsterdam, Afdeel. Natuurkunde (1st Sect.), Deel 17, No. 2, 1939 pp. 1-53
2. J. M. Burgers, "A Mathematical Model Illustrating the Theory of Turbulence," Advances in Applied Mechanics, Vol. 1, New York, Academic Press, 1948, pp. 171-199
3. E. Hopf, "A Mathematical Example Displaying Features of Turbulence," Comm. on App. Math., Vol. 1, No. 4, 1948, pp. 303-322
4. R. J. Dickson, "Bounds for Solutions of Some Parabolic Problems," Unpublished Ph.D. Dissertation, California Institute of Technology, 1954