NUMERICAL ALGORITHMS FOR CONTROLLABILITY AND EIGENVALUE ALLOCATION



George S. Miminis
School of Computer Science
McGill University
May 1981

bу

A thesis submited to the Faculty of Graduate Studies and Research, in partial fulfillment of the requirements for the degree of Master of Science.

Algorithms for Controllability and Eigenvalue Allocation

Contents

Abstract

Resumé

Acknowledgements

CHAPTER 1

- (I) · Introduction
 - (1) Overview
 - (2) Notation

CHAPTER -2

Jo J

- (I) Introduction
 - (1) Dynamic/systems
 - (2) Controllability of dynamic systems
- (II) The algebraic problem of controllability
- (III) The computational problem of controllability
- (IV) Distance from an uncontrollable system

CHAPTER 3

- (I) Introduction
 - (1) Open-loop (Nonfeedback) and closed-loop (Feedback) control systems
 - (2) Stability problem in dynamical systems
 - (a) Stability for discrete-time systems
 - (b) stability for continuous-time systems
- (II) On Pole assignment in single-input controllable linear systems
 - (1) The case of real eigenvalues
 - (a) Explicitly shifted method
 - (b) Implicitly shifted method
 - (2) The case of complex eigenvalues

CHAPTER 4

(I) Conclusion

APPENDIX

REFERENCES

Abstract

As part of this thesis we are concerned with the numerical properties of some methods for computing controllability of a linear dynamic system with constant coefficients. We show, using counterexamples, why some well known methods for determining the controllability of such a system are often numerically unstable or computationally inefficient. We also show that recently developed algorithms by C. Paige [2] and others [17], [18], [19], work satisfactorily for the same examples.

We also give a different algorithm which determines whether a system is controllable or not. The same algorithm with some changes can be used to compute the distance of the given system from the nearest uncontrollable one, but it is pointed out why this algorithm is computationally inefficient.

We also present an algorithm for solving the problem of eigenvalue allocation for a single-input controllable, linear dynamic system with constant coefficients. This algorithm is based on numerically stable transformations.

These latter two algorithms have been developed jointly with C. Paige.

Resumé

Dans cette thèse, nous examinons les propriétés numériques de certaines méthodes pour le calcul de la commandabilité d'un système linéaire dynamique avec coefficients constants. Nous montrons, utilisant des contre-exemples, pourquoi certaines méthodes très connues pour déterminer la commandabilité sont souvent instables numériquement ou inefficaces à calculer. Nous montrons aussi que des algorithmes recemment developpés par C. Paige [2] et autres [17], [18], [19], sont satisfaisant pour les mêmes exemples.

Nous indiquons aussi un algorithme différent qui détermine si un système est commandable ou non. Ce même algorithme avec quelques changements peut être utilisé pour calculer la distance du système donné au système non commandable le plus rapproché, mais on note pourquoi cet algorithme est inefficace.

De même, nous indiquons un algorithme pour résoudre le problème d'allocation des valeurs propres pour un système linéaire dynamique commandable monovariable avec coefficients constants. Cet algorithme est basé sur des transformations numériquement stables.

Ces deux derniers algorithmes ont été developés conjointement avec .

C. Paige.

Acknowledgements

I would like to thank Professor C. Paige, not only for his invaluable supervision of the work reported here but also for all that I have learned from him about computational linear algebra. It is always a pleasure to work with him and I am grateful to him for his guidance and encouragement during the entire period of the work.

I also thank Diane Chan for typing this thesis.

(I) Introduction

(1) Overview

C. Paige [2] points out that a measure of the distance of a given dynamic system from the nearest uncontrollable one would be quite useful because it would not only determine whether a system is controllable or not, but it would also indicate, in the case of a controllable system, if the given system is quite controllable or it is nearly uncontrollable. In the same paper it is mentioned that the above problem appears to be an open problem.

This problem and the fact that apparently there are no numerically reliable and efficient algorithms for the problem of eigenvalue allocation, gave us the impetus to write this thesis. The outline of the thesis is as follows.

In the second section of this chapter we establish our notational conventions. In Chapter 2 we first prove three widely used mathematically equivalent theorems. We also prove another theorem; suggested by C. Paige in [2], using a reduction of a matrix to a block upper Hessemberg matrix, also suggested by C. Paige in [2] as a very useful tool. Then we show, using counterexamples, why three well known methods often fail to determine whether a system is controllable or not, while the new algorithm in [2] succeeds in determining the controllability of the systems in the same examples. Finally in the same chapter we give an algorithm which uses the notion of the distance of a system from the nearest uncontrollable one to determine the controllability of the system. As we men-

tioned in the abstract, the same algorithm with some changes can be used for the computation of the distance of a system from the nearest uncontrollable one when the system is controllable, but it is pointed out why this algorithm is not computationally efficient.

In Chapter 3 we present some new algorithms for the eigenvalue allocation problem. First we describe how we can assign a certain desired set of real engenvalues to a real upper Hessenberg matrix by determining its first row, using an explicit method related to the QR algorithm. Then we do the same thing but using an implicit version of the method. A good description of the QR algorithm is given in [12]. Finally we describe how to assign a certain desired set of eigenvalues which occur in complex conjugate pairs to a real upper Hessenberg matrix by determining its first row, using an implicitly shifted method.

In Chapter 4 we comment on the methods which are developed in this thesis.

In the Appendix we present the subroutine EVA which performs the eigenvalue allocation described in Chapter 3. We also present four examples, which indicate that the methods described in Chapter 3 are very efficient and reliable.

(2) Notation

a) Matrices will be represented by upper case Roman letters. Vectors
will be represented by lower case Roman. Scalars will be represented by lower case Greek letters, and indices by lower case Roman.

- b) A^T will represent the transpose of the matrix A, while A^H will represent the complex conjugate transpose of A.
- c) R represents the set of real numbers, C the set of complex numbers and N the set of natural numbers.
- d) $\forall \alpha \in A$ means, for every α element of the set A.
- e) e with k & N represents a column vector, the elements of which are all zero except the kth one which is unity, while e will represent a column vector with all its elements equal to one.
- f) $A = (\alpha_{ij})$ means that the matrix A has elements α_{ij} .
- g) i(j)h will represent the arithmetical progression with first element i, step j and final element h. j can be either positive or negative, so i can be less than or greater than h.
- h) $A \cup B$ means the union of the sets A and B_1 .
- i) Δ_n , where $n \in \mathbb{N}$ will represent the set of integers in the interval '[1, n], that is, $i \in \Delta_n$ if and only if $i \in \mathbb{N}$ and $1 \le i \le n$.
- j) Let $x = \begin{cases} \xi_1 \\ \xi_2 \end{cases}$ be an n-vector. The numbers $\xi_1, \xi_2, \dots, \xi_n$ are

called the components or the elements of x.

- k) Rⁿ represents the set of all n-vectors with real components. Similarly Cⁿ represents the set of all n-vectors with complex components.
- 1) $\mathbb{R}^{m \times n}$ (respectively $\mathbb{C}^{m \times n}$) represents the set of all matrices of dimensions m, n and real (respectively complex) elements.

- m) Let $x = (\xi_1, \xi_2, \dots, \xi_n)^T$, then the 2-norm of x is represented by $||x||_2$ and defined by $||x||_2 = \sqrt{\sum_{i=1}^n |\xi_i|^2}$.
- Let A be an m×n matrix with elements (α_{ij}) , then the 2-norm (respectively the F-norm) of A is represented by $||A||_2$ (respectively $||A||_F$) and is defined by $||A||_2 = \max_{\substack{||A||_2 \neq 0 \\ ||x||_2 \neq 0}} \frac{||Ax||_2}{||x||_2 \neq 0}$ (respectively $||A||_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |\alpha_{ij}|^2}$).
- p) A subspace R is called trivial if and only if $R = \{0\}$.
- q) A permutation matrix is the unity matrix with rearranged columns or equivalently, rows.

e.g.
$$P = \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{array} \right)$$

r) Let $A \in \mathbb{R}^{m \times n}$ then R(A) will represent the subspace spanned by the columns of A.

CHAPTER 2

(I) Introduction

(1) Dynamic systems

A system is called dynamic when it changes states with respect to time and the characteristics of the state at one time are related with those at other times. The dynamical systems considered here can be represented mathematically by either difference or differential equations. The choice of difference or differential equation corresponds to whether the system is observed in discrete or continuous time respectively.

(2) Controllability of dynamic systems

The concept of controllability of a dynamical system is introduced by the following question:

can a dynamical system be transferred from any given initial state to any desired state in finite time by some control action?>>

If we answer <<yes>> then the system is controllable.

(II) The algebraic problem of controllability

Consider the continuous-time system

 $\dot{x}(t) = Ax(t) + Bu(t)$

(1)

where $t \in \mathbb{R}$ $(t \ge 0)$ is the time; $x \in \mathbb{R}^n$ is the state of the system; $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$ and $u \in \mathbb{R}^m$ is the control vector.

Note: Throughout this thesis n will represent the order of matrix A and m the number of columns of B, unless it is mentioned that they represent something else.

The set of all $\,x\,$ is a real n-dimensional vector space, called the state space of the system (1) and it will be denoted by $\,R\,$.

<u>Definition (I)</u>: The system (1) is said to be controllable if for any given pair of states $(x_I, x_F) \in \mathbb{R}^n \times \mathbb{R}^n$ there exists a time t > 0 and a control u defined on [0,t] such that the solution of (1) which corresponds to the initial value $x(0) \equiv x_I$ gives us $x(t) \equiv x_F$.

The state space R can be decomposed into two parts: the controllable subspace R_1 and the uncontrollable subspace R_2 which is the orthogonal complement of R_1 [16]. So, R is generated by $R_1 \cup R_2$ and we write

$$R = R_1 \oplus R_2 \tag{1a}$$

From the definition (I) we can observe the following: if the system (1) is controllable then $R_2 = \{0\}$. So $R = R_1$. If (1) is uncontrollable then $R_2 \neq \{0\}$ but R_1 need not be trivial.

We shall now define the Discrete-time system

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k$$
 , $k \in \mathbb{N}$ (2)

where the symbols have meaning similar to the previously defined continuous-

Definition (II): The system (2) is called controllable if for any given pair of states $(x_I^-, x_F^-) \in \mathbb{R}^n \times \mathbb{R}^n$ there exists a positive integer q and a sequence of controls u_1, u_2, \ldots, u_q^- such that using (2) we can transfer state $x_1^- \equiv x_1^-$ to state $x_{q+1}^- \equiv x_F^-$ in exactly q steps.

Definition (III): Define

$$\widetilde{\mu} = \min_{\delta A, \delta B} ||_{\zeta}, \zeta \in \{2, F\}, \delta A \in \mathbb{R}^{n \times n}, \delta B \in \mathbb{R}^{n \times m}$$
(3)

such that the systems

$$\dot{\hat{\mathbf{x}}} = (\mathbf{A} + \delta \mathbf{A})\mathbf{x} + (\mathbf{B} + \delta \mathbf{B})\mathbf{u} \tag{4}$$

and

$$\mathbf{x}_{k+1} = (\mathbf{A} + \delta \mathbf{A})\mathbf{x}_{k} + (\mathbf{B} + \delta \mathbf{B})\mathbf{u}_{k}$$
 (4\alpha)

are uncontrollable, [2], (another approach is used by Moore in [15], [16]).

- Remark (I): μ is a useful measure of the distance of the system (1) (respectively (2)) from the nearest uncontrollable system (4) (respectively (4 α)).
- Remark (II): Since the controllability of a problem depends only on the matrices A and B (as we will see in the following), we conclude that the distance between the systems (1) and (4) and the distance between the systems (2) and (40) are identical.
- Remark (III): If $\mu = 0$ then (1) and (2) are uncontrollable, otherwise $\mu > 0$ and systems (1) and (2) are controllable.

8

Note: As remark (II) indicates and as we will see in the following the conditions on A and B for controllability are the same for the discrete and continuous system. So in the following if we want to say that system (1) or (2) is controllable or not we will just say that the pair (A,B) is controllable or not. In general, if we want to derive a result which can be applied to both (1) and (2) we can work with the pair (A,B).

Proposition (I): Applying orthogonal transformations to (1) leaves μ unchanged.

Proof: After applying the orthogonal transformations to the system (1) .

we get:

$$Q^{T}\dot{x} = (Q^{T}AQ)(Q^{T}x) + (Q^{T}BP)(P^{T}u)$$
(48)

where Q \in R^{n×n}, P \in R^{m×m} are orthogonal matrices. Then from definition (III) of μ let

$$\widetilde{\mu} = \min \left| \left| \left(Q^{T} \delta A Q, Q^{T} \delta B P \right) \right| \right|_{\zeta}, \quad \zeta \in \{2, F\}$$
 (5)

such that the system

$$Q^{T}\dot{x} = [Q^{T}(A + \delta A)Q](Q^{T}x) + [Q^{T}(B + \delta B)P](P^{T}u)$$

is uncontrollable. Now we can see that

(5) =>
$$\tilde{\mu}$$
 = min $\left| \left| Q^{T}(\delta A, \delta B) \left(\begin{array}{cc} Q & O \\ O & P \end{array} \right) \right| \right|_{\zeta}$
= min $\left| \left| \left(\delta A, \delta B \right) \right| \right|_{\zeta}$

• µ

where μ is given in (3).

Corollary (I): The system (1) is controllable if and only if the system

$$\widetilde{\widetilde{\mathbf{x}}} = \widetilde{\mathbf{A}}\widetilde{\mathbf{x}} + \widetilde{\mathbf{B}}\widetilde{\mathbf{u}} \tag{5a}$$

is controllable, where $\tilde{x} = Q^T x$, $\tilde{u} = P^T u$, $\tilde{A} = Q^T A Q$, $\tilde{B} = Q^T B P$ with $Q \in \mathbb{R}^{n \times n}$ and $P \in \mathbb{R}^{m \times m}$ orthogonal matrices. In other words orthogonal transformations do not alter the controllability of (1).

Proof: Obviously (5α) is the same as (4β). So from proposition (I) since the distance of (5α) from the nearest uncontrollable system is the same as the distance of (1) from the nearest uncontrollable system, this implies that corollary (I) is true.

Note: The arguments of proposition (I) and corollary (I) can also be applied for discrete-time systems.

Before We mention theorems about controllability, we will transform the pair (A,B) by orthogonal transformations to another simpler pair $(\widetilde{A},\widetilde{B})$. This transformation will be very useful in the following when we will use it to prove theorems about controllability.

It can be shown [2] that (A,B) can be transformed, using orthogonal transformations, to $(\widetilde{A},\widetilde{B})$ where for some positive integer $k \le n$ and $n_0 = m$

where $A_{ij} \in \mathbb{R}^{n_1 \times n_j}$, rank $(A_{i,i-1}) = n_i$ for i = l(1)k-1 and $A_{k,k-1}$ is either zero or has rank n_k . By corollary (I) (A,B) is controllable if and only if $(\widetilde{A},\widetilde{B})$ is controllable. Thus instead of looking for the controllability of (A,B) we will look for the controllability of $(\widetilde{A},\widetilde{B})$.

Matrix transformations like (6) form some of the tools of the numerical Analyst's trade, and related algorithms are used by Wilkinson [20], [21] and Van Dooren [22] in analyzing the generalized eigenvalue problem, and by Van Dooren [23] in analyzing general linear systems.

There are four mathematically equivalent theorems on controllability for the pair $(\widetilde{A},\widetilde{B})$:

Theorem (I) [2] If $A_{k,k-1} = 0$ (in (6)) then $(\widetilde{A},\widetilde{B})$ is uncontrollable, if $\operatorname{rank}(A_{k,k-1}) = n_k$ then $(\widetilde{A},\widetilde{B})$ is controllable.

Theorem (II) [4], [5] The pair $(\widetilde{A}, \widetilde{B})$ is controllable if and only if rank $(\widetilde{B}, \widetilde{AB}, \ldots, \widetilde{A}^{n-1}\widetilde{B}) = n$

Definition (IV): Let $A \in \mathbb{R}^{n \times n}$ and λ_i , i = 1(1)n be its eigenvalues, then we define $\lambda(A)$ to be the set which contains the eigenvalues of A only, $(\lambda(A) = \{\lambda_i | i=1(1)n\})$.

Theorem (III) [1], [6] The pair $(\widetilde{A}, \widetilde{B})$ is controllable if and only if

 $rank[\widetilde{B}, \widetilde{A} - \lambda_{i}] = n$ for i = 1(1)n

where λ_i , i = 1(1)n are the eigenvalues of A or \widetilde{A} (since $\lambda(A) = \lambda(\widetilde{A})$, [12]).

there exists at least one matrix $F \in \mathbb{R}^{m \times n}$ such that

$$\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) = \emptyset$$

<u>Proof</u>: We show the equivalence of, and then prove, these four theorems

- (I) If $A_{k,k-1} = 0$ then
 - (1) $\operatorname{rank}[\widetilde{B},\widetilde{AB},\ldots,\widetilde{A}^{n-1}\widetilde{B}] < n$
 - (2) $\operatorname{rank}[\widetilde{B},\widetilde{A}-\lambda_{\underline{1}}I] < n$ for at least one integer i with $1 \le i \le n$.
 - (3) For every $F \in \mathbb{R}^{m \times n}$, $\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) \neq \emptyset$.
- (II) If $rank(A_{k,k-1}) = n_k$ then
 - (1) $\operatorname{rank}[\widetilde{B}, \widetilde{AB}, \dots, \widetilde{A}^{n-1}\widetilde{B}] = n$
 - (2) $\operatorname{rank}[\widetilde{B}, \widetilde{A} \lambda_i I] = n \text{ for } i = 1(1)n$
 - (3) There exists at least one $F \in \mathbb{R}^{m \times n}$ such that $\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) = \emptyset$.

and then proving

(III) (1) The system (1) is controllable if and only if

$$rank[B,AB,...,A^{n-1}B] = n$$

(2) The system (2) is controllable if and only if

$$rank[B,AB,...,A^{n-1}B] = n$$

- (I) Let $A_{k,k-1} = 0$ then
 - (1) We want to prove that $\operatorname{rank}[\widetilde{B},\widetilde{AB},\ldots,\widetilde{A}^{n-1}\widetilde{B}] < n$. Since $A_{k,k-1} = 0$ the last n_k rows of the matrix \widetilde{A}^k , $\ell = 0(1)n-1$ are going to be $(0,0,\ldots,0,A_{k,k}^k)$. So the last n_k rows of the matrix $\widetilde{A}^k\widetilde{B}$, $\ell = 0(1)n-1$ are zero. Thus for every non-zero vector z, with $z \in \mathbb{R}^n$ there exists a non-zero vector $y^T = (0,0,\ldots,0,z^T)$, $y \in \mathbb{R}^n$ such that

$$y^{T}[\widetilde{B},\widetilde{AB},...,\widetilde{A}^{n-1}\widetilde{B}] = 0 \Rightarrow rank[\widetilde{B},\widetilde{AB},...,\widetilde{A}^{n-1}\widetilde{B}] < n$$
.

(2) We want to prove that for some integer i with $1 \le i \le n$ and $\lambda_i \in \lambda(\widetilde{A})$, rank $[\widetilde{B}, \widetilde{A} - \lambda_i I] < n$. So we have to find a vector $y \in \mathbb{R}^n$ which satisfies the following conditions:

$$y^{T}\widetilde{B} = 0$$

$$y^{T}(\widetilde{A} - \lambda_{4}I) = 0$$
(7)

If we choose λ_i to be a solution of the eigenproblem $z^T A_{kk} = \lambda z^T$ and $z \in \mathbb{R}^n$ be a left eigenvector of A_{kk} corresponding to λ_i then the vector $y^T = (0,0,\ldots,0,z^T)$ satisfies all the conditions (7).

(3) We want to prove that for every $F \in \mathbb{R}^{m \times n}$, $\lambda(A) \cap \lambda(\widetilde{A} + \widetilde{B}F) \neq \emptyset$. Let $F \in \mathbb{R}^{m \times n}$ be a matrix such that $F = (F_1, F_2, \dots, F_k)$ where $F_i \in \mathbb{R}$ for i = 1(1)k. Then

If D is the matrix which is derived from $\tilde{A} + \tilde{B}F$ by omitting its last k rows and its, last k columns, then

$$\det[(\widetilde{A} + \widetilde{BF}) - \lambda I_n] = \det(A_{kk} - \lambda I_{n_k}) \det(D - \lambda I_{n-n_k})$$

So $\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) \supseteq \lambda(A_{kk})$ for every $F \in \mathbb{R}^{m \times n}$

(II) Let rank $(A_{k,k-1}) = n_k$

We want to prove that $rank[\widetilde{B}, \widetilde{AB}, ..., \widetilde{A}^{n-1}\widetilde{B}] = n$. We know that (1)

we want to prove that
$$\operatorname{rank}[B,AB,\ldots,A^{-1}B] = n$$
. We know that
$$\widetilde{B} = \begin{pmatrix} A_{10} \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix}, \widetilde{AB} = \begin{pmatrix} A_{11}A_{10} \\ A_{21}A_{10} \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix}, \ldots, \widetilde{A}^{k-1}\widetilde{B} = \begin{pmatrix} C_1 \\ C_2 \\ \vdots \\ C_{k-1} \\ A_k, k-1, k-2, \ldots, A_{10} \end{pmatrix}$$

where C_i , i = 1(1)k-1 are appropriate $n_i \times n_0$ matrices. Besides we know that $rank(A_{i,i-1}) = n_i$ for i = 1(1)k. That is, the matrices $A_{1,i-1}$, i = 1(1)k are all of full row rank. So the matrices $A_{i,i-1} A_{i-1,i+1} ... A_{10}$, i = 1(1)k are of full row rank. Thus there is no non-zero left null vector for the matrix $[\widetilde{B}, \widetilde{AB}, \ldots, \widetilde{A}^{n-1}\widetilde{B}]$.

- We want to prove that $\operatorname{rank}[\widetilde{B},\widetilde{A}-\lambda_{1}I]=n$ for every i with i=1(1)n and $\lambda_{1}\in\lambda(\widetilde{A})$ for i=1(1)n. It is clear that, since $A_{1,i-1}$ is of full row rank for i=1(1)k then $[\widetilde{B},\widetilde{A}-\lambda_{1}I]$ has full row rank for every λ_{1} .
- (3) We want to prove that there is at least one $F^{(1)} \in \mathbb{R}^{m \times n}$ such that $\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) = \emptyset$. We know that

and the matrices $A_{j,j-1}$ are of full row rank, that is, rank $(A_{j,j-1}) = n_j$, j = 1(1)k.

We are going to give two proofs, the first will be for the simple case $m=n_0=1$, that is, when \widetilde{B} is a vector and the second will be for the general case.

First Proof: $m = n_0 = 1$, so let $\tilde{B} = b = (\alpha_{10}, 0, ..., 0)^T$ and

$$\widetilde{A} = \begin{cases} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1,n-1} & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & \cdots & \alpha_{2,n-1} & \alpha_{2n} \\ & & & & & & & \\ & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & & \\ & & & \\ & & & & \\ & & & & \\ & &$$

Let λ_i , $i=1(1)^\ell$ with $\ell \leq n$, be the ℓ distinct eigenvalues of \widetilde{A} , and x_i , $i=1(1)^\ell$ the eigenvectors of \widetilde{A} corresponding to λ_i . Let also $F\equiv f^T$, $f\in \mathbb{R}^n$. Then, if for some λ_i , $(\widetilde{A}+bf^T)-\lambda_i I$ is singular there exists $y_f\in \mathbb{R}^n-\{0\}$, where $y_f^T=(\psi_{1f},y_{2f}^T)$ and $\psi_{1f}\in \mathbb{R}$, such that

$$y_{f}^{T}[(\tilde{A} + bf^{T}) - \lambda_{1}I] = 0 \implies (8)$$

$$y_{f}^{T}bf^{T} + y_{f}^{T}(\tilde{A} - \lambda_{1}I) = 0 \implies$$

$$y_{f}^{T}bf^{T}x_{1} + y_{f}^{T}(\tilde{A} - \lambda_{1}I)x_{1} = 0 \implies$$

$$y_{f}^{T}bf^{T}x_{1} = 0 \implies \psi_{1f}\alpha_{10}f^{T}x_{1} = 0 \implies$$

But $\psi_{1f} \neq 0$, because if $\psi_{1f} = 0$ then from (8) and because of the fact $\alpha_{j,j-1} \neq 0$, j = 2(1)n, y_{2f}^T would be zero too, so $y_f^T = 0$, thus $\psi_{1f} \neq 0$ and from $\psi_{1f}\alpha_{10}f^Tx_1 = 0$ => $f^Tx_1 = 0$ since $\alpha_{10} \neq 0$. But if we choose,

$$f^{T} - e^{T}(X^{T}X)^{-1}X$$
) (8a)

where the columns of X are the eigenvectors of \widetilde{A} , then $f^Tx_i \neq 0$, $i = 1(1)\ell$, which is a contradiction. So for f^T given in (8a) $\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + bf^T) = \emptyset$.

Second Proof: For $m \ge 1$. Here we shall prove that for almost all $m \times n$ matrices F we have

$$\lambda(\widetilde{A}) \cap \lambda(\widetilde{A} + \widetilde{B}F) = \emptyset$$
 (9)

and we shall construct an F for which (9) will be true.

Fact: For any integer p ≥ 1 , almost all p×p matrices are nonsingular.
By this is meant that if the elements of a matrix are independently chosen from a continuous uniform distribution, then with probability 1 -the matrix will be nonsingular.

Suppose \widetilde{A} has eigenvalues λ_i , i=1(1)n, then we perform the following construction: (Note that $A_{j,j-1}$ has n_j linearly independent columns) for given i, with i=1(1)n and for j=2(1)k choose n_j columns of $\widetilde{A}-\lambda_j I$ corresponding to n_j linearly independent columns of $A_{j,j-1}$ and move them to the front of the matrix, that is to positions

$$1 + \sum_{i=2}^{j-1} n_i$$
 to $\sum_{i=2}^{j} n_i$, (where $\sum_{i=2}^{1}$ is ignored).

For example the matrix

where stars represent the n_2 and n_3 linearly independent columns of A_{21} and A_{32} respectively. We may then write the result as, with permutation matrix P,

$$(\widetilde{A}-\lambda_{1}^{1})P = \underbrace{\begin{bmatrix} A_{11}^{(1)} & A_{12}^{(1)} \\ A_{21}^{(1)} & A_{22}^{(1)} \end{bmatrix}}_{A_{22}^{(1)}} \cdot n_{1}$$

$$(9a)$$

where $A_{21}^{(i)}$ is block upper triangular and nonsingular since it has nonsingular $n_j \times n_j$ blocks on its diagonal, j = 2(1)k. (These blocks being the same for all choices of λ_i ,

i = 1(1)n). Now write
$$\widetilde{B}FP = \begin{bmatrix} F_{11} & F_{12} \\ \hline 0 & 0 \end{bmatrix}_{n-n_1}^{n_1}$$
 so $n-n_1$

$$[(\widetilde{A} + \widetilde{B}F) - \lambda_{1}I]P = \begin{bmatrix} A_{11}^{(1)} + F_{11} & A_{12}^{(1)} + F_{12} \\ A_{21}^{(1)} & A_{22}^{(1)} \end{bmatrix}$$

For this to be singular there exists $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} n - n_1$, $x \neq 0$

such that $[(\widetilde{A} + \widetilde{B}F) - \lambda_{1}I]Px = 0$ so first

$$A_{21}^{(1)}x_1 + A_{22}^{(1)}x_2 = 0 \implies$$

$$\mathbf{x}_{1} = -\left[\mathbf{A}_{21}^{(1)}\right]^{-1} \mathbf{A}_{22}^{(1)} \mathbf{x}_{2} \tag{98}$$

From (9 β) we can see that $x_2 = 0$ implies $x_1 = 0$ and therefore x = 0, thus $x_2 \neq 0$. Next

$$(A_{11}^{(i)} + F_{11})x_1 + (A_{12}^{(i)} + F_{12})x_2 = 0 \Rightarrow$$

$$[A_{12}^{(i)} + F_{12} - (A_{11}^{(i)} + F_{11})[A_{21}^{(i)}]^{-1}A_{22}^{(i)}]x_2 = 0 \qquad (9\tilde{\gamma})$$

where $x_2 \neq 0$.

Let
$$C^{(1)} = A_{12}^{(1)} + F_{12} - (A_{11}^{(1)} + F_{11})[A_{21}^{(1)}]^{-1} A_{22}^{(1)}$$
.

From (9 γ) we conclude that $C^{(1)}$ must be singular. But if the elements of F are chosen randomly from the uniform distribution [0,1], then ${}_{i}F_{11}$ and F_{12} will be independent, so given F_{11} , for almost all F_{12} $C^{(1)}$ will be nonsingular. For a given F there are at most f different f so for almost all f, f and f are f such that f and f are f and f and f and f are f such that f and f are chosen f and f and f and f are chosen f so that f and f are chosen f are chosen f are chosen f and f are f and f are chosen f and f are chosen f and f are f and f are chosen f are chosen f and f are chosen f are chosen f and f are chosen f and f are chosen f are chosen f and f are chosen f are chosen f are chosen f are chosen f and f are chosen f are chosen f are chosen f are chosen f an

that is let $F = A_{10}^T (A_{10} A_{10}^T)^{-1} (0, -\gamma I) P^T$. Note that $(A_{10} A_{10}^T)^{-1}$ exists since A_{10} has full row rank. Also note that F is independent of the choice of λ_1 . Then

$$[(\widetilde{A} + \widetilde{B}F) - \lambda_{1}I]P = \begin{bmatrix} A_{11}^{(1)} & A_{12}^{(1)} \\ A_{21}^{(1)} & A_{22}^{(1)} \end{bmatrix}$$

4,

So for this choice of F the corresponding relation to (9γ) will be:

$$[(A_{12}^{(1)} - A_{11}^{(1)}[A_{21}^{(1)}]^{-1}A_{22}^{(1)}) - \gamma I]x_2 = 0 , x_2 \neq 0$$
 (98)

So γ must be an eigenvalue of the $n_1 \times n_1$ matrix

$$A^{(i)} = A_{12}^{(i)} - A_{11}^{(i)} [A_{21}^{(i)}]^{-1} A_{22}^{(i)}$$

 $A^{(1)}$ has at most n_1 distince eigenvalues, and as there are at most n distinct eigenvalues of \widetilde{A} , there are at most n distinct $A^{(1)}$, so there are at most nn_1 distinct values of γ for which $\widetilde{A}+\widetilde{B}F$ has any eigenvalue equal to any eigenvalue of A. Thus for almost all F of even this most restricted form, \widetilde{A} and $\widetilde{A}+\widetilde{B}F$ have completely distince eigenvalues.

(III) For this part the simplest proof was found in [3] (p. 277-281, 282-283) and it is as follows.

Lemma (I): Let $M = [B,AB,...,A^{n-1}B]$ then for any $q \ge n$ we have

$$rank[B,AB,...,A^{q-1}B] = rank(M)$$

Proof: As k increases by one unit the rank of the matrix

M_k = [B,AB,...,A^{k-1}B] either increases (by at least 1) or remains constant. Suppose that k is an integer such that the rank of M_{k+1} is equal to the rank of M_k. That means that the m columns comprising A^kB are each linearly

dependent on the (previous) columns in $\begin{subarray}{c} M_k \end{subarray}$. That is, there is a relation of the form

$$\mathbf{A}^{\mathbf{k}}\mathbf{B} = \mathbf{B}\mathbf{D}_0 + \mathbf{A}\mathbf{B}\mathbf{D}_1 + \ldots + \mathbf{A}^{\mathbf{k}-1}\mathbf{B}\mathbf{D}_{\mathbf{k}-1}$$

where each $\mathbf{D}_{\mathbf{i}}$ is an m×m matrix. Now multiplication of this relation by A leads to the new relation

$$A^{k+1}_{B} = ABD_{0} + A^{2}BD_{1} + ... + A^{k}BD_{k-1}$$

which shows that the columns comprising $A^{k+1}B$ are linearly dependent on the columns in M_{k+1} . Therefore, the rank of M_{k+2} is the same as the rank of M_{k+1} . By continuing this afgument, we see that for all j > k the rank of M_j is equal to that of M_k . Thus we have shown that, in the progression of M_k 's, once the rank fails to increase, it will remain constant even as additional columns are adjoined. In view of the above, the rank of M_k increases by at least 1 at each increment of k until it attains its maximum rank. Since the maximum rank is at most n, the maximum rank is attained within n steps (that is, by M_n).

(1) Suppose first that the rank condition does not hold. For any t₁ > 0 and any integrable control u(t) defined on [0, t₁] we have

$$x(t_1) = \int_0^{t_1} e^{A(t_1-t)} B u(t) dt \Rightarrow$$
 (10)

$$\mathbf{x}(t_1) = \mathbf{B} \int_0^{t_1} \mathbf{u}(t) dt + \mathbf{AB} \int_0^{t_1} (t_1 - t) \mathbf{u}(t) dt + \dots$$

*

when evaluated, the integrals in the above expression are simply constant m-dimensional vectors. Therefore, the expression shows that $x(t_1)$ is a linear combination of the columns of B,AB,... By the earlier lemma, if the rank of M is less than n, then even the infinite set of vectors B,AB,A²B,... does not contain a full basis for the entire n-dimensional space. Thus, there is a vector x_1 that is linearly independent of all these vectors, and therefore cannot be attained.

Now suppose that the rank condition does hold. We will show that the system is completely controllable and that in fact the state can be transferred from zero to an arbitrary point x_1 within an arbitrary short period of time.

We first show that for any $t_1 > 0$ the matrix

$$K = \int_{0}^{t_1} e^{-At_{BB}T} e^{-A^{T}t} dt$$

is nonsingular. To prove this, suppose there is a vector a such that Ka = 0. Then

$$a^{T}Ka = 0$$

or more explicitly

$$\int_{0}^{t_{1}} a^{T} e^{-At} B B^{T} e^{-A^{T}} t a dt = 0$$
 (11)

The integrand above has the form $c(t)^T c(t)$, where $c(t) = B^T e^{-A^T t} a$. It follows that the integrand is always nonnegative. For the integral (11) to vanish, it follows that the integrant must vanish identically for $0 \le t \le t_1$. Therefore

$$a^{T}e^{-At}B = 0$$
 for all t, $0 \le t \le t$,

Evaluation of this expression, and its successive derivatives, with respect to t, at t = 0 leads to the following sequence of equations:

$$\mathbf{a}^{\mathsf{T}}\mathbf{B} = 0$$

$$\mathbf{a}^{\mathsf{T}}\mathbf{A}\mathbf{B} = 0$$

$$\vdots$$

$$\mathbf{a}^{\mathsf{T}}\mathbf{A}^{\mathsf{n}-1}\mathbf{B} = 0$$

This means that the vector a must be orthogonal to all columns of the matrix M . Since it is assumed that this matrix has rank n , it must follow that a=0. Therefore K is non-singular. Now, given \mathbf{x}_1 , select any $\mathbf{t}_1>0$ and set

$$u(t) = B^{T} e^{-A^{T}} t K^{-1} e^{-At_{1}} x_{k}$$
 (12)

Then from (10) $x(t_1) = \int_0^{t_1} e^{A(t_1-t)} BB^T e^{-A^T t} K^{-1} e^{-At_1} x_1 dt \Rightarrow x(t_1) = e^{At_1} KK^{-1} e^{-At_1} x_1 = x_1$

Therefore the control (12) transfers the state from zero to x_1 , and the system is completely controllable.

Suppose a sequence of controls $u_1, u_2, ..., u_q$ is applied to the system (2) with $x_1 = 0$ it follows that

$$x_{q+1} = A^{q}B u_1 + A^{q-1}B u_2 + ... + B u_q$$

From this formula we see that points in state space can be reached if and only if they can be expressed as linear combinations of powers of A times B. Thus the issue of complete controllability rests on whether the infinite sequence B,AB,A²B,... has a finite number of columns that span the entire n-dimensional space. By the earlier lemma, however, these span the full n-dimensional space if and only if

$$rank[B,AB,...,A^{n-1}B] = n .$$

Remark (IV) From step (III) we can see that (1) is controllable if and only if (2) is controllable. So it is the matrices A and B which determine controllability and not the type of the system (continuous or discrete).

(III) The computational problem of controllability

In the previous section we presented four mathematically equivalent theorems on controllability. These theorems can lead to several different computational approaches, four of which are the following.

- C1: Involves transforming A and B to the form in (6) and is described in [2].
- C2: Form the matrix [B,AB,...,Aⁿ⁻¹B] and then compute its rank.
- C3: Compute the eigenvalues λ_i , i = 1(1)n of A and then compute the ranks of the matrices $[B,A-\lambda_i I]$ for i = 1(1)n.
- C4: For a random matrix $F \in \mathbb{R}^{m \times n}$ form the matrix A + BF and then compute the eigenvalues of A + BF and A.

In this section we will prove that the algorithms which can be made by using C2, C3, and C4 are poor. This will be done by using counterexamples. We will also show that the algorithm which is based on C1 works perfectly for the same examples.

We will show the weakness of C2 first. Before this we give two definitions.

Definition (V) [12], [14]: Let $D \in \mathbb{C}^{n \times m}$ with $n \ge m$, then the singular value decomposition (SVD) of D is $D = USV^H$ where $U \in \mathbb{C}^{n \times n}$, $V \in \mathbb{C}^{m \times m}$, $S = \begin{bmatrix} S_1 \\ 0 \end{bmatrix}$ with $S_1 = \operatorname{diag} \left(\sigma_1, \ldots, \sigma_m\right)$, $\sigma_1 \ge \sigma_2 \ge \ldots \ge \sigma_m \ge 0$, U and V are unitary matrices. Let $V_1 = V$ and $U = (U_1, U_2)$ where $U_1 \in \mathbb{C}^{n \times m}$. Then $D = U_1 S_1 V_1^H$. The real scalars σ_1 , i = 1(1)m are called the singular values of D and the columns of V_1 and U_1 are called the right and left singular vectors, respectively.

<u>Definition (VI)</u>: Let $D \in \mathbb{R}^{n \times n}$, then the number

$$\chi(D) \stackrel{\Delta}{=} \frac{\sigma_1(D)}{\sigma_n(D)} \ge 1$$

is called the condition number of D for solution of equations, where $\sigma_1(D)$ and $\sigma_n(D)$ are the largest and the smallest singular values of D. If $\chi(D)$ is large then the problem of the solution of the equation Dy = d, where $d \in \mathbb{R}^n$, $y \in \mathbb{R}^n$, is ill-conditioned, for more details about condition numbers see [12, page 192]. Consider now the clearly controllable system

$$\dot{x} = A_{n}x + B_{n}u , \qquad (13)$$
where $A_{n} = \text{diag } (1, 2^{-1}, \dots, 2^{1-n}) , B_{n} = (\underbrace{1, 1, \dots, 1}_{n})^{T} .$

In this case the (i,j) element of the matrix $C_n = [B,AB,...,A^{n-1}B]$ is $2^{(i-1)(1-j)}$ which can be formed and stored with full accuracy on most computers. Even so with $n \in \{8,9,10\}$ we have the following results:

TABLE 1

RESULTS OF COMPUTATIONAL TESTS ON (13) USING C2

n	σ ₁ (C _n)	σ ₈ (C _n)	σ ₉ (c _n)	σ ₁₀ (C _n)	χ(A _n)	χ(C _n)
8	3.3033	2.0688×10 ⁻⁸	-	-	128	0.1596×10 ⁹
9	3.4704	4.7752×10 ⁻⁸	1.5846×10 ⁻¹⁰	-	256	0.219×10 ¹¹
10	3.6298	7.1246×10 ⁻⁸	3.6402×10 ⁻¹⁰	6.1288×10 ⁻¹³	512	0.5922×10 ¹³

where $\sigma_{\bf i}(C_{\bf n})$ is the ith singular value of $C_{\bf n}$. So the numerical rank of $C_{\bf n}$ when ${\bf n}$ ϵ {8,9,10} will be less than ${\bf n}$ using the SVD on a computer with relative precision no smaller than 10^{-7} . Obviously, in

these cases C2 would fail, indicating that the controllable system (13) is uncontrollable when $n \in \{8,9,10\}$.

This happens because if we transform a problem by matrix multiplication and the transforming matrix is ill-conditioned with respect to solution of equations, that is, it has large condition number, then the transformed problem will usually be more sensitive to changes in data, than the original one. As we can see from table I $\chi(A_n) >> 1$ when $n \in \{8,9,10\}$ and as a consequence of this the resulting matrix C_n has much larger condition number $\chi(C_n)$ than A_n . So C_n is more sensitive to changes in data than A_n .

In order to compare Cl with C2 we applied C1 to (13) for $n \in \{8,9,10\}$. Since in (13) B is a column vector then in (6) k is equal to 8, 9 or 10 when n is equal to 8, 9 or 10, respectively, and $A_{i,i-1}$ is scalar for i = 1(1)k. Theoretically Cl should give

$$A_{i,i-1} \neq 0$$
 , $i = 1(1)k$.

The results which were taken from the actual performance of C1 are presented in table II.

TABLE 11

RESULTS OF COMPUTATIONAL TESTS ON (13) USING C1

n	A ₁₀	i = 2(1) k-1	A _{k,k-1}
8	-2.828426	0.024< A _{1,1-1} <0.323	-0.010119
9	-3	0.012< A _{1,1-1} <0.314	0.005113
10	-3.162277	0.006< A _{i,i-1} <0.305	. _, -0. 002570

So by observing the results presented in table II we can conclude that (13) is controllable when $n \in \{8,9,10\}$, and this conclusion would be reached using a relative precision as low as 10^{-4} . Thus Cl succeeded where C2 failed.

In the rest of this section we will show the weakness of C3 and C4. We will do this by the same way we showed the weakness of C2, that is, by considering a counterexample, because we know that a good example is sufficient to condemn a poor algorithm. The approaches in C3 and C4 depend on finding eigenvalues. Unlike singular values the eigenvalues of some matrices can be very ill-conditioned, that is, very sensitive to small changes in the matrix. So we choose a matrix which has already been presented, by Wilkinson [13, p. 90], as an extremely ill-conditioned case regarding the eigenvalues. The matrix Wilkinson considered was the following:

$$\hat{A} = \begin{bmatrix}
20 & 20 & & & & \\
& 19 & 20 & & & \\
& & 18 & 20 & & \\
& & & & 2 & 20 \\
& & & & & 1
\end{bmatrix}$$

The eigenvalues of \hat{A} are $\lambda_1 = i$, i = 1(1)20. Consider now the uncontrollable system

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \tag{14}$$

with
$$A = Q^{T} \hat{A} Q$$
, $B^{T} = (\underbrace{1,1,\ldots,1}_{19},0)Q$

where Q is a random orthogonal matrix found by finding the SVD of a square matrix whose elements are uniform random numbers on (-1,1) obtained by the VNI(O) random number generator. To test C4 we calculated the eigenvalues of A and A + BF. The sets $\tilde{\lambda}(A)$ of the computed eigenvalues of A and $\tilde{\lambda}(A+BF)$ of A + BF were obtained using single precision (6 significant hexidecimal digits) on the AMDAHL 470/V7 computer at McGill University. The elements of F were random numbers on (-1,1) from the VNI(O) random number generator. To test C3 we calculated the values

$$\rho(B,A-\widetilde{\lambda}_{1}I) \stackrel{\triangle}{=} ratio \ of \ smallest \ to \ largest \ singular$$

$$value \ of \ (B,A-\widetilde{\lambda}_{1}I) \ .$$
 where $\widetilde{\lambda}_{1} \in \widetilde{\lambda}(A)$ for $i=1(1)20$.

The ratio $\rho(B,A-\widetilde{\lambda}_1^{-1}I)$ was the same for both eigenvalues of a complex conjugate pair. The results are given in table III.

TABLE III

RESULTS OF COMPUTATIONAL TESTS ON (14)

ρ (Β,Α-λ̃ 1)	$\widetilde{\lambda}(A)$	$\widetilde{\lambda}(A + BF)$
2.024 × 10 ⁻³	-0.32985±j1.06242	0.99999
3.979×10^{-3}	0.92191±j3.13716	-8.95872±j 3.73260
7.3203×10 ⁻³	3.00339±j4.80414	-5.11682±j 9.54329
1.187 × 10 ⁻²	5.40114±j6.17864	-0.75203±j14.148167
1.846 × 10 ⁻²	8.43769±j7.24713	5.77659±j15.58436
2.5905×10 ⁻²	11.82747±j7.47463	11.42828±j14.28694
3.221 × 10 ⁻²	15.10917±j6.90721	13.30227±j12.90197
4.036 × 10 ⁻²	18.06886±j5.66313	18.59961±j14.34739
5.179 × 10 ⁻²	20.49720±j3.81950	23.94877±j11.80677
6.436 × 10 ⁻²	22.06287±j1.38948	28.45618±j 8.45907
		32.68478

Theoretically one of the eigenvalues of A should be the same as an eigenvalue of A+BF, this one being unity. But as we can see from table III the computed eigenvalues of A are almost unrelated with the true eigenvalues of A and the computed eigenvalues of A+BF. So in this case C4 would tell us a clearly uncontrollable system is controllable. As we can conclude from the results of table III the failure of C4 is due to the ill-conditioned eigenvalues of \hat{A} . The approach C3 failed too, because in theory one value of $\hat{\rho}$ should be zero, but no value of $\hat{\rho}$ can be considered as zero within the previously mentioned accuracy. So C3 would have indicated that the uncontrollable system (14) was controllable. C3 was also carried out with the true eigenvalue, unity, of A, and this gave

$$\rho(B,A-I) = 5.293 \times 10^{-8}$$

indicating that it is the eigenvalue computation that caused the failure of C3, just as for C4.

In order to compare the ability of approach C1 we applied it to system (14). Since, as in the previous counterexample, B is a column vector then in (6) k = 20 and $A_{i,i-1}$ is scalar for i = 1(1)20. Theoretically, C1 should give

$$A_{i,i-1} \neq 0$$
, $i = 1(1)19$ and $A_{20,19} = 0$.

In practice C1 gave $A_{1,0} = 4.35887$; $A_{2,1} = 8.30008$; $19 < |A_{1,i-1}| \le 22^*$, i = 3(1)19 and $A_{20,19} = 0.000006914$, showing that the system is uncontrollable to this precision of computation. To see if this was just a

coincidence the algorithm was applied to (14) with 50 different random orthogonal matrices Q. In all cases similar good results were obtained.

Some of the ideas and results of this section have already been presented in [2].

(IV) DISTANCE FROM AN UNCONTROLLABLE SYSTEM

As we saw in equation (3) of section (II), μ is defined to be the distance between the system defined by the matrix (A,B) and the nearest uncontrollable system defined by the matrix (A + δ A , B + δ B). A knowledge of μ would generally be more useful than just knowing if the system is controllable or not. In this section an algorithm is given which will make use of remark (III) to determine whether a system is controllable or not. The approach could be extended to compute μ , but this is not done here as it appears to be inefficient.

We will relate now $\,\mu\,$ with another quantity which can be calculated, that is, we are going to prove that

$$\mu = \min \sigma_{\mathbf{n}}(B, \mathbf{A} - \lambda \mathbf{I}) , \lambda \in \mathbb{C}$$
 (16)

where $\sigma_n(B, A-\lambda I)^n$ is the smallest singular value of $(B, A-\lambda I)$.

Note: In the following we will represent the singular values of $(B,A-\lambda I)$ by either $\sigma_1(B,A-\lambda I)$ or $\sigma_1(\lambda)$ meaning exactly the same thing, that is, $\sigma_1(\lambda) \equiv \sigma_1(B,A-\lambda I)$, i=1(1)n and $\forall \lambda \in C$.

Before we go ahead to prove (16) we shall first prove a lemma which will help us in the proof of (16).

Lemma (II): Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{m \times n}$, where $m \ge n$ and let $A = U_A \Sigma_A V_A^H \quad , \quad B = U_B \Sigma_B V_B^H \quad \text{be the singular value decompositions of } A$ and $B \cdot L$ Let also $\alpha_1 \ge \alpha_2 \ge \dots \ge \alpha_n$ and $\beta_1 \ge \beta_2 \ge \dots \ge \beta_n$ be the singular values of A and B respectively. Then for given A we shall prove that:

minimum
$$|A - B|_2 = \alpha_{k+1}$$
, $k \le n$ and $\alpha_{n+1} = 0$ (17) rank(B) $\le k$

Proof: We know ([12], p. 231, th. 6.6) that

$$||A - B||_{2} \ge |\alpha_{1} - \beta_{1}|$$
, $i = 1(1)n$

so obviously $||A - B||_2 \ge \max |\alpha_i - \beta_i|$, $i \in \Delta_n$. Now if rank(B) $\le k$ with $k \le n$ we have that $\beta_i \ge 0$, i = 1(1)k and $\beta_i = 0$, i = k+1(1)n.

Let α be the maximum of $|\alpha_i - \beta_i|$ with $i \in \Delta_k$, that is, $\alpha = \max_i |\alpha_i - \beta_i|$, $i \in \Delta_k$. Then we have that

$$||A - B||_2 \ge \max(\alpha, \alpha_{k+1})$$

So when $/ \operatorname{rank}(B) \le k$ we have $||A - B||_2 \ge \alpha_{k+1}$ for every $B \in \mathbb{R}^{m \times n}$. Now when $\operatorname{rank}(B) \le k$ there is a B such that $B = U_A \Sigma_B V_A^H$ and $\beta_1 = \alpha_1$, i = 1(1)k. For this B we obviously have $||A - B||_2 = \alpha_{k+1}$. So indeed $\min_{rank} ||A - B||_2 = \alpha_{k+1}$.

Now we are ready to prove (16).

<u>Proof:</u> Let (A,B) define a system and $[(A + \delta A), (B + \delta B)]$ define an uncontrollable system. Also let $D(\lambda) = [(B + \delta B), (A + \delta A) - \lambda I]$, $\lambda \in C$, and

be the singular value decomposition of (B, A- λ I) with $\sigma_1(\lambda) \geq \sigma_2(\lambda) \geq \ldots \geq \sigma_n(\lambda)$.

From lemma (I) for $k \le n$ we get ($\sigma_{n+1}(\lambda)$ will be considered zero)

minimum | | (B,A- λ I) - [(B + δ B), (A + δ A) - λ I] | |₂ = $\sigma_{k+1}(\lambda)$ => rank[D(λ)] \leq k

minimum .
$$||(\delta A, \delta B)||_2 = \sigma_{k+1}(\lambda)$$
 / (19) $||(\delta A, \delta B)||_2 = \sigma_{k+1}(\lambda)$

The smaller the 2-norm of $(\delta A, \delta B)$ is, the nearer the uncontrollable system $[(A + \delta A), (B + \delta B)]$ is to the system (A,B). Since $[(A + \delta A), (B + \delta B)]$ is uncontrollable obviously $k \le n-1$. We can see that in (19) there are two parameters λ and k. Regarding k, $(\delta A, \delta B)|_{2}$ takes its minimum value and at the same time keeps $[(A + \delta A), (B + \delta B)]$ uncontrollable, if and only if k = n-1. So we have that:

$$\mu = \min_{\lambda} \left[\min_{\substack{\text{minimum} \\ \text{rank}[D(\lambda)] \le n-1}} ||(\delta A, \delta B)||_2 \right] = \min_{\lambda} \sigma_n(\lambda)$$

Thus $\mu = \min_{\lambda} \sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda \mathbf{I})$, $\lambda \in \mathbb{C}$ (20)

Remark (V): An important property of the singular values is that they are not very sensitive to changes in the matrix. In fact if $\sigma_1 \geq \ldots \geq \sigma_n \quad \text{are the singular values of the matrix } \quad D \quad \text{and} \quad \widetilde{\sigma}_1 \geq \ldots \geq \widetilde{\sigma}_n \quad \text{the singular values of the perturbed matrix } \quad D + \delta D \quad \text{then}$

$$|\sigma_{i} - \tilde{\sigma}_{i}| \le ||\delta D||_{2}$$
 , $i = 1(1)n$. (21)

For more information see [12, page 321] or [14, page 24]. Thus, from (21), we conclude that the singular values are well-conditioned with respect to perturbations in the matrix. So, from (20) we can also conclude that μ is well-conditioned with respect to perturbations in the matrix [2, page 137].

Remark (VI): We can easily observe that $\sigma_n(B, A-\lambda I)$ is a function with domain $\mathcal{D} \subseteq \mathfrak{C}^{n \times (n+m)}$ and range $\mathcal{P} \subseteq \mathbb{R}$, that is,

$$\sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A}-\lambda\mathbf{I}): \mathcal{D} \longrightarrow \mathcal{P}$$
 (21a)

We are going to present two algorithms, one for $\lambda \in \mathbb{R}$ and one for $\lambda \in \mathbb{C}$. But before we present the algorithms we will make some useful observations about the function $\sigma_n(B, A-\lambda I)$

Remark (VII): Let $\sigma_1' \geq \sigma_2' \geq \ldots \geq \sigma_n'$ be the singular values of (B,A), then using (21) we can see that (0, λ I) shifts the singular values of (B,A) by at most $|\lambda|$, because

$$|\sigma_{i}(\lambda) - \sigma_{i}^{\dagger}| \le ||(B, A-\lambda I) - (B,A)||_{2} = ||(0,\lambda I)||_{2} = |\lambda|$$
for $i = 1(1)n$.

Since we are interested only in the smallest singular value of (B, $A-\lambda I$) , we have

$$|\sigma_{\mathbf{n}}(\lambda) - \sigma_{\mathbf{n}}^{\dagger}| \le |\lambda| \Longrightarrow$$

$$\sigma_{\mathbf{n}}^{\dagger} - |\lambda| \le \sigma_{\mathbf{n}}(\lambda) \qquad (22)$$

Similarly, (B,A) can shift the singular values of (0, λ I) by at most $||(B,A)||_2$, because

$$|\sigma_{1}(\lambda) - |\lambda|| \le ||(B, A-\lambda I) - (0, -\lambda I)||_{2} = ||(B,A)||_{2} = \sigma_{1}^{*},$$

$$i = 1(1)n.$$

For i = n we have that

$$|\lambda| - |\lambda|| \le \sigma_1' = \frac{1}{|\lambda|} - \sigma_1' \le \sigma_1(\lambda) . \tag{23}$$

From (22) (respectively from (23)) we can easily observe that if $|\lambda| < \sigma_n' \quad \text{(respectively} \quad |\lambda| > \sigma_1' \quad \text{) then} \quad \sigma_n(\lambda) > 0 \quad \text{. Since we are eigenstated in finding if} \quad \sigma_n(\lambda) \quad \text{is ever zero we should search in the ring} \quad \sigma_n' \leq |\lambda| \leq \sigma_1' \quad \text{(see fig. 1)}.$

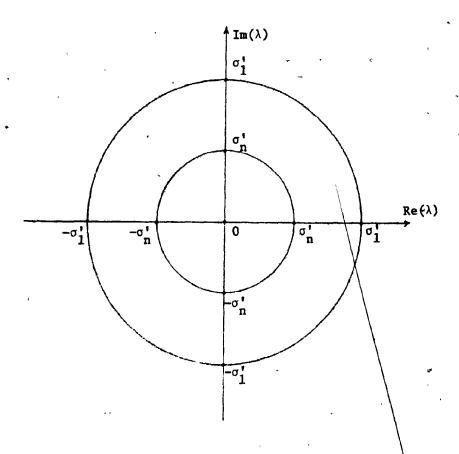
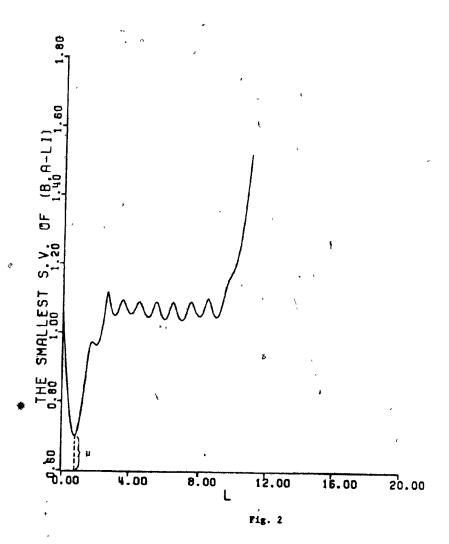


Fig. 1

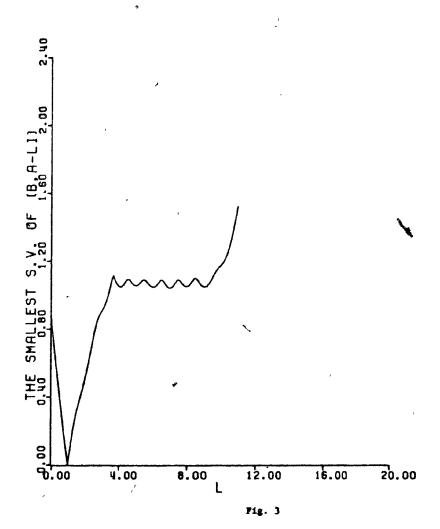
where $Re(\lambda)$ and $Im(\lambda)$ is the real and the imaginary part of λ , respectively.

Remark (VIII): For A having real eigenvalues we ran some programs in which we plotted the values of $\sigma_n(B, A-\lambda I)$ for λ in the real interval Δ where every eigenvalue λ_i , i=1(1)n of A was in Δ too, that is, $\lambda_i \in \Delta$, i=1(1)n. The plots we obtained had the general form of fig. 2 when $\mu > 0$ (system controllable) and fig. 3 when $\mu = 0$ (system uncontrollable).

CONTROLLABLE . CASE



UNCONTROLLABLE CASE



- Note: As we can see from fig. 3 if the system is uncontrollable then the function $\sigma_n(B, A-\lambda I)$ is almost linear in the neighbourhood of λ' , where λ' is such that $\sigma_n(B, A-\lambda' I) = 0$.
- Remark (IX): Now (21) will help us bound the rate of change in $\sigma_n(B,A-\lambda I) \quad \text{with respect to} \quad \lambda \quad \text{Let} \quad \lambda_1,\lambda_2 \quad \text{be two distinct complex numbers, then}$

$$|\sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda_{1}\mathbf{I}) - \sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda_{2}\mathbf{I})| \leq ||(\mathbf{B}, \mathbf{A} - \lambda_{2}\mathbf{I}) - (\mathbf{B}, \mathbf{A} - \lambda_{2}\mathbf{I})||_{2}$$

$$\leq ||[0, (\lambda_{2} - \lambda_{1})\mathbf{I}]||_{2}$$

$$= |\lambda_{1} - \lambda_{2}| \Rightarrow$$

$$\frac{|\sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda_{1}\mathbf{I}) - \sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda_{2}\mathbf{I})|}{|\lambda_{1} - \lambda_{2}|} \leq 1$$
(24)

This means that $|\sigma_n(B,A-\lambda I)| \le 1$ wherever the derivative is defined, in other words, (24) means that at every point of the curve $\sigma_n(B,A-\lambda I)$ the tangent (if it can be defined) makes an angle ω with the horizontal axis such that

$$\omega \in \left[0, \frac{\pi}{4}\right] \cup \left[\frac{3}{4}, \pi\right] . \tag{25}$$

Remark (X): In remark (VII) we mentioned that if $\sigma_n(B,A-\lambda I)$ has a zero, then this zero will be in the ring of fig. 1. So, we should search all the ring in order to find the possible root. Now we will show that we do not have to search all the ring because the plot of $\sigma_n(B,A-\lambda I)$ above the axis $Re(\lambda)$ is completely identical with the

corresponding plot below the axis $Re(\lambda)$. So we should prove that $\sigma_n(B,A-\lambda I)=\sigma_n(B,A-\bar{\lambda}I)$ for every $\lambda\in\mathbb{C}$, where $\bar{\lambda}$ is the complex conjugate of λ . We know that for a matrix D such that $D\in\mathbb{C}^{n\times n}$ we have $\lambda(D^T)=\lambda(D)$ [12, page 267]. Also we know that the singular values of the matrices $(B,A-\lambda I)$ and $(B,A-\bar{\lambda}I)$ are the eigenvalues of the matrices $(B,A-\lambda I)^H(B,A-\bar{\lambda}I)$ and $(B,A-\bar{\lambda}I)^H(B,A-\bar{\lambda}I)$, respectively. But

$$[(B,A-\lambda I)^{H}(B,A-\lambda I)]^{T} = (B,A-\lambda I)^{T}[(B,A-\lambda I)^{H}]^{T}$$

$$= (B,A-\lambda I)^{T}(B,A-\overline{\lambda}I)$$

$$= (B,A-\overline{\lambda}I)^{H}(B,A-\overline{\lambda}I)$$

since A and B are real. Thus

$$\sigma_{\mathbf{i}}(B,A-\lambda I) = \sigma_{\mathbf{i}}(B,A-\overline{\lambda}I)$$
 , $i = 1(1)n$. (26)

Remark (XI): In this remark we will calculate the $\sigma_n(B,A-\lambda I)$, which will be useful for the esitmation of the derivative of the function $\sigma_n(B,A-\lambda I)$ wherever it can be defined. In general let $D \in \mathbb{C}^{m \times n}$ and $D = U \Sigma V^H$ be the singular value decomposition of D, then we have the following.

$$DV = U\Sigma^{l} \Rightarrow \qquad \text{(when these derivatives exist)}$$

$$\dot{D}V + D\dot{V} = \dot{U}\Sigma + U\dot{\Sigma} \Rightarrow \Rightarrow$$

$$U^{H}\dot{D}V + U^{H}D\dot{V} = U^{H}\dot{U}\Sigma + U^{H}U\dot{\Sigma} \Rightarrow \Rightarrow$$

$$U^{H}\dot{D}V + \Sigma V^{H}\dot{V} = U^{H}\dot{U}\Sigma + \dot{\Sigma} \qquad (27)$$

But we know that if a vector $\mathbf{x}(t)$ has constant magnitude then it is orthogonal to its derivative, that is $\dot{\mathbf{x}}^H\mathbf{x} = 0$. So in (27) the diagonal elements of the matrices $\Sigma \mathbf{V}^H\dot{\mathbf{v}}$ and $\mathbf{U}^H\dot{\mathbf{U}}\Sigma$ are zero. So, from (27), we have

$$\dot{\sigma}_{1} = \mathbf{u}_{1}^{H} \dot{\mathbf{D}} \mathbf{v}_{1} \tag{28}$$

Thus, in our case since $D = (B,A-\lambda I)$ we have from (28)

$$\dot{\sigma}_{1}(\lambda) = -u_{1}^{H}(0, I_{n})v_{1}, \quad i = 1(1)n$$
 (29)

where u_{i} , v_{i} are the columns of U,V respectively.

We will first present an algorithm for real λ and then we will extend it for complex λ . Acutally our problem is the solution of the equation $\sigma_n(B,A-\lambda I)=0$ with one unknown, where in the real case the root, if there exists any, will be in the interval

$$\Delta = [-\sigma_1', -\sigma_n'] \cup [\sigma_n', \sigma_1']$$

$$-\sigma_1' \dots -\sigma_{n-1}' \quad -\sigma_n' \quad 0 \quad \sigma_n' \quad \sigma_{n-1}' \quad \dots \quad \sigma_1'$$

$$(30)$$

Fig. 4

Of course we do not have any formula for this equation but we can find as many values of the function $\sigma_n(B,A-\lambda I)$ as we want, and numerically in many cases that is what really counts. A first thought is to use Newton's method, but using Newton's method we may face the following problem.

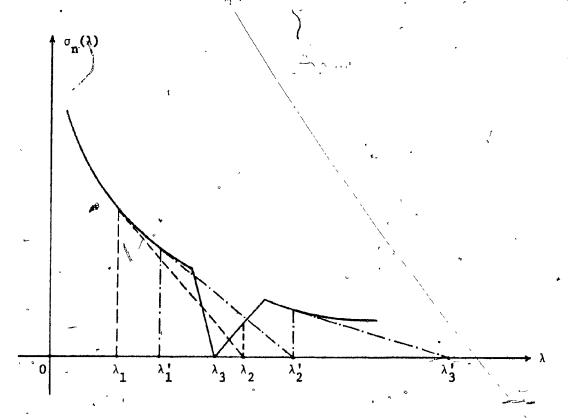


Fig. 5

As we can see from fig. 5, starting from λ_1 we will get the root, but starting from λ_1' we will miss it. So we cannot use Newton's method. We will try an iterative method of the following form $\lambda_{k+1} = \lambda_k + \text{step}_k$ where $\lambda_1 = -\sigma_1'$ and $\text{step}_k = \sigma_n(B, A - \lambda_k I)$, that is, the step, in order to go from λ_k to λ_{k+1} , is the value of the function, the root of which we are looking for, at λ_k . Using this step size and because of remark (IX) we will never miss the root. This can also be seen from fig. 6.

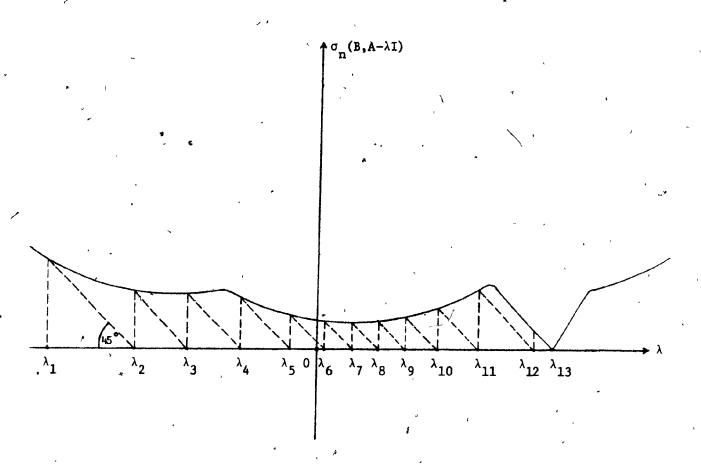


Fig. 6

Since we are looking for a possible root in the interval Δ , this algorithm can be quite fast if the system is quite controllable, that is, if μ is quite greater than zero, this can be observed from fig. 7.

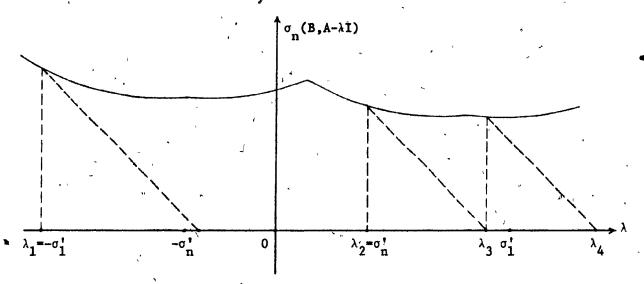


Fig. 7

Z.

But the same algorithm can be inefficient if the system is uncontrollable and we have reached the point where the function $\sigma_n(B,A-\lambda I)$ is almost linear, this can be seen from fig. 8.

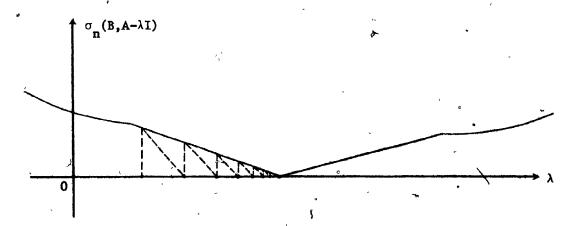


Fig. 8

In order to improve this situation we should take the fact $\dot{\sigma}(\lambda) \equiv \dot{\sigma}(B,A-\lambda I) \equiv \frac{d\sigma(B,A-\lambda I)}{d\lambda} \quad \text{into account.} \quad \text{We can use the quantity}$ $\frac{\sigma(\lambda_{k+1}) - \sigma(\lambda_{k})}{\sigma(\lambda_{k+1}) - \sigma(\lambda_{k})} \quad \text{as an estimate for } \dot{\sigma}(\lambda_{k}) \quad \text{This quantity will be quite}$

accurate for the calculation of the derivative when we have reached the point where the function $\sigma_n(B,A-\lambda I)$ is almost linear. At this point two consecutive points will have almost the same derivative as we can see from fig. 9. So this will be the criterion for detecting if we have reached the point where the function $\sigma_n(B,A-\lambda I)$ is almost linear.

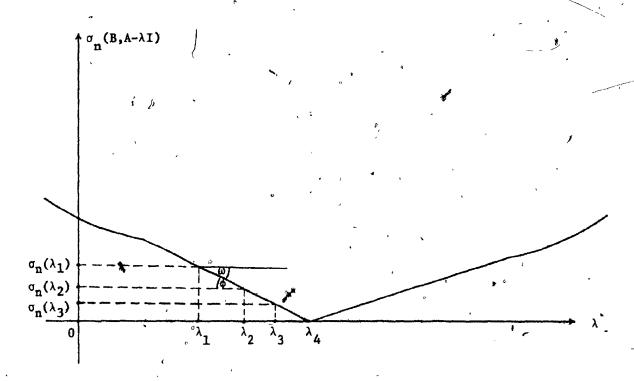


Fig. 9

From fig. 9 we can obviously see that $|\dot{\sigma}_n(\lambda_1) - \dot{\sigma}_n(\lambda_2)| < \delta$ where δ is an appropriate, small enough, positive real number. So from now on, until we reach the root we will use Newton's method, for example in fig. 9 we can have the following:

$$\lambda_4 = \lambda_2 - \frac{\sigma(\lambda_2)}{\dot{\sigma}(\lambda_2)}$$

So the algorithm for $\lambda \in \mathbb{R}$ is as follows:

Part 1 (Initialization)

Al: Calculate the singular values of (B,A) and let them be
$$\sigma_1' \geq \sigma_2' \geq \ldots \geq \sigma_n'$$
;

A2: Set
$$\lambda := -\sigma_1'$$
; $\lambda_{old} := \lambda$;

A3: Evaluate
$$\sigma_n(B,A-\lambda I)$$
; and set $\sigma_{old} := \sigma_n(B,A-\lambda I)$;

A4: If
$$\sigma_{\text{old}} < \varepsilon \sigma_1'$$
, where ε is a small positive real number, then system uncontrollable, stop;

A5:
$$\lambda := \lambda + \sigma_{old}$$
; $\lambda_{new} := \lambda$;

•A6: If
$$\lambda > -\sigma_n^*$$
 then go to Al77

A7: Evaluate
$$\sigma_{n}(B,A-\lambda I)$$
; $\sigma_{new} := \sigma_{n}(B,A-\lambda I)$;

A8: If
$$\sigma_{\text{new}} < \epsilon \sigma_1^{\dagger}$$
 then system uncontrollable, stop;

A9:
$$\sigma_{\text{old}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\lambda_{\text{new}} - \lambda_{\text{old}}}$$
; $\sigma_{\text{old}} := \sigma_{\text{new}}$; $\lambda_{\text{old}} := \lambda_{\text{new}}$;

Part 2 (We search the interval $[-\sigma_1', -\sigma_1']$)

Alo:
$$\lambda := \lambda + \sigma'_{old}$$
; $\lambda_{new} := \lambda$;

All: If
$$\lambda > -\sigma_n^*$$
 then go to Al7

Al2: Evaluate
$$\sigma_{n}(B,A-\lambda I)$$
; $\sigma_{new} = \sigma_{n}(B,A-\lambda I)$;

Al3: If
$$\sigma_{\mathbf{n}}(\mathbf{B},\mathbf{A}-\lambda\mathbf{I}) < \varepsilon \sigma_{\mathbf{1}}^{t}$$
 then system uncontrollable, stop;

Al4: $\sigma_{\mathbf{new}} := \frac{\sigma_{\mathbf{new}} - \sigma_{\mathbf{old}}}{\lambda_{\mathbf{new}} - \lambda_{\mathbf{old}}}$;

A14:
$$\sigma_{\text{new}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\lambda_{\text{new}} - \lambda_{\text{old}}}$$

Al5: If
$$|\dot{\sigma}_{new} - \dot{\sigma}_{old}| < \delta$$
 then go to Al8 else

Al6:
$$\sigma_{\text{old}} := \sigma_{\text{new}}$$
; $\lambda_{\text{old}} := \lambda_{\text{new}}$; $\sigma_{\text{old}} := \sigma_{\text{new}}$; go to Al0;

4

Part 3 (We search the interval $[\sigma_1^*, \sigma_1^*]$)

Al7: If $\lambda < \sigma_n^t$ then $\lambda := \sigma_n^t$; $\lambda_{old} := \lambda$; and go to A3 but instead of steps A6 and A11 we shall have the following step in both cases:

If $\lambda > \sigma_1^t$ then system controllable, stop; else go to A10 but instead of step A11 we shall have the above mentioned step.

Part 4 (Newton's method)

A18:
$$\lambda := \lambda + \frac{\sigma}{\sigma_{\text{new}}}$$

Al9: If
$$|\lambda - \lambda_{new}| > M^{(*)}$$
 then go to Al6;

A20:
$$\lambda_{\text{old}} := \lambda_{\text{new}}$$
; $\lambda_{\text{new}} := \lambda$;

A21: Evaluate
$$\sigma_n(B,A-\lambda I)$$
; $\sigma_{old} := \sigma_{new}$; $\sigma_{new} := \sigma_n(B,A-\lambda I)$;

A22: If
$$\sigma_n(B,A-\lambda I) < \varepsilon \sigma_1^{\gamma}$$
 then system uncontrollable, stop

else
$$\sigma_{\text{new}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\lambda_{\text{new}} - \lambda_{\text{old}}}$$
; go to A18;

Now we will extend the algorithm A for $\lambda \in \mathbb{C}$. As we saw in remark (X), we do not have to search all the ring but only the part of the anulus which is above the real axis. In order to be able to apply an algorithm similar to algorithm A we will first do the following:

^(*) There might be some case where $|\dot{\sigma}_{new} - \dot{\sigma}_{old}| < \delta$ in Al5 but we are not in a neighborhood of the possible root, so we must have a bound and a check built in for safety (see fig. 5).

Let us imagine a plane P which contains the axis $Re(\lambda)$ and is vertical to the complex plane. Now we take the projection of the surface which is above the upper half and is defined by the function $\sigma_n(B,A-\lambda I)$ onto the plane P. The curve we get on plane P is similar with the curve we had for $\lambda \in \mathbb{R}$ (see fig. 10 and fig. 11). The most important thing is that remark (IX) can be applied in this case too. So algorithm A can also be applied in this case with the following changes.

Let us assume that $\lambda_1 = \alpha_1 + j\beta_1$, where $j^2 = -1$, then we start with $\alpha_1 = -\sigma_1^i$ and $\beta_1 = 0$, so $\lambda_1 = -\sigma_1^i$. The corresponding value of the function is $\sigma_n^{(1)} = \sigma_n(B,A-\lambda_1I)$. At the next step we have the following: $\alpha_2 = \alpha_1 + \sigma_n^{(1)}$, but we do not know β_2 . So we will find $\sigma_n^{(2)}$ as the minimum of the curve, which lies on the plane perpendicular to the complex plane and contains the line segment $\alpha_2 k$ (see fig. 10), considering α_2 and k as boundaries of the curve, the minimum of which we are looking for. Thus for calculating $\sigma_n^{(2)}$ we have the following:

$$\sigma_{n}^{(2)} = \min_{\beta_{2}} \sigma_{n}^{[B,A-(\alpha_{2}+j\beta_{2})I]}$$
, where $\beta_{2} \in [0, \sqrt{(\sigma_{1}^{'})^{2\zeta}-\alpha_{2}^{2}}]$ see triange α_{2} 0k fig. 10

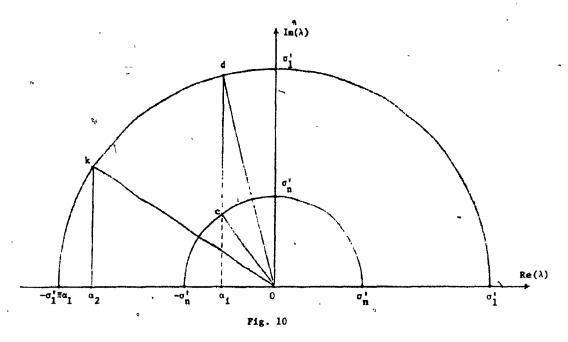
In general we have the following:

$$\sigma_{n}^{(i)} = \min_{\beta_{i}} \sigma_{n}[B, A-(\alpha_{i} + j\beta_{i})I], \qquad (31)$$

where
$$\beta_{i} \in \begin{cases} \left[0, \sqrt{(\sigma_{1}^{i})^{2} - \alpha_{i}^{2}}\right] & \text{if } \alpha_{i} \in \left[-\sigma_{1}^{i}, -\sigma_{n}^{i}\right] \cup \left[\sigma_{n}^{i}, \sigma_{1}^{i}\right] \\ \left[\sqrt{(\sigma_{n}^{i})^{2} - \alpha_{i}^{2}}, \sqrt{(\sigma_{1}^{i})^{2} - \alpha_{i}^{2}}\right] & \text{if } \alpha_{i} \in \left(-\sigma_{n}^{i}, \sigma_{n}^{i}\right) \end{cases}$$

$$(32)$$

(see fig. 10).



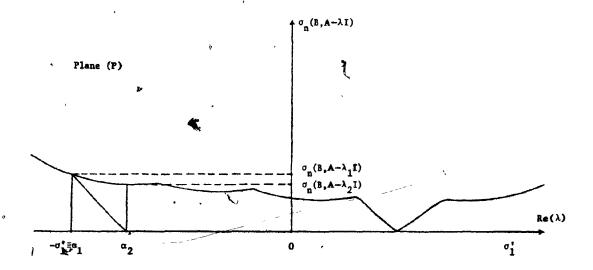


Fig. 11 (uncontrollable case)

So the algorithm for $\lambda \in \mathbb{C}$ is as follows:

Algorithm B

Part 1 (Initialization)

B1: Calculate the singular values of (B,A) and let them be $\sigma_1' \geq \sigma_2' \geq \ldots \geq \sigma_n'$;

B2: Set $\alpha := -\sigma_1^{\dagger}$; $\alpha_{\text{old}} := -\sigma_1^{\dagger}$;

B3: Evaluate $\sigma_n(B,A-\alpha I)$; $\sigma_{old} := \sigma_n(B,A-\alpha I)$;

B4: If $\sigma_{\text{old}} < \epsilon \sigma_1^{\dagger}$ then system uncontrollable, stop else

B5: $\alpha := \alpha + \sigma_{\text{old}}$; $\alpha_{\text{new}} := \alpha$;

B6: If $\alpha > \sigma_1^*$ then system controllable, stop else

B7: Evaluate $\min_{\beta} \sigma_{n}[B, A-(\alpha+j\beta)I]$ with $b_{1} \le \beta \le b_{2}$, where b_{1}, b_{2} are appropriate bounds (see (32));

Set $\sigma_{\text{new}} := \min_{\beta} \sigma_{\text{n}}[B, A-(\alpha+j\beta)I]$

B8: If $\sigma_{\text{new}} < \varepsilon \sigma_1^{\prime}$ then system uncontrollable, stop

B9: $\sigma_{\text{old}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\sigma_{\text{new}} - \sigma_{\text{old}}}; \sigma_{\text{old}} := \sigma_{\text{new}}; \alpha_{\text{old}} := \alpha_{\text{new}};$

Part 2 (We search in the interval $[-\sigma_1', \sigma_1']$)

B10: $\alpha := \alpha + \sigma_{\text{old}}$; $\alpha_{\text{new}} := \alpha$;

Bll: If $\alpha > \sigma_1^*$ then system controllable, stop else

B12: Evaluate $\min_{\beta} \sigma_{n}[B, A-(\alpha+j\beta)I]$ with $b_{1} \leq \beta \leq b_{2}$, where b_{1}, b_{2} are appropriate bounds (see (32))

Set $\sigma_{\text{new}} := \min_{\beta} \sigma_{\text{n}}[B, A-(\alpha+j\beta)I]$;

B13: If $\sigma_{\text{new}} < \varepsilon \sigma_1'$ then system uncontrollable, stop

B14:
$$\dot{\sigma}_{\text{new}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\sigma_{\text{new}} - \sigma_{\text{old}}}$$
;

B15: If
$$|\dot{\sigma}_{\text{new}}^{*} - \dot{\sigma}_{\text{old}}| < \delta$$
 then go to B17

B16:
$$\sigma_{\text{old}} := \sigma_{\text{new}}$$
; $\alpha_{\text{old}} := \alpha_{\text{new}}$; $\dot{\sigma}_{\text{old}} := \dot{\sigma}_{\text{new}}$; \dot{g}_{o} to B10;

Part 3 (Newton's method)

B17:
$$\alpha := \alpha + \frac{\sigma_{new}}{\sigma_{new}}$$
;

B17: $\alpha := \alpha + \frac{\sigma_{new}}{\dot{\sigma}_{new}}$;

B18: If $|\alpha - \alpha_{new}| > M$ then go to B16:

B19:
$$\alpha_{\text{old}} := \alpha_{\text{new}}$$
; $\alpha_{\text{new}} := \alpha$;

Evaluate $min[B, A-(\alpha+j\beta)I]$ with $b_1 \le \beta \le b_2$;

$$\sigma_{\text{old}} := \sigma_{\text{new}} ; \sigma_{\text{new}} := \min_{\beta} [B, A - (\alpha + j\beta)I] ;$$

B21: If $\sigma_{\text{new}} < \epsilon \sigma_1'$ then system uncontrollable, stop

else
$$\sigma_{\text{new}} := \frac{\sigma_{\text{new}} - \sigma_{\text{old}}}{\sigma_{\text{new}} - \sigma_{\text{old}}}$$
; go B17

Note: Obviously the algorithms A and B are quite expensive because we need to evaluate the function $\sigma_n(B,A-\lambda I)$ many times. We feel there may be a much more simple approach to this problem, but have not yet found one.

In page 32 it is pointed out that in order to decide whether a system is controllable or not it is adequate to see whether the equation $\sigma_n(B,A-\lambda I) = 0$ has any root or not, with respect to $\lambda \in \mathbb{C}$. But instead of finding the roots of $\sigma_n(B,A-\lambda I)=0$, if any, we can solve the optimization problem

$$\min_{\lambda} \sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda \mathbf{I}) , |\lambda| \in [\sigma_{\mathbf{n}}', \sigma_{\mathbf{I}}']$$
 (32\alpha)

which will not only tell us if the system is controllable or not, but in the case of a controllable system it will also tell us how far the system is from the nearest uncontrollable one, that is, it will give us the distance of the system from the nearest uncontrollable one. But we have already pointed out that even the solution of the equation

$$\sigma_{\mathbf{n}}(\mathbf{B}, \mathbf{A} - \lambda \mathbf{I}) = 0 \tag{328}$$

is quite expensive, so the minimization problem (32α) will be more expensive, at least, with the methods we have so far.

CHAPTER 3

(I) Introduction

(1) Open-loop (Nonfeedback) and closed-loop (Feedback) control systems.

A control action is required to take a controllable dynamic system from a given initial state to any other desired state. The two basic ways for carrying out such a control are the following:

(a) Open-loop control or nonfeedback:

When the determination of the control action is independent of output measurements on the system.

(b) Closed-loop control or feedback:

When the determination of the control action is based on the output behavior of the system.

Closed-loop control is often preferable to open-loop control. The following example will illustrate the difference between closed-loop and open-loop controls and will also illustrate some of the advantages of closed-loop control.

Example: Consider a home heating system. Also consider two different control systems. The first (open-loop) turns the heating on every 30 minutes and keeps it on for 20 minutes. The second (closed-loop) has a thermostat which turns on the heating when the temperature of the home drops under 20°C and turns it off when the temperature exceeds 25°C.

The advantages of the closed-loop system here over the open-loop system are obvious. While the closed-loop system will keep the temperature between

20°C - 25°C the open-loop system may make the home too hot or too cold. For more details and examples see [3], [7].

(2) Stability problem in dynamical systems

C

Stability in dynamical systems is defined with respect to equilibrium points, so we first need to give the definition of the equilibrium point.

<u>Definition (VII)</u>: A vector $\tilde{\mathbf{x}}$ is an equilibrium point of a dynamical system when it has the following property; when the state of the system reaches $\tilde{\mathbf{x}}$ it remains equal to $\tilde{\mathbf{x}}$ forever (as time approaches infinity), when there is no control input.

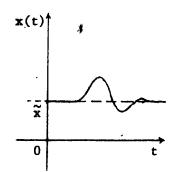
An equilibrium point is stable if when the state vector is moved slightly away from that point, it tends to return to it, or at least does not keep moving further away (marginal stability). The following example will illustrate this point.

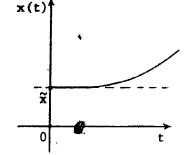
Example: Let a stick be perfectly aligned with the vertical. Then the dynamical system which describes the movements of the stick is in an equilibrium point. If the stick is balanced on its bottom end, then the slighest impulse input will destroy the balance and the stick not only is not going to return to the same equilibrium point but it will keep moving further away. So if the stick is balanced on its bottom end, it is in an unstable equilibrium point. If on the other hand it is hanging from a support at the top end and an impulse input moves it from this equilibrium point then the stick tends to return to it. So in this case the stick is in a stable equilibrium point.

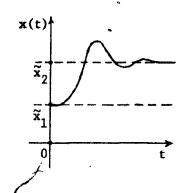
Now we will define when an equilibrium point is stable, marginally stable, or unstable.

<u>Definition (VIII)</u>: An equilibrium point $\tilde{\mathbf{x}}$ of a dynamical system is called stable when in response to an impulse input, if the system is at the equilibrium point $\tilde{\mathbf{x}}$, the output of the system tends to return to $\tilde{\mathbf{x}}$ as time increases. Otherwise $\tilde{\mathbf{x}}$ is called either an unstable or marginally stable equilibrium point according to whether the output of the system goes to infinity or settles at another equilibrium point, see [7].

In fig. 12 we give one example for each case.







a. Stable equilibrium point

c. Marginally stable equilibrium point

Fig. 12

Remark (XII): Consider the systems

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b} \tag{33}$$

$$\mathbf{x}_{\mathsf{L},\mathsf{L}} = \mathbf{A}\mathbf{x}_{\mathsf{L}} + \mathbf{b} \tag{34}$$

and let $\tilde{\mathbf{x}}$ be an equilibrium point for these two systems, then from (33) we have

from (34), we have

Setting $z = x - \tilde{x}$ and $y_k = x_k - \tilde{x}$, the equations (35) and (36) give, respectively

$$\dot{z} = Az \tag{37}$$

and

C

$$y_{k+1} = Ay_k . (38)$$

So, it is clear that the conditions

$$\lim_{t \to +\infty} x(t) = \tilde{x} \qquad \text{and} \qquad \lim_{k \to +\infty} x_k = \tilde{x}^{-1}$$

are equivalent to the conditions

$$\lim_{t \to +\infty} z(t) = 0 \qquad \text{and} \qquad \lim_{k \to +\infty} y_k = 0$$

respectively, that is, the condition for x(t) and x_k to tend to \tilde{x} in (33) and (34), respectively, is equivalent to the condition for z(t) and y_k to tend to zero in the homogeneous sytems (37), (38), respectively.

So, instead of looking for the stability of (33) and (34), we will look for the stability of (37) and (38) respectively.

(a) Stability for discrete-time systems

Consider the system

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k \tag{39}$$

then $\mathbf{x}_{k+1} = \mathbf{A}^k \mathbf{x}_1$ for the given initial state $\mathbf{x}_1 \in \mathbb{R}^n$. Therefore lim $\mathbf{x}_{k+1} = 0 \iff \lim_{k \to +\infty} \mathbf{A}^k \mathbf{x}_1 = 0 \iff |\lambda_1| < 1$ for i = 1(1)n where $\lambda_1 \in \lambda(A)$, i = 1(1)n and for every initial state $\mathbf{x}_1 \in \mathbb{R}^n$ [3, page 155]. Thus system (34) is stable if and only if the magnitude of each of the eigenvalues of A is less than one (see figure 13).

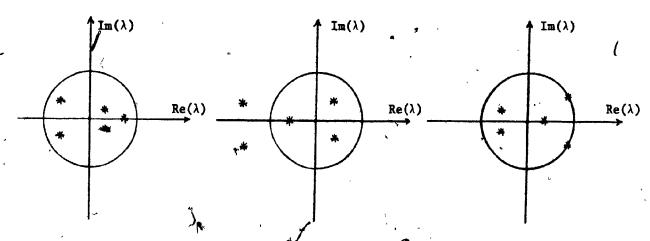
(b) Stability for continuous-time systems

Consider the system

$$\dot{\mathbf{x}}(\mathsf{t}) = \mathsf{A}\mathbf{x}(\mathsf{t}) \tag{40}$$

then $x(t) = e^{At}x(0)$ for the given initial stage $x(0) \in \mathbb{R}^n$. So, $\lim_{t \to +\infty} x(t) = 0 \iff \lim_{t \to +\infty} e^{At}x(0) = 0 \iff \operatorname{Re}(\lambda_1) < 0$ for i = 1(1)n where $\lambda_1 \in \lambda(A)$, i = 1(1)n and for every initial state $x(0) \in \mathbb{R}^n$ [3, page 157]. Therefore system (33) is stable if and only if the real parts of each of the eigenvalues of A are less than zero (see figure 14).

G



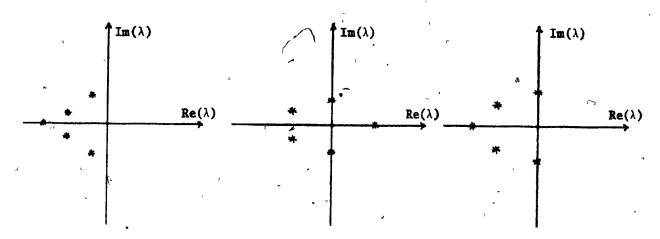
Stable system

(

b. Unstable system c. Marginally stable

system

Fig. 13 Discrete-time case



Stable system

b. Uństable syst**em** ,

c. Marginally stable

Fig. 14 Continuous-time case

いってはいれば、からからのははないのはのはないであるとなっています。ないできないということできないということできないというということ

(II) On pole assignment in single-input controllable linear systems

The problem of pole assignment or eigenvalue allocation for controllable linear systems arises as follows: consider the controllable system

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \tag{41}$$

which is made into a closed-loop system by defining

$$\mathbf{u} = \mathbf{F}\mathbf{x} \ , \tag{42}$$

where F o Rm×n. Then (41), in view of (42), gives

$$\dot{\mathbf{x}} = (\mathbf{A} + \mathbf{B}\mathbf{F})\mathbf{x} . \tag{43}$$

We want to choose F such that the system (43) is stable at every equilibrium point $\tilde{\mathbf{x}}$. So we should choose F such that the real parts of each of the eigenvalues of the matrix $\mathbf{A} + \mathbf{BF}$ are less than zero.

Note: Instead of a time-continuous system, we can also have a discrete-time system. The only difference from continuous-time systems is that in the discrete case we will choose F such that the magnitude of all the eigenvalues of the matrix A + BF is less than one.

So the pole assignment or eigenvalue allocation probelm is as follows:

Given $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$ find $F \in \mathbb{R}^{m \times n}$ such that the matrix A + BF has a desired set of eigenvalues.

Luenberger in [3] gives a theorem (p. 299) based on a controllability canonical form also discussed in the same text which can be used to derive an

algorithm to solve the above problem when m = 1, that is, when we have single-input controllable linear system. But this algorithm is based on the calculation of the coefficients of the characteristic polynomial of A, which can be achieved by calculating the eigenvalues of A ' But the calculation of the eigenvalues of a matrix A, when A is t symmetric, can be an ill-conditioned problem [13], [12], [2]. The same problem (single-input) is also considered in [32]. The generalization of the canonical forms, discussed in [3], to multi-input linear dynamical systems are discussed by Luenberger in [10]. One of these forms might be useful in the eigenvalue allocation problem (e.a.p.) for the multi-input case, but it should be pointed out that it is quite difficult to make an algorithm which solves the e.a.p. using this canonical form, and one of the reasons is that in general, there is not a unique way to bring the system to such a canonical form. One particular derivation of a canonical form of a multi-input dynamical system is given in [11] where a different approach from [10] is used. For a more explicit way of bringing a system to its canonical form see Wonham [9], where the fundamental result of pole-assignability was presented. The fact that the e.a.p. is quite difficult for the multi-input case, is emphasized in [24] where a geometric viewpoint of the problem is given. An algorithm which partially solved the e.a.p., was given by Davison [27], where he proved the following:

Let the system

$$\begin{array}{c} x = Ax + Bu \\ y = Cx \end{array}$$

where x, u, A and B are as we have already described, and $C \in \mathbb{R}^{r \times n}$ $y \in \mathbb{R}^r$, $(r \le n)$. Then if (43a) is controllable and rank(C) = r, there

is a matrix $F \in \mathbb{R}^{m \times r}$ such that r eigenvalues can be assigned to the matrix A + BFC . An algorithm is given. This result was further extended by Davison and Chatterjee [28] and by Sridhar and Lindorff [29] as follows: Under the same conditions of [27] max(r,m) eigenvalues can be assigned to A + BFC . Numerical difficulties of the algorithm given in [27], when n >> 10 led Davison and Chow to extend the result of [27] to systems where $n \gg 10$ [31]. An algorithm is given. The results of [27], [28], [29] are further extended by Davison and Wang [30] where they show that under the same conditions of [27] min(n, m+r-1) eigenvalues can be assigned to A + BFC . An algorithm is also given. A similar extension to [27], [28], [29] is given by Kimura [25] where he proves that under the same conditions of [27] an $n \le m + r - 1$, an arbitrary set of distinct eigenvalues can be assigned to A + BFC . In this paper, different from the conventional (up to that time) approach using the characteristic equation, an approach based on the properties of the eigenvectors is used. This result is further extended by Kimura in [26], where an algorithm is also given.

Although the above mentioned works are important theoretically, the resulting algorithms have usually not been designed with as wide an understanding of numerical difficulties as we have today. The main aim of this thesis is to present algorithms which are reliable numerically, as well as theoretically, and to do this we start with the most simple form of the problem.

Note: In this section we will deal with single-input controllable systems, that is, m=1 or $B\equiv b$, $b\in \mathbb{R}^n$. So that $F\equiv f^T$, $f\in \mathbb{R}^n$. The algorithm we will present will calculate an n-dimensional real vector f such that given an $n\times n$ matrix A and an n-dimensional real vector b the matrix $A+bf^T$ has a desired set of eigenvalues $\lambda_1,\lambda_2,\ldots,\lambda_n$.

Remark (XIII): For a controllable system we can always find an $F \in \mathbb{R}^{m \times n}$ such that the matrix A + BF has any desired set of eigenvalues. A proof for this can be found in [9], [10] or [11] (see also earlier here, page 14).

Remark (XIV): The problem of finding a matrix F such that the matrix A + BF has any desired set of eigenvalues is equivalent to the problem of finding a matrix \widetilde{F} such that the matrix $\widetilde{A} + \widetilde{BF}$ has the same desired set of eigenvalues as A + BF, where $\widetilde{A} = Q^TAQ$, $\widetilde{B} = Q^TBP$, $\widetilde{F} = P^TFQ$ with P and Q orthogonal matrices. The above mentioned property holds because

$$\widetilde{A} + \widetilde{BF} = Q^{T}AQ^{3} + (Q^{T}BP)(P^{T}FQ) = Q^{T}(A + BF)Q^{T}$$

So,

$$\lambda(\widetilde{A} + \widetilde{BF}) = \lambda(A + BF)$$
.

If P, Q are chosen such that \widetilde{A} and \widetilde{B} are of the form (6), then the problem is much simpler, for example in our case where m=1 we have the following:

where $\alpha_{i,i-1} \neq 0$ for i = 1(1)n, because the system defined by $(\widetilde{A}, \widetilde{b})$ is considered to be controllable. So,

$$\widetilde{A} + \widetilde{b}\widetilde{f}^{T} =
\begin{cases}
\alpha_{11}^{+\alpha_{10}\phi_{1}} & \alpha_{12}^{+\alpha_{10}\phi_{2}} & \cdots & \alpha_{1,n-1}^{+\alpha_{10}\phi_{n-1}} & \alpha_{1n}^{+\alpha_{10}\phi_{n}} \\
\alpha_{21}^{-\alpha_{22}} & \cdots & \alpha_{2,n-1}^{-\alpha_{2n}} & \alpha_{2n}^{-\alpha_{2n}} \\
& \alpha_{32}^{-\alpha_{32}} & \cdots & \alpha_{2,n-1}^{-\alpha_{2n}} & \alpha_{3n}^{-\alpha_{2n}} \\
& \vdots & \vdots & \vdots & \vdots \\
& \alpha_{n,n-1}^{-\alpha_{nn}} & \alpha_{nn}^{-\alpha_{nn}}
\end{cases}$$

Thus the original problem can be stated as follows:

Given an upper Hessenberg matrix $A \in \mathbb{R}^{n \times n}$ with elements $\alpha_{i,j}$, where i=1(1)n, j=1(1)n, and $\alpha_{i,j-1} \neq 0$, i=2(1)n, we are concerned in choosing real α_{lj} , j=1(1)n (that is, we are concerned in choosing the first row of A) in order to give A a certain desired set of eigenvalues which may occur in complex conjugate pairs. Thus from now on we will talk about determining the first row of A, rather than determining \widetilde{f} in $\widetilde{A} + \widetilde{bf}^T$.

We shall now describe a method which solves this problem. The description will be divided into two parts: as an easy introduction the first part will describe how to deal with real eigenvalues, while the second will deal with the more difficult case of eigenvalues which occur in complex conjugate pairs. The first part also will be divided into two other parts, the first part will describe an explicitly shifted method, while in the second part we shall show how the same thing can be done implicitly. As in methods for

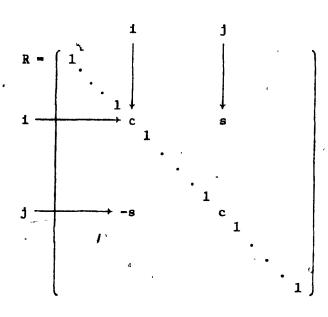
さいない、これできますしているのではないのでは、

solving the eigenproblem, the implicitly shifted approach will not only be more accurate in some cases, but will also allow the complex case to be handled in real arithmetic.

From now until the end of this chapter we will use a notation which we describe immediately.

1) Sometimes the elements of the matrices will be represented by either * , x , 0 or blank, meaning that the corresponding element is a variable, a known scalar, or in the last two cases is zero respectively. By "variable" we mean an element that has not yet been described. For example usually the first row of our upper Hessenberg matrix A is "variable".

Definition (IX) [12], [14]: Det $a = (\alpha_1, \dots, \alpha_i, \dots, \alpha_j, \dots, \alpha_n)^T$ be an n-dimensional vector with $\alpha_i \neq 0$ let also



be an orthogonal matrix such that the vector

Ra =
$$\begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_1 \\ \vdots \\ \alpha_j \\ \vdots \\ \alpha_n \end{bmatrix}$$

has its jth element zero. Thus c,s should be chosen such that if $\rho = \sqrt{\alpha_1^2 + \alpha_j^2} \quad \text{then } c = \frac{\alpha_1}{\rho} \;, \; s = \frac{\alpha_j}{\rho} \;. \quad \text{Then matrix } R \; \text{ is a rotation matrix}$ in the 1,j plane.

2) The following notation for a 3×3 matrix

C

$$\begin{pmatrix}
* & * & * \\
x & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
x & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
x & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
x & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & x & x
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & * \\
* & * & *
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & *
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & *
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & *
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & *
\end{pmatrix}$$

$$\begin{pmatrix}
* & * & * \\
* & * & *
\end{pmatrix}$$

means that, first we multiply our matrix from the right by a rotation matrix which combines the second and third columns without eliminating any element, then we apply the transpose of the same rotation from the left combining the second and third row, so a non-zero element is introduced at the place (3,1). After this the element (3,1) is eliminated into the element (3,2) by applying another rotation matrix from the right, and

finally the transpose of the same rotation matrix is applied from the left and it introduces variable elements at the second row of the matrix.

Note: All the methods we will describe in the following have two main parts, the forward sweep and the backward sweep.

(1) The case of real eigenvalues

(a) Explicitly shifted method:

from the right, where Q_1 is a product of appropriate rotation matrices determined such that $(A_1^{-\lambda}_1I)Q_1 = R_1$ where R_1 is upper triangular, that is,

$$(A_1^{-\lambda_1}I)Q_1 \equiv \begin{pmatrix} \star & \star & \star \\ \hline & \star & \star \\ \hline & & \times \\$$

Since $\alpha_{1,i-1} \neq 0$, i = 2(1)n, the (i,i) elements of (M1) will be nonzero, i = 2(1)n. Thus the element (1,1) of the matrix (M1) must be zero for (M1) to be singular; but since this element is variable we set it equal to zero, thereby using some of our freedom in choosing the first row of A.

Then we multiply $(A_1^{-\lambda}_1^{I})Q_1$ from the left by Q_1^T and we get

$$Q_{1}^{T}(A_{1}-\lambda_{1}I)Q_{1} \equiv \begin{pmatrix} 0 & \star & \star \\ 0 & x & x \end{pmatrix} 2^{t} \\ 0 & x & x \end{pmatrix} 1^{t} \rightarrow \begin{pmatrix} 0 & \star & \star \\ 0 & \star & \star \\ 0 & x & x \end{pmatrix} \tag{M2}$$

We set the row vector which is formed by the elements (1,2), (1,3) of (M2) equal to b_1^T and the matrix which is formed by the elements (2,2), (2,3), (3,2), (3,3) equal to $A_2-\lambda_1 I$. So we have the following:

$$Q_{1}^{T}(A_{1}-\lambda_{1}I)Q_{1} = \begin{pmatrix} 0 & b_{1}^{T} \\ 0 & A_{2}-\lambda_{1}I \end{pmatrix} = \begin{pmatrix} \lambda_{1} & b_{1}^{T} \\ 0 & A_{2} \end{pmatrix} (44)$$

From relation (44) we can obviously see that $\lambda(A_1) = \{\lambda_1\} \cup \lambda(A_2)$. Note that the rotations in (M1) were nontrivial since $\alpha_{1,i-1} \neq 0$, i = 2(1)n-1 and the (2,2) element of (M1) is nonzero, as a result the (3,2) element in (M2) is nonzero and A_2 has the same form as A_1 (the general case will be proved in remark (XV)).

In the next step we will work with A_2 , the first row of which can be chosen arbitrarily and we will choose an orthogonal $(n-1)\times(n-1)$ matrix Q_2 such that the matrix $Q_2^T(A_2-\lambda_2I)Q_2$ is singular. At the kth step, we choose an orthogonal $(n-k+1)\times(n-k+1)$ matrix, Q_k , such that the matrix $Q_k^T(A_k-\lambda_kI)Q_k$ is singular, that is,

$$Q_{\mathbf{k}}^{\mathbf{T}}(\mathbf{A}_{\mathbf{k}}^{-\lambda}\mathbf{k}^{\mathbf{I}})Q_{\mathbf{k}} = \begin{pmatrix} 0 & \mathbf{b}_{\mathbf{k}}^{\mathbf{T}} \\ 0 & \mathbf{b}_{\mathbf{k}}^{\mathbf{T}} - \lambda_{\mathbf{k}}^{\mathbf{I}} \end{pmatrix} \Rightarrow Q_{\mathbf{k}}^{\mathbf{T}}\mathbf{A}_{\mathbf{k}}Q_{\mathbf{k}} = \begin{pmatrix} \lambda_{\mathbf{k}} & \mathbf{b}_{\mathbf{k}}^{\mathbf{T}} \\ 0 & \mathbf{A}_{\mathbf{k}+1}^{\mathbf{I}} - \lambda_{\mathbf{k}}^{\mathbf{I}} \end{pmatrix}, \quad \mathbf{k}=1(1)\mathbf{n}-1$$

$$(45)$$

so that $\lambda(A_k) = \{\lambda_k\} \cup {}_{b}\lambda(A_{k+1})$ for k = 1(1)n-1.

€.

At the final step, when $\,k$ = n-1 , we choose an orthogonal 2×2 matrix Q_{n-1} such that

$$Q_{n-1}^{T}(A_{n-1}^{-1}-\lambda_{n-1}^{-1})Q_{n-1} = \begin{pmatrix} 0 & \beta_{n-1} \\ 0 & A_{n}^{-1} \end{pmatrix} = \begin{pmatrix} 0 & \beta_{n-1} \\ 0 & A_{n-1}^{-1} \end{pmatrix}$$

$$Q_{n-1}^{T} A_{n-1} Q_{n-1} = \begin{pmatrix} \lambda_{n-1} & \beta_{n-1} \\ 0 & A_{n} \end{pmatrix}$$

$$(46)$$

But the matrix A_n is well known because it is a scalar, that is, $A_n \in \mathbb{R}$ and since it is a variable we set it equal to λ_n .

Now we will go backwards and we will be calculating one by one the first row of the matrices A_k for k=n-1(-1)1. Finally we will calculate the first row of A, since $A\equiv A_1$.

Now we will describe the kth step of the backward sweep. At the kth step of the backward sweep we have from (45)

$$Q_{\mathbf{k}}^{\mathbf{T}}(\mathbf{A}_{\mathbf{k}}^{-\lambda} \mathbf{A}_{\mathbf{k}}^{\mathbf{I}}) Q_{\mathbf{k}} = \begin{bmatrix} 0 & \mathbf{b}_{\mathbf{k}}^{\mathbf{T}} \\ 0 & \mathbf{A}_{\mathbf{k}+1}^{-\lambda} \mathbf{A}_{\mathbf{k}}^{\mathbf{I}} \end{bmatrix}$$

$$(47)$$

From the previous. k-1 steps in the backward sweep we will know the first row of A_{k+1} and so the first row a_{k+1}^T of $A_{k+1}^{-\lambda}{}_kI$ and we want to evaluate the first row a_k^T of $A_k^{-\lambda}{}_kI$ and from that the first row of A_k . Let

$$Q_k = P_1 P_2 \dots P_{n-k}$$
 (48)

where P_i , i=1(1)n-k are the known rotation matrices (see Definition (IX)). We know that first we applied the transformation Q_k from the right, and since this transformation combines the columns of $A_k^{-\lambda}{}_k I$, we see that every one of the rotation matrices P_i , i=1(1)n-k affects the first row of $A_k^{-\lambda}{}_k I$. But the same thing does not happen with Q_k^T , which combines the rows of $A_k^{-\lambda}{}_k I$. As we can see, the only rotation matrix of Q_k^T which affects the first row of $A_k^{-\lambda}{}_k I$ is the last one, that is P_{n-k}^T , because P_{n-k}^T combines the first with the second row of (the altered) $A_k^{-\lambda}{}_k I$. Let

$$\begin{pmatrix}
\mathbf{e_1^T} \\
\mathbf{e_2^T}
\end{pmatrix}
\mathbf{P_{n-k-1}^T} \cdot \cdot \cdot \mathbf{P_1^T} (\mathbf{A_k^{-\lambda_k I}}) \mathbf{Q_k} = \begin{pmatrix}
0 & \mathbf{\bar{a}_k^T} \\
0 & \mathbf{\bar{a}_k^T}
\end{pmatrix} (49\alpha)$$

Note that $a \stackrel{T}{k}$ will be available because at the corresponding stage of the forward sweep we will not combine it with the first row which contains variable elements.

Let

$$P_{n-k}^{T}$$
 then we have (498)

$$\begin{pmatrix}
\mathbf{c}_{k} & \mathbf{s}_{k} \\
-\mathbf{s}_{k} & \mathbf{c}_{k}
\end{pmatrix}
\begin{pmatrix}
\mathbf{0} & \bar{\mathbf{a}}_{k}^{T} \\
\mathbf{0} & \mathbf{a}_{k}^{T}
\end{pmatrix}
=
\begin{pmatrix}
\mathbf{0} & \mathbf{b}_{k}^{T} \\
\mathbf{0} & \mathbf{a}_{k+1}^{T}
\end{pmatrix}
=
\begin{pmatrix}
\mathbf{e}_{1}^{T} \\
\mathbf{e}_{2}^{T}
\end{pmatrix}
\begin{pmatrix}
\mathbf{0} & \mathbf{b}_{k}^{T} \\
\mathbf{e}_{2}^{T}
\end{pmatrix}
=> (50)$$

$$c_k a_k^T - s_k \overline{a}_k^T = a_{k+1}^T$$

$$\overline{a}_{k}^{T} = \frac{c_{k}a_{k}^{T} - a_{k+1}^{T}}{s_{k}}, \text{ since } s_{k} \neq 0, \qquad (51)$$

which determines the "variable row \overline{a}_{k}^{T} . But from (49a) and the form of the P_{1} ,

$$\mathbf{e}_{1}^{T}(\mathbf{A}_{k}-\lambda_{k}\mathbf{I})\mathbf{Q}_{k} = \mathbf{a}_{k}^{T}\mathbf{Q}_{k} = (0, \overline{\mathbf{a}}_{k}^{T})$$

So

$$a_k^T = (0, \bar{a}_k^T) Q_k^T$$
 (52)

From (52) the first row of A_k is

$$\mathbf{e}_{1}^{T}\mathbf{A}_{k} = (0, \overline{\mathbf{a}}_{k}^{T})\mathbf{Q}_{k}^{T} + \lambda_{k}\mathbf{e}_{1}^{T}$$
 (52a)

Thus from knowing the first row of A_{k+1} we can find the first row of A_k . But we know the first (only) row of A_n in (46), and so we can go backwards to finally produce the first row of $A \equiv A_1$, our desired result.

<u>Comment</u>: From (51) we can see that if s_k is relatively small then we may get into numerical trouble. So from now until the end of this section we will comment of this possible trouble.

<u>Definition (X)</u>: Let $A \in \mathbb{C}^{n \times n}$ be an upper Hessenberg matrix with elements α_{ij} , i = 1(1)n, j = 1(1)n. Then we say A is unreduced if $\alpha_{i+1,i} \neq 0$ for i = 1(1)n-1.

Remark (XV): If we write the result of the kth step of the forward sweep of the explicitly shifted method as

$$Q_{\mathbf{k}}^{\mathbf{T}} \mathbf{A}_{\mathbf{k}} Q_{\mathbf{k}} = \left(\begin{array}{cc} \lambda_{\mathbf{k}} & \mathbf{b}_{\mathbf{k}}^{\mathbf{T}} \\ 0 & A_{\mathbf{k+1}} \end{array} \right)$$

then if A_k is unreduced, A_{k+1} will also be unreduced. This result follows because a zero (i,i-1) element of A_{k+1} implies that either the corresponding (i,i-1) element of A_k is zero or the corresponding (i+1,i) element of A_k is zero. This can be seen from the following example of a 3×3 matrix. In the forward sweep of the explicitly shifted method, one step is described by

$$\begin{bmatrix}
\begin{pmatrix} * & * & * \\
x & x & x \\
\hline
& x & x
\end{bmatrix}$$

$$\begin{bmatrix}
\begin{pmatrix} * & * & * \\
x & x & x \\
\hline
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
& & & \\
&$$

Let α_{ij}^k represent the (i,j)-element of the kth matrix in this example. Suppose the first matrix is unreduced, so neither rotation is trivial, and suppose $\alpha_{32}^5 = 0$. Then we have either

$$\alpha_{32}^5 = 0$$
 => $\alpha_{32}^4 = 0$ => $\alpha_{22}^3 = 0$ => $\alpha_{21}^2 = 0$ => $\alpha_{21}^1 = 0$,

or

$$\alpha_{32}^5 = 0$$
 => $\alpha_{32}^1 = 0$ since rotation 1 will not exist.

But in both cases we have a contradiction, so $\alpha_{32}^5 \neq 0$. More formally, after some trivial calculations we get: First, from rotation 1

$$\alpha_{32}^{1} c_{1} + \alpha_{33}^{1} s_{1} = 0 , \frac{c_{1}^{2}}{s_{1}^{2}} = \frac{(\alpha_{33}^{1})^{2}}{(\alpha_{32}^{1})^{2}} , \text{ so if } \rho_{1}^{2} = (\alpha_{33}^{1})^{2} + (\alpha_{33}^{1})^{2} , s_{1} = -\frac{\alpha_{32}^{1}}{\rho_{1}^{2}}$$

Similarly from rotation 2

$$\rho_2^2 = (\alpha_{21}^2)^2 + (\alpha_{22}^2)^2 = (\alpha_{21}^1)^2 + (\alpha_{22}^2)^2$$

Now we can calculate a_{32}^5 as follows:

$$\alpha_{32}^{5} = \alpha_{32}^{4}$$

$$= -s_{1}\alpha_{22}^{3}$$

$$= \frac{\alpha_{32}^{1}}{\rho_{1}} \rho_{2}$$

$$= \frac{\alpha_{32}^{1}}{\rho_{1}} \left[(\alpha_{21}^{1})^{2} + (\alpha_{22}^{2})^{2} \right]^{1/2}$$
(528)

were ρ_1, ρ_2 are well known numbers associated with rotations 1 and 2 respectively (see definition (IX)), and

$$\alpha_{22}^2 = \frac{1}{\rho_1} \det \begin{bmatrix} \alpha_{22}^1 & \alpha_{23}^1 \\ \alpha_{32}^1 & \alpha_{33}^1 \end{bmatrix}$$
 (53)

From (528) we can see that if, $\alpha_{32}^5 = 0$ then either $\alpha_{32}^1 = 0$ or both $\alpha_{21}^1 = 0$, $\alpha_{22}^2 = 0$.

So if A_k is unreduced then $A_k^{-\lambda}{}_kI$ is unreduced which implies that $A_{k+1}^{-\lambda}{}_kI$ and so A_{k+1} are unreduced too. But we know that A is unreduced because of the controllability of the system $\dot{x} = Ax + bu$ and theorem (I) therefore A_1, A_2, \dots, A_{n-1} are all unreduced.

Now since each matrix $A_k - \lambda_k I$ is unreduced its (2,1)-element is non zero, so the rotation P_{n-k} which is made to eliminate this element is non trivial, thus $s_k \neq 0$. But s_k can be relatively small if some of the subdiagonal

elements of $A_k - \lambda_k I$ are small (this can be seen from (528), if α_{32}^1 is small in comparison with α_{33}^1 the α_{32}^5 may be small and so will be the corresponding s used to eliminate α_{32}^5 if α_{32}^5 is not small, since $a_{32}^5 = \frac{\alpha_{32}^5}{\sqrt{(\alpha_{32}^5)^2 + (\alpha_{33}^5)^2}}$.

Now from (51) we have

$$s_{k}^{T} = c_{k}^{T} = c_{k$$

Let a_k^T be the first row of A_k , k = 1(1)n. Then from (52a) we know that

$$(0, \bar{\mathbf{a}}_{\mathbf{k}}^{\mathrm{T}})\mathbf{Q}_{\mathbf{k}}^{\mathrm{T}} = \tilde{\mathbf{a}}_{\mathbf{k}}^{\mathrm{T}} - \lambda_{\mathbf{k}}\mathbf{e}_{\mathbf{l}}^{\mathrm{T}}$$
(54a)

We also know that from the second row of (47)

$$(0, \mathbf{a}_{k+1}^{T}) = (0, \tilde{\mathbf{a}}_{k+1}^{T}) - \lambda_{k} \mathbf{e}_{2}^{T}$$
 (548)

Now (54) in view of (54 α) and (54 β) gives us

$$s_{k}(\widetilde{a}_{k}^{T} - \lambda_{k}e_{1}^{T}) + (0, \ \widetilde{a}_{k+1}^{T})Q_{k}^{T} - \lambda_{k}e_{2}^{T}Q_{k}^{T} = c_{k}(0, \ a^{T}_{k})Q_{k}^{T} < \Longrightarrow$$

$$s_{k}\widetilde{a}_{k}^{T} + (0, \ \widetilde{a}_{k+1}^{T})Q_{k}^{T} = c_{k}(0, \ a^{T}_{k})Q_{k}^{T} + s_{k}\lambda_{k}e_{1}^{T} + \lambda_{k}e_{2}^{T}Q_{k}^{T}$$
(54 γ)

Let

$$c_{k}(0, a_{k}^{T})Q_{k}^{T} + a_{k}\lambda_{k}e_{1}^{T} + \lambda_{k}e_{2}^{T}Q_{k}^{T} = d_{k}^{T},$$
 (546)

which we see is wholly available from the forward sweep then (54γ) in view of (54δ) gives us

$$a_{k} \tilde{a}_{k}^{T} + (0, \tilde{a}_{k+1}^{T})Q_{k}^{T} = d_{k}^{T}$$
 $k = n-1(-1)1$ (54\vec{\varepsilon}{\varepsilon})

Now let
$$\hat{\mathbf{a}}_{k}^{T} = (\underbrace{0, \dots, 0}_{k-1}, \tilde{\mathbf{a}}_{k}^{T})$$
, $\hat{\mathbf{d}}_{k}^{T} = (\underbrace{0, \dots, 0}_{k-1}, \mathbf{d}_{k}^{T})$, $\hat{\mathbf{Q}}_{k} = \begin{bmatrix} \mathbf{I}_{k-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_{k} \end{bmatrix}$

for k = l(1)n, so \hat{a}_k^T , \hat{d}_k^T are n-dimensional row vectors and \hat{Q}_k is an $n \times n$ orthogonal matrix, then instead of (54 ϵ) we have the following:

$$s_{k}\hat{a}_{k}^{T} + \hat{a}_{k+1}^{T} \hat{Q}_{k}^{T} = \hat{d}_{k}^{T} <=> -$$

$$s_{k}\hat{a}_{k}^{T}\hat{Q}_{k} + \hat{a}_{k+1}^{T} = \hat{d}_{k}^{T} \hat{Q}_{k} , \quad k = n-1(-1)1$$
(545)

Now from (545) we can observe the following:

For
$$k = n-1$$
 we have $\begin{pmatrix} s_{n-1} \\ 1 \end{pmatrix}$, $\begin{pmatrix} \hat{a}_{n-1}^T \\ \hat{a}_{n-1}^T \end{pmatrix} = \hat{d}_{n-1}^T \hat{Q}_{n-1}$

For
$$k = n-2$$
 we have $(s_{n-2}, 1)$ $\left\{ \begin{array}{l} \hat{a}_{n-2}^T & \hat{Q}_{n-2} & \hat{Q}_{n-1} \\ \\ \hat{a}_{n-1}^T & \hat{Q}_{n-1} \end{array} \right\} = \hat{d}_{n-2}^T \hat{Q}_{n-2} \hat{Q}_{n-1}$

So finally we get the system:

$$\begin{bmatrix} \mathbf{s}_1 & \mathbf{1} & & & & & \\ & \mathbf{s}_2 & \mathbf{1} & & & & \\ & & \mathbf{x}_2^T & & & & & \\ & & \mathbf{x}_2^T & & & & & \\ & & & \mathbf{x}_{n-2}^T & & & & \\ & & & & & \mathbf{x}_{n-1}^T & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ \end{bmatrix} \qquad \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ \end{bmatrix} \qquad \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ & & \\ &$$

where
$$\mathbf{x}_{i}^{T} = \hat{\mathbf{a}}_{i}^{T} \hat{\mathbf{Q}}_{i} \hat{\mathbf{Q}}_{i+1} \dots \hat{\mathbf{Q}}_{n-1} \hat{\mathbf{Q}}_{n}$$
 and $\mathbf{g}_{i}^{T} = \hat{\mathbf{d}}_{i}^{T} \hat{\mathbf{Q}}_{i} \hat{\mathbf{Q}}_{i+1} \dots \hat{\mathbf{Q}}_{n-1} \hat{\mathbf{Q}}_{n}$, $i = 1(1)n$ with $\mathbf{Q}_{n} = 1$, (540)

and since $\widetilde{\mathbf{a}}_n^T = \lambda_n$ (see (46)), the first (and only) row of A_n , we have taken $s_n = 1$, $g_n^T = (0, \dots, 0, \lambda_n)$, giving $\mathbf{x}_n^T = g_n^T$. Note that $\widehat{\mathbf{a}}_1^T$ is the desired first row, and knowing \mathbf{x}_1^T we know $\widehat{\mathbf{a}}_1^T = \mathbf{x}_1^T \ \widehat{\mathbf{Q}}_n^T \ \widehat{\mathbf{Q}}_{n-1}^T$. In the system of equations (54 η) the unknowns are effectively the rows $\widetilde{\mathbf{a}}_k^T$, k = 1(1)n-1. So system (54 η) can be solved from $\widetilde{\mathbf{a}}_n^T$ to $\widetilde{\mathbf{a}}_1^T = \widehat{\mathbf{a}}_1^T$. It should be pointed out that this effective solution of equations appears to be as late in the computation as possible, and probably reffects the ill-condition of the problem. Thus there is some reason to believe that this algorithm will be numerically stable, since there is no other solution of equations, and all other computations are nice stable rotations. The ratio of the largest to smallest singular values of the matrix in (54 η) may well give a measure of the sensitivity of the eigenvalue allocation problem.

(b) Implicitly shifted method:

Here we will describe a method in which every step of the forward sweep is achieved without having to subtract the eigenvalue from the diagonal and then restore it. We first describe and prove a theorem which will help us develop the implicitly shifted method.

Theorem (V): Let A, Q and H be real $n \times n$ matrices with $Q^TQ = I$, and H upper Hessenberg with elements $n_{i+1,i} > 0$ for i = m+1(1)n-1 where $1 \le m \le n-2$.

3

where e_1 and e_m have n-m and m elements respectively. If AQ = QH then Q_2 , H_2 and $R(Q_1)$ are uniquely determined by A and the last column of Q_2 . If $n_{m+1,m} = 0$ then $R(Q_1)$ is an eigensubspace of A. If m = 1 then apart from the sign of the first column of Q, Q and H are uniquely determined.

<u>Proof</u>: Let $Q = (q_1, q_2, ..., q_n)$ be partitioned into columns. Suppose that we already know q_n . Then we have

$$HQ^{T} = Q^{T}A \tag{55}$$

Suppose we have already calculated the last n-k columns of Q, that is, we have calculated the q_1 for i=k+1(1)n. Also suppose that we have already calculated the last n-k-1 rows of H. Now we will calculate the (k+1)th row of H and the kth column of Q. From (55) and the fact that H is upper Hessneberg, we obtain, with $H = (n_{11})$

$$\eta_{k+1,k}q_k^T + \eta_{k+1,k+1}q_{k+1}^T + \dots + \eta_{k+1,n}q_n^T = q_{k+1}^T A$$
, $k = n-1(-1)1$. (56)

Since Q is orthogonal, we multiply (56) by q from the right to obtain

$$n_{k+1,i} = q_{k+1}^T A q_i$$
, $i = n(-1)k+1$. (57)

Now we will calculate q_k and $n_{k+1,k}$. From (56), we obtain

$$q_k^T = \eta_{k+1,k}^{-1} (q_{k+1}^T A - \sum_{i=k+1}^n \eta_{k+1,i} q_i^T)^*, k = n-1(-1)m+1$$
 (58)

Here (58) and the requirements $||\mathbf{q}_k||_2 = 1$ and $\eta_{k+1,k} > 0$ uniquely determine \mathbf{q}_k and $\eta_{k+1,k}$. This way we obtain all but the first m columns of Q and the first m+1 rows of H. The rest of \mathbf{H}_2 is given by $\mathbf{H}_2 = \mathbf{Q}_2^T \mathbf{A} \mathbf{Q}_2$, and $\mathcal{R}(\mathbf{Q}_1)$ is that space orthogonal to \mathbf{Q}_2 . Now

$$AQ = QH (59)$$

so if $\eta_{m+1,m} = 0$ we have $AQ_1 = Q_1H_1$, $Q_1^TQ_1 = I$, and $R(Q_1)$ is an eigensubspace of A. If m = 1 then $Q_1 = q_1$ is uniquely determined but for its sign, so $H = Q^TAQ$ is uniquely determined to this extent. If m=1 and $\eta_{2,1} = 0$ then

$$Aq_1 = q_1 \eta_{1,1},$$
 (59a)

so q_1 is an eigenvector of A corresponding to the eigenvalue $\lambda = \eta_{11}$.

Remark (XVI): The requirement that $n_{k+1,k}$ be positive in the statement of the theorem was necessary only to tie down uniqueness. Actually $n_{k+1,k}$, k=m+1(1)n-1 and q_k , k=m+1(1)n are determined up to a constant factor of absolute value unity simply by the requirement that $n_{k+1,k} \neq 0$, k=m+1(1)n-1. It is this essential uniqueness of H and Q that we shall use.

. To see how theorem (V) can be applied to improve the explicitly shifted method, suppose that one step of the forward sweep with shift λ has been applied to the unreduced upper Hessenberg matrix A to yield an upper

Hessenberg matrix H such that H = $\begin{pmatrix} \lambda & b^T \\ 0 & \widetilde{A} \end{pmatrix}$ with \widetilde{A} an unreduced

upper Hessenberg matrix. Then,

 $H = Q^{T}AQ$, and Q is lower Hessenberg.

The following algorithm is an alternate way of computing $\mbox{\em H}$.

1) Find an orthogonal matrix P such that

$$Pe_{n} = Qe_{n} \tag{60}$$

that is, P and Q have the same last column. To do this choose the rotation P^T such that $P^TQe_n = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = e_n$, more will be said on this later.

- 3) Find orthogonal matrices U_i , i = 1(1)n-2 to reduce P^TAP to upper Hessenberg form H', that is

$$\mathbf{U}_{n-2}^{\mathsf{T}} \dots \mathbf{U}_{1}^{\mathsf{T}} \mathbf{P}^{\mathsf{T}} \mathbf{A} \mathbf{P} \mathbf{U}_{1} \dots \mathbf{U}_{n-2} = \mathbf{H}'$$
 (61)

₹

where continuing the same example H' would have the form

$$\left(\begin{array}{cccc}
\mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\
\mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\
2 & \mathbf{x} & \mathbf{x} & \mathbf{x}
\end{array}\right)$$

If we set $Q' = PU_1 \cdot \cdot \cdot U_{n-2}$, then from (61) we obtain

$$Q^T'AQ' = H'$$

and we see from the form of the U_1 that $U_1e_n=e_n$, so $Q^*e_n\stackrel{*}{,}Pe_n=Qe_n$ and from Theorem (V) we have effectively $Q^*=Q$, $H^*=H$. As a result we set the (unknown) (2,1) element of H^* to zero, and (1,1) element to λ .

To determine P we should find the last row of Q^T . We know that $(A-\lambda I)Q=R$ where R is upper triangular matrix with elements ρ_{ij} , i=1(1)n, j=1(1)n. Now

$$(A-\lambda I)Q = R =>$$

$$A - \lambda I = RQ^{T} =>$$

$$e_{n}^{T}(A - \lambda I) = \rho_{nn}^{T} q_{n}^{T}, \qquad (62)$$

where $Q = (q_1, q_2, \dots, q_n)$ is a multiple of the last row a^T of $A = \lambda I$. Since A is upper Hessenberg we know that

$$a^T = (0, \dots, 0, \alpha_{n,n-1}, \alpha_{nn} - \lambda)$$

,

Then we choose P such that

$$\mathbf{a}^{\mathrm{T}}\mathbf{P} = \pm ||\mathbf{a}||_{2} \mathbf{e}_{\mathbf{n}}^{\mathrm{T}}$$
 (62a)

(a P like this can always be found, see [12], page 232, theorem 3.4). So from (62α) we get

$$\mathbf{a}^{\mathrm{T}} = \pm ||\mathbf{a}||_{2} \mathbf{e}_{n}^{\mathrm{T}} \mathbf{P}^{\mathrm{T}} \implies$$

$$\mathbf{e}_{n}^{\mathrm{T}} \mathbf{P}^{\mathrm{T}} = \pm \frac{\mathbf{a}^{\mathrm{T}}}{||\mathbf{a}||_{2}}$$
(63)

So using the P of (63), the last row of which is a multiple of the last row of $A - \lambda I$, we accomplish our aim, to find a matrix P such that $P = Qe_{n}$.

Now we will give an example with a 3×3 matrix to demonstrate how the forward sweep of the implicitly shifted method works.

we apply it to the matrix $A - \lambda I$ from the right it will eliminate the (3,2)-element of $A - \lambda I$ into the (3,3)-element of the same matrix. Of course, P will not do the same to the matrix A unless $\lambda = 0$. So we will have the following:

$$\begin{bmatrix}
1 & * & * & * \\
x & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
x & x & x
\end{bmatrix}
\xrightarrow{1} \begin{bmatrix}
* & * & * \\
x & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
x & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & * & *
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{2} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{3} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x
\end{bmatrix}
\xrightarrow{4} \begin{bmatrix}
* & * & * \\
* & x & x$$

Now we set the (1,1)-element of the 5th matrix equal to λ , so the (2,1)-element of the same matrix becomes zero in view of the theorem (V). Finally, we have

(

$$Q^{T}AQ = \left(\begin{array}{cc} \lambda & b^{T} \\ 0 & \widetilde{A} \end{array}\right)$$

That is what exactly we would have, using the explicitly shifted method, but without subtracting the shift from the diagonal of A and then restore it back.

Now we will describe the kth step of the backward weep of the implicitly shifted method, which as we will see is almost e same as the backward sweep of the explicitly shifted method. At the kth step, we have

$$Q_{k}^{T} A_{k} Q_{k} = \begin{pmatrix} \lambda_{k} & b_{k}^{T} \\ 0 & A_{k+1} \end{pmatrix} . \tag{63a}$$

As in the explicitly shifted method we know the first row a_{k+1}^T of A_{k+1} and we want to find the first row a_k^T of A_k . Let $(\bar{\lambda}_k$, $\bar{a}_k^T)$ be the first row of $A_k^TQ_k$ before the effect of applying the orthogonal matrix Q_k^T on the left, or, as in the explicitly shifted method, before the effect of the rotation matrix P_{n-k}^T as this is the only rotation matrix of the matrix Q_k^T that will affect the first row of $A_k^TQ_k$, then we have

$$\begin{pmatrix}
c_{k} & s_{k} \\
-s_{k} & c_{k}
\end{pmatrix} \qquad
\begin{pmatrix}
\overline{\lambda}_{k} & \overline{a}_{k}^{T} \\
\gamma & a_{k}^{T}
\end{pmatrix} \qquad
\begin{pmatrix}
\lambda_{k} & b_{k}^{T} \\
0 & a_{k+1}^{T}
\end{pmatrix} \qquad (63\beta)$$

from (638) and since (γ, a, T) is well known, we obtain

$$c_{k}^{\gamma} - s_{k}^{\overline{\lambda}}_{k} = 0 \iff$$

$$\dot{\lambda}_{k} = \frac{c_{k}^{\gamma}}{s_{k}} \tag{64}$$

and

$$c_k a_k^T - s_k a_k^T = a_{k+1}^T \iff$$

$$\bar{\mathbf{a}}_{k}^{T} = \frac{c_{k} \mathbf{a}^{t}_{k}^{T} - \mathbf{a}_{k+1}^{T}}{s_{k}}$$
 (65)

So now we know the first row $(\bar{\lambda}_k^{}$, $\bar{a}_k^T)$ of $A_k^{}Q_k^{}$, that is

$$\mathbf{a}_{k}^{T}\mathbf{Q}_{k} = (\overline{\lambda}_{k}, \overline{\mathbf{a}}_{k}^{T}) \quad \Longleftrightarrow$$

$$\cdot \mathbf{a}_{k}^{T} = (\overline{\lambda}_{k}, \overline{\mathbf{a}}_{k}^{T}) \mathbf{Q}_{k}^{T}. \tag{66}$$

Thus from knowing the first row of A_{k+1} we can find the first row of A_k . But we know the first (only) row of $A_n = \lambda_n$ in (63 α) with k = n - 1 and so we can go backwards to finally produce the first row of $A = A_1$, our desired result.

We now make one remark and prove a theorem based on the remark. This theorem will help us develop an implicitly shifted method for the case of complex eigenvalues.

Remark (XVII): It may have been noticed that the explicitly and the implicitly shifted methods which have been described here are quite similar with the QR explicitly and implicitly shifted algorithms, respectively. The basic difference is that in our methods instead of trying to get an upper triangular matrix R by multiplying a matrix $A^2 \mu I$, $\mu \in \mathbb{R}$ from the left by an orthogonal matrix Q^T , we multiply $A - \mu I$ from the right, where μ is a well known shift. So the kth step of an algorithm based on our explicitly shifted method, to calculate the eigenvalues of A, would be as follows:

- (1) Form $A_k \mu_k I$
- (2) Factor the above matrix in the form

$$(A_{k} - \mu_{k}I)Q_{k} = R_{k} \qquad \stackrel{\text{\tiny (67)}}{=}$$

$$A_{k} - \mu_{k} I = R_{k} Q_{k}^{T} \tag{67a}$$

(3) Finally compute A_{k+1} as

$$A_{k+1} = Q_k^T R_k + \mu_k I$$
 (68)

From (67) and (68) we obtain:

$$\mathbf{A}_{k+1} = \mathbf{Q}_k^{\mathrm{T}} \mathbf{A}_k \mathbf{Q}_k . \tag{68a}$$

Theorem (VI): Let $\tilde{Q}_m = Q_1 Q_2 \dots Q_m$ and $\tilde{R}_m = R_1 R_2 \dots R_m$, where Q_1, R_1 with i = 1(1)m are the orthogonal and upper triangular matrices of the ith step of the above mentioned algorithm applied for m steps, then

$$\widetilde{R}_{m}\widetilde{Q}_{m}^{T} = (A - \mu_{1}I) (A - \mu_{2}I) \dots (A - \mu_{m}I) .$$
 (69)

3

<u>Proof</u>: Let $A = A_1$. Then obviously, we obtain

(.)

$$A_{k+1} = \widetilde{Q}_k^T \wedge \widetilde{Q}_k . \tag{70}$$

we will prove (69) by induction. For m = 1 (69) gives

$$\widetilde{R}_{1}\widetilde{Q}_{1}^{T} = A_{1} - \mu_{1}^{T}$$

$$R_1 Q_1^T = A_1 - \mu_1 I$$

which is true if in (67α) we set k = 1. Suppose (69) is true for m = k-1, that is, suppose that the following relation is true.

$$\tilde{R}_{k-1}\tilde{Q}_{k-1}^{T} = (A - \mu_1 I) (A - \mu_2 I) . . . (A - \mu_{k-1} I)$$
 (71)

Now we will prove that (69) is true for m * k , that is, we will prove that the relation

$$\tilde{R}_{k}\tilde{Q}_{k}^{T} = (A - \mu_{1}I) (A - \mu_{2}I) \dots (A - \mu_{k-1}I) (A - \mu_{k}I)$$
 (72)

is true. (68) because of (70) gives us

$$Q_{\mathbf{k}}^{\mathbf{T}} R_{\mathbf{k}} = \widetilde{Q}_{\mathbf{k}}^{\mathbf{T}} A \widetilde{Q}_{\mathbf{k}} - \mu_{\mathbf{k}} \mathbf{I} \iff$$

$$Q_{\mathbf{k}}^{\mathbf{T}} R_{\mathbf{k}} = \widetilde{Q}_{\mathbf{k}}^{\mathbf{T}} (A - \mu_{\mathbf{k}} \mathbf{I}) \widetilde{Q}_{\mathbf{k}} \iff$$

$$R_{\mathbf{k}} = \widetilde{Q}_{\mathbf{k}-1}^{\mathbf{T}} (A - \mu_{\mathbf{k}} \mathbf{I}) \widetilde{Q}_{\mathbf{k}} \implies$$

$$\widetilde{R}_{\mathbf{k}} = \widetilde{R}_{\mathbf{k}-1} \widetilde{Q}_{\mathbf{k}-1}^{\mathbf{T}} (A - \mu_{\mathbf{k}} \mathbf{I}) \widetilde{Q}_{\mathbf{k}} \iff$$

$$\widetilde{R}_{\mathbf{k}} \widetilde{Q}_{\mathbf{k}}^{\mathbf{T}} = \widetilde{R}_{\mathbf{k}-1} \widetilde{Q}_{\mathbf{k}-1}^{\mathbf{T}} (A - \mu_{\mathbf{k}} \mathbf{I}) \tag{72a}$$

6)

(72a) becuase of (71) gives us

$$\widetilde{R}_{k}\widetilde{Q}_{k}^{T} = (A - \mu_{1}I) (A - \mu_{2}I) \dots (A - \mu_{k}I) (A - \mu_{k}I)$$

which is (72).

(2) The case of complex eigenvalues:

Here we will describe how to accomplish two steps of the forward sweep of the implicifly shifted method in one pass, and we will show how we can use this double step to assign a pair of complex conjugate eigenvalues to a matrix A. Let λ_1 , λ_2 be the two eigenvalues and let $A_1 \equiv A$, then from (68a) for k=1 and k=2 we have respectively:

$$A_{2} = Q_{1}^{T} A_{1} Q_{1} \quad \text{and}$$

$$A_{3} = Q_{2}^{T} A_{2} Q_{2} \quad \Longrightarrow$$

$$A_{3} = Q_{2}^{T} Q_{1}^{T} A_{1} Q_{1} Q_{2}$$
(73)

where A_1 , A_2 and A_3 are upper Hessenberg. In order to accomplish these two steps implicitly, we first find an orthogonal matrix P_1 such that $P_1e_n=Q_1Q_2e_n$. To do this we should find the last row of $Q_2^TQ_1^T$. From theorem (VI) we have

$$R_1 R_2 Q_2^T Q_1^T = (A - \lambda_1 I) (A - \lambda_2 I)$$
, (74)

where the upper triangular matrices R_1 , R_2 have elements $\rho_{ij}^{(1)}$, $\rho_{ij}^{(2)}$, respectively. From (74) we can see that if q_n^T is the last row of $Q_2^TQ_1^T$ then

$$(\rho_{nn}^{(2)} \ \rho_{nn}^{(1)}) q_n^T = e_n^T (A - \lambda_1 I) (A - \lambda_2 I)$$
 (75)

So q_n^T is a multiple of the last row a^T of the matrix $(A - \lambda_1 I)(A - \lambda_2 I)$. But A is upper Hessenberg so a^T has only three non-zero elements, $a^T = (0, \dots, 0, \alpha_{0, n-2}, \alpha_{0, n-1}, \alpha_{0n})$ say. Then after some trivial calculations, we find that

$$\alpha_{0,n-2} = \alpha_{n-1,n-2} \alpha_{n,n-1}$$

$$\alpha_{0,n-1} = \alpha_{n,n-1} [\alpha_{n-1,n-1} + \alpha_{nn} - (\lambda_1 + \lambda_2)]$$

$$\alpha_{0,n} = \alpha_{nn}^2 + \alpha_{n,n-1} \alpha_{n-1,n} - \alpha_{nn} (\lambda_1 + \lambda_2) + \lambda_1 \lambda_2$$
(76)

Now choose P₁ such that

$$a^{T}P_{1} = \pm ||a||_{2} e_{n}^{T} \iff$$

$$a^{T} = \pm ||a||_{2} e_{n}^{T} P_{1}^{T} \iff$$

$$e_{n}^{T} P_{1}^{T} = \pm \frac{a^{T}}{||a||_{2}} \qquad (77)$$

Thus

$$P_1e_n = \pm \frac{a}{\|a\|_2} = \pm \frac{\rho_{nn}^{(2)} \rho_{nn}^{(1)}}{\|a\|_2} q_n = \pm Q_1Q_2e_n$$

and we have accomplished our aim.

From (76) we can see that if $\lambda_1=\lambda$ and $\lambda_2=\overline{\lambda}$, where $\lambda\in\mathbb{C}$ and $\overline{\lambda}$ is the complex conjugate of λ , P_1 is real, and can be found by using real arithmetic, since $\lambda+\overline{\lambda}$ and $\lambda\overline{\lambda}$ are reals even when $\lambda\in\mathbb{C}$. Now, after having found P_1 , the next step is to find orthogonal matrices U_1 , i=1(1)n each with $U_1e_n=e_n$, such that the matrix A^* is upper Hessenberg in

$$A' = Q^T A_1 Q$$
, where $Q = P_1 U_1 \dots U_n$

Then we will have $Qe_n = P_1e_n = \pm Q_1Q_2e_n$, and since A_3 has nonzero (i+1, i) elements, i = 3(1)n-1, the last n-2 columns of Q and Q_1Q_2 are essentially the same, as are the last $(n-2)\times(n-2)$ diagonal blocks of A' and A_3 . This follows from Theorem (V). We will give an example with a 6×6 matrix to demonstrate how the forward sweep of the double step method works. Let

then we find an orthogonal matrix P such that if we apply it to the matrix $(A - \lambda_1^{-1}I)(A - \lambda_2^{-1}I)$ from the right it will first eliminate the (6,4)-element of $(A - \lambda_1^{-1}I)(A - \lambda_2^{-1}I)$ into the (6,5)-element of the same matrix, and then the resulting (6,5)-element of $(A - \lambda_1^{-1}I)(A - \lambda_2^{-1}I)$ into the (6,6)-element of the same matrix. So when we first form the matrix P^TAP , and then we reduce it to upper Hessenberg form, we will have the following:

Note: The element 3' cannot be eliminated because after the multiplication of the rotation matrix 6', this element becomes unknown, as we can see from matrix 3. Since from Theorem (V) the first two columns of Q span the eigenspace of A corresponding to λ_1 and λ_2 , we see that the elements (3,1) and (3,2) of matrix 3 must be equal to zero, and the elements (1,1), (1,2), (2,1), (2,2) of the same matrix will be forced to be equal to γ , δ , ∞ , ∞ , ε , respectively, where $\gamma + \zeta = 2\alpha$ and $\det \begin{bmatrix} \gamma & \delta \\ \varepsilon & \zeta \end{bmatrix} = \alpha^2 + \beta^2$, $\lambda_1 = \alpha + j\beta$ and $\lambda_2 = \alpha - j\beta$ with $j^2 = -1$. In other words, γ , δ , ε , ζ are such that λ_1 and λ_2 are the two eigenvalues of the matrix $\begin{pmatrix} \gamma & \delta \\ \varepsilon & \zeta \end{pmatrix}$. Finally, we have after this first double step of the forward sweep

$$\begin{array}{ccccc}
Q^T & A & Q & \stackrel{\bullet}{-} & \begin{pmatrix} \gamma & \delta & b_1^T \\ \epsilon & \epsilon & b_2^T \\ 0 & 0 & \widetilde{A} \end{pmatrix} & ,
\end{array}$$

where \mathbf{b}_1 , \mathbf{b}_2 and the first row of $\widetilde{\mathbf{A}}$ are as yet unknown.

Let us now examine the backward sweep when we know the first row of $$A_{\bf k+2}$$ and we want to find the first row of $$A_{\bf k}$$.

$$Q_{k}^{T} A_{k} Q_{k} = \begin{bmatrix} \gamma & \delta & b_{k}^{T} \\ \epsilon & \zeta & b_{k+1}^{T} \\ 0 & 0 & A_{k+2} \end{bmatrix}$$

Let a_{k+2}^T and a_k^T be the first rows of A_{k+2} and A_k , respectively. From the above example, we can see that the first row of A was affected only by rotation 5', coming from the matrix Q^T . But in order to calculate the first row of A we should take into consideration rotation 6' also. So the first three rows of the last two rotation matrices from the left were:

$$\begin{bmatrix} 1 & & & & & \\ & & & & \\ & & & c_{k} & & \\ & & -s_{k} & & c_{k} \end{bmatrix} \begin{bmatrix} c_{k-1} & s_{k-1} & & \\ -s_{k-1} & c_{k-1} & & \\ & & & 1 \end{bmatrix} = \begin{bmatrix} c_{k-1} & s_{k-1} & 0 \\ -c_{k}s_{k-1} & c_{k}c_{k-1} & s_{k} \\ & & & \\ s_{k}s_{k-1} & -s_{k}c_{k-1} & c_{k} \end{bmatrix}$$

Let $(\widetilde{\alpha},\widetilde{\beta},\widetilde{\widetilde{b}}^T)$ be the first row of (A_kQ_k) . Then

$$\begin{bmatrix}
c_{k-1} & s_{k-1} & 0 \\
-c_{k}s_{k-1} & c_{k}c_{k-1} & s_{k} \\
s_{k}s_{k-1} & -s_{k}c_{k-1} & c_{k}
\end{bmatrix}
\begin{bmatrix}
\widetilde{\alpha} & \widetilde{\beta} & \widetilde{b}^{T} \\
\widetilde{\mu} & v & f^{T} \\
\xi & \pi & g^{T}
\end{bmatrix}
= \begin{bmatrix}
\gamma & \delta & b_{k}^{T} \\
\varepsilon & \zeta & b_{k+1}^{T} \\
0 & 0 & a_{k+2}^{T}
\end{bmatrix}$$
(78)

Then, from (78) and since the rows (μ, ν, f^T) and (ζ, π, g^T) are computed in the forward sweep, we obtain:

$$(\mathbf{s_k s_{k-1}}, -\mathbf{s_k c_{k-1}}, \mathbf{c_k}) \begin{pmatrix} \widetilde{\alpha} & \widetilde{\beta} & \widetilde{b}^T \\ \mu & \nu & \mathbf{f}^T \end{pmatrix} = (0, 0, \mathbf{a_{k+2}}^T) \iff$$

So the first row of the matrix $A_k Q_k$ is $(\widetilde{\alpha}$, $\widetilde{\beta}$, $\widetilde{b}^T)$ that is

Thus

$$\mathbf{a}_{\mathbf{K}}^{\mathbf{T}} = (\widetilde{\alpha}, \widetilde{\beta}, \widetilde{\mathbf{b}}^{\mathbf{T}}) \mathbf{Q}_{\mathbf{k}}^{\mathbf{T}}$$
.

So knowing the first row of A_{k+2} we can find the first row of A_k , when we are using such a double step.

Now if we want to assign a set of eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, where the complex conjugate pairs appear together, to an unreduced upper Hessenberg matrix A, we will use a combination of double steps or implicitly shifted

single steps, according to what type of eigenvalue we have to assign at each step; a complex conjugate pair or a single real eigenvalue. So knowing the first row of A_n or A_{n-1} , according to whether the last eigenvalue to be assigned is a real one or a complex conjugate pair respectively, we can go backwards to finally produce the first row of $A_n \equiv A$, the desired result.

In the following and until the end of this section we shall compare the two previously described methods of assigning eigenvalues, that is, the double and the single step.

Note that four multiplications are involved when any rotation $\, R \,$, is applied to any vector $\, x \,$.

Single step

(a) Forward sweep:

6

Suppose we want to assign one eigenvalue to a $k \times k$ unreduced upper Hessenberg matrix with k = 4, then we have the following

from which we see that rotations π and π' involve 4×4 and 4×3 multiplications respectively. Similarly rotations 1, 1', 2 and 2' involve 4×4 , 4×4 , 4×3 and 4×4 multiplications. Thus the number of multiplications which are involved in one forward step of a 4×4 matrix can be found as follows:

$$4 \times 4 + 4 \times 3$$

$$4 \times 4 + 4 \times 4$$

$$4 \times 3 + 4 \times 4 + 4$$

$$4 \times 3 + 4 \times 4 + 4$$

$$4 \times 3 + 4 \times 4 + 4$$

$$4 \times 3 + 4 \times 4 + 4$$

$$4 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times 4 + 4$$

$$1 \times 3 + 4 \times$$

' So in general we have the following:

Hence in the single step method one forward step involves multiplications of the order $4k^2$ when it is applied to a $k \times k$ matrix. So the order of the number of multiplications involved in the forward sweep is the following:

$$4 \sum_{k=1}^{n} k^{2} = \frac{4}{3} n(n + \frac{1}{2})(n + 1)$$

$$\approx \frac{4}{3} n^{3}, \qquad (81)$$

where n is the order of the matrix we deal with.

(b) Backward sweep:

Suppose we are at the kth step of the backward sweep, that is, we are to find the first row of the k×k submatrix knowing the first row of the (k-1)×(k-1) submatrix, then from (64), (65) and since at this step a_k^T and a_{k+1}^T are (k-1)-vectors, we need k multiplications to calculate the first row of the k×k submatrix. So the order of the number of multiplications involved in the backward sweep is $\sum_{k=1}^n k^2 \frac{1}{2} n^2$.

Double step

(a) Forward sweep:

Suppose we want to assign two eigenvalues to a $k \times k$ unreduced upper Hessenberg matrix with k = 6, then we have the following:

from which we see that rotations π_1 , π_2 , π_1' , π_2' involve 4×6 , 4×6 , 4×4 and 4×4 multiplications respectively. Similarly we can find the number of multiplications involved in rotations 1, 1', 2, 2', 3, 3', 4, 4', 5, 5', 6, 6'. Thus the number of multiplications involved in one forward step of the double step method of a 6×6 matrix, can be found as follows:

$$4 \times 6 + 4 \times 6 + 4 \times 4 + 4 \times 4$$

$$4 \times 6 + 4 \times 6 + 4 \times 5 + 4 \times 5$$

$$4 \times 5 + 4 \times 5 + 4 \times 6 + 4 \times 6$$

$$4 \times 4 + 4 \times 4 + 4 \times 6 + 4 \times 6 + 4$$

$$4 \times 6 + 4 \sum_{i=4}^{6} 1_{i} = 336$$

. So in general we have the following:

$$4k + 4k + 4 \times 4 + 4 \times 4$$

$$4k + 4k + 4 \times 5 + 4 \times 5$$

$$4(k-1) + 4(k-1) + 4 \times 6 + 4 \times 6$$

$$\vdots \qquad \vdots \qquad \vdots$$

$$4 \times 5 + 4 \times 5 + 4k + 4k$$

$$4 \times 4 + 4 \times 4 + 4k + 4k + 4k$$

$$4(4k + 4 \sum_{k=4}^{k} 1) = 16k + 16 \sum_{k=1}^{k} 1 - 96$$

$$= 8k^{2} + 24k - 96$$

$$\approx 8k^{2}$$
(82)

Hence in the double step method one forward step/involves multiplications of the order $8k^2$ when it is applied to a k×k matrix. So the order of the number of multiplications involved in the forward sweep, when n is even, is the following:

$$8 \sum_{k=1}^{n/2} (2k)^{2} = 32 \sum_{k=1}^{n/2} k^{2}$$

$$\approx \frac{32}{3} (\frac{n}{2})^{3}$$

$$= \frac{4}{3} n^{3}.$$
(83)

ľ

If n is odd then obviously we will have to use at least one single step.

Note that the forward sweep of the double step method can be carried out using elementary reflectors [12] instead of rotations.

(b) Backward sweep:

Suppose we are at the kth step of the backward sweep, then from (79) since f, g and a_{k+2}^T are (k-2)-vectors we need multiplications of order 2k to calculate the first row of the k×k matrix. So the order of the number of multiplications involved in the backward sweep is $\frac{n^2}{2}$.

Hence from the above results we can observe that one double step involves almost as many multiplications as are involved in two single steps, thus the double step, with rotations, and the single step are equivalent in solving the eigenvalue allocation problem, in the sense that they involve the same order of number of multiplications.

CHAPTER 4

(I) Conclusion

In the second Chapter of this thesis we considered the most widely, so far, used algorithms in computing controllability and we compared them with a recently developed algorithm by C. Paige in [2]. The comparison was done by using appropriate counterexamples and it showed us that there are cases where the algorithms which were used so far fail, while the new algorithm gives the correct results. Finally in the same Chapter we presented an algorithm which using the notion of the distance of a system from the nearest uncontrollable one determines whether the system is controllable or not, the same algorithm with some changes (see note at the end of Chapte 2) can be used to calculate the distance of the given system from the nearest uncontrollable one, but it is pointed out why this algorithm is not efficient. Thus for determining the controllability of a system we have a reasonable algorithm, but for a controllable system we do not yet have an efficient algorithm to measure its distance from the nearest uncontrollable system, but we have the theoretical approach developed in section (III) of Chapter 2 of this thesis which may prove useful in developing an efficient algorithm. So the efficient computation of the distance of a system from the nearest uncontrollable one remains an open problem.

In Chapter 3 we have presented a very efficient, fast and we hope numerically stable algorithm to solve the eigenvalue allocation problem which can be used to determine the feedback such that a single input closed

loop system is stable. An algorithm for the same problem which was suggested by Luenberger in [3] appears to be unstable because it makes use of the coefficients of the characteristic polynomial of A which can be found by computing the eigenvalues of A, but we know that the problem of calculating the eigenvalues is in general ill-conditioned (if A is not symmetric), however Luenberger's algorithm can be extended to multi-input dynamic systems [10]. The method given here has also been designed with extension to multi-input systems in mind, but this has not yet been done.

APPENDIX

Computational results based on the implicitly shifted and double step algorithms. Subroutine EVA.

Based on the implicitly shifted and the double step algorithms, which were described in Chapter 3 we made a subroutine named EVA (Eigenvalue Allocation) in FORTRAN IV. This subroutine takes a matrix A without its first row and a set L of eigenvalues and evaluates the first row of A such that $\lambda(A) = L$. With this subroutine we checked some cases, where the reliability of the results was checked as follows: First the full matrix A was passed to the IMSL routine EQRH3F, in order to calculate its eigenvalues. Then the eigenvalues computed by the EQRH3F routine together with A (without its first row) were given to subroutine EVA and the first row of A was computed, and finally the full matrix A (with its computed first row) was given again to the IMSL routine EQRH3F, in order to compute its eigenvalues. The results were obtained using single precision arithmetic on the AMDAHL 470/V7 at McGill University. We also give an example in which we point out why using the above way to check EVA we may obtain results which differ significantly from those expected. This is not necessarily a fault of the algorithm (that is, it does not mean that the algorithm is necessarily numerically unstable). It may just be a result of the poor conditioning of the specific eigenproblem.

Example (I)

n = 5

The matrix A (with the original first row)

5.279000	9.125000	4.433000	6.297000	5.687000
38.345000	39.492000	3.605000	5.987000	7.770000
ļ. ·	-5.564000	6.396000	6.492000	5.889000
**		3.564000	9.539000	6.364000
	,	7	-5.977000	4.796000,

The eigenvalues of A using EQRH3F

46.726480

-3.995152

3.844520 '

• 9.463070 + j 3.235559

9.463070 - j 3.235559 1

The first row of A using EVA

•

5.278953 9.124928

4.433008

6.296913

5.687075

The eigenvalues of A, with the computed first row, using EQRH3F

46.726427

-3.995148

3.844508

9.463083 + j 3.235580

9.463083 - j 3.235580

Example (II)

n = 6

The matrix A with given first row

9.452000	-4.279000	5.126000	6.433000	3.297000	4.687000	1
6.474000	-8.345000	79.490000	`3 7.605000	0.987000	8.770000	
a a	4.657000	5.564000	7.396000	7.492000	7.890000	
-		-0.998000	4.564000	9.540000	9.364000	١,
		¥.	-7.463000	-7.977000	3.796000	
				-9.897000	8.697000	

The eigenvalues of A using EQRH3F

-21.014220

17.953000

-2.881598 + j 9.934898

-2.881598 - 19.934898

9.028457

11.750950

The first row of A using EVA

9.452018 -4.278798 5.125249

6.429889

3.297468

688825

The eigenvalues of A , with the computed first row, using EQRH3F

-21.014265

17.953045

-2.881592 + j 9.934899

-2.881592 - j 9.934899

9.028496

11.750926

The matrix A with given first row

```
8.657000 5.254000 6.254000 -5.565000 6.874000 0.657000 7.364000 4.655000 -0.763000
3.546000 -6.356000 5.567000 -6.534900 5.578000 5.546000 4.356000 0.768000
                                                                            0.658000
         6.865000 7.456000 6.457000 6.845000 -0.176000 9.763000 5.765000 -7.765000
                  7.645000 -6.656000 -6.645000 5.657000 7.456000 5.765000
                                                                           -5.655000
                           7.565000 5.343000 0.678000 0.567000 6.678000
                                                                            8.568000
                                    -8.876000 5.684000 8.563000 7.456000
                                                                           4.357000
                                              7.536000 -7.657000 7.456000 -6.672000
                                                      -4.453000 -6.366000 -6.366000
                                                                 5.464000 -4.366000
```

The eigenvalues of A using EQRH3F

```
12.579090
```

9.194291 + j 2.895247 -8.440416 + j 6.0173679.194291 - j 2.895247 · -8.440416 - j 6.017367 -0.525684_+ j 6.229814 -12.636600 -0.525684 - j 6.229814 -4.659868

The first row of A using EVA

8.656885 5.254084 6.254014 -5.565294 6.874562 0.657171 7.363661 4.657894 -0.763001

The eigenvalues of A , with the computed first row, using EQRH3F

12.579074 -8.440431 + j 6.0173749.194260 + j 2.895261 -8.440431 - j 6.017374 9.194260 - j 2.895261 -12.636604 -0.525680 + j 6.229808 -4.659882 -0.525680 - 1 6.229808

Example (IV)

n = 10.

The matrix A with given first row

The eigenvalues of . A using EQRH3F

17.525430 -8.598625 0.776537 + j 10.545620 2.859068 + j 2.695472 0.776537 - j 10.545620 2.859068 - j 2.695472 8.885132 + j 5.104329 -1.585043 8.885132 - j 5.104329 4.878755

The first row of A using EVA

 9.450891
 8.320937
 9.788930
 5.789191
 -9.535405

 4.651105
 0.547856
 8.618833
 6.374817
 -0.460253

The eigenvalues of A, with the computed first row, using EQRH3F

-8.598647

0.776526 + j 10.545624

0.776526 - j 10.545624

2.859084 + j 2.695479

2.859084 - j 2.695479

8.885115 + j 5.104358

-1.585052

8.885115 - j 5.104358

4.878776

Remark (XVIII): Using the previous way of checking subroutine EVA we may get some poor results due to some ill-conditioned eigenvalues.

Example (V):

Let the 10×10 matrix A without its first row be

Now using EVA we want to calculate the first row of A such that A will be assigned 10 eigenvalues each equal to 7. The first row of A computed by EVA is as follows:

21.000100 -196.998600 646.753400 -1994.128000 -873.761200 -3130.766000 -4675.128000 -8100.312000 -11546.110000 -11095.480000

The magnitude of this first row suggests this was a poorly conditioned problem. Now giving the full matrix A to EQRH3F we get the following eigenvalues:

 4.653267 + j
 1.855709
 8.417780 + j
 2.677467

 4.653267 - j
 1.855709
 8.417780 - j
 2.677467

 6.345307 + j
 2.945413
 9.807051 + j
 1.054002

 6.345307 - j
 2.945413
 9.807051 - j
 1.054002

 7.531523
 4.021682

As we can see no one of the above eigenvalues is 7. The above example does not necessarily mean that the algorithm gave an inaccurate first row for A. The difficulty probably arises because unreduced upper Hessenberg matrices are always defective when they have multiple eigenvalues and this happens because if λ is an eigenvalue of A then $\operatorname{rank}(A-\lambda I)=n-1$; so there is one and only one linearly independent eigenvector x such that $(A-\lambda I)x=0$, if λ is a multiple eigenvalue then it has only one linearly independent eigenvector and so A is defective. In this example the Jordan canonical form of A is a 10×10 Jordan block, the eigenvalues of which are very ill-conditioned [12, p. 301], [13], and so cannot be computed accurately.

(A+N+N1+TESTL+RSTACK+SSTACK+N2+QSTACK)

EVA IS A SUBROUTINE WHICH BY CALLING TWO OTHER SUBROUTINES FORW AND BACKW CALCULATES AND RETURNS THE FIRST ROW OF A STHOLE FRECISION NXN MATRIX A IN ORDER TO GIVE CERTAIN DECIFED FIGENVALUES 11/12+1... IN WHICH MAY OFCUR IN COMFITY CONTUGATE PAIRS.

UN ENTRY

SINGLE PERCISION MYN ARRAY WHICH COMTAINS THE MARRY WHOSE ETRSE YOU WE SILL CALCULATE. THE FIRST COW OF A 15 TONORED AND OVERWRITTEN IN THE COMPUTATION.

N

INTEGER.

N IS THE ORDER OF A

11

MI IS THE NUMBER OF ELEMENTS OF THE ARRAY TESTL

WHERE NI=74N.

TESTL

SINGLE PRECISION ONE DIMENSINAL ARRAY OF DEDER INT WHICH CONTAINS THE CIGENVALUES WHICH ARE GOING TO BE THE-EIGENVALUES OF A . EVERY EIGENVALUE OCCUPTES TWO CONSCIUTIVE FLACES OF THE ARRAY TESTL. THE FIRST FOR ITS FEAL PART

AND THE SECOND FOR ITS IMAGINARY FART. COMPLEX CONJUGATE PATES MUST AFFEAR TOGETHER.

RSTACK

SINGLE FRECISION HYN APPAY. RSTACK IS A WOPE ARRAY.

SSTACK

INTEGER ONE DIMENSIONAL ARRAY OF ORDER SSTACK IS A WORK ARRAY.

N2

INTEGEF.

N2 IS THE NUMBER OF ROWS OF THE ARRAY

WHERE N2=(N*(N-1)/2) 1

QSTACK

SINGLE PRECISION NOX3 ARRAY. DOTACK IS A WORK ARRAY.

ON RETURN

SEE THE PREVIOUS DISCUSSION ABOUT A ON RETURN WE HAVE THE DESIRED FIRST ROW OF THE REST OF . IS AS ON ENTRY.

INTEGER N. SSTACK (N) . INDEX, ISTEP, RSTEP REAL A(N.N), TESTL(N1), QSTACN (N2,3), RSTACK(N,N)

WE REARRANGE THE EIGENVALUES SUCH THAT THE LAST TWO FIGENVALUES CONSIST EITHER A COMPLEX CONJUGATE FAIR OR A PAIR OF TWO REAL

С C

C C C

```
CIDENVALUES.
C
          IF (ARS(TESTE(N1)).GT.O.SE-S.OR.ANS(TESTE(N1)).FE.O.SE S.AN
              ABS(TESTL(N1-2)).LE.0.5F-5) 00 TO 5
          C = IE STL(N1-1)
          TESTI (N1-1) = TESTI (N1-5)
          TEST! (N1-5)=C
          TESTI (N1) #TESTI (N1 4)
          IFSTL(N1-4)=0.
          CALL FORW (N.Nt. INDEC. ISTEP-ESTEP, A. TESTE, QSTACK, RSTACK, SSTACK,
                     NO)
          CALL BACKW (N.HI. &NDEY, TOTTE, RSTEE, A. TESTE, OSTACK, ESTACK, HO,
                       SSTACES
       RETURN
      THD
      SUBROUTING FORM (NyN1+INDEX+ISTED+RSTED+A+HOTE+O516E+E516C
                         PSSTACKINED -
C
٠,
      FORM SUBROUTING FERFORMES THE FORWARD CHIEF.
C
C
         COMMON /RI/ C.S.R.H
         INTEGER RSTEP: FY: ISTEP: SW: M: DROP: THUE X:N: EINE: DROP1: SSTACK (N)
         REAL TICISIA (NIN) ITESTE (NE) IRSTACE (N2.3) INSTACE (NIN) IN ANIHALI
               AN2
C
C
L
      INITIALIZATION
C
C
         INDEX=1
         ISTEF = 1
         RSTET=1
         EV=1
      WE DETERMINE THE FORWARD SWEFF, STARTING FROM N UNTIL
  30
         M=N-RSTEP+1
         DROF=N-M
     THE LAST TWO EIGENVALUES FOR SIMPLICITY ARE TREATED FORETHER.
         IF (M.EQ.2) GO TO 10
     WE CHECK IF THE COMING EIGENVALUE IS COMPLEX OR REAL.
```

0000

C

ε C

С \mathbf{C} C

C C

C C

```
IF (ARS(TESTL(EV+1)).QT.0.5E-5) ON TO 5
 C'
      IN THIS PROGRAM SEGMENT WE DEAL WITH REAL FIGRNIVALUES.
 C
           SSTACK(ISTER)-1
           ISTEP#ISTEP#1
       WE DETERMINE THE TRANSFORMATION MATRIX WHICH WILL HELD HS MOVE ONE FORWARD STEP. BY PALCULATING ITS LAST COLUMN.
 C
 C
 C;
 C
           SS=A(N.N-1)
           CC=A(N,N) TISTL(EV)
 C
 נ
נ
       WE AFFLY THE TRANSFORMATION.
          CALL RROT (OSTACKIN 1:N:99:CC:A:INDEX:N:H2):
          EU=EU+2
          SW=1
          'CALL IROT (A.N-1,N,SW,N)
          LINE=DROPH 4
          I=LINE-1
          ÍF(M.EQ.3) OO TO 15
          DO 20 II=LINE,N
             -I=N+LINT II
SS=A(I,I 2)
             CC=A(I,I-1)
             CALL RROT (QSTACK+I 2:1-1:55:00:A:INDEY:N:N2)
             CALL LROT (A,I-2,I-1,SW,N)
   20
          CONTINUE
          I = I - 1
          SS=A(I,I-2)
   15
          CC=A(I,I-1)
          CALL REDT (OSTACK) 1 2-1-1-SS-CC/A-INDEY-N-H2)
          DROP1 = DROF +1
          DO 25 JJ=DROP1.N
             LL-1409044=
             RSTACK(RSTEF+J)=A(I-1+J)
   25
          CONTINUE
          RSTEP=RSTEF+1
          GO TO 30
0000
      IN THIS PROGRAM SEGMENT WE DEAL WITH TWO EIGENVALUES WHICH ARE
     · COMPLEX CONJUGATES.
      SO WE DETERMINE THE TRANSFORMATION MATRIX WHICH WILL HELF US HOVE
```

TWO FORWARD STEPS BY CALCULATING ITS LAST COLUMN.

5

C

AN2=A(N-1;N-2)*A(N;N-1) AN1=A(N,N-1)*(A(N-1,N-1)+A(N,N)-2,*TESTL(FU)) AN=A(N+N)##2+A(N+N-1)#A(N-1+N)-A(N+N)#2.#TESTL(EV)+TESTL(EV)##2

```
TITESTL (EVI1)**2
          EV=FU+4
          SSTACK (ISTER) ="
          ISTEP= ISTEP+1
C
C
       WE APPLY THE TRANSFORMATION
C
c
C
          CALL FROT (ROTACHIN DIN 1/AND-ANT-A-INDEY-N-ND)
          SW = 1
          CALL IROT (A.N 2.N 1.SW.N)
          T-r
          CALL REDT (OSTACKIN-1:N:T:AN:A:1MDEX:N:N2)
          SW=2
          CALL LEGT (A+H t+N+SW+H)
          I INC = TOKOF +5 .
          I=LINE 1
          IF, (M.EQ.4) ON TO 37 NO 40 II-LINE.N
             I=NILINC-II
             SS=A(I:I 3)
             CC-A(I, [-2)
             CALL PROT (GSTACK) I 3.1 2.50.00.0.THDEX.H.H?)
             SM= 1
             CALL LROT (A) I 3/I 2/SW/ID
             SS=A(I,I 2)
             CC=A(I,I 1)
             CALL FROT (QSTACK, T 2) T-1,55,60,6,THUTY, N, M2,
             SW=2
             CALL LROT (A.T-2.I-1.5W.N)
          CONTINUE
   40
          I = I - 1
   35
          SS=A(I,I-3)
          CC=A(I,I-2)
          CALL RROT (GSTACK, I-3, I-2, SS, CC, A, INDEX, N, N2)
          SS=A(I:I-2)
          CC=A(I,I-1)
          CALL RROT (RSTACK, I-2, I-1, SS, CC, A, INDEX, N, N2)
          DROP1=DROP+1
          DO 45 JJ=DROP1+N
             LL-190904N=L
             RSTACK(RSTEP+J)=A(T-2+J)
             RSTACK(RSTEP+1,J)=A(I-1,J)
   45
          CONTINUE
          RSTEP=RSTEP+2
          GO TO 30
С
     " SPECIAL TREATMENT FOR THE 2X2 MATRIX.
      INSTEAD OF MOVING ONE FORWARD STEF AND THEN ONE BACKWARD STEP. IN
C
      ORDER TO CALCULATE THE FIRST ROW OF THE 2X2 MATRIX, WE DO THIS
C
      DIRECTLY USING TWO FORMULAS.
C
C
C
   10 -
         A(N-1,N) = -(A(N,N)**2-(TESTL(EU)+TESTL(EU+2))*A(N,N)+TESTI(LU)*
                   TESTL(EV+2)+TESTI(EV+1)**2)/A(N+N-1)
          A(N-1;N-1)=TESTL(EV)+TESTL(EV+2)-A(N;N)
```

1-1.

,

4

PAGE

RETURN '

```
SUBROUTINE RROT (RSTACK+J1,J2,51,C1,A,IMDEX,N+N2)
RROT IS A SUBROUTING WHICH MULTIFLIES THE MATRIX A FROM THE RIGHT BY A ROTATION MATRIX, COMBINING THE 11.17 COLUMNS OF A
    COMMON* /R1/ C.S.R.M.
    INTEGER MANAIIAUIAUS
    REAL C.S.R.C1.S1.A(N.N).OSTACK(N2.3).A1.A2
    C = C1
    S=S1
    R=SQBT(S***+C**2)
    I1=N-M+?
NO 5 I=I1+N
       A1=(A(I, J1)*C-A(T, J2)*S1/R
       AZ=(A(I,J1)*S+A(I,J2)*C)/R
A(I,J1)=A1
       A(I,J2)=A2
   CONTINUE
    QSTACK(INDEX:1)=C
  , RSTACK(INDEX,2)=S
   RSTACK(INDEX:3)=R
   INDEX=INDEX+1
RETURN
END
```

SUBROUTINE LROT (A.11.12.5W.N)

LROT IS A SUBROUTINE WHICH MULTIPLIES THE MATRIX. A FROM THE LEFT BY A ROTATION MATRIX COMBINING THE $11.12~{
m ROWS}$ OF A .

00000

C

PAGE &

```
SUBROUTINE BACKW (N.N1. INDEX, PSTEP, RSTEP, A, TESTI, OSTACK, ESTACK, N.
                  SSTACK)
 ¢
 C
       BACKW IS A SUBROUTINE WHICH PEFFORMES THE BACKWARD SWEEP.
 C
          REAL A(N.N), TESTL(N1), RSTACK(N2,3), RSTACK(N,N), C,S,A1,A2,A5,
               C1,S1
      1
          INTEGER SSTACK(N), NO.LINE, RSTEF
          INDEX=INDEX-1
          R=RSTACK(INDEX:3)
          C=RSTACK(INDEX+1)/R
          S=QSTACK(INDEX,2)/R
   20
         ISTEF'=ISTEP-1
         RSTEP=RSTEP-SSTACK (ISTEP)
         M=N-RSTEP +1
         LINE=RSTEP
      WE DETERMINE IF THE CURRENT BACKWARD STEP IS ASSOCIATED WITH A
      REAL EIGENVALUE OR WITH A COMPLEX CONJUGATE PATH OF LIBERALUES.
         IF (SSTACK(ISTEP).EG.2) GO TO 5
      HERE WE MOVE ONE BACKWARD STEP BECAUSE THE LURRENT STEP IS
      ASSOCIATED WITH A REAL EIGENVALUE.
         NO=M-1
      FIRST WE FIND HOW THE FIRST ROW OF A SUBMATRIX OF A
                                                              *ASSOCIATED
      WITH THE CURRENT' STEP, WAS, BEFORE THE EFFECT OF THE
      TRANSFORMATION MATRIX FROM THE LEFT+
         A(LINE, LINE) = - C*RSTACK (RSTEP, RSTEP)/S
         LINE1=LINE+1
         DO 10 J=LINE1.N
            A(LINE, J)=(-C*RSTACK(RSTEP, J)+A(LINE+1, J))/S
   10
         CONTINUE
C
   PHERE WE FIND THE CORRECT FIRST ROW OF THE SAME SUBMATRIX OF
¢
č
         DO 15 J=1.NO
            A1=A(LINE+LINE+J-1)*C+A(LINE+LINE+J)*S
            A2=-A(LINE+LINE+J-1)*S+A(LINE+LINE+J)*C
            A(LINE,LINE+J-1)=A1
            A(LINE, LINE+J) =A2
            INDEX=INDEX-1
           IF (INDEX.EQ.O) GO TO 35
           R=OSTACK(INDEX,3)
            C=QSTACK(INDEX,1)/R
```

```
ድለፀር 7
             S=QSTACK(INDEX:2)/R
    15
          CONTINUE
          GD TO 20
C
       HERE WE MOVE TWO BACKWARD STEPS DECAUSE THE CURRENT STEP 15
Ĉ
       ASSOCIATED WITH A COMPLEX CONJUGATE PAIR OF CIDENVALUES.
         R1≈QSTACK(INDEX-1,3)
C1=QSTACK(INDEX-1,1)/F1
         S1=OSTACK(INDEX=1.2)/E1
          NO=2*M -4
C
000000
     FIRST WE FIND HOW THE FIRST ROW OF A SURMAIRIX OF A WITH THE CURRENT STEP, WAS, REFORE THE EFFECT OF THE
                                                              *ASSOCIATED
      TRANSFORMATION MATRIX FROM THE LEFT.
         )/(S#S1)
     1
         A(LINE,LINE(1)=(-5*C1*RSTACN(RSTEP,RSTEP+1)-C*RSTACK(RSTEP+1,
     1
                         RSTEP+1)>)/(9#91)
         LINE2=LINE+2
         110 25 J=LINE2,N
          A(LINE, J)=(-S*C1*PSTACK(RSTEF, ))-C*RSTACK(RSTCC+1, 1)+
                      A(LINE+2,J))/(5*51)
   25
         CONTINUE
         NO=NO/2
000
     HERE WE FIND THE CORRECT FIRST ROW OF THE SAME SURMATRIX OF
C
         DO 30 J=1,NO
            A1=A(LINE,LINE+J-1)*C1+A(LINE,LINE+J)*C*S1
     1
               +A(LINE+LINE+J+1)*S*S1
            A2=-A(LINE,LINE+J-1)*S1+A(LINE,LINE+J)*C*C1+
               A(LINE,LINE+J+1)#S#C1
            AJ=-A(LINE,LINE+J)*S+A(LINE,LINE+J+1)*C
            A(LINE,LINE+J-1)=A1
            A(LINE+LINE+J)=A2
            A(LINE+LINE+J+1)=A3
            INDEX=INDEX-2
            IF (INDEX.ED.O) GO TO 35
            R=QSTACK(INDEX,3)
            C=QSTACK(INDEX,1)/R
            S=QSTACK(INDEX,2)/R
           R1=QSTACK(INDEX-1,3)
```

35 RETURN

CONTINUE GO TO 20

30

C1=QSTACK(INDEX-1,1)/R1 S1=GSTACK(INDEX-1,2)/R1

REFERENCES

- [1] H.H. Rosenbrock, "State-space and Multivariable Theory". London:
 Nelson, 1970.
- [2] C.C. Paige, "Properties of Numerical Algorithms Related to computing controllability", IEEE Transactions on Automatic Control, Feb. 1981, Vol. AC-26, No. 1, pp. 130-138.
- [3] D.G. Luenberger, "Introduction to Dynamic Systems, Theory, Models and Applications", New York: John Wiley, 1979.
- [5] R.E. Kalman, Y.C. Ho and K.S. Narendra, "Controllability of linear dynamical systems". Contrib. Diff. Equations 1, No. 2, pp. 189-213 (1963).
- [6] M.L.J. Hautus, "Controllability and observability conditions of linear autonomous systems" in Proc. Kon. Akad. Wetensh. Ser. A., Vol. 72, pp. 443-448, 1969.
- [7] R.N. Clark, "Introduction to Automatic control systems", Wiley, New York, 1962.
- [8] E.J. Davison, W. Gesing, and S.H. Wang, "An algorithm for obtaining the minimal realization of a linear time-invariant system and determining if a system is stabilized-detectable", IEEE Trans.

 Automat., Contr., Vol. AC-23, pp. 1048-1054, Dec. 1978.
- [9] W.M. Wonham, "On Pole Assignment in Multi-input controllable linear systems", IEEE Trans. Automatic Contror, Vol. AC-12, December (1967), pp. 660-665.

- [10] D.G. Luenberger, "Canonical Forms for Linear Multivariable Systems",

 IEEE Transactions on Automatic Control, Vol. AC-12, No. 3, June (1967),

 pp. 290-293.
- [11] D.G. Luenberger, "Observers for Multivariable systems" IEEE Transactions on Automatic Control, Vol. AC-11, No. 2, April (1966), pp. 190-197.
- [12] G.W. Stewart, "Introduction to Matrix Computations", New York:
 Academic Press, 1973.
- [13] J.H. Wilkinson, "The Algebraic Eigenvalue Problem," London: Oxford University Press, 1965.
- [14] C.L. Lawson and R.J. Hanson, "Solving Least Squares Problems", Englewood Cliffs, N.J.: Prentice-Hall, 1974.
- [15] B.C. Moore, "Singular value analysis of linear systems", in Proc. 1978, IEEE Conference on Decision and Control, San Diego, CA, pp. 66-73, Jan. 1979.
- [16] B.C. Moore, "Computational problems with modal analysis", Sixteenth Annual Allerton Conf. on Communication, Control, and Computing, Oct. 1978.
- [17] P. Van Dooren, "The Generalized Eigenstructure Problem in Linear System Theory", IEEE Transactions on Automatic Control, Vol. AC-26, No. 1, pp. 111-129, Feb. 1981.
- [18] P. Van Dooren, A. Emami-Naeini, and L. Silverman, "Stable extraction of the Kronecker structure of pencils," in Proc. 17th IEEE Conf.

 Decision Contr., Jan. 1979, pp. 521-524.
- [19] R.V. Patel, "Computation of matrix fraction descriptions of linear time-invariant systems," in Proc. 18th IEEE Conf. Decision Contr.,

 Dec. 1979, pp. 325-328.

- [20] J.H. Wilkinson, "Linear differential equations and Kronecker's canonical form," in Recent advances in Numerical Analysis,

 C. de Boor and G.H. Golub, Eds., New York: Academic, 1978.
- [21] J.H. Wilkinson, "Kronecker's canonical form and the QZ algorithm",
 Nat. Phys. Lab., Teddington, Middlesex, England, Rep. DNACS 10/78,
 Nov. 1978.
- [22] P. Van Dooren, "The Computation of Kronecker's canonical form of a singular pencil," Linear Algebra and Its Applications, Vol. 27, pp. 103-140, 1979.
- [23] P. Van Dooren, "The generalized eigenstructure problem applications in linear systems theory," Ph.D. dissertation, Katholieke Univ.

 Leuven, Leuven, Belgium, May 1979.
- [24] R.W. Brockett and C.I. Byrnes, "Multivariable Nyquist Criteria, Root Loci, and Pole Placement: A Geometric Viewpoint,", IEEE Transactions on Automatic Control, Vol. AC-26, pp. 271-284, Feb. 1981.
- [25] H. Kimura, "Pole assignment by gain output feedback," IEEE Trans.

 Automat. Contr. Vol. AC-20, pp. 509-516, 1975.
- [26] H. Kimura, "A further result in the problem of pole assignment by output feedback," IEEE Trans. Autom. Contr., Vol. AC-22, pp. 458-463, 1977.
- [27] E.J. Davison, "On pole assignment in linear systems with incomplete state feedback," IEEE Trans. Autom. Contr. (Short Papers), Vol. AC-15, pp. 348-351, June 1970.
- [28] E.J. Davison and R. Chatterjee, "A note on pole assignment in linear system with incomplete state feedback," IEEE Trans. Autom. Contr. (Corresp.), Vol. AC-16, pp. 98-99, Feb. 1971.

- [29] B. Sridhar and D.P. Lindorff, "Pole-placement with constant gain output feedback," Int. J. Contr. Vol. 18, pp. 993-1003, Nov. 1973.
- [30] E.J. Davison and S.H. Wang, "On Pole Assignment in Linear Multi-variable Systems Using Output Feedback," IEEE Trans. Autom. Contr., Vol. AC-20, pp. 516-518, Aug. 1975.
- [31] E.J. Davison and S.G. Chow, "An algorithm for the assignment of closed-loop poles using output feedback in large linear multivariable systems," IEEE Trans. Automat. Contr., Vol. AC-18, pp. 74-75, Feb. 1973.
- [32] D.Q. Mayne, P. Murdoch, "Modal Control of Linear Time Invariant"

 Systems," Int. J. Control, Vol. 11, pp. 223-227, 1970.