

NUMERICAL ANALYSIS OF THE LYAPUNOV EQUATION
WITH
APPLICATION TO INTERCONNECTED POWER SYSTEMS

by

Thomas Mac Athay

SB, Michigan Technological University
(1974)

SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF
MASTER OF SCIENCE
and
ELECTRICAL ENGINEER
at the
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June, 1976

Signature of Author *Thomas M. Athay*
Department of Electrical Engineering and Computer Science, May 14, 1976

Certified by *Mil. B. Sandell*
Thesis Supervisor

Accepted by
Chairman, Departmental Committee on Graduate Students

NUMERICAL ANALYSIS OF THE LYAPUNOV EQUATION

WITH

APPLICATION TO INTERCONNECTED POWER SYSTEMS

by

Thomas Mac Athay

Submitted to the Department of Electrical Engineering and Computer Science on May 14, 1976, in partial fulfillment of the requirements for the degree of Master of Science and Electrical Engineer at the Massachusetts Institute of Technology.

ABSTRACT

The Lyapunov equation is fundamental to control theory. A number of numerical solution methods are compared, with special emphasis placed on applicability to large scale system matrices of a general sparse structure. An iterative decoupling algorithm is developed to exploit this special form, and a computer program that realizes this method is reported. An introduction to linear error analysis is included, and a number of results are developed and extended to several Lyapunov equation solution techniques.

Thesis Supervisor: Nils R. Sandell, Jr.

Title: Assistant Professor of Electrical Engineering and Computer Science.

ACKNOWLEDGMENT

I would like to express my gratitude to Professor Nils Sandell, Jr., for supervising this thesis. He stimulated my original interest in the subject and provided valuable intellectual influence and moral support throughout the course of this study.

My association with the Electronic Systems Laboratory and resulting interaction with other members, particularly fellow graduate students, has been a valuable and pleasant experience. I am also grateful for the financial support I received during this period.

Final thanks go to my wife and partner, Roberta Ann, mostly for her love, but also for typing this thesis.

TABLE OF CONTENTS

Abstract	1
Acknowledgement	2
I. Introduction	
1.1 Use of Lyapunov Equation	4
1.2 Summary of Thesis	5
II. Solution Methods	
2.1 Two Special Criteria	7
2.2 Direct Solution Methods	10
2.3 Transformation Solution Methods	14
2.4 Iterative Solution Methods	21
2.5 Iterative Decoupling Method	23
III. Error Analysis	
3.1 Introduction to Error Analysis	29
3.2 Conditioning	32
3.3 Perturbation Bounds	53
3.4 Summary of Results	83
IV. Applications	
4.1 Iterative Decoupling Program.	94
4.2 Power Systems Application.	96
4.3 Conclusions	99
Appendix	
A. Canonical Forms101
B. Round-off Errors105
References109

CHAPTER I
INTRODUCTION

1.1 USE OF LYAPUNOV EQUATION

Modern control theory has in the past two decades experienced a rapid development that is both exciting and significant. Its engineering application has understandably been a slower process, although recently it appears that the scope and diversity of application efforts are increasing. Important and rapidly changing limits to real world applications are computational requirements that arise in both the analysis and design of control systems. Because the theory has much to offer, efforts to understand and perhaps ease computational requirements are well motivated.

The Lyapunov equation arises in many aspects of both the analysis and design of linear control systems. Its solution is important in stability analysis of linear continuous systems [25], in pole assignment [17], when evaluating quadratic integrals which are often used as cost functions in optimal control [19, 25], and when evaluating covariance matrices in filtering and estimation for continuous systems. In addition, the algebraic Riccati equation, which occurs in some important filtering and optimal control problems, can be solved iteratively where each iteration is the solution of a Lyapunov equation [9, 11]. Another situation where the Lyapunov equation arises is in the design of decentralized control systems. Current research in large scale systems and decentralized control in the Electronic Systems

Laboratory is being directed toward physical systems that, although of large dimension, have sparse system matrices with particular structural forms. Examples of such research are the decentralized control of a freeway traffic corridor [12] and of large scale interconnected power systems [14]. In these studies the Lyapunov equation plays an important role, both in analysis and in design.

1.2 SUMMARY OF THESIS

There are many solution methods of the Lyapunov equation. In Chapter Two a number of them are introduced and compared, with special emphasis being placed on two special criteria. These are a method's potential for exploiting a general, sparse system matrix form and an algorithm's efficiency for re-solution. Basically, three properties of an algorithm are used to quantify the analysis: computational speed, storage requirements, and accuracy. The first two are included in Chapter Two, while the latter is covered in Chapter Three. In the final section of Chapter Two an iterative decoupling algorithm is developed that is specifically designed to meet the two criteria of special interest in this thesis.

Chapter Three is concerned with analyzing the error properties of several important solution methods. In order to do this, a general introduction to the method of backward error analysis is presented. This type of analysis is basically comprised of two steps, one of which involves the concept of numerical conditioning while the other

requires obtaining perturbation bounds that represent the effects of round-off error in the computer. Each of these steps are discussed in a separate section. The final section of the chapter is primarily a summary of the results obtained and final remarks regarding the comparison of the solution methods of interest.

Chapter Four contains a description of a computer program that was written to implement the iterative decoupling algorithm mentioned previously. The results of several small examples are briefly presented. In the next section, an unsuccessful application of the iterative decoupling algorithm to a power system example is reported. Finally, some conclusions and suggestions for future research are presented.

In addition, two appendices are included that are referenced at appropriate points in the main text. One is a summary of some pertinent facts of linear algebra, while the other briefly outlines some elementary tools of error analysis that are used in the analyses of Chapter Three.

CHAPTER II
SOLUTION METHODS

2.1 TWO SPECIAL CRITERIA

The Lyapunov equation is the steady state solution, $P(\infty)$, to the linear matrix equation

$$\frac{d}{dt} P(t) = A^T P(t) + P(t)A + Q \quad A, P, Q \text{ } n \times n \quad (2.1.1)$$

If the eigenvalues of A are such that $\operatorname{re} \lambda_j < 0, j = 1, 2, \dots, n$, then the steady state solution $P(\infty) = P$ of

$$0 = A^T P + PA + Q \quad (2.1.2)$$

exists and is unique, and is given by the convergent integral

$$P = \int_0^{\infty} e^{A^T t} Q e^{At} dt \quad (2.1.3)$$

Furthermore, if Q is symmetric and positive definite, then P will also be symmetric and positive definite. It is occasionally convenient to represent equation (2.1.2) in the form

$$L_A : R^{n \times n} \rightarrow R^{n \times n} \quad P \rightarrow A^T P + PA$$

A great deal of attention has been given to the numerical solution of the Lyapunov equation. A useful classification of the variety of solution techniques are the groupings of direct, transformation, and iterative methods. The purpose of this chapter is to summarize those methods that are (at least partly) favorable numerically, and to especially consider their application to sparse equations that must be solved many times. The following chapter will analyze in greater detail some of the algorithms introduced in this chapter. In particular, accuracy is not considered in this chapter.

The general sparse structure considered is system matrices of the form

$$A = \begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{1N} \\ A_{21} & A_{22} & A_{23} & A_{2N} \\ A_{31} & A_{32} & A_{33} & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ A_{N1} & A_{N2} & \cdot \cdot \cdot & A_{NN} \end{bmatrix} \quad \begin{array}{l} A_{ii} \quad n_i \times n_i \\ A_{ij} \quad n_i \times n_j \end{array}$$

In the modelling of some physical systems of current interest [12, 14] the diagonal blocks represent physical subsystems and may or may not be sparse. In the dynamic power system problem, for instance, the diagonal blocks model the incremental dynamics of a primemover and electrical generator pair and associated mechanical and electrical

(excitation) controls, while the off-diagonal blocks represent the coupling of machine angle and voltage state variables through transmission network. For even a moderate number of interconnected subsystems (say $N = 5$ to 10) a solution method that efficiently exploits this sparse structure would be valuable.

The other special criteria that is important is a method's economy in solving $L_A(P) = -Q$ many times with the same A matrix but different Q matrices. This is an important consideration in a number of applications. A particular example is an approach, outlined by Sandell and Athans [14], to the design of decentralized control systems with fixed communication structures. In this approach, the optimal (infinite time, quadratic cost criteria) design is the solution to the constrained parameter optimization problem

$$\text{Min}\{\text{tr}Q(z)P\} \quad (2.1.4)$$

$$\text{subject to} \quad A(z)P + PA^T(z) = -R(z) \quad (2.1.5)$$

where $Z^T = (z_1, z_2, \dots, z_p)$ is the vector of parameters that characterize the design. The problem is to search for z^* such that $\text{tr}Q(z^*)P(z^*) \leq \text{tr}Q(z)P(z)$. The point here is that the gradient of (2.1.4) with respect to z at $z = z^k$ (k^{th} step in search) is evaluated by resolving a Lyapunov equation p times with the same A matrix but different driving terms, i.e., with $L_A(P) = -Q$, different Q matrices.

2.2 DIRECT SOLUTION METHODS

The matrix equation

$$A^T P + PA = -Q \quad (2.2.1)$$

or $L_A(P) = -Q$, is an equation in $R^{n \times n}$. It can be conveniently rewritten using Kronecker product notation [25] as an equation in R^{n^2} . Let q and p be vectors that correspond to the n^2 elements of Q and P , taken by rows, respectively. Then (2.2.1) becomes

$$(A^T \otimes I + I \otimes A^T)p = -q \quad (2.2.2)$$

or
$$K_A p = -q \quad K_A: R^{n^2} \rightarrow R^{n^2} \quad (2.2.3)$$

Bellman [19] discusses the basic properties of the Kronecker product, a number of which are also contained in section 3.2. Equation (2.2.3) represents n^2 equations in n^2 unknowns, but in most applications $Q = Q^T$, therefore $P = P^T$ and we can rewrite equation (2.2.3) as

$$Cp = -q \quad (2.2.4)$$

where
$$C \text{ is } \frac{n(n+1)}{2} \times \frac{n(n+1)}{2}$$

Two simple algorithms for forming the matrix C are presented by Chen and Shieh[7] and Bingulac [8]. Some easy subscript arithmetic yields

a simpler algorithm, however. Let $s = (2n-i)(i-1)/2 + j$. The s^{th} row in (2.2.4), using (2.2.1) is

$$\sum_{k=1}^n (a_{ki} p_{kj} + p_{ik} a_{kj}) = -q_s$$

The elements a_{ki} and a_{kj} belong in the r^{th} column of C, with r given by

$$p_{kj}: r = \begin{cases} (2n-k)(k-1)/2 + j & k < j \\ (2n-j)(j-1)/2 + k & k \geq j \end{cases}$$

$$p_{ik}: r = \begin{cases} (2n-i)(i-1)/2 + k & k \geq j \\ (2n-k)(k-1)/2 + i & k < i \end{cases}$$

Once C is formed, any standard algorithm for solving linear equations can be used. Because one often wishes to resolve (2.2.3) with different q vectors, LU decomposition of C is an efficient approach [21]. With this approach, ignoring sparsity, solving equation (2.2.4) for p requires operations (multiplications and divisions) of the order of $n^6/24$, where A is $n \times n$. Once the LU decomposition of C is accomplished, however, computing p given q requires operations of the order of $n^4/4$. The memory requirement is approximately $n^4/4$ words.

The above operation counts are pessimistic. Even if A is full (n^2 elements), C will have a known sparse structure and $n^2 + [n(n-1)/2][2n-1]$ elements. Let α_c be the number of elements in C and define a sparsity ratio $\beta_c: \beta_c = \alpha_c/n_c^2, 0 \leq \beta_c \leq 1$. Then

for a given $A(n \times n)$ with a sparsity ratio β_A , $\beta_C \approx 4\beta_A/n$. This strongly suggests that sparse matrix techniques should be considered when using the direct method. Experience in using a sparse matrix package that includes an optimal ordering algorithm for minimizing the fill-in of the triangular factors of C has resulted in an operation count for the direct method of

$$\text{ops} = \frac{1}{3}\beta_{\text{Lu } C}^2 n^3 + \frac{3}{2}\beta_{\text{Lu } C} n^2 \approx 6\beta_A^2 n^4 + \frac{9}{2}\beta_A n^3$$

Resolving (2.2.4) for new q requires approximately $3\beta_A n^3$ operations. Although significantly faster, a sparse direct method algorithm is considerably more complex to code. An advantage of either direct method is that a solution can be iteratively refined, with each iteration similar to resolving for new q , to a specified accuracy if the Kronecker matrix is at least moderately well conditioned. This issue will be covered in more detail in the following chapter.

An idea that is standard in the solution of large, sparse sets of linear equations and that exploits the same general sparse structure of interest was studied as a potential solution method for the Lyapunov equation. Again, consider solving $Ax = b$ for x when A has the form:

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdot & \cdot & A_{1N} \\ A_{21} & A_{22} & & & \cdot \\ \cdot & \cdot & & & \cdot \\ \cdot & \cdot & & & \cdot \\ A_{N1} & \cdot & \cdot & \cdot & A_{NN} \end{bmatrix}$$

Each subsystem (A_{ii}) is a stability matrix of dimension n . Let the total number of interconnections between the subsystems equal K . Then we can split A into two parts:

$$A = B + UV^T \quad \text{where } B = \begin{bmatrix} A_{11} & & & & \\ & A_{22} & & & \\ & & \cdot & & \\ & & & \cdot & \\ & & & & A_{NN} \end{bmatrix}$$

$U \text{ (} nN \times K \text{), } V \text{ (} nN \times K \text{)}$

then [28] $A^{-1} = B^{-1} - Y \Lambda Z^T$ (2.2.5)

$$\text{where } Y = B^{-1}U$$

$$Z = V^T B^{-1}$$

$$\Lambda = (I + V^T B^{-1}U)^{-1}$$

Note that Λ is ($K \times K$) matrix; the dimension of Λ will depend on the geometry of the interconnections, but $\dim \Lambda \leq K$. This idea can be reformulated into an algorithm that involves an LU decomposition of B (i.e., the subsystems individually) and of Λ , and several steps of forward and back substitution. The higher order terms in the operation count are

$$Nn^3/3 + kNn^2 + k^3/3.$$

Although this count is reasonable for the solution of $Ax = b$, when the algorithm is extended to the Lyapunov equation, the corresponding operation count (for the definitions of N , n , and k) is roughly

$$N^2 n^6 / 24 + kN^3 (n^4 / 4 + 4k^2 n^3)$$

which is considerably worse than the sparsity-exploiting direct method.

2.3 TRANSFORMATION SOLUTION METHODS

Transformation solution methods of the equation

$$A^T P + PA = -Q \tag{2.3.1}$$

generally involve introducing a similarity transformation on the state space such that the new system matrix $A = N^{-1}AN$ is in a canonical form and the transformed equation

$$\tilde{A}^T \tilde{P} + \tilde{P} \tilde{A} = -\tilde{Q} \tag{2.3.2}$$

is more easily solvable. This section will consider two solution methods that utilize the companion and Schur canonical forms.

The use of the companion form actually characterizes a number of different solution methods. A number of methods view the problem in the frequency domain. The solution to equation (2.3.1) represents

the variance of the system $x(t) = A^T x(t) + W(t)$ in steady state, where $W(t)$ is stationary, zero mean white noise with spectral density matrix Q . The variance is also given by an integral of the form

$$P = \frac{1}{2\pi j} \int_{-j\infty}^{j\infty} H(s)H(-s)ds \quad (2.3.3)$$

where $H(s)$ is the system transfer function matrix. Astrom gives a very compact solution for the variance of the output of a single input, single output system using results from spectral factorization theory [32]. Hagander has made the extension to multivariable systems [1]. Others have implicitly used the companion form to solve equation (2.3.1), e.g., Muller [3].

Molinari [4] has developed an algorithm that uses the companion form explicitly. Consider:

$$TAT^{-1} = C_A, \quad S = T^{-T}PT^{-1}, \quad R = -T^{-T}QT^{-1}$$

$$\det [sI - A] = S^n + a_n S^{n-1} + a_{n-1} S^{n-2} + \dots + a_2 S + a_1$$

$$\text{Then } A^T P + PA = -Q \text{ becomes } C_A^T S + SC_A = R \quad (2.3.3)$$

C_A is of the form:

$$\begin{bmatrix} 0 & & & & & & & & & \\ & 1 & & & & & & & & \\ & & 1 & & & & & & & \\ & & & \cdot & & & & & & \\ & & & & \cdot & & & & & \\ & & & & & \cdot & & & & \\ & & & & & & 1 & & & \\ -a_1 & -a_2 & -a_3 & \cdot & \cdot & \cdot & \cdot & \cdot & -a_n & \\ & & & & & & & & & \end{bmatrix}$$

Letting $C_A^T S = U$ and $SC_A = V$, equation (2.3.3) becomes $U + V = R$, where

$$\begin{aligned} u_{ij} &= -a_i s_{nj} & i=1 & \quad j=1, 2, \dots, n \\ u_{ij} &= s_{i-1, j} - a_i s_{nj} & i=2, 3, \dots, n & \\ & & j=1, 2, \dots, n & \\ v_{ij} &= -a_i s_{in} & i=1 & \quad i=1, 2, \dots, n \\ v_{ij} &= s_{i, j-1} - a_j s_{in} & j=2, 3, \dots, n & \\ & & i=1, 2, \dots, n & \end{aligned}$$

This set of equations can be rewritten by forming a new equation that results from defining an alternating summation along a diagonal $i + j = \text{constant}$. That is

$$\sum_{j=\sigma}^{\theta} (-1)^{j+1} [u_{jk} + v_{jk}] = \sum_{j=\sigma}^{\theta} (-1)^{j+1} r_{jk} \quad i=1, 2, \dots, (2n-1)$$

$$\sigma = \begin{bmatrix} 1 & i \leq n \\ i+1-n & i > n \end{bmatrix} \quad \theta = \begin{bmatrix} i & i \leq n \\ n & i > n \end{bmatrix} \quad (2.3.4)$$

Molinari gives a total operation count of $5n^3$, with storage requirements of the order of $4n^2 + 4n$. Aspects of this solution approach will be analyzed in the following chapter; in general, however, this algorithm is superior in terms of speed and storage requirements. Applying this algorithm directly when the system matrix A has the general sparse structure of interest (2.1) may be difficult, however. First, the companion form only exists when A is non-derogatory (discussed in Appendix A). Some other difficulties, which are numerical rather than theoretical, will be discussed in the next chapter. Roughly 3/5 of the total operation counts involve solving $R = T^{-T}QT^{-1}$ for R and $P = T^TST$ for P . In addition, although T will generally reflect the sparseness of A and be sparse, T^{-1} tends to become full. To resolve for new Q matrices with the Molinari algorithm requires $\sim 3n^3$ operations; only the transformation of A is not redone. It is interesting to note that resolving equations (2.3.5) and (2.3.6) for S with different R matrices requires less than $2n^2$ operations.

Another transformation, that to a real Schur form, has been applied to the solution of the Lyapunov equation by Bartels and Stewart [2]. As in the companion form approach, a similarity transformation is applied to the system matrix A (i.e., on the state space), but in this case the transformation is orthogonal as well. The real Schur form immediately yields the eigenvalues of A , and the algorithms that affect the transformation are often used for this purpose. Some of the implications of this fact are important in error analysis, and hence will be considered in the next chapter.

Again, the equation

$$A^T P + PA = Q \quad (2.3.7)$$

is transformed to the equation

$$B^T X + XB = C \quad (2.3.8)$$

where $B = U^T A U$ $X = U^T P U$ $C = -C^T Q U$.

Note that $U^T U = U U^T = I$, which for real U is the definition of an orthogonal matrix. The matrix B (real Schur) is of the form:

$$B = \begin{bmatrix} B_{11} & B_{12} & \cdots & B_{1n} \\ & B_{22} & & \cdot \\ & & \cdot & \cdot \\ & & & \cdot \\ & & & B_{n-1,n} \\ & & & & B_{nn} \end{bmatrix} \quad \begin{array}{l} \text{each } B_{KK} \quad K = 1, \dots, n \\ \text{is of, at most, dimension two.} \end{array}$$

By partitioning X and C conformally with B , equation (2.3.8) can be solved recursively using the equation

$$B_{ii}^T X_{ij} + X_{ij} B_{jj} = C_{ij} - \sum_{k=1}^{i-1} B_{ki}^T X_{kj} - \sum_{k=1}^{j-1} X_{ik} B_{kj} \quad (2.3.9)$$

$$i = 1, 2, \dots, n \quad j = 1, 2, \dots, n$$

X_{ij} is, at most, a 2×2 matrix and can be easily found using the direct method (2.2).

The transformation to real Schur form involves two distinct steps. First A is transformed to the Hessenberg matrix H ,

$$N^T AN = H = \begin{bmatrix} h_{11} & h_{12} & h_{13} & \cdot & \cdot & \cdot & h_{1n} \\ h_{21} & h_{22} & h_{23} & & & & \cdot \\ 0 & h_{32} & h_{33} & & & & \cdot \\ & 0 & h_{43} & h_{44} & & & \cdot \\ & & 0 & \cdot & \cdot & & \cdot \\ & & & & \cdot & \cdot & h_{n-1,n} \\ & & & & 0 & h_{n,n-1} & h_{nn} \end{bmatrix}$$

The Hessenberg form is important in a variety of numerical problems and a number of algorithms exist to form it [18, 22, 24]. Bartels and Stewart use Householder's method, which is an "exact" transformation with an operation count of the order of $2n^3$. The real Schur form matrix, B , is obtained from H using the QR algorithm, which is an iterative (and inexact) procedure where each iteration is an orthogonal similarity transformation. Because this step is in effect solving the eigenvalue problem, it can be a numerically difficult one. The operation count is of the order of $4\sigma n^3$, with σ the number of iterations before a specified tolerance is met. The rest of the solution is straight forward. Evaluating $C = -U^T Q U$ and $P = U X U^T$ takes $3n^3$ operations (less

if A has known sparsity). The count on the recursion (2.3.9) depends on the number of complex eigenvalues of A; assuming one half are complex the count is n^3 . An average total count using the real Schur form then, is $(4\sigma + 2)n^3 + (4n^3 + 15n^2)$ while storage requirements are of the order of $3n^2$. Of the total operation count, the second term in brackets is the amount required to resolve the problem for different Q matrices.

The transformation of A to real Schur form is clearly the key to this solution method. Because this step solves the eigenvalue problem, the method has little or no potential for directly exploiting the general sparse structure of interest. In addition, the solution of the eigenvalue problem is typically limited to maximum dimensions on the order of 150-200 for numerical reasons that will become more clear in the next chapter. Nonetheless, it is an interesting general solution technique, and of course particularly attractive if the eigenvalues of A are also sought.

2.4 ITERATIVE SOLUTION METHODS

The solution of $A^T P + PA = -Q$ for stable A can be written as

$$P = \int_0^{\infty} e^{A^T t} Q e^{At} dt \quad (2.4.1)$$

This can be computed using the approximation

$$P(t+\Delta) = e^{A^T \Delta} P(t) e^{A \Delta} + \int_0^{\Delta} e^{A^T t} Q e^{At} dt \triangleq \Phi^T P(t) \Phi + \Gamma$$

$$\text{then } P = \sum_{k=0}^{\infty} (\Phi^T)^k \Gamma \Phi^k \quad (2.4.2)$$

An accelerated version of equation (2.4.2) generates the series much faster, that is

$$P_{k+1} = (\Phi^T)^{2^k} P_k \Phi^{2^k} + P_k \quad P_0 = \Gamma \quad (2.4.3)$$

Kleinman characterizes this as a safe approach; the operation count is proportional to n^4 , which must be repeated for new Q matrices. Davidson and Man [5] applied the Crank-Nicolson numerical integration method using a formula similar to equation (2.4.3), but with different approximations to Φ and Γ . They report operation counts of the order of $(3\sigma + 4)n^3$, with σ the number of iteration steps of equation (2.4.3). The iteration must be redone for new Q .

Another approximation for Φ results from introducing the bilinear transformation [16],

$$A \rightarrow \Phi = -(A + aI)(A - aI)^{-1}$$

$$Q \rightarrow \Gamma = (A^T - aI)^{-1} Q (A - aI)^{-1} \frac{1}{2a}$$

The solution can then be generated using the accelerated iteration equation (2.4.3). The operation count is similar to the Davidson and Man method. Smith [6] used this approach and the previous one for large ($n \leq 146$), lightly damped systems and found the bilinear

approach superior in terms of accuracy. Hagander [1] also favors this approach over the other iterative methods. However, the bilinear transformation algorithm needs to be redone for new Q matrices, and in terms of the other special requirement, exploitation of sparsity, it is not particularly favorable. This is because Φ will almost never be sparse even when A is.

2.5 ITERATIVE DECOUPLING METHOD

The algorithm discussed in the last section are essentially general purpose solution methods. An iterative decoupling approach naturally suggests itself when the general sparse system matrix structure is directly exploited. Consider a partitioned matrix of the form previously mentioned (i.e., off-diagonal submatrices very sparse)

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \begin{array}{l} A_{11} (n_1 \times n_1), A_{22} (n_2 \times n_2) \\ N = 2 \end{array}$$

The K_A matrix that is $A^T \times I + I \times A^T$, does not have the same general structure of A unless the P and Q matrices are similarly partitioned and K_A is formed in the following way:

$$P = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix} \quad Q = \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix}$$

let p_{ij} = vector of $n_i n_j$ elements of P_{ij} taken row-wise q_{ij} is similarly defined

and p = vector of N^2 P_{ij} sub-vectors that correspond to taking the P_{ij} blocks row-wise, $p \in R^{n^2}$

Then $A^T P + P A = -Q$ becomes

$$(A^T \times I)p + (I \times A^T)p = -q \quad (2.5.1)$$

The first term of equation (2.5.1) is of the form:

$$\begin{bmatrix} A_{11}^T \times I_1 & & A_{21}^T \times I_1 & \\ & A_{11}^T \times I_2 & & A_{21}^T \times I_2 \\ A_{12}^T \times I_1 & & A_{22}^T \times I_1 & \\ & A_{12}^T \times I_2 & & A_{22}^T \times I_2 \end{bmatrix} \cdot \begin{bmatrix} P_{11} \\ P_{12} \\ P_{21} \\ P_{22} \end{bmatrix} \quad (2.5.2)$$

The second term of equation (2.5.1) is of the form:

$$\begin{bmatrix} I_1 \times A_{11}^T & I_1 \times A_{21}^T & & \\ I_1 \times A_{12}^T & I_1 \times A_{22}^T & & \\ & & I_2 \times A_{11}^T & I_2 \times A_{21}^T \\ & & I_2 \times A_{12}^T & I_2 \times A_{22}^T \end{bmatrix} \cdot \begin{bmatrix} P_{11} \\ P_{12} \\ P_{21} \\ P_{22} \end{bmatrix} \quad (2.5.3)$$

Written in this way, the diagonal blocks of K_A include only the diagonal blocks of A , and the same correspondence holds for the off-diagonal blocks. Assuming that Q is symmetric, the equations corresponding to the off-diagonal blocks are redundant, i.e., $P_{ij} = P_{ji}^T$, $Q_{ij} = Q_{ji}^T$. Decomposing the A matrix, $A = A_0 + A_1$

$$A_0 = \begin{bmatrix} A_{11} & & & & \\ & A_{22} & & & \\ & & \cdot & & \\ & & & \cdot & \\ & & & & \cdot \\ & & & & & A_{NN} \end{bmatrix} \quad A_{ii} \quad n_i \times n_i$$

Then $L_A(P) = L_{A_0}(P) + L_{A_1}(P) = -Q$. The term $L_{A_0}(P)$ corresponds to the diagonal sub-blocks of equations (2.5.2, 2.5.3) and consists of $N(N+1)/2$ uncoupled equations for the block elements of P ,

$$A_{ii}^T P_{ij} + P_{ij} A_{jj} = -Q_{ij} \quad \begin{array}{l} i=1, 2, \dots, N \\ j=i, i+1, \dots, N \end{array} \quad (2.5.4)$$

The idea, then is to consider the sequence of solutions P^k , where

$$L_{A_0}(P^k) = -Q - L_{A_1}(P^{k-1}) \quad (2.5.5)$$

Laub discusses this iteration for some general classes of linear operators [29]; for square matrices this method is well-known [23].

Consider solving $Ax = b$ $x \in \mathbb{R}^n$, $b \in \mathbb{R}^n$. Let $A = L + D + U$, where L is strictly lower, U is strictly upper and D is diagonal (this decomposition can be in an element or block-partitioned sense). The general iteration is

$$\begin{aligned} A_0 x^k &= b - A_1 x^{k-1} & k = 1, 2, \dots & \quad (2.5.6) \\ x^0 &= 0 \end{aligned}$$

If equation (2.5.6) converges, $\lim x^k = x = (A_0 + A_1)^{-1} b$ and assuming A^{-1} and A_0^{-1} exist,

$$\begin{aligned} (A_0 + A_1)x &= b & x &= (A_0 + A_1)^{-1} b \\ & & &= (I + A_0^{-1} A_1)^{-1} A_0^{-1} b. \end{aligned}$$

If $\rho(A_0^{-1} A_1) < 1$, that is $A_0^{-1} A_1$ is a contraction, then

$$x = (I - A_0^{-1} A_1 + (A_0^{-1} A_1)^2 - \dots) A_0^{-1} b$$

or

$$A_0 x = \sum_{i=0}^{\infty} (-1)^i [A_0^{-1} A_1]^i b.$$

Define

$$x^k = A_0^{-1} \sum_{i=0}^{k-1} [A_0^{-1} A_1]^i b = A_0^{-1} [b - A_1 x^{k-1}]$$

is simply a rewriting of equation (2.5.5). The power series approach, therefore, has the same conditions for convergence, namely that A_0 is stable and $A_0^{-1}A_1$ is a contraction.

In order to obtain an approximate operation count using the Jacobi iterative method, assume that there are N subsystems, each of dimension n , and that an off-diagonal block of A has at most γ elements. Let the maximum number of elements in an A_{ii} be α , so $\beta = \alpha/n^2$. If the reduced order equations are solved using a transformation method, the total operation count is of the order of

$$10Nn^3 + \sigma [N^2n^3 + \gamma(N^3n + \frac{N^2n^2}{2})]$$

where σ is the number of iterations. At each iteration, cubic terms result from going in and out of the transform frame, raising the intriguing possibility of reduction of the operation count by a power of n for problems that can be set up to make this transformation-retransformation unnecessary. Storage requirements are approximately $N^2n^2 + 7/2 Nn^2$. If a sparse direct method is used, the count is approximately

$$6\beta^2 N^2 n^4 + \sigma [N^3 n \gamma + 3\beta n^3 N^2].$$

The storage requirements are roughly $\beta(6N^2n^3 + 2Nn^2) + 5n^2N^2$.

CHAPTER III
ERROR ANALYSIS

3.1 INTRODUCTION TO ERROR ANALYSIS

Any solution method for the Lyapunov equation can be analyzed in terms of computational speed, storage requirements, and accuracy. The first two are relatively straight forward to estimate and have been covered in the previous chapter. This chapter considers the third issue, error analysis.

Whenever computations are performed on a digital computer, one fundamental problem is to assess the numerical significance of the solutions. Four basic sources of error can be identified for a typical computation in scientific or engineering work:

- i) Modelling errors occur whenever mathematical models are used to represent physical processes.
- ii) Measurement errors represent the difference between ideal and computed parameters of the mathematical model as well as whatever errors exist in data that is input to the computer.
- iii) Truncation errors represent approximations made during the computation to functions that are not simple algebraic relationships. For example, if e^x is approximated by a finite series, then the neglected higher order terms are the truncation error.

- iv) Rounding errors represent the affects of the finite word length of the computer.

By definition, the first two types of error are not a property of the algorithm to be analyzed and will not be considered in this chapter. In addition, truncation errors will not be considered, simply because most of the algorithms analyzed here consist of a finite number of steps and therefore do not involve such errors. Two exceptions are the iterative decoupling method and the QR algorithm, which are infinite (iterative) procedures terminated with some appropriate tolerance test. Although the resulting truncation error will not be considered explicitly, the nature of the analysis used will nonetheless yield results that are consistent with those of the "exact" algorithms and allow useful comparisons to be made. For the remainder of this chapter then, the term error will refer to rounding error only.

There are two fundamental approaches to error analysis, namely forward and backward. The forward approach attempts to carry errors made at each computational step forward and compare the final result with the ideal one. The backward approach views the computed result as the exact solution of a perturbed equation. The latter approach is more modern and usually easier to perform and interpret. Its development is primarily due to Wilkinson [18], a well-known authority in several areas of error analysis. Although he has apparently not studied the Lyapunov equation, this chapter relies heavily on his work [18, 24]. The backward approach actually involves two steps,

representing the inherent conditioning of the problem and obtaining bounds for the perturbations of the ideal equation. The formulation of these two steps is naturally dependent on the problem being analyzed, but a simple description illustrates the basic philosophy.

$$\begin{array}{l}
 \text{data} \\
 \text{parameters}
 \end{array}
 \} \rightarrow \text{computation} \rightarrow \text{answer}$$

$$u_1, u_2, \dots, u_m \quad t \text{ digit approximation} \quad x_1, x_2, \dots, x_n$$

The conditioning of the problem is defined to be the sensitivity of relative errors in x to perturbations in the inputs u . If the sensitivity is high, then the problem is ill-conditioned. At each step in the computation, rounding errors are treated as effective perturbations of the inputs. In general, the conditioning depends on the parameters and the general problem being solved, but not on the specific algorithm used or t (computer word length). The effective perturbations, on the other hand, depend strongly on the specific algorithm, the parameters, t , and possibly x .

The remainder of this chapter is comprised of three sections. The first two consider the two steps of backwards analysis, conditioning and perturbation bounds, respectively. The third summarizes and compares the results obtained for the direct, transformation, and iterative decoupling Lyapunov equation solution methods.

3.2 CONDITIONING

In order to discuss the numerical conditioning of a linear set of equations, some basic tools of linear algebra are needed. The necessary ones are primarily basic and commonly known. Error analysis particularly relies on the use and manipulation of vector and matrix norms however, so a brief collection of definitions and relationships that will be important later are included here.

A convex set K , $K \in \mathbb{R}^n(\mathbb{C}^n)$, is a set of points such that $x, y \in K$, $0 \leq \lambda \leq 1 \rightarrow \lambda x + (1-\lambda)y \in K$. A convex body is a closed, bounded, convex set with interior points. K is an equilibrated convex body if

$$k \in K, \quad |w| \leq 1 \rightarrow wk \in K.$$

Notice that in this case the origin is interior to K . The equilibrated convex body of particular interest here is the unit sphere S , i.e.,

$$S = \{x \in \mathbb{C}^n \mid x^H x \leq 1\}.$$

Now let K be an equilibrated convex body. Then the norm of the vector $x \in \mathbb{R}^n(\mathbb{C}^n)$ with respect to K is defined [22]

$$\|x\|_K = \inf\{v \mid v \geq 0, x \in vK\} \quad (3.2.1)$$

and the (least upper) bound of the matrix $A \in \mathbb{R}^{n \times n}(\mathbb{C}^{n \times n})$ with respect to K is

$$\text{lub}_K(A) = \inf(\alpha \mid \alpha \geq 0, AK \subset \alpha K). \quad (3.2.2)$$

These definitions satisfy the following properties:

$$\begin{array}{ll}
 \text{i)} \quad \|x\| > 0 \text{ if } x \neq 0 & \text{v)} \quad \text{lub}(A) > 0 \text{ if } A \neq 0 \\
 \text{ii)} \quad \|\alpha x\| = |\alpha| \|x\| & \text{vi)} \quad \text{lub}(\alpha A) = |\alpha| \text{lub}(A) \\
 \text{iii)} \quad \|x+y\| \leq \|x\| + \|y\| & \text{vii)} \quad \text{lub}(A+B) \leq \text{lub}(A) + \text{lub}(B) \\
 \text{iv)} \quad \|Ax\| \leq \|x\| \text{lub}(A) & \text{viii)} \quad \text{lub}(AB) \leq \text{lub}(A) \text{lub}(B)
 \end{array} \quad (3.2.3)$$

where the subscript K has been omitted for notational simplicity only.

For any matrix A , there is at least one $x \neq 0$ such that

$$\|Ax\| = \|x\| \text{lub}(A).$$

A matrix norm is usually defined as any real valued function of the elements of A such that properties v) - viii) are satisfied with $\|\cdot\|$ replacing $\text{lub}(\cdot)$. The most commonly used vector norms are given by

$$\|x\|_p = (|x_1|^p + \dots + |x_n|^p)^{1/p} \quad (3.2.4)$$

$$p = 1, 2, \infty$$

where $\|x\|_p = \max(x_i)$. Analogous matrix norms that are subordinate to a vector norm satisfy

$$\|A\| = \max \|Ax\|, \quad \|x\| = 1$$

and are computed as

$$\|A\|_1 = \max_j \sum_i |a_{ij}|$$

$$\|A\|_\infty = \max_i \sum_j |a_{ij}|$$

$$\|A\|_2 = (\max \lambda(A^H A))^{1/2}. \quad (\text{spectral norm of } A)$$

Norms are said to be consistent or compatible if

$$\|Ax\| \leq \|A\| \|x\|.$$

Notice that subordinate norms must be consistent, but consistent norms are not necessarily subordinate. The useful euclidian norm $\|A\|_E$ is consistent with $\|x\|_2$, where

$$\|A\|_E = (\text{tr}(A^H A))^{1/2} = (\sum_{ij} |a_{ij}|^2)^{1/2}.$$

The eigenvalues of $A^H A$ are called the singular values of A , i.e., $\lambda(A^H A) = \sigma^2$. The eigenvalues of $A^H A$ are real and positive, so we can write

$$\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2 > 0$$

from which

$$\|A\|_2 = \sigma_1.$$

Notice that if S is the unit sphere, then

$$\text{lub}_S(A) = \sigma_1 = \|A\|_2. \quad (3.2.5)$$

A (greatest lower) bound for A with respect to the equilibrated convex body K is defined by

$$\text{glb}_K(A) = \sup\{\alpha \mid \alpha K \subset AK\}. \quad (3.2.6)$$

Now $\text{glb}_K(A) = 1/\text{lub}_K(A^{-1})$, and in particular for the unit sphere

$$\text{glbs}(A) = \alpha_n = 1/\|A^{-1}\|_2. \quad (3.2.7)$$

A few important relationships between the different matrix norms that are frequently useful in error analysis are summarized below. First, from

$$Ax = \lambda x$$

we see that

$$\|A\| \|X\| \geq \|AX\| = \|\lambda X\| = |\lambda| \|X\|,$$

so $|\lambda| \leq \|A\|$ for any norm.

Using this fact, a number of inequalities can be deduced, i.e.,

$$\|A\|_2^2 = \max \lambda(A^H A) \leq \|A^H A\|_E \leq \|A\|_E^2 = \text{tr}(A^H A) = \sum \lambda(A^H A) \leq n^{1/2} \|A\|_2,$$

or $\|A\|_2 \leq \|A\|_E \leq n^{1/2} \|A\|_2$.

Also $\|A\|_2^2 \leq \|A^H A\|_\infty \leq \|A^H\|_\infty \|A\|_\infty = \|A\|_1 \|A\|_\infty$.

Denoting $|A|$ as the matrix whose elements are $|a_{ij}|$, notice that

$\| |A| \| = \|A\|$ for all norms except $\|\cdot\|_2$, while

$\|I\| = 1$ for 1, 2, ∞ norms

$= n^{1/2}$ for $\|\cdot\|_E$.

For a linear set of equations the concept of numerical conditioning is expressed with some simple inequalities and quantified by a single number called the condition number. Consider the sensitivity of the solution of the set of equations

$$Ax = b \tag{3.2.8}$$

to variations in A and b . Let h and k be perturbations in x and b , respectively, where

$$A(x + h) = b + k, \quad (3.2.9)$$

or
$$h = A^{-1}k,$$

so
$$\|h\| = \|A^{-1}k\| \leq \|A^{-1}\| \|k\|. \quad (3.2.10)$$

From (3.2.9),

$$\|b\| \leq \|A\| \|x\|, \text{ or } \|x\| \geq \|b\| \|A\|^{-1},$$

and combining with (3.2.10) yields

$$\|h\| / \|x\| \leq \|A\| \|A^{-1}\| \|k\| / \|b\|. \quad (3.2.11)$$

Now let E be a perturbation of A and consider

$$(A + E)(x + h) = b \quad (3.2.12)$$

or
$$(A + E)h = -Ex.$$

Now
$$(A + E) = A(I + A^{-1}E), \text{ so}$$

$$(A + E)^{-1} \text{ exists if } \lambda(A^{-1}E) < 1.$$

In addition,

$$\begin{aligned}
 (I + X)^{-1} &= I - X + X^2 - X^3 + \dots \quad \text{if } \|X\| < 1, \quad . \\
 \|(I + X)^{-1}\| &= \|I - X + \dots\| \leq \|I\| + \|X\| + 1/2\|X^2\| + \dots \\
 &= \frac{1}{1 - \|X\|} .
 \end{aligned}$$

Assuming $\|A^{-1}E\| < 1$, and using the above inequality, equation (3.2.12) can be manipulated to yield

$$\frac{\|h\|}{\|X\|} \leq \frac{\|A\| \|A^{-1}\| \|E\| / \|A\|}{1 - \|A\| \|A^{-1}\| \|E\| / \|A\|} . \quad (3.2.13)$$

Now equations (3.2.11) and (3.2.13) are true for the spectral norm (a fortiori for the euclidian norm), and in each case the critical quantity relating the perturbations is defined to be the spectral condition number,

$$k(A) = \|A\|_2 \|A^{-1}\|_2 .$$

The relationship of (3.2.13) illustrates clearly the basic idea of backward error analysis applied to the linear matrix equation (3.2.8). The spectral condition number of A relates the sensitivity of the solution accuracy to the "size" of relative perturbations in A. It is a property of A and independent of solution methods, while the

matrix E, which represents the effects of round off errors, will depend explicitly on the specific algorithm used and computer word length. Note that the condition number of A is defined relative to the problem being solved. For instance, if X is the matrix whose columns are the eigenvectors of A (including pseudo eigenvectors if A has repeated roots) then the condition number of A with respect to the eigenvalue problem is defined to be $\|X\|_2 \|X^{-1}\|_2$ [24]. In this thesis, conditioning will always mean conditioning with respect to the solution of sets of linear equations.

As an interesting application of $K(A)$, suppose that $Ax = b$ is solved in some way and no rounding errors occur. In this case, each element of A is correct to t digits, so $|e_{ij}| \leq 2^{-t} |a_{ij}|$, or

$$\|E\|_E \leq 2^{-t} \|A\|_E, \quad \|E\|_2 \leq n^{1/2} 2^{-t} \|A\|_2,$$

and equation (3.2.13) yields

$$\frac{\|h\|_2}{\|x\|_2} = \frac{k(A) n^{1/2} 2^{-t}}{1 - k(A) n^{1/2} 2^{-t}}. \quad (3.2.14)$$

This shows that unless $n^{1/2} k(A) 2^{-t} \ll 1$ the computed solution will be numerically insignificant. A simple calculation illustrates the use of equation (3.2.14). Suppose that $k(A) = 100$, $n = 100$, and single precision is used on an IBM-370. In this case, the word length is

32 bits, of which 24 are given to the mantissa in floating point, i.e.,
 $t = 24$. Then

$$\begin{aligned} (K(A)n^{1/2})2^{-t} &\approx (K(A)n^{1/2})10^{-t} \quad (.30) \\ &= (10^3)10^{-7.2} = 10^{-4.2}, \end{aligned}$$

so we can expect only 4 or 5 significant figures in our solution!

A condition number for the Lyapunov equation can be similarly defined. The perturbed equation

$$(A + E)^T (P + H) + (P + H) (A + E) = -Q \quad (3.2.15)$$

or

$$(K_A + K_E) (P_V + H_V) = -Q_V,$$

leads to

$$\frac{\|H_V\|}{\|P_V\|} < \frac{\|K_A\| \|K_A^{-1}\| \|K_E\| / \|K_A\|}{1 - \|K_A\| \|K_A^{-1}\| \|K_E\| / \|K_A\|}, \quad (3.2.16)$$

and the condition number for the Lyapunov equation is

$$k(L_A) = \|K_A\|_2 \|K_A^{-1}\|_2, \quad K_A = A^T \times I + I \times A^T$$

In the remainder of this section, two topics are covered. First, some properties of the condition number are presented and discussed. The simpler notation of $Ax = b$ will be used, realizing that the

results apply to the Lyapunov equation analogously. Secondly, some properties of the Kronecker matrix are investigated in an attempt to obtain a useful relationship between its condition number and that of A.

The spectral condition number depends on the singular values of A, for combining equations (3.2.5) and (3.2.7) we have

$$k(A) = \|A\|_2 \|A^{-1}\|_2 = \text{lubs}(a)/\text{glbs}(A) = \sigma_1/\sigma_n, \quad (3.2.17)$$

where

$$\sigma_1^2 = \max \lambda(A^H A), \text{ and } \sigma_n^2 = \min \lambda(A^H A).$$

The computation of $k(A)$ then, is a major task and for this reason other norm measures are often used in practice to bound $k(A)$, usually the euclidian norm, i.e.,

$$k(A) = \|A\|_2 \|A^{-1}\|_2 \leq \|A\|_E \|A^{-1}\|_E.$$

In this case the obvious difficulty is computing A^{-1} . Suppose, however, that an approximation of A^{-1} is obtained such that

$$I - AC = R, \quad C \approx A^{-1},$$

or
$$A^{-1}(I - R) = C,$$

$$A^{-1} = C(I - R)^{-1},$$

then

$$\|A^{-1}\| \leq \frac{\|C\|}{1-\|R\|} \quad \text{if } \|R\| < 1. \quad (3.2.18)$$

This use of an approximate inverse is apparently a common way in practice to assess the accuracy of a computed solution x [17, 18, 21, 23, 33], but assuming that C is computed when x is, it does not provide an a priori estimate. Now recall the notation introduced with the iterative decoupling solution methods of the previous chapter. (The notation $Ax = b$ is used here, but the results, again, are analogous to $L_A(P) = -Q$.) We had

$$A = A_0 + A_1, \quad A_0 \text{ (block) diagonal,}$$

and the condition for convergence of the methods was

$$P(A_0^{-1}A_1) < 1,$$

which is certainly true if $\|A_0^{-1}A_1\| < 1$. This is the same condition on R in (3.2.18), however, so

$$\|A^{-1}\| \leq \frac{\|A_0^{-1}\|}{1-\|A_0^{-1}A_1\|}, \quad \|A_0^{-1}A_1\| < 1, \quad (3.2.19)$$

showing that if any of the iterative methods are used to solve $Ax = b$ then (3.2.19), together with the easily computed $\|A\|_E$, provide a useful approximation to the condition number.

Intuitively one may be tempted to use the ratio $\lambda_{\max}(A)/\lambda_{\min}(A)$ as a measure of the conditioning of A with respect to $Ax = b$. This can be a very misleading approximation. Certainly

$$\|A\|_2 \geq |\lambda_{\max}(A)| = |\lambda_1| ,$$

$$\|A^{-1}\| \geq |\lambda_{\max}(A^{-1})| = 1/|\lambda_{\min}(A^{-1})| = 1/|\lambda_n| ,$$

and
$$k(A) = \sigma_1/\sigma_n \geq |\lambda_1|/|\lambda_n| , \quad (3.2.20)$$

with equality being obtained only for matrices of special form, e.g., symmetric and anti-symmetric matrices. In order to illustrate why an ill-conditioned matrix may not have a small eigenvalue, consider the following:

normalize A such that $\|A\|_2 = 1$, and

let P be the unitary matrix (which always exists) such that $P^H A P = \text{diag}(\lambda_i) + T$,

$= D + T$, T strictly upper triangular.

Now
$$P^H A^{-1} P = (D + T)^{-1} = (I - R + R^2 - \dots - (-1)^{n-1} R^{n-1}) D^{-1} ,$$

where
$$R = D^{-1} T .$$

We have

$$\|T\|_2 \leq \|T\|_E \leq \|D + T\|_E = \|A\|_E \leq n^{1/2} \|A\|_2 = n^{1/2},$$

so
$$\|R\|_2 \leq \|D^{-1}\|_2 \|T\|_2 \leq n^{1/2} / |\lambda_n|,$$

and
$$\|A^{-1}\|_2 \leq (1 + n^{1/2} / |\lambda_n| + (n^{1/2} / |\lambda_n|)^2 + \dots + n^{1/2} / |\lambda_n|^{n-1}) 1 / |\lambda_n|$$

from which

$$|\lambda_n| \leq \frac{n^{1/2(n+1)}}{\|A^{-1}\|_2} \quad 1/n \quad \doteq \quad \frac{n^{1/2}}{\|A^{-1}\|_2^{1/n}} \quad \cdot \quad (3.2.21)$$

This shows that for fixed n , as $\|A^{-1}\|_2$ approaches infinity, $|\lambda_n|$ approaches zero, but very slowly. Because of the direction of the inequality (3.2.20), we can only conclude $|\lambda_1| / |\lambda_n|$ large implies that ill-conditioning is to be expected.

Because the spectral condition number is an important quantity in backwards error analysis, it is unfortunate that it is difficult to compute. The use of an approximate inverse has been mentioned, but this too requires considerable computation. If an a posteriori estimate of a solution's accuracy is sufficient, the most practical procedure probably involves the use of iterative refinement, and this is illustrated in the next section. If a linear equation solution is a step in a more complicated computation, however, neither approach may be practical and an a priori estimate of accuracy would be very useful. For example, with the iterative decoupling solution method for the Lyapunov equation, a simple error analysis (Section 3.4) relates the final solution error to those that occur at each step, so

an a priori estimate of conditioning, combined with a perturbation bound, could be used to choose a convergence tolerance. Additional research into the nature of the condition number, however, has led to the conclusion that a simple, easily computed approximation for it is unlikely to be found [22, 24, 35, 36, 37, 38, 39]. A few theoretical properties of the condition number are included here that illustrate the practical difficulty, but aid understanding.

First, Wehl [35] has proven a number of interesting inequalities that relate the singular values and eigenvalues of A. Let

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$$

$$\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2 > 0,$$

then the Wehl inequalities are

$$|\lambda_1| |\lambda_2| \dots |\lambda_k| \leq \sigma_1 \sigma_2 \dots \sigma_k \quad k=1, 2, \dots, n-1,$$

$$|\lambda_1| |\lambda_2| \dots |\lambda_n| = \sigma_1 \sigma_2 \dots \sigma_n \quad (3.2.22)$$

and

$$|\lambda_1|^s + |\lambda_2|^s + \dots + |\lambda_k|^s \leq \sigma_1^s + \sigma_2^s + \dots + \sigma_k^s$$

$s > 0$, real.

Manipulating (3.2.22) yields a bound for $R(A)$

$$K(A) \leq \sigma_1^n / \det(A) \leq \|A\|_E^n / \det(A). \quad (3.2.23)$$

This is probably the simplest bound not involving A^{-1} which has been found. A geometrical insight into the spectral condition number results by recognizing that A maps the unit sphere onto an ellipse whose axes lie along the eigenvectors of $A^H A$. Let S be the unit sphere, and y_k an eigenvector of $A^H A$. So

$$A^H A y_1 = \sigma_1^2 y_1 \quad A^H A y_n = \sigma_n^2 y_n$$

while from the definitions of matrix bounds (3.2.2)

$$\max_{x \in S} \|Ax\| \equiv \|Ax_1^*\| = \sigma_1, \quad \min_{x \in S} \|Ax\| \equiv \|Ax_n^*\| = \sigma_n$$

then $Ax_1^* = \sigma_1 y_1, \quad Ax_n^* = \sigma_n y_n. \quad (3.2.24)$

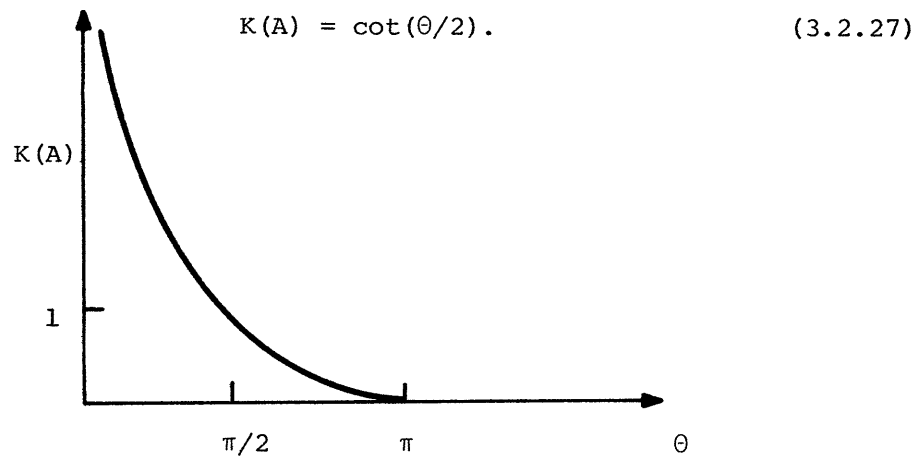
Again, $K(A) = \sigma_1 / \sigma_n$, which is the ratio of the lengths of the major and minor axes of the ellipse AS . Another geometrical representation of the condition number formalizes the intuitive notion that ill-conditioning is related to the distortion of the ellipse (notice that as $\det(A) \rightarrow 0$, the ellipse degenerates to a plane in at least one dimension). If the angle θ is defined by

$$(k(A) - k^{-1}(A))/(k(A) + k^{-1}(A)) = \cos \theta \quad (3.2.25)$$

then the inequality of Wielandt [22] is given by

$$\frac{|\mathbf{x}^H \mathbf{A} \mathbf{A} y|^2}{(\mathbf{x}^H \mathbf{A} \mathbf{A} \mathbf{x})(y^H \mathbf{A} \mathbf{A} y)} \leq \cos^2 \theta, \quad \mathbf{x}, \mathbf{y} \text{ any orthogonal pair of vectors} \quad (3.2.26)$$

The geometrical interpretation is that θ is the minimal angle between $\mathbf{A}x$ and $\mathbf{A}y$, for all orthogonal pairs x and y . Applying a standard trigonometric identity to (3.2.21), we obtain



This discussion of the condition number is concluded with the following simple example. The mapping of the unit sphere is illustrated, along with the relationships between the singular values, eigenvalues, and norms of A .

$$A = \begin{bmatrix} -1 & a \\ 0 & -2 \end{bmatrix} \quad K^2(A) = \frac{5+a^2 + ((a^2+1)(a^2+9))^{1/2}}{5+a^2 - ((a^2+1)(a^2+9))^{1/2}}$$

It is easy to verify that $K(A)$ is an increasing function of a (for $a > 0$).

Let $a = 3$, so

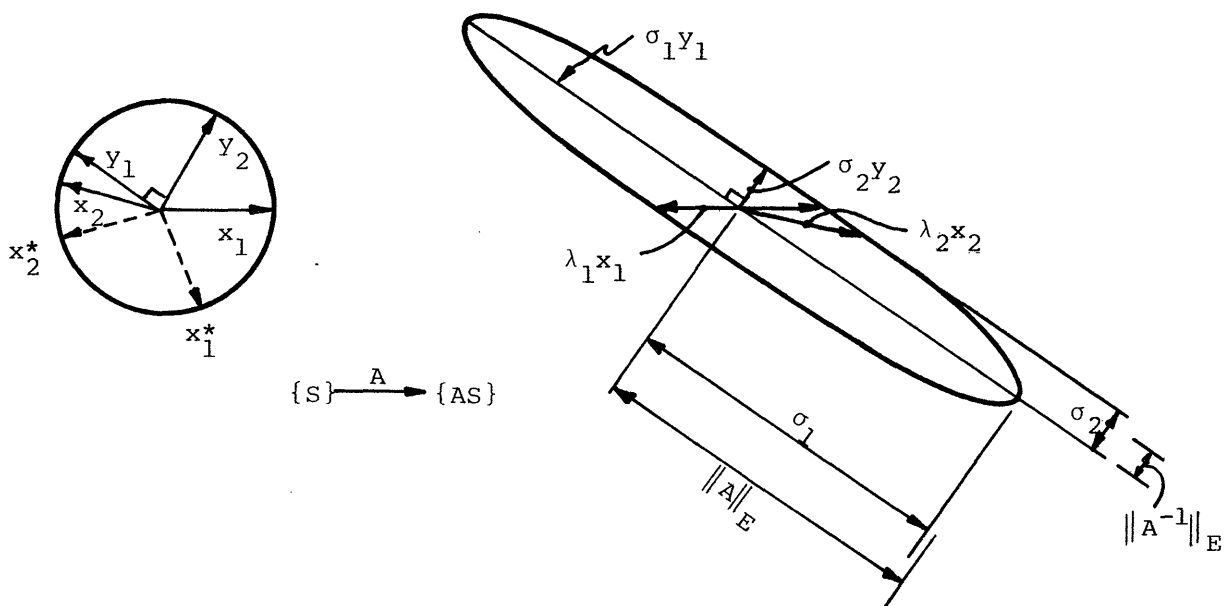
$$A^H A = \begin{bmatrix} 10 & -6 \\ -6 & 4 \end{bmatrix} \quad \sigma_1^2 = (3.70)^2, \quad y_1 = \begin{bmatrix} -.851 \\ .525 \end{bmatrix}$$

$$K(A) = \frac{3.70}{.54} = 6.85 \quad \sigma_2^2 = (0.54)^2, \quad y_2 = \begin{bmatrix} .525 \\ .851 \end{bmatrix} .$$

Similarly, the eigenvalues and eigenvectors of A are

$$\lambda_1 = -1 \quad x_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \lambda_2 = -2 \quad x_2 = \begin{bmatrix} -.950 \\ .316 \end{bmatrix}$$

and $\|A\|_E = 3.74 \quad \|A^{-1}\|_E = 0.53$



Again, it is emphasized that although the preceding discussion of conditioning was phrased in terms of the equation $Ax = b$, analogous remarks apply for the Lyapunov equation, whose condition number was defined in equation (3.2.16) to be

$$k(L_A) = \|K_A\|_2 \|K_A^{-1}\|_2, \quad K_A = A^T \times I + I \times A^T.$$

Because the Kronecker matrix K_A is $n^2 \times n^2$ however, the practical difficulties in either computing or approximating its condition number are even more severe than in the case of A itself which is $n \times n$. For this reason the properties of the Kronecker matrix were investigated to see if its condition number is related in any simple way to that of A . The conclusion reached is that no such relationship exists, and the purpose of the remainder of this section is to briefly summarize why this is true.

First, the properties of the Kronecker matrix are naturally very dependent on those of the Kronecker product $A \times B$. If A is $n \times m$ and B is $l \times p$, then the matrix $A \times B$ is $(nl) \times (mp)$ and partitioned into the (mn) blocks $(a_{ij}B)$. Many interesting properties of the Kronecker product are not developed here, but are given by Bellman [19] and Barnett and Storey [30]. The eigenvalues of $A \times B$ are important here, however, so consider

$$\begin{aligned} A & n \times n, & Ax^i &= \lambda_i x^i, \\ B & m \times m, & By^j &= \mu_j y^j, \end{aligned}$$

then $A \times B$ has eigenvalues $\lambda_i \mu_j$ ($i=1, 2, \dots, n; j=1, 2, \dots, m$) and eigenvectors $z_{ij} \in R^{nm}$, where

$$z_{ij} = \begin{bmatrix} x_1^i y^j \\ x_2^i y^j \\ \cdot \\ \cdot \\ x_n^i y^j \end{bmatrix} \quad (3.2.28)$$

This result is easy to verify by expanding the defining relationship $(A \times B)z_{ij} = \lambda_i \mu_j z_{ij}$, using the block structure of $(A \times B)$. We can use this result to easily find the eigenvalues and eigenvectors of the Kronecker matrix by using a theorem of Frobenius [19], i.e., if the roots of A are $\lambda(A)$, then the roots of the polynomial function of A , $f(n)$, are $\lambda(f(A)) = f(\lambda)$. Consider

$$(I_n + \epsilon A) \times (I_m + \epsilon B) = I_n \times I_m + \epsilon (A \times I_m + I_n \times B) + \epsilon^2 (A \times B).$$

From the above we have

$$\lambda[(I_n + \epsilon A) \times (I_m + \epsilon B)] = (1 + \epsilon \lambda_i)(1 + \epsilon \mu_j) = 1 + \epsilon(\lambda_i + \mu_j) + \epsilon^2 \lambda_i \mu_j.$$

and the matrix $(A \times I_m + I_n \times B)$, therefore has eigenvalues $(\lambda_i + \mu_j)$ with associated eigenvectors z_{ij} . By letting $A = A^T$ and $B = A^T$, this result applies to the Kronecker matrix directly.

Now the singular values of $A \times B$ may be found using the relations (3.2.28), i.e.,

$$(A \times B) = [\lambda (A \times B)^H (A \times B)]^{1/2} = [\lambda (A^H A \times B^H B)]^{1/2} = \sigma(A) \sigma(B).$$

The approach used to obtain the eigenvalues of $(A \times I_m + I_n \times B)$ does not extend to its singular values, however. The essential difficulty may be seen by considering

$$(A \times I_m + I_n \times B)^H (A \times I_m + I_n \times B) = (A^H A \times I_m + I_n \times B^H B) + A^H \times B + A \times B^H \quad (3.2.29)$$

If the last two terms were absent, then the singular values of $(A \times I_m + I_n \times B)$ would be $\sigma_i(A) + \sigma_j(B)$ as was the case for the eigenvalues. Although the eigenvalues of each term on the right hand side of (3.2.29) are products or sums of the eigenvalues of A and B, no expression has been found for the eigenvalues of the complete expression. This is basically unsurprising, for the roots of a sum of matrices are essentially unrelated to the roots of each single matrix, except for matrices of very special form [20].

Some insight may be gained from considering the inequality given previously,

$$k(L_A) \geq \frac{|\max \lambda(K_A)|}{|\min \lambda(K_A)|} = \frac{\max |\lambda_i + \lambda_j|}{\min |\lambda_i + \lambda_j|} . \quad (3.2.30)$$

Suppose that A is 2x2, with a complex pair of eigenvalues $\alpha \pm jw$, $\alpha \ll w$. Then for A $|\lambda_{\max}|/|\lambda_{\min}| = k(A) = 1$, while from equation (3.2.30)

$$k(L_A) \geq 2|\alpha + jw|/\alpha \doteq 2w/\alpha \gg 1.$$

This is significant because lightly damped poles are common in many engineering systems, and this simple example shows that the condition number of the Lyapunov equation can become very large although the A matrix itself is very well-conditioned!

Finally, consider the commonly used approximation for the condition number,

$$k(L_A) \leq \|K_A\|_E \|K_A^{-1}\|_E.$$

Now it is easy to see that

$$\|K_A\|_E \leq 2n^{1/2} \|A\|_E ,$$

but again no reasonable bound for $\|K_A^{-1}\|_E$ in terms of A can be obtained. Barnett and Storey [30] give an explicit expression for the inverse of the Kronecker matrix K_A , from which a bound can be obtained, but the result is very pessimistic and has no practical value.

3.3 PERTURBATION BOUNDS

In the previous section the first step in the backward error analysis of the linear matrix equation was developed and extended to the Lyapunov equation. This section pursues the second and more difficult step, which is to assess specific algorithms and obtain bounds for the matrix that represents the equivalent perturbations in the elements of A or L_A . In order to obtain a perturbation bound, an algorithm must be analyzed in great detail, i.e., broken down to the level of each addition and multiplication (division). At this level, the effect of a finite word length is assumed to be such that $\bar{x}(\cdot)\bar{y} \equiv (x(\cdot)y)(1+e)$, $|e| \leq 2^{-t}$, where the bar denotes the computed quantity and (\cdot) can be any of the four basic arithmetic operations. With this assumption, a number of intermediate error bounds, e.g., for inner products, can be obtained. These results have been placed in Appendix B, and are the building blocks of the bounds obtained here. It is generally true that this process of successively bounding elementary operations yields conservative results, and this fact motivates using a statistical approach (see [34] for example), but the bounding approach is preferred by Wilkinson and used here. The purpose of this section, which relies heavily on the works of Wilkinson, is to describe simply the algorithms of interest and the essential steps of their error analysis, omitting many of the details. Once the fundamentals of Appendix B and the basic concepts of backward analysis are grasped, the intermediate results follow in a simple, but somewhat tedious, manner.

Although several of the transformation solution methods of the Lyapunov equation are quite different from the direct method, their error analyses are similar and useful comparisons can be made. Both types of methods rely on a series of elementary similarity transformations at some point in the algorithm, so it is natural to consider these elementary operations separately before proceeding. There are two basic types of similarity transformations, unitary and non-unitary or elementary [22, 24]. The basic formulation of a similarity transformation using elementary matrices is

$$X A X^{-1} = B$$

where X is an elementary matrix. There are many types of elementary matrices, two of which will be used later:

- i) The matrices I_{ij} , equal to the identity matrix except in rows and columns i and j , which are of the form

$$\begin{array}{cc} \left[\begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array} \right] & \begin{array}{l} \text{row } i \\ \text{row } j \end{array} \\ \text{col. } i & \text{col. } j \end{array}$$

Note that $I_{ij} I_{ij} = I$, i.e., it is orthogonal. Premultiplication by I_{ij} interchanges rows i and j and postmultiplication interchanges columns i and j . Any permutation matrix is a product of matrices of this type.

- ii) The matrices N_i , equal to the identity matrix except for the i^{th} column, which is (written as a row vector)

$$(0, 0, \dots, 0, 1, -n_{i+1,i}, -n_{i+2,i}, \dots, -n_{ni}).$$

The inverse of N_i is obtained by changing the signs of the $n_{k,i}$ elements. An important property is that the product $N_1 N_2 \dots N_{n-1}$ is a lower triangular matrix with unit diagonal elements, while the $i-j^{\text{th}}$ element is $-n_{ij}$ ($i > j$).

The practical use of elementary matrices in computations is illustrated by the fact that there always exists an X that is the product of elementary matrices such that

$$Xx = ke_1 \quad \text{with: } \begin{array}{l} e_1 \text{ unit vector} \\ k \text{ real constant} \\ x \text{ arbitrary vector} \end{array} \quad (3.3.1)$$

The other type of similarity transformations is based on unitary matrices. Two kinds of elementary unitary matrices are used in practice. The first kind is called a plane rotation; the matrix $R(p,q)$ is defined by

$$\begin{array}{ll} r_{pp} = e^{j\alpha} \cos\theta & r_{pq} = e^{j\beta} \sin\theta \\ r_{qp} = -e^{-j\beta} \sin\theta & r_{qq} = e^{-j\alpha} \cos\theta \\ r_{ij} = \delta_{ij} \text{ otherwise} & \alpha, \beta, \theta \text{ real} \end{array}$$

With $\alpha=\beta=0$, the transformation is equivalent to a rotation in the p,q plane through the angle θ . The four non-trivial elements can be expressed also by

$$\begin{aligned} r_{pp} &= \bar{x}/r & r_{qq} &= x/r & r_{pq} &= \bar{y}/r & r_{qp} &= -y/r \\ r^2 &= |x|^2 + |y|^2 \end{aligned} \quad (3.3.3)$$

The pre (post)-multiplication of A with $R(p,q)$ affects (independently) only rows (columns) p and q . Let $z = R(p,q)x$ and construct $R(p,q)$ using (3.3.3) with x_p for x and x_q for y . Then $z_p = r$, $z_q = 0$, and $z_i = x_i$ otherwise.

The second kind of elementary unitary matrix is called an elementary Hermitian matrix. These are of the form

$$P(r) = I - 2w(r)w^H(r), \quad \text{where } \|w(r)\|_2 = 1 \quad (3.3.4)$$

and
$$w^T(r) = (0, 0, \dots, 0, w(r)_{r+1}, w(r)_{r+2}, \dots, w(r)_n).$$

It is easy to verify that $P(r)$ is Hermitian and unitary. When $w(r)$ is real, $P(r)$ is real, symmetric and orthogonal, Pre (post)-multiplication of A by $P(r)$ affects only rows (columns) $r+1$ to n and treats each column (row) of A independently. Given a vector x , we can choose $w(r)$ such that $P(r)x$ has zero elements in positions $r+2, r+3, \dots, n$. To do this consider

$$S^2 = |x_{r+1}|^2 = \dots + |x_n|^2,$$

$$T = (|x_{r+1}|^2 S^2)^{1/2}, \quad H = S^2 + T \quad (3.3.5)$$

$$w^T(r) = (0, 0, \dots, w_{r+1} (1+S^2/T), x_{r+2}, \dots, x_n).$$

When x is real (3.3.5) can be expressed more simply by

$$S^2 = x_{r+1}^2 + \dots + x_n^2 \quad H = S^2 + |x_{r+1}| S,$$

$$w^T(r) = (0, 0, \dots, w_{r+1} \pm S, x_{r+2}, \dots, x_n), \quad (3.3.6)$$

where $\pm S$ is chosen to have the same signs as x_{r+1} . Notice that the above transformation $P(r)x$ is identical to the succession of transformations $R(r+1, r+2), R(r+1, r+3), \dots, R(r+1, n)$ applied to x . Since each column of A is treated independently in forming $P(r)A$, we can reduce elements $r+2, \dots, n$ of any column of A to zero without affecting the first r elements.

In matrix problems similarity transformations are a major part of many numerical algorithms [24]. The two basic types differ in an important way that may affect the numerical stability of the algorithm. The difference is that unitary similarity transformations preserve the conditioning of the original matrix while elementary transformations, in general, do not. Let U be unitary and consider

$$U A U^H = B. \quad U^H U = I$$

Now

$$\begin{aligned} \|B\|_2^2 &= \|UAU^H\|_2^2 = \max \lambda (UAU^H)^H (UAU^H) \\ &= \max \lambda (UA^H AU^H) = \max \lambda (A^H A) = \|A\|_2^2, \end{aligned} \quad (3.3.7)$$

so the spectral norm is invariant under unitary similarity transformations. The difficulty with elementary similarity transformations can be illustrated simply. Equation (3.3.6) shows how to construct a $P(r)$ such that for

$$y = P(r)x, \quad y_k = 0 \quad k = r+2, \dots, n.$$

Suppose we attempt the same transformation with an elementary matrix N_r , i.e.,

$$y = N_r x \quad (3.3.8)$$

Clearly the appropriate elements of N_r must be

$$n_{k,r} = x_k/x_r, \quad k > r,$$

and the transformation breaks down if $x_r = 0$ and $x_k \neq 0$ for some $k > r$. Generally, if x_r is small relative to some x_k , then the rounding errors become large. (The details of why this is true are outlined in Appendix B.) This difficulty of numerical instability can be greatly

reduced using the familiar operation of pivoting, which is affected using the orthogonal matrices I_{ij} mentioned previously. To illustrate this, and to complete the transformation (equation 3.3.8), let q be the smallest integer such that

$$|x_q| = \max_{i \geq r} |x_i|$$

Then let $z = I_{rq} x$ and again choose N_r so that $N_r z$ has zero elements in positions $r+1, \dots, n$. We have

$$n_{ir} = z_i/z_r, \quad |n_{ir}| \leq 1$$

and the difficulty is avoided. The combination I_{ra} and N_r is called a stabilized elementary matrix.

Theoretically then, unitary similarity transformations are superior with respect to numerical stability. Wilkinson (and others, e.g. [17, 22]) states that in practice stabilized elementary transformations are almost as stable as unitary transformations, provided the pivoting strategy is successful. The latter fact complicates the a priori error analysis, however.

The basic properties of elementary similarity transformation have been introduced. In the remainder of this section algorithms are formulated and perturbation bounds are given for the following problems:

- i) $Ax = b$
- ii) Transformation to Hessenberg form
- iii) Transformation to real Schur form
- iv) Transformation to Companion form.

The results of these analyses will then be extended to direct and transformation Lyapunov solution methods in the final section of this chapter.

If x is a solution of

$$Ax = b \quad (3.3.9)$$

and S is a square, nonsingular matrix, then it is also a solution of

$$SAX = Sb \quad (3.3.10)$$

The basic idea is to construct a simple matrix S such that SA is upper triangular. This is accomplished in a series of steps, each of which produces a set of equations equivalent to equation (3.3.9). At the $(r-1)^{\text{th}}$ step we have

$$A_{r-1}x = b_{r-1} \quad (3.3.11)$$

where A_{r-1} is upper triangular in the first $r-1$ columns. The r^{th} step consists of finding an elementary matrix of the form N_r such that $N_r A_{r-1}$

is upper triangular in the first r columns and is identical to A_{r-1} in the first $(r-1)$ rows and columns. Thus, the r^{th} step is precisely the problem covered previously (e.g., equation 3.3.8). If we ignore the pivoting issue for now and denote the elements of A at the r^{th} step a_{ij}^r , then the critical elements of N_r are given by

$$n_{ir} = a_{ir}^{r-1} / a_{rr}^{r-1} \quad , \quad (3.3.12)$$

and
$$A_r = N_r A_{r-1} \quad , \quad b_r = N_r b_{r-1} \quad (3.3.13)$$

Combining equations (3.3.13) for $r=1$ to $n-1$, we have

$$N_{n-1} \dots N_2 N_1 A_0 = A_{n-1} \quad , \quad N_{n-1} \dots N_2 N_1 b_0 = b_{n-1} \quad ,$$

where $A_0 = A$ is the original matrix. Recalling the property of the product $N_{n-1} \dots N_2 N_1$ mentioned previously, define

$$L = N_1^{-1} N_2^{-1} \dots N_{n-1}^{-1} \quad (\text{unit lower triangular matrix})$$

where
$$l_{ij} = n_{ij} \quad i > j \quad (3.3.14)$$

$$= \delta_{ij} \quad \text{otherwise.}$$

The arguments made earlier regarding pivoting definitely apply here, but it can be readily verified that the inclusion of a matrix I_{rq} at the r^{th} step does not affect the essential result, which is

$$A_0 x = L A_{n-1} x \equiv LUx = b. \quad (3.3.15)$$

L and U are lower and upper triangular matrices, respectively, so x may be found by performing simple forward and back substitution. Of course rounding errors occur at each step, so the computed L and U correspond to the exact triangularization of $A_0 + E$, i.e.,

$$LU = A_0 + E. \quad (3.3.16)$$

For $n = 3$,

$$L = \begin{bmatrix} 1 & & & \\ n_{21} & 1 & & \\ n_{31} & n_{32} & 1 & \end{bmatrix} \quad U = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ & a_{22} & a_{23} \\ & & a_{33} \end{bmatrix}$$

This basic bounds outlined in Appendix B can be applied for each step in the process to construct the perturbation matrix E. If the pivoting strategy is successful, $|n_{ij}| \leq 1$. If the maximum element of A_r is denoted by g, then

$$|E| \leq 2g2^{-t} \begin{bmatrix} 0 & 0 & 0 & . & . & . & 0 & 0 \\ 1 & 1 & 1 & . & . & . & 1 & 1 \\ 1 & 2 & 2 & . & . & . & 2 & 2 \\ 1 & 2 & 3 & . & . & . & 3 & 3 \\ . & . & . & . & . & . & . & . \\ 1 & 2 & 3 & . & . & . & n-1 & n+1 \end{bmatrix} \quad (3.3.17)$$

Similar analyses can be applied to the solution of the triangular sets of equations

$$(L + \delta L)y = b \qquad (U + \delta U)x = y \qquad (3.3.18)$$

Before collecting these results, however, it is worthwhile to repeat an argument made by Wilkinson [18] regarding the solution of triangular sets of equations that has important practical consequences. To illustrate this consider solving $(L + \delta L)y = b$. Assume for now that $\|\delta L\| \leq n2^{-t}\|L\|$. Equation (3.2.13) then provides the estimate

$$\frac{\|y - L^{-1}b\|}{\|L^{-1}b\|} \leq \frac{n2^{-t}\|L\|\|L^{-1}\|}{1 - n2^{-t}\|L\|\|L^{-1}\|} \qquad (3.3.19)$$

This implies that if $\|L\|\|L^{-1}\|$ is large, then the accuracy of the computed solution will deteriorate. However, in practice, large errors do not occur in the solution of triangular sets of equations and Wilkinson has found that a realistic estimate is

$$\frac{\|y - L^{-1}b\|}{\|L^{-1}b\|} \leq f(n)2^{-t} \quad , \qquad (3.3.20)$$

where $f(n)$ is a linear function of n and the accuracy is not affected by the condition number of the matrix L . This contrasts greatly with the general case of solving $Ax = b$, and it is important to realize that this is a special property of triangular sets of equations [18]. The

practical consequence for the general case is that when b is such that the corresponding solution reveals the ill-conditioning of A , it is the errors due to the LU factorization which limit the accuracy of the solution of $Ax = b$.

Combining equations (3.3.16, 3.3.18), we have

$$(A + E + L\delta U + U\delta L + \delta L\delta U)x = b,$$

$$\text{or} \quad (A + K)x = b. \quad (3.3.21)$$

Taking the norm of equation (3.3.17) and of analogous results for the other perturbation bounds yields

$$\|E\|_{\infty} \leq 2g2^{-t}(n/2 + 1)(n - 1)$$

$$\|\delta L\|_{\infty} \leq 1/2(n^2 + n + 2)2^{-t}$$

$$\|\delta U\|_{\infty} \leq g/2(n^2 + n + 2)2^{-t}$$

$$\|L\|_{\infty} \leq n$$

$$\|U\|_{\infty} \leq gn,$$

so
$$\|K\|_{\infty} \leq 2^{-t}g(2n^2 + n^3).$$

Again $g = \max|a_{ij}|$, but Wilkinson has found that $\max|a_{ij}^r| \leq 8\max|a_{ij}|$ for almost all matrices, when pivoting is used. So if A is initially normalized such that $|a_{ij}| \leq 1/8$, then g is essentially 1. In addition he states

that this bound is pessimistic, which is to be expected as statistical variations have been neglected throughout. For practical purposes, the expected error is

$$\|K\|_{\infty} \leq gn2^{-t} \quad (3.2.22)$$

with inner product accumulation. If accumulation is not available, the factor n becomes $n^{3/2}$. Before beginning the analysis of the reduction to Hessenberg form, a few concluding remarks are appropriate regarding the solution of $Ax = b$. First, the factorization algorithm presented was chosen to emphasize the use of similarity transformations because they conceptually link all the analyses of this section. Algorithms which carry out this process in a computationally different and more efficient manner are well known [21]; in addition, if accumulation of inner products is available, a direct LU factorization scheme provides a better bound for E , namely $\|E\|_{\infty} \leq g \cdot n2^{-t}$. The result of equation (3.2.22) still applies, however. Secondly, it is important to clearly recognize the difference between residuals and accuracy. In order to illustrate this, call the computed solution x' , so that the residual vector is

$$r' = b - Ax' \quad (3.3.23)$$

Combining equations (3.3.21, 3.3.22) then

$$\|r'\|_{\infty} = \|b - Ax'\|_{\infty} \leq gn2^{-t}\|x'\| .$$

Thus r' is bound to be small relative to x' regardless of the accuracy of x' . For an a priori accuracy estimate we use equation (3.2.13). Of course one conclusion of the last section was that finding the actual condition number is a major computation, and that a reliable a priori estimate of it is not generally available (except for the use of an approximate inverse). Another approach, which is more practical, employs iterative refinement of the solution.

Define a sequence of back-substitutions using the computed L and U by

$$\begin{aligned} r^S &= b - Ax^S \\ x^{S+1} &= x^S + (LU)^{-1}r^S . \end{aligned} \tag{3.3.24}$$

Wilkinson proves that if

$$\|A^{-1}\| \|E\| < 2^{-p} , \quad p > 1$$

then

$$\|x^{S+1} - x\| \leq \frac{2^{-p}}{p-2^{-p}} \|x^S - s\| . \tag{3.3.25}$$

For $p \geq 2$, at least p significant bits are gained per iteration. The additional work is $O(1/n)$ times the original factorization, but it is essential that the residuals are accumulated in double precision. For most A matrices a single refinement is sufficient, in which case output

of the quantity $\|x^2 - x^1\|_\infty / \|x^2\|_\infty$ provides a reliable assessment of the accuracy of x^2 .

The second problem for which perturbation bounds are required is that of transforming a general matrix A to Hessenberg form, H. Recall that a Hessenberg matrix is of the form

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} & \cdot & \cdot & \cdot & h_{1n} \\ h_{21} & h_{22} & h_{23} & & & & \cdot \\ 0 & h_{32} & h_{33} & & & & \cdot \\ & 0 & h_{43} & h_{44} & & & \cdot \\ & & 0 & & \cdot & & \cdot \\ & & & \cdot & & & h_{n-1,n} \\ & & & & \cdot & & \\ & & & 0 & h_{n,n-1} & h_{nn} & \cdot \end{bmatrix}$$

For symmetric matrices a tridiagonal matrix is analogous to H, and methods that affect the transformation are conceptually identical. The transformation is carried out in a series of n-1 steps (A n×n); at the beginning of the rth step we have

$$A_{r-1} = \begin{bmatrix} H_{r-1} & \begin{array}{c} | \\ | \\ | \end{array} & C_{r-1} \\ \hline 0 & \begin{array}{c} | \\ b_{r-1} \\ | \end{array} & B_{r-1} \end{bmatrix} \begin{matrix} r \\ n-r \end{matrix} \quad (3.3.26)$$

b_{r-1} is $(n - r \times 1)$,

and the problem is to construct a matrix X_r such that

$$A_r = X_r A_{r-1} X_r^{-1} \quad (3.3.27)$$

Any of the transformation matrices introduced at the beginning of this section can be used. The use of plane rotation matrices is referred to as Given's method, while Householder's method [22, 24] employs the unitary Hermitian transformation matrices. Householder's method will be used to illustrate the basic process, so X_r is replaced by P_r , i.e.,

$$P_r = \begin{bmatrix} I & & & \\ & \vdots & & \\ & & \vdots & \\ & & & Q_r \end{bmatrix} \begin{matrix} r \\ \\ \\ n-r \end{matrix}$$

where $Q_r = I - 2v_r v_r^T$ and we assume that the elements of A are real (P_r is orthogonal). The result of the multiplications (3.3.27) with P_r is

$$\begin{aligned} H_r &= H_{r-1} \quad , \quad C_r = C_{r-1} Q_r \\ b_r &= Q_r b_{r-1} \quad , \quad B_r = Q_r B_{r-1} Q_r \end{aligned} \quad (3.3.29)$$

The method relies on constructing a Q_r such that $b_r = Q_r b_{r-1}$ has zero elements in positions $r+2, r+3, \dots, n$. This construction was given

previously by equation (3.3.6). Thus, proceeding from $r=1$ to $r=n-1$, we have

$$P_{n-1} P_{n-2} \dots P_1 A P_1 \dots P_{n-2} P_{n-1} = P A P = H \quad (3.3.30)$$

Again, rounding errors occur at each step and the aim of backward error analysis is to derive perturbation bounds, viewing the computed H as the exact transformation of a perturbed A matrix, i.e.,

$$H = P(A + E)P. \quad (3.3.31)$$

Wilkinson [24] does this in great detail for a number of variations of the basic algorithm. The following development, which is valid for any unitary transformation, is both illustrative of the common features of his approach and useful for obtaining a number of results because of the generality. The matrices E , X , Y , and Z will denote perturbation matrices; other matrices are assumed exact unless over-lined with a bar, which symbolizes a computed quantity. Finally, although the analysis is for unitary similarity transformations, we assume that the P_r matrices are orthogonal merely to simplify the notation.

So, $A_r = P_r A_{r-1} P_r$ is the exact transformation at the r^{th} step, and

$$\bar{A}_r = \bar{P}_r \bar{A}_{r-1} \bar{P}_r + E_r \quad (3.3.32)$$

where

$$\bar{P}_r = P_r + X_r.$$

For any of the different algorithms, we can compute a constant a that depends on the specific arithmetic operations used in forming P_r such that $\|X_r\|_2 \leq a2^{-t}$. Now

$$\bar{A}_r = (P_r + X_r)\bar{A}_{r-1}(P_r + X_r) + E_r \quad (3.3.33)$$

$$= P_r \bar{A}_{r-1} P_r + Y_r$$

where
$$Y_r = X_r \bar{A}_{r-1} P_r + P_r \bar{A}_{r-1} X_r + X_r \bar{A}_{r-1} X_r + E_r . \quad (3.3.34)$$

Let $r = 1, 2, \dots, n$

$$\bar{A}_n = G_1 A_0 G_1 + G_2 Y_1 G_2 + \dots + G_n Y_{n-1} G_n + Y_n \quad (3.3.35)$$

where
$$G_r = P_r P_{r-1} \dots P_1 ,$$

or
$$\bar{A}_n = G_1 A_0 G_1 + Y ,$$

or
$$\bar{A}_n = G_1 (A_0 + Z) G_1 \quad (3.3.36)$$

where
$$Z = L_n Y_n L_n + L_{n-1} Y_{n-1} L_{n-1} + \dots + L_1 Y_1 L_1$$

$$L_r = P_1 P_2 \dots P_r .$$

Certainly

$$\|Y\|_2 \leq \|Y_n\|_2 + \|Y_{n-1}\|_2 + \dots + \|Y_1\|_2 ,$$

and using $\|X_r\|_2 \leq a2^{-t}$ as well as the invariance of the spectral norm under unitary similarity transformations, equation (3.3.34) implies

$$\|Y_r\|_2 \leq \|\bar{A}_{r-1}\|_2 (2a2^{-t} + a^2 2^{-2t}) + \|E_r\|_2 \quad (3.3.37)$$

while (3.3.33) yields

$$\|\bar{A}_r\|_2 \leq \|\bar{A}_{r-1}\|_2 + \|Y_r\|_2 \leq (1+a2^{-t})^2 \|\bar{A}_{r-1}\|_2 + \|E_r\|_2 \quad (3.3.38)$$

Considering equation (3.3.32), we can analyze the specific algorithm of interest to find a bound of the form

$$\|E_r\|_2 \leq f(r, n) 2^{-t} \|\bar{A}_{r-1}\|_2$$

where $f(r, n)$ is a simple function of r and n in general, although in some cases it is a constant. Using this, equations (3.3.37, 3.3.38) become

$$\begin{aligned} \|Y_r\|_2 &\leq (2a2^{-t} + a^2 2^{-2t} + f(p, n) 2^{-t}) \|\bar{A}_{r-1}\|_2 \\ \|\bar{A}_r\|_2 &\leq ((1+a2^{-t})^2 + f(p, n) 2^{-t}) \|\bar{A}_{r-1}\|_2 . \end{aligned}$$

Finally, for

$$\begin{aligned} \bar{A}_n &= G_1 (A_0 + Z) G , \\ \|Z\|_2 &\leq 2^{-t} \|A_0\|_2 \sum_{r=1}^n \{ [2a+a^2 2^{-t} + f(r, n)] \prod_{i=1}^r [(1+a2^{-t})^2 + f(i, n) 2^{-t}] \} . \end{aligned}$$

Applying the basic results of Appendix B to the Householder method one can find that $f(i, n)$ should be of the form $f(i, n) = k(n-i)$. Shoving this into the above result yields a complicated series which blows up rapidly as n gets very large, but in the useful range the result is (for a computer without accumulation of inner products)

$$\|Z\|_2 \leq \|A_0\|_2 (k_1 n^2 + k_2 n^4 2^{-t}) 2^{-t} . \quad (3.3.39)$$

If accumulation of inner products is available, replace n^2 with n and n^4 with n^2 .

The third problem for which perturbation bounds are required is that of transforming a matrix in Hessenberg form to real Schur form. Recall that the real Schur form is a generalization of the triangular form wherein 2×2 blocks along the diagonal correspond to complex conjugate pairs of eigenvalues. Of course, this transformation solves the eigenvalue problem, which is of considerable importance and difficulty. As a result, a thorough understanding and analysis of the Bartels and Stewart Lyapunov equation solution method requires a preliminary study of the eigenvalue problem, which would be in itself a substantial thesis topic.

There exist many ways to obtain the Schur form which are variations of two basic methods, the LR and QR algorithms [24]. Both are inexact procedures in which an infinite sequence of similarity transformations are successively applied to a general matrix A ; the former

relies on elementary matrices to do this while the latter employs unitary transformations. Neither algorithm actually requires the Hessenberg form, but the computational advantages of obtaining it for the general A first are great. As was mentioned previously, elementary transformations are almost as numerically stable as unitary ones if they can be stabilized. A stabilized process is one in which the elements of the transformation matrices are strictly bounded. For the LR algorithm numerical stability and speed of convergence are conflicting properties because preserving the Hessenberg form at each step (least number of elements below the diagonal, i.e., fewest operations) eliminates pivoting options and vice versa. The QR algorithm is, therefore, generally "better", and will be used to illustrate the basic reduction. It is defined by the simple recursion

$$A_s = Q_s R_s, \quad A_{s+1} = Q_s^H A_s Q_s = R_s Q_s, \quad (3.3.40)$$

where Q_s is unitary and R_s is upper triangular. Manipulating (3.3.40),

$$A_{s+1} = (Q_s^H Q_{s-1}^H \cdots Q_1^H) A_1 (Q_1 Q_2 \cdots Q_s),$$

or

$$(Q_1 Q_2 \cdots Q_s) A_{s+1} = A_1 (Q_1 Q_2 \cdots Q_s). \quad (3.3.41)$$

Denoting $(Q_1 Q_2 \cdots Q_s) = P_s$ and $(R_s R_{s-1} \cdots R_1) = U_s$, we have

$$\begin{aligned}
P_S U_S &= (Q_1 \dots Q_{S-1}) (Q_S R_S) (R_{S-1} \dots R_1) \\
&= (Q_1 \dots Q_{S-1}) A_S (R_{S-1} \dots R_1) \\
&= A_1 (Q_1 \dots Q_{S-1}) (R_{S-1} \dots R_1) \\
&= A_1 P_{S-1} U_{S-1} .
\end{aligned}$$

Repeating for $P_{S-1} U_{S-1}$, $P_{S-2} U_{S-2}$, etc. yields

$$P_S U_S = A_1^S .$$

Therefore $P_S U_S$ is the corresponding factorization of A_1^S , and in either case Wilkinson [24] shows that the factorization is unique if the diagonal elements of the R_S are taken to be positive. Computationally, each iteration (3.3.40) involves two essential steps. The first step is to construct an orthogonal Q_S such that (assume A real)

$$Q_S^T A_S = R_S , \tag{3.3.42}$$

where R_S is upper triangular. Notice that this step is analogous to the triangular factorization covered previously (equation 3.3.15).

Now it is easy to verify that if A_1 is in Hessenberg form, so is each A_S . Therefore, Q_S^T can be constructed as a product of $(n-1)$ plane rotations in the $(1, 2)$, $(2, 3)$, ...k $(n-1, n)$ planes, i.e.,

$$Q_S^T = R(n-1, n)R(n-2, n-1) \dots R(1, 2) \quad (3.3.43)$$

where the $R(p, q)$ are defined by equation (3.3.3). The second step is to then successively post-multiply R_S with the transposes of $R(p, q)$,

$$R_S (R^T(1, 2)R^T(2, 3) \dots R^T(n-1, n)) = R_S Q_S = A_{S+1} \quad (3.3.44)$$

The A_S tend in the limit to a matrix of the form

$$\begin{bmatrix} x_1 & x & \cdot & \cdot & \cdot & x \\ & x_2 & x & \cdot & \cdot & x \\ & & x_3 & & & \cdot \\ & & & \cdot & & \cdot \\ & & & & \cdot & \cdot \\ & & & & & x_{p-1} & x \\ & & & & & & x_p \end{bmatrix} \quad x_i (n_i \times n_i)$$

The dimension of each X_i is equal to the number of distinct eigenvalues of equal modulus. An interesting property of the convergence is that the eigenvalues of the X_i converge to eigenvalues of A , while elements above the (block) diagonal do not tend to a strict limit, but may change from iteration to iteration by a factor ultimately of modulus unity [22, 24].

The analysis developed previously, equations (3.3.32, 3.3.38), can be applied to the QR algorithm in order to obtain perturbation

bounds. In terms of the notation used there, however, the summation was over columns of A, but is now taken over steps of the QR algorithm. In this case $f(s, n)$ must also be computed with the same procedure, realizing that the successive A_s are in Hessenberg form. Accomplishing this, we find

$$f(s, n) \doteq k_1 n \quad (3.3.45)$$

$$a \doteq k_2 n$$

from which the final perturbation bound is

$$\begin{aligned} \bar{A}_s &= P_s^T (A_0 + Z_s) P_s, \quad A_0 \text{ in Hessenberg form,} \\ \|Z_s\| &\leq \|A_0\| [n \cdot s (K_3 + K_2^2 n 2^{-t}) 2^{-t} + n^2 s^2 (K_3^2 + K_2^2 K_3 n 2^{-t}) 2^{-2t}], \\ & \quad (3.3.46) \\ K_3 &= K_1 + 2K_2. \end{aligned}$$

Note that this bound neglects terms that go to infinity as the product ns gets very large, but is applicable in the useful range of ns .

The final problem of this section is that of transforming a general matrix A to Companion (or Frobenius) form. As usual, several methods may be used. A well-known algorithm is that of Danilewski [31], which uses elementary matrices to affect the transformation. The difficulty with this method can be seen by considering the r^{th} step:

$$\begin{bmatrix} m_{11} & m_{12} & \cdot & \cdot & m_{15} \\ & m_{22} & \cdot & \cdot & \cdot \\ & & m_{33} & \cdot & \cdot \\ & & & \cdot & \cdot \\ & & & & m_{55} \end{bmatrix} \begin{bmatrix} h_{11} & h_{12} & \cdot & \cdot & h_{15} \\ h_{21} & h_{22} & \cdot & \cdot & h_{25} \\ & h_{32} & \cdot & \cdot & h_{35} \\ & & h_{43} & \cdot & h_{45} \\ & & & h_{54} & h_{55} \end{bmatrix} = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \\ p_5 \end{bmatrix} \begin{bmatrix} m_{11} & m_{12} & \cdot & \cdot & m_{15} \\ & m_{22} & \cdot & \cdot & \cdot \\ & & m_{33} & \cdot & \cdot \\ & & & \cdot & \cdot \\ & & & & m_{55} \end{bmatrix}$$

(3.3.48)

The rows of M may be determined successively ($m_{11} = 1$, free choice),
 $r = 1, 2, \dots, n$ with

$$m_{r,j+1} = (m_{k-1,j} - \sum_{k=r}^j m_{r,k} h_{k,j})^{1/h_{j+1,j}} \quad (3.3.49)$$

$$j = r-1, \dots, n-1.$$

Then $p_r = (m_{r-1,n} - \sum_{k=r}^n m_{r,k} h_{k,n})^{1/m_{n,n}}$.

If any of the $h_{j+1,j}$ are zero or very small, the $m_{r,j+1}$ become very large and large rounding errors will result. Some criteria should be used to set those elements to zero and therefore decompose H into a direct sum of block Hessenberg matrices, each of which may then be transformed individually. Unfortunately, perturbation bounds can not be obtained for this second step in the reduction. An interesting idea [24] which could further improve the numerical stability of this step is to modify the first step such that the sub-diagonal elements

are either 0 or 1. Theoretically, this is possible as the modified Hessenberg form must be exactly similar to the corresponding Frobenius canonical form. Note that if no zeros appear on the subdiagonal, then the matrix must be non-derrogatory and the Frobenius form degenerates to the Companion form. The following algorithm, which is a modification of an algorithm developed in [24], illustrates an approach that may solve this problem. Stabilized elementary transformations are used to form a matrix N such that

$$AN = NH \quad (3.3.50)$$

where N is lower triangular except in the first column, which is e_1 (unit vector). If the r^{th} column of the products in (3.3.50) are equated, we have

$$\left. \begin{aligned} \sum_{k=r}^n a_{ij} n_{kr} &= \sum_{k=2}^i n_{ik} h_{kr} & i=2, 3, \dots, r \\ &= \sum_{k=2}^r n_{ik} h_{kr} + n_{i,r+1} & i=r+1, \dots, n \end{aligned} \right\} r=2, 3, \dots, n$$

and at the r^{th} step the r^{th} column of H and the $(r+1)^{\text{th}}$ column of N can be obtained.

Consider the situation at the r^{th} step for $n=5$ and $r=3$, where critical elements are represented as they might be stored in the computer,

$$\begin{bmatrix} h_{11} & h_{12} & a_{13} & a_{14} & a_{15} \\ n_{22} & h_{22} & a_{23} & a_{24} & a_{25} \\ n_{32} & n_{33} & a_{33} & a_{34} & a_{35} \\ n_{42} & n_{43} & a_{43} & a_{44} & a_{45} \\ n_{52} & n_{53} & a_{53} & a_{54} & a_{55} \end{bmatrix}$$

↑
rth step

i) determine rth column of H

$$\left. \begin{aligned} h_{ir} &= \left(\sum_{k=r}^n a_{ik} n_{kr} \right)^{1/n_{ii}} \\ &\text{store in } a_{ir} \end{aligned} \right\} i=2, 3, \dots, r$$

ii) compute

(3.3.51)

$$\left. \begin{aligned} n_{i,r+1} &= \sum_{k=r}^n a_{ik} n_{kr} - \sum_{k=2}^r n_{ik} h_{kr} \\ &\text{store in } a_{ir} \end{aligned} \right\} i=r+1, \dots, n-1$$

iii) let (r+1)' be the first integer such that

$$|n_{i(r+1)',}| = \max_{i>r} |n_{i,r+1}|$$

exchange rows and columns (r+1) and (r+1)'

- iv) if $|n_{r+1',r+1'}| \leq \epsilon$ (some tolerance), then take the r^{th} subdiagonal of H to be zero and set $n_{(r+1)',(r+1)'}$ (which is now in position $(r+1, r)$) to 1.

Assuming that the Hessenberg form is constructed with subdiagonal elements zero or unity, then the Frobenius form can be obtained as a direct sum of Companion forms. Each companion form can then be obtained with a slight modification of equations (3.3.48, 3.3.49), where now both the diagonal elements of M and the subdiagonal elements of each reduced order Hessenberg matrix are unity. Numerically, the modification is significant because no divisions are required for the second step. An error analysis of this modified procedure yields a perturbation bound of the form

$$C = M(H + E)M^{-1},$$

$$|E| \leq 2g2^{-t} \begin{bmatrix} 1 & 2 & \cdot & \cdot & \cdot & \cdot & n-1 & n \\ 0 & 1 & 2 & \cdot & \cdot & \cdot & n-2 & n-1 \\ & 0 & & \cdot & & & \cdot & \cdot \\ & & & \cdot & & & \cdot & \cdot \\ & & & & \cdot & & \cdot & \cdot \\ & & & & & \cdot & 0 & 1 \end{bmatrix} \quad (3.3.52)$$

where $g = \max_{r,i,j} |h_{i,j}^r|$.

Notice that this result is similar to that for LU decomposition except for the factor g. This factor complicates the a priori analysis

and reflects the fact that, at one point or another, reduction of a general matrix to companion form requires a series of unstabilized elementary transformations. In order to obtain this bound, it was assumed that the modifications suggested above can be successfully carried out. If the additional assumption is made that the perturbation bound of the reduction of A to Hessenberg form using stabilized elementary transformations is approximately the same as the bound previously obtained for unitary transformations, then the bounds for each step of the two step reduction of A to companion form can be combined. From equation (3.3.39)

$$H = G(A + Z)G,$$

$$\|Z\|_2 \leq \|A\|_2 (K_1 n^2 + k_2 n^4 2^{-t}) 2^{-t},$$

and from equation (3.3.52)

$$C = M(H + E)M^{-1},$$

$$\|E\|_1 = \|E\|_\infty \leq gn^2 2^{-t}.$$

Using the fact that the spectral norm is invariant under unitary similarity transformations and that for any matrix B

$$\|B\|_2^2 \leq \|B^H B\|_\infty \leq \|B^H\|_\infty \|B\|_\infty = \|B\|_1 \|B\|_\infty,$$

then

$$C = MG(A + X)GM^{-1} = D(A + X)D^{-1} ,$$

$$\|x\|_2 \leq (\|A\|_2^{K_1} n^2 + gn^2 + \|A\|_2^{K_2} n^4 2^{-t}) 2^{-t} , \quad (3.3.53)$$

where

$$g = \max_{r,i,j} |h_{ij}^r| .$$

3.4 SUMMARY OF RESULTS

The purpose of this section is to apply the main results of the analyses of the previous sections to the following Lyapunov equation solution techniques: direct method, transformation methods based on the Schur and Companion canonical forms, and the iterative decoupling method. It is intended to be essentially self-contained, but it is important to realize that error analysis does not yield exact results, and that a number of assumptions and qualifying remarks made previously are not repeated here.

Recall that the technique of backward error analysis assumes that the computed solution exactly satisfies the perturbed equation

$$(A^T + E^T)(P + \delta P) + (P + \delta P)(A + E) = -Q \quad (3.4.1)$$

where the matrix E represents the affects of round-off errors that occur during each step of the particular algorithm of interest. The matrix δP represents the resulting solution error, and its relative size is given by equation (3.2.16).

Now from that equation, we see that unless

$$\|K_A\| \|K_A^{-1}\| \|K_E\| / \|K_A\| \ll 1$$

the solution will be very inaccurate, so assuming that this is true, equation (3.2.16) becomes

$$\|\delta P\| / \|P\| \leq k(L_A) \|K_E\| / \|K_A\| , \quad (3.4.2)$$

where $K_A = A^T \times I + I \times A^T$ and $\|K_A\|_2 \|K_A^{-1}\|_2 \equiv k(L_A)$, the spectral condition number of the Kronecker matrix. Technically, equation (3.4.2) is valid only for the 2-norm, due to the definition of the condition number, but we will use whatever norm is most convenient. This is reasonable because it is the ratio on the right hand side of (3.4.2) that is of primary interest here. In addition, some solution methods do not involve the Kronecker matrix, so the relationships

$$\begin{aligned} n^{1/2} \|x\|_E &\leq \|K_X\|_E \leq 2n^{1/2} \|\bar{x}\|_E , \\ \|x\|_2 &\leq \|K_X\|_2 \leq 2\|x\|_2 \end{aligned} \quad (3.4.3)$$

will be used. Notice that in this case, (3.4.2) becomes

$$\|\delta P\| / \|P\| \leq 2k(L_A) \|E\| / \|A\| . \quad (3.4.4)$$

Finally, the results will be developed assuming that the computer used does not accumulate inner products (e.g., IBM-370), although the appropriate modifications are given if accumulation is available. (See Appendix B for explanation of inner product accumulation.)

Direct Method

First, re-write equation (3.4.1) using the Kronecker notation,

$$(K_A + K_E)(p + \delta p) = -q.$$

Now, from equation (3.2.22), a bound for K_E that accounts for errors made in the LU factorization and in solving the resulting sets of triangular equations is

$$\|K_E\|_E \leq g(n^2/2)^{3/2} 2^{-t} \leq \frac{g}{2} n^3 2^{-t}. \quad (3.4.5)$$

Using (3.4.2, 3.4.3), we have

$$\|\delta p\|_E / \|p\|_E \leq k(L_A) \left(\frac{g}{2} n^{5/2} 2^{-t} \right) / \|A\|_E, \quad (3.4.6)$$

where $g = \max | (K_A)_{ij}^r | \leq 8 \max | (K_A)_{ij}^0 | \leq 16 \max | a_{ij} |$. The first inequality reflects the modest growth of the elements of the matrix being factored when pivoting is used (which is essential) and the second follows easily from the formulation of the Kronecker matrix. If accumulation is used, replace $n^{5/2}$ by $n^{3/2}$. It is important to

realize that this bound actually includes a statistical factor suggested by Wilkinson [18] and explained more fully in the previous section.

Equation (3.4.6) is the main result, but an interesting extension, although not entirely consistent with the bounding approach used here, facilitates a comparison with the Bartels and Stewart method analysis that will follow. The basic idea is that if an assumption is made on a statistical distribution of the magnitudes of the elements of A, then the factor g can be related to the euclidian norm and eliminated in equation (3.4.6). To illustrate, let $g = 16 \max |a_{ij}|$ and each element of A be such that $|a_{ij}| = x_{ij} \frac{g}{16}$, where the x_{ij} are independent random variables equally distributed on (0, 1). This formulation is technically incorrect (i.e., g is a random variable now, etc.) but informally at least the expected value of the euclidian norm of A is

$$E[|A|_E] = ng/16\sqrt{3}$$

and using this in equation (3.4.6) we have

$$|\delta P|_E / |P|_E \leq k(L_A) 8\sqrt{3} n^{3/2} 2^{-t} . \quad (3.4.7)$$

Bartels and Stewart Method

Letting R denote the product of orthogonal matrices that affected the transformation of A to Hessenberg to real Schur form, the exact Lyapunov equation becomes

$$(RA^T R^T) (RPR^T) + (RPR^T) (RAR^T) = -RQR^T \quad (3.4.8)$$

or
$$A_S^T Y + YA_S = -C$$

collecting the lower order terms of equations (3.3.39, 3.3.46), we have that the computed Schur matrix \bar{A}_s is exactly similar to the perturbed equation

$$\begin{aligned} \bar{A}_s &= R(A + Z_s)R^T, \\ \|Z_s\|_2 &\leq \|A\|_2 (k_1 n^2 + k_2 ns) a^{-t} \end{aligned} \quad (3.4.9)$$

where s is the number of iterations of the QR algorithm. Comparing equations (3.2.11, 3.2.13), we see that as far as relative perturbations in the solution are concerned, the errors that occur in forming RQR^T can be added to those that result from transforming A . In addition, assuming that the transformations are applied to Q at each step of the algorithm, the resulting errors are similar to those of (3.4.9), i.e.,

$$\begin{aligned} \bar{C} &= R(Q + \delta Q)R^T, \\ \|\delta Q\|_2 &\leq \|Q\|_2 (k_1 n^2 + k_2 ns) 2^{-t}. \end{aligned} \quad (3.4.10)$$

The next source of error occurs in solving for Y . This step is essentially that of solving a block triangular set of equations,

where each block is at most 2×2 . Consider ignoring these errors temporarily. Then, combining equations (3.4.4, 3.4.9, 3.4.10)

$$\|\delta Y\|_2 / \|Y\|_2 \leq 4k(L_A) (k_1 n^2 + k_2 ns) 2^{-t} \quad (3.4.11)$$

Now in the previous chapter we found that the solution of a triangular set of equations produces a low relative error that in practice does not depend on the condition number of the triangular matrix, e.g., from equation (3.3.20) the solution of $Lx = b$ yields the very satisfactory bound $\|\delta X\| \leq Kn\|X\|$. Because the 2×2 blocks can be solved explicitly, it seems reasonable for this analysis to assume that we effectively obtain Y from the solution of a sparse set of triangular equations of dimension $n^2/2$ with a relative error bound (ignoring sparsity) on the order of $\|\delta Y\| \leq \frac{k}{2} n^2 \|Y\|$. For even moderate $k(L_A)$ in (3.4.11), this term is relatively unimportant.

The final step in the algorithm is to compute

$$\bar{P} = \bar{R}^T (Y + \delta Y) \bar{R} = P + \delta P,$$

where the first order perturbation terms in δP are

$$\delta P = \delta R^T Y R + R^T Y \delta R + R^T \delta Y R. \quad (3.4.12)$$

In taking the norm of equation (3.4.12), we find that the contribution due to the first two terms is small compared to the last one, so the

final result is essentially that of (3.4.11), i.e.,

$$\|\delta Y\|_2 / \|Y\|_2 \leq 4\kappa(L_A) (k_1 n^2 + k_2 ns) 2^{-t} \quad (3.4.13)$$

A final comment on this result is that we expect s , the number of QR iterations, to be some linear function of n , and that accumulation of inner products reduces the $k_1 n^2$ term to $k_3 n$ if Householder's method is used to affect the transformation to Hessenberg form.

Companion Form Methods

Unfortunately, a complete error analysis of Lyapunov equation solution methods that rely on the Companion form cannot be obtained. There are several reasons for this. First, the various comments made in the previous section apply here, and they all basically reflect the fact that at some point in transforming A to Companion form unstabilized elementary transformations must be used. Molinari's algorithm was criticized in particular for ignoring this difficulty entirely. Several improvements were suggested, however, and they led to the intermediate result (3.3.53)

$$\begin{aligned} \bar{C} &= D(A + X)D^{-1} , \\ \|X\|_2 &\leq (k_1 n^2 \|A\|_2 + gn^2) 2^{-t} , \\ g &= \max |h_{ij}^r| , \end{aligned} \quad (3.4.14)$$

where H is the Hessenberg matrix and r denotes the r^{th} step in transforming H to C. Now the next step in this solution method is to form the Hurwitz matrix W, which is constructed from the computed elements of C, i.e., the coefficients of the characteristic polynomial of A. W is then factored, from which the last row of $Y = D^{-T}PD^{-1}$ can be obtained. The difficulty is that the method of backward analysis tells us that \bar{C} is exactly similar to a perturbed A matrix (i.e., (3.4.14)), but this gives no indication of the relative accuracy of the elements of \bar{C} . The perturbations in these elements affect not only the analysis of the factorization of W and solution of $Y_r(n)$ (last row), but that of the recursion used to obtain the remaining rows $Y_r(j)$, $j=n-1, n-2, \dots$, as well.

Although several analyses of these steps were performed using a combination of forward and backward techniques, the results are not very consistent and are therefore not reported here. The subjective conclusion reached, however, is that this Lyapunov equation solution method is probably less accurate than the two analyzed previously.

Iterative Decoupling Algorithm

For this analysis, the notation $Ax = b$ will be used to illustrate the approach because it is simpler. The result is then easily extended to the Lyapunov equation.

Again, using the backward approach, the k^{th} iteration can be viewed as

$$(A_0 + E)(x^k + \delta x^k) = b - A_1(x^{k-1} + \delta x^{k-1}) \quad (3.4.15)$$

Consider the first step, i.e.,

$$(A_0 + E)(x' + \delta x') = b,$$

or
$$\delta x' = -(I + A_0^{-1}E)^{-1}A_0^{-1}Ex'.$$

From equation (3.2.13)

$$\|\delta x'\| \leq k(A_0) \frac{\|E\|}{\|A_0\|} \|x'\|$$

where we assume that $k(A_0) \frac{\|E\|}{\|A_0\|} \ll 1$. (Note that we can effectively force this by using iterative refinement if necessary.) Now let $r = k(A_0) \frac{\|E\|}{\|A_0\|}$, $s = \|A_0^{-1}A_1\|$, and consider the second step,

$$(A_0 + E)(x^2 + \delta x^2) = b - A_1(x' + \delta x'),$$

or
$$\delta x^2 = -(I + A_0^{-1}E)^{-1}A_0^{-1}Ex^2 - (I + A_0^{-1}E)^{-1}A_0^{-1}A_1\delta x'.$$

and
$$\|\delta x^2\| \leq r\|x^2\| + sr\|x'\|.$$

Continuing in this fashion, we find

$$\|\delta x^k\| \leq r\|x^k\| + sr\|x^{k-1}\| + \dots + s^{k-1}r\|x'\|$$

or
$$\|\delta x^k\| \leq r\|x^*\| (1 + s + s^2 + \dots + s^{k-1}) = r\|x^*\| \frac{1-s^k}{1-s},$$

$$s < 1, \text{ where } \|x^*\| = \max_k \|x^k\|.$$

The desired result is

$$\|\delta x^k\| / \|x^*\| \leq \frac{k(A_0) \|E\| / \|A_0\|}{(1 - \|A_0^{-1} A_1\|)} \quad (3.4.16)$$

Now the condition $s = \|A_0^{-1} A_1\| < 1$ is sufficient for the convergence of the algorithm, so expressing the series $1 + s + s^2 + \dots + s^{k-1}$ in the above form is reasonable, but the bound (3.4.16) may be very pessimistic if s is near 1. Equation (3.4.16) is valid for the 2-norm, and the matrices E and A_0 are assumed to be block diagonal, i.e.,

$$E = \text{diag}(E_i), \quad A_0 = \text{diag}(A_i) \quad i=1, 2, \dots, N$$

$$A_i \quad n_i \times n_i$$

so

$$\|E\|_2 = \max_i \|E_i\|_2, \quad \|A_0\|_2 = \max_i \|A_i\|_2 \equiv \|A_i^*\|_2$$

and

$$k(A_0) \leq \frac{\max_i \|A_i\|_2}{\min_j \|A_j^{-1}\|_2} \quad (3.4.17)$$

Extending this to the Lyapunov equation, we see that the bound will be similar to that obtained for the direct method previously, except that the numerator in (3.4.16) depends only on the reduced order equations. Suppose that equality is obtained in (3.4.17), so that the i^{th} subsystem ($i=j$) has the largest condition number, $k(L_{A_i})$. Let $n_m = \max_i n_i$. From (3.4.5)

$$\|K_E\|_2 \leq \|K_E\|_E \leq n_m^{1/2} \|K_E\|_2 \leq \frac{g}{2} n_m^{7/2} 2^{-t}$$

so combining (3.4.6, 3.4.16), the final result is

$$\|\delta P\|_2 / \|P\|_2 \leq \frac{k(L_{A_i}) (\frac{g}{2} n_m^3 2^{-t}) / \|A_i^*\|_2}{(1 - \|A_0^{-1} A_1\|)} \quad (3.4.18)$$

Although the notation is somewhat cumbersome, the essential point here is that for the iterative decoupling algorithm, it is the errors made in solving the reduced order equations, along with the contraction condition, that limit the accuracy of the final solution.

CHAPTER IV
APPLICATION

4.1 ITERATIVE DECOUPLING PROGRAM

A Fortran computer program that realizes the iterative decoupling algorithm has been written. A general purpose subroutine for solving the Lyapunov equation by the direct method is used to solve the reduced order equations. Sparsity coding techniques have been employed; in particular, a commercial sparse matrix package is used to perform optimal ordering of the Kronecker matrices to make the LU factorizations as sparse as possible, within certain constraints on the relative magnitude of pivotal elements that are important for numerical stability. Although $N(N+1)/2$ Kronecker matrices must be constructed and factored, the LU factorization is first done symbolically and then numerically. This is useful when some of the diagonal blocks of A share a similar structure because the symbolic factorization need not be repeated for those Kronecker matrices that involve the similar, but not necessarily identical subsystems.

The algorithm has been tested initially with relatively small system matrices, so that the centralized solution can be computed and iteratively refined to a specified accuracy to provide a reliable check. The same accuracy tolerance is used to terminate the decoupling iteration. Specifically, the iteration terminates when each element of the diagonal blocks of P^k , i.e., P_{ii}^k , $i=1, \dots, N$, changes in value from that of the previous step by less than TOL. Note that this is an accuracy tolerance.

In the following example $TOL = 10^{-4}$ and the solution is the steady state covariance of the system $\dot{x}(t) = Ax(t) + \omega(t)$, $E\{\omega(t)\omega'(J)\} = I\delta(t-J)$.

$$A = \left[\begin{array}{ccc|ccc|cc} -3 & & & & & \alpha_1 & & & & & \\ & -2 & -2 & & & & & & & & \\ 1 & 2 & -2 & & & & & & & & \\ \hline & & & -4 & & 1 & & & & & -\alpha_2 \\ \alpha_1 & & & & -3 & -2 & & & & & \\ \hline & & & & 2 & -3 & & & & & \\ & & & & & & & & & & \\ & & & -\alpha_1 & & & & & -2 & & \\ & & & & & & & & & -1 & -0.5 \\ & & & & & & & & & 0.5 & -1 \end{array} \right]$$

The coupling elements α_1, α_2 were varied as $\alpha_1 = \frac{1}{2}k, \alpha_2 = k$ for $k = 0, 1, 2, \dots, 7$. As k increases, the elements in the diagonal blocks of P move increasingly away from the initial decentralized solution ($k=0$), and the number of iterations required for convergence naturally grows accordingly. This range of coupling elements was sufficient to vary some of the solution values by two order of magnitude. The number of iterations varied from 2 for $k=1$ to 9 for $k=7$. The first solution, $k=1$, had an execution time of approximately .6 seconds, most of which is forming the Kronecker matrices and factoring them. Each additional solution executed in less than .06 second. The computer used was an IBM-370/168.

4.2 POWER SYSTEM APPLICATION

In this section, another numerical experiment using the iterative decoupling algorithm is reported. Unfortunately the algorithm failed to converge for the system matrices used, so the purpose of this section is to briefly describe what was attempted and why it did not work.

A particular problem of current interest in the study of power system dynamic behavior is that of obtaining transient stability equivalents. In one approach to this problem [42], an important preliminary step is to identify coherent groups of generators, i.e., machines that tend to swing together under the influence of a severe network disturbance. A reliable, but time-consuming, method of making this identification is to simply run a number of large transient stability programs, and visually compare plots of the rotor angle responses of all the generators of interest. Now in some cases simple, linearized machine models may be sufficient for the purpose of identifying coherent groups, and in this case, the solution of a Lyapunov equation provides valuable information. For example, let x_i and x_j be the rotor angles of machines i and j , and suppose the matrix Q is null except for $q_{ij}=q_{ji}=1$. Then

$$M = \int_0^{\infty} x^T Q x dt = \text{tr}(X_0 P) , \quad (4.2.1)$$

where P is the solution of the Lyapunov equation, provides a measure of the coherency of machines i and j .

In this experiment, a simple three machine-infinite bus system was used, where each machine was represented by a constant voltage behind transient reactance, i.e., two state swing equation model. When linearized, these equations are of the form:

$$\begin{aligned}
 M_1 \Delta \dot{w}_1 + \frac{\Delta w_1}{R_1 w_0} &= -Y_{11} \delta_1 + Y_{12} \delta_2 + Y_{13} \delta_3 \\
 \dot{\delta}_1 &= \Delta w_1 \\
 M_2 \Delta \dot{w}_2 + \frac{\Delta w_2}{R_2 w_0} &= Y_{12} \delta_1 - Y_{22} \delta_2 + Y_{23} \delta_3 \quad (4.2.2) \\
 \dot{\delta}_2 &= \Delta w_2 \\
 M_3 \Delta \dot{w}_3 + \frac{\Delta w_3}{R_3 w_0} &= Y_{13} \delta_1 + Y_{23} \delta_2 - Y_{33} \delta_3
 \end{aligned}$$

where δ_i = perturbation of machine i's rotor angle from operating point

R_i = droop of machine i

Y_{ij} = transfer admittance between machines i and j

Y_{ii} = self admittance of machine i

Lee [41] studied this same system, and the per unit values used here were taken from his work. He did not include the damping term, but this is necessary in order for the Lyapunov equation solution to exist. A typical set of values used was:

$$\begin{aligned}
M_1 &= .1326 & Y_{11} &= 2.20 & Y_{12} &= 1.0 \\
M_2 &= .1592 & Y_{22} &= 2.60 & Y_{13} &= .90 \\
M_3 &= .1194 & Y_{33} &= 2.30 & Y_{23} &= 1.2 \\
R_1 &= R_2 = R_3 & &= .01
\end{aligned}
\tag{4.2.3}$$

With parameters on the order of (4.2.3), the iterative decoupling algorithm did not converge, but slowly moved away from the initial decentralized solution. The reason is that the necessary condition for convergence of the algorithm, i.e., $\rho(L_{A_0}^{-1}L_{A_1}) < 1$, is not satisfied for these typical values. It is interesting to note, however, that $\rho(A_0^{-1}A_1)$ is less than one. In order to see if any simple normalization of the elements of A might help, $L_{A_0}^{-1}L_{A_1}$ was computed symbolically. To illustrate the difficulty, consider that the row norm of this product is of the form:

$$\|L_{A_0}^{-1}L_{A_1}\|_{\infty} = \sum_{ij} Y_{ij}/M_i (1 + w_0 R_i M_i + 1/w_0 R_i Y_{ii}) .$$

We can see that normalizing A is useless and, although only a sufficient condition for convergence, that the values of (4.2.3) must be drastically changed to make this quantity less than 1.

4.3 CONCLUSIONS

The Lyapunov equation is both theoretically interesting and practically useful. Although commonly associated with stability theory, the various physical interpretations of its solution and relationship to the evaluation of quadratic integrals make it a basic tool in a number of areas of control theory. Many different methods can be used to solve it and a number of these were discussed in Chapter Two. The iterative decoupling algorithm developed in that chapter is basically an original contribution of this thesis, although the idea upon which it is based is not new. It is a special purpose solution method with several desirable properties that requires more development in order to assess its real potential. A suggestion for future work here is to extend the algorithm to the over-relaxation scheme, as this would add flexibility to the method.

The error analyses of Chapter Three are based heavily on the works of others, primarily Wilkinson, although several of the results are original in their specific application and extension to the Lyapunov equation. One important conclusion here is that the bounding approach is primarily useful for comparing different algorithms and the results obtained should not be interpreted too literally. This is especially true regarding a priori accuracy estimates. There is little doubt that theoretical analysis and numerical experience are both necessary in order to perform useful error analyses. For this reason, an obvious

suggestion for future research would be a systematic, well-organized set of numerical experiments designed to correlate and refine some of the bounds obtained here.

APPENDIX A

The purpose of this appendix is to summarize a few properties of, and relationships between, the canonical forms used in the main text, with special emphasis on the existence conditions of the Companion canonical form. This material is standard and may be found in most linear algebra texts, and for this reason, the account is brief and factual.

Let the matrix A be an element of $C^{n \times n}$. A scalar $\lambda \in C$ is called an eigenvalue of A if there exists a non-zero vector $x \in C^n$ such that $Ax = \lambda x$, and the vector x is called an eigenvector of A associated with the eigenvalue λ . The eigenvalues are the roots of the characteristic equation of A , which is a polynomial of degree n given by $\det(\lambda I - A) = 0$.

If the eigenvalues of A are distinct, then the n eigenvectors of A are linearly independent and form a basis for C^n . In this case, the matrix P , whose columns are the eigenvectors, induces a similarity transformation on C^n such that $P^{-1}AP = \text{diag}(\lambda_i)$.

Suppose that A has r distinct eigenvalues, $\lambda_1, \lambda_2, \dots, \lambda_r$ with multiplicities m_1, m_2, \dots, m_r , where $\sum_{k=1}^r m_k = n$. Then the generalization of the diagonal form is the Jordan canonical form, i.e., there exists a matrix H such that $H^{-1}AH = J$, where J is in Jordan form and is the direct sum of p ($p \geq r$) Jordan sub-blocks $J_k(\lambda_i)$. As an example, for $n=6$, $p=4$, $m_1=4$, $m_2=1$, $m_3=1$:

$$J = \begin{bmatrix} J_2(\lambda_1) & & & \\ & J_2(\lambda_1) & & \\ & & J_1(\lambda_2) & \\ & & & \lambda_1(\lambda_3) \end{bmatrix} \quad (A.1)$$

The general form of a Jordan sub-block $J_k(\lambda_i)$ is

$$\begin{bmatrix} \lambda_i & 1 & & & & \\ & \lambda_i & 1 & & & \\ & & \cdot & \cdot & & \\ & & & \cdot & \cdot & \\ & & & & \cdot & 1 \\ & & & & & \lambda_i \end{bmatrix} \quad (k \times k \text{ matrix})$$

Now, the number of independent eigenvectors of A is equal to the number of Jordan sub-blocks. For the above example, the eigenvalues of J are $e_1, e_3, e_5,$ and e_6 , while those of A are $He_1, He_3, He_5,$ and He_6 . Note that, apart from the ordering of the sub-blocks, the transformation is unique.

The elementary divisors of A are the p polynomials $\det(J_k(\lambda_i - \lambda))$, and the minimal polynomial of A is the product of those elementary divisors that correspond to the Jordan blocks of largest dimension of each distinct eigenvalue, i.e., for the above example the minimal polynomial is

$$m(\lambda) = \det(J_2(\lambda_1 - \lambda)) \det(J_1(\lambda_2 - \lambda)) \det(J_1(\lambda_3 - \lambda)).$$

Now a matrix with distinct eigenvalues must necessarily have linear elementary divisors, while one with one or more non-linear elementary divisors is called defective. If there is more than one Jordan sub-block for any distinct eigenvalue, then the matrix is called derogatory, and in this case the degree of the minimal polynomial is less than n .

If the characteristic polynomial of A is

$$\det(\lambda I - A) = \lambda^n + a_n \lambda^{n-1} + a_{n-1} \lambda^{n-2} + \dots + a_2 \lambda + a_1,$$

then the Companion canonical form C_A is

$$\begin{bmatrix} 0 & 1 & & & & & \\ & \cdot & \cdot & & & & \\ & & \cdot & \cdot & & & \\ & & & \cdot & \cdot & & \\ & & & & \cdot & & \\ & & & & & 1 & \\ & & & & & \cdot & \\ & & & & & 0 & 1 \\ -a_n & -a_{n-1} & \cdot & \cdot & \cdot & -a_2 & -a_1 \end{bmatrix} \quad (\text{A.2})$$

A matrix A is similar to a matrix in Companion form only if it is non-derogatory. Such a matrix is also said to be cyclic of period n . For those interested in control theory, an equivalent statement is that the matrix A is non-derogatory (and hence, a similar companion

form exists) iff there exists a vector b such that the pair (A, b) is completely controllable.

The generalization of the Companion form for derogatory matrices is the Frobenius (or Rational) canonical form, which is the direct sum of m sub-blocks of dimension n_i , $i=1, 2, \dots, m$, and each sub-block is of the form of (A.2). For the example (A.1) $m=2$, $n_1=4$, and $n_2=2$. Any matrix A is similar to a matrix in Frobenius form, and in this case, m is the smallest integer such that there exists a B ($n \times m$) such that the pair (A, B) is completely controllable.

APPENDIX B

Round-off errors in algebraic computations occur because real numbers must be represented with a finite number of digits or bits. Matrix computations, however complicated, are performed by a series of elementary algebraic operations. In this appendix, some basic results for the fundamental arithmetic operations as performed on a digital computer are given [18]. Only the case of floating point arithmetic is considered.

In floating point, the real number x is represented as

$$x = 2^b(a); \quad b \text{ integer, } -1 \leq a \leq -1/2 \text{ or } 1/2 \leq a \leq 1.$$

Consider the addition of two scalars, x_1 and x_2 . Define:

$fl(x_1 + x_2)$ computed quantity

$x_1 + x_2$ exact quantity

t number of digits assigned to mantissa.

In the bounding approach of error analysis [18], it is assumed that the rounding errors of the elementary operations are such that

$$fl(x_1 + x_2) \equiv (x_1 + x_2)(1 + e)$$

$$fl(x_1 x_2) \equiv x_1 x_2 (1 + e) \quad |e| \leq 2^{-t}$$

$$fl(x_1/x_2) \equiv x_1/x_2 (1 + e)$$

Using these assumptions, similar results can be obtained which will be useful in later sections. In order to illustrate the procedure, consider the computation

$$s_n = fl(x_1 + x_2 + \dots + x_n)$$

let $s_1 = fl(x_1)$

$$s_r = fl(s_{r-1} + x_r) = (s_{r-1} + x_r)(1 + e) \quad r=2, 3, \dots, n$$

then $s_n = fl(x_1 + x_2 + \dots + x_n) = x_1(1+e) + x_2(1+e) + \dots + x_n(1+e_n)$

where $(1 - 2^{-t})^{n-r+1} \leq (1 + e_r) \leq (1 + 2^{-t})^{n-r+1}$.

Now, a bound of the form $(1 - 2^{-t})^r \leq (1 + e) \leq (1 + 2^{-t})^r$ arises frequently and is somewhat inconvenient. With the very reasonable assumption that $r2^{-t} < 0.1$, it can be replaced with $|e| \leq r2^{-t}1$, where $2^{-t}1 = 1.06 2^{-t}$. So, for the above result, we have $|e_r| \leq (n-r+1)2^{-t}1$.

Notice that the bound depends on the order of summation; the best procedure is to sum the smallest terms first.

In a similar manner we may obtain the following

$$\begin{aligned}
 p_n &= fl(x_1 x_2 \dots x_n) = x_1 x_2 \dots x_n (1+e_2)(1+e_3) \dots (1+e_n) \\
 &= x_1 x_2 \dots x_n (1+E) \quad |E| < (n-1)2^{-t_1}
 \end{aligned}$$

$$\begin{aligned}
 s_n &= fl(x_1 y_1 + x_2 y_2 + \dots + x_n y_n) = x_1 y_1 (1+e_1) + x_2 y_2 (1+e_2) + \dots + x_n y_n (1+e_n) \\
 & \quad |e_r| < (n-r+2)2^{-t_1}
 \end{aligned}$$

The results for the extended sum and inner product assumed that the machine does not accumulate with a $2t$ -digit mantissa. As far as round-off errors are concerned, accumulation is definitely an advantage. In a machine with this feature, intermediate results in a series of elementary operations are not rounded to t -digits, i.e., the working registers that contain the intermediate results carry a $2t$ -digit mantissa. For a machine that accumulates the operations are denoted $fl_2(\cdot)$, and comparing the following bounds with those given previously illustrates the significance of accumulation. (Note that higher level languages, e.g., Fortran, on the IBM-370, do not have this capability.)

$$\begin{aligned}
 fl_2(x_1 + x_2 + \dots + x_n) &= [x_1(1+e_1) + x_2(1+e_2) + \dots + x_n(1+e_n)](1+e) \\
 |e| &\leq 2^{-t} \quad |e_r| < \frac{3}{2} (n+1-r) a^{-2t_2}
 \end{aligned}$$

where $2^{-2t_2} = 1.06 2^{-2t}$

$$fl_2(x_1 y_1 + x_2 y_2 + \dots + x_n y_n) = [x_1 y_1 (1+e_1) + x_2 y_2 (1+e_2) + \dots + x_n y_n (1+e_n)] (1+e)$$

$$|e| \leq 2^{-t} \quad |e_r| < \frac{3}{2} (n+2-r) 2^{-2t_2}$$

Some other useful results:

$$B = fl(kA) \quad \|B - kA\|_E \leq |k| 2^{-t} \|A\|_E \quad k\text{-scaler}$$

$$y = fl(Ax) = Ax + e \quad \|e\|_E \leq 2^{-t} n \|A\|_E \|x\|_2$$

$$C = fl(AB) = AB + E \quad \|E\|_E \leq 2^{-t} n \|A\|_E \|B\|_E$$

$$C = fl_2(AB) = AB + E \quad \|E\|_E \leq 2^{-t} \|AB\|_E + \frac{3}{2} n \|A\|_E \|B\|_E$$

REFERENCES

1. P. Hagander, "Numerical Solution of $A^T S + SA + Q = 0$," Information Sciences, 4, 1972.
2. R.H. Bartels and G.W. Stewart, "Solution of the Matrix Equation $AX + XB = C$," Communications of the ACM, Vol. 15, No. 9, September, 1972.
3. P.C. Muller, "Solution of the Matrix Equations $AX + XB = -Q$ and $S^T X + XS = -Q$," SIAMJ Appl. Math., Vol. 18, No. 3, May, 1970.
4. B.P. Molinari, "Algebraic solution of Matrix Linear Equations in Control Theory," Proc. IEEE, Vol. 116, No. 10, October, 1969.
5. E.J. Davison and F.T. Man, "The Numerical Solution of $A'Q + QA = -C$," IEEE Trans. AC-13, August, 1968.
6. P.G. Smith, "Numerical Solution of the Matrix Equation $AX + XA^T + B = 0$," IEEE Trans. AC-16, June, 1971.
7. C.F. Chen and L.S. Shieh, "A Note on Expanding $PA + A^T P = -Q$," IEEE Trans. AC-13, February, 1968.
8. S.P. Bingulac, "An Alternative Approach to Expanding $PA + A^T P = -Q$," IEEE Trans. AC-15, February, 1970.
9. D.L. Kleinman, "On An Iterative Technique for Riccati Equation Computations," IEEE Trans. AC-13, February, 1968.
10. N.R. Sandell, Jr., P. Varaiya, and M. Athans, "A Survey of Decentralized Methods for Large Scale Systems," to be published.
11. N.R. Sandell, Jr., "On Newton's Method for Riccati Equation Solution," IEEE Trans. AC-19, June, 1974.
12. D. Looze, Decentralized Control of a Freeway Traffic Corridor, M.S. Thesis, M.I.T., August, 1975.
13. P.V. Kokotovic, "Feedback Design of Large Linear Systems," in Feedback Systems, edited by J.B. Cruz, Jr., McGraw-Hill, N.Y., 1972.
14. N.R. Sandell, Jr., and M. Athans, "A Think Piece for Three-Year Research Plan for the Development of Decentralized Strategies for the Control of Interconnected Power Systems," Electronic Systems Laboratory Report, M.I.T., January, 1975.

15. M. Athans, "The Role and Use of the Stochastic Linear - Quadratic - Gaussian Problem in Control System Design," IEEE Trans. AC-16, December, 1971.
16. R.A. Smith, "Matrix Calculations for Lyapunov Quadratic Forms," J. Differential Equations, 2, April, 1966, pp. 208-17.
17. Rall, Louis B. (Ed.), Error in Digital Computation, Vol, I. John Wiley & Sons, 1965.
18. Wilkinson, J.H., Rounding Errors in Algebraic Processes. Prentice-Hall, 1963.
19. Bellman, Richard, Introduction to Matrix Analysis. McGraw-Hill, 1960.
20. Schneider, Hans (Ed.), Recent Advances in Matrix Theory. University of Wisconsin Press.
21. Forsythe, George and Moler, Cleve, Computer Solution of Linear Algebraic Systems. Prentice-Hall, 1967.
22. Householder, A.S., The Theory of Matrices in Numerical Analysis. Blaisdell Publishing Company, 1965.
23. Varga, R.S., Matrix Iterative Analysis. Prentice-Hall, 1962.
24. Wilkinson, J.H., The Algebraic Eigenvalue Problem. Clarendon Press, 1965.
25. Brockett, Roger W., Finite Dimensional Linear Systems. John Wiley and Sons, Inc., 1970.
26. Byerly, R.T. and Kimbark, E.W. (Eds.), Stability of Large Electric Power Systems. IEEE Press, New York, 1974.
27. Elgerd, O.I., Electric Energy Systems Theory: An Introduction. McGraw-Hill.
28. Selby, S.M. (Ed.), Standard Mathematical Tables. CRC Press, Inc., 1974.
29. Laub, Alan J., Decentralization of Optimization Problems by Iterative Coordination. Center for Comparative Studies in Technological Development and Social Change, University of Minnesota, 1974.

30. Barnett, S. and Storey C., Matrix Methods in Stability Theory. Thomas Nelson & Sons, 1970.
31. Faddeeva, V.N., Computational Methods of Linear Algebra. Dover Publications, Inc., 1959.
32. Astrom, K.J., Introduction to Stochastic Control Theory. Academic Press, 1970.
33. Gastinel, Noel, Linear Numerical Analysis. Academic Press, 1970.
34. Oppenheim, Alan V., Papers on Digital Signal Processing. M.I.T. Press, 1973.
35. Wehl, H., "Inequalities Between Two Kinds of Eigenvalues of a Linear Transformation," Proc. Nat. Acad. Sci., Vol. 35, pp. 408-11, 1949.
36. Von Neumann, J. and Goldstine, H.H., "Numerical Inverting of Matrices of High Order," Bull. Amer. Math. Soc., Vol. 53, pp. 1021-1099, 1947.
37. Fan, K., "On a Theorem of Wehl," Proc. Nat. Acad. Sci., Vol. 36, pp. 760-66, 1951.
38. Marcus, M., "Remark on a Norm Inequality for Square Matrices," Proc. Amer. Math. Soc., Vol. 6, pp. 117-19, 1955.
39. Kato, T., "·", Journal of Phys. Soc. of Japan, Vol. 4, pp. 334-339, 1949.
40. Jameson, A., "Solution of the Equation $AX + XB = C$ by the Inversion of an $M \times M$ or $N \times N$ Matrix," SIAMJ Appl. Math., Vol. 16, pp. 1020-1023, 1968.
41. Lee, Stephen, Transient Stability Equivalents for Power System Planning, EPSEL Report No. 38, M.I.T., June, 1972.
42. deMello, R., Podmore, R., Stanton, K.N., Coherency Based Dynamic Equivalents for Transient Stability Studies, Electric Power Research Institute, Final Report EPRI 904, January, 1975.