

UNIVERSIDADE ESTADUAL DE CAMPINAS
INST. DE MATEMÁTICA, ESTATÍSTICA E CIÊNCIA DA COMPUTAÇÃO

SOBRE ERROS DE ARREDONDAMENTO
EM PROCESSOS ALGÉBRICOS

DISSERTAÇÃO DE MESTRADO

Márcia Ferrari Orsi
Orientador: Prof. Dr. Odelar Leite Linhares

UNICAMP
BIBLIOTECA CENTRAL

CAMPINAS
1973

Este trabalho foi realizado graças ao apoio financeiro da Fundação de Amparo à Pesquisa do Estado de São Paulo, através da Bolsa de Aperfeiçoamento que nos outorgou no período de março de 1971 a julho de 1973.

Aos meus pais, Cida e Arthur Ferrari, aos quais nunca serei suficientemente grata pelo carinho e abnegação com que se dedicaram à minha formação.

Do meu marido, Armando Geraldo Orsi, cujo apoio, incentivo e orientação foram fundamentais para a realização deste trabalho e de todo o programa de pós-graduação.

Sinceros agradecimentos

Ao Prof. Dr. Odemar Leite Binhares, nosso orientador e caro mestre, pela dedicada atenção que nos dispensou durante a realização não só deste trabalho, mas também de todo o nosso programa de Mestrado.

Ao Prof. Dr. Ivan de Queiroz Barros, cujas valiosas sugestões, apresentadas durante a discussão de inúmeros tópicos deste trabalho, foram de fundamental importância para a conclusão do mesmo.

Ao Prof. Dr. Waldemar W. Setzer, pela investida colaboração que nos dispensou, permitindo-nos a utilização de todos os recursos humanos e materiais do Departamento de Matemática Aplicada do IME da USP.

Ao meu marido, Armando Geraldo Orsi, que, com paciência e dedicação, colaborou de forma decisiva na elaboração geral deste trabalho.

Aos professores e funcionários do IMECC da UNICAMP, do IME da USP e do ICM de São Carlos, que em inúmeras ocasiões nos auxiliaram nas diversas atividades do nosso programa de pós-graduação.

À Srta. Maria Alice F. Carneira, pela primorosa organização das referências bibliográficas.

Resumo

Neste trabalho são estudados os erros de arredondamento em processos algébricos, com ênfase no método regressivo de análise dos mesmos. Inicialmente se efetua um estudo dos erros de arredondamento nas aritméticas de vírgula fixa e flutuante, recebendo o segundo caso um tratamento mais detalhado. A seguir, são descritas as análises progressiva e regressiva de erros, sendo a última estudada de forma mais completa. Uma breve evolução histórica também se encontra incluída, com a finalidade de situar cronologicamente o desenvolvimento da análise de erros. Finalmente, é efetuada uma análise regressiva de erros do método de eliminação de Gauss, assim como de alguns problemas que podem ser resolvidos através do mesmo.

Abstract

In this work we study the rounding errors occurring in algebraic processes. Emphasis is given to the backward method of error analysis. We make at first a study of rounding errors in fixed and floating-point arithmetic; the second case receives a more detailed treatment. Forward and backward error analysis are then developed, with more attention given to the latter. A brief historical evolution of the field development is included, so that one can place chronologically the development of error analysis. We conclude by making a backward error analysis of Gaussian elimination as well as of some problems which can be solved with this method.

Índice

	<u>página</u>
<u>Introdução</u> -----	1
<u>1. Pré-requisitos</u> -----	4
1.1. Espaços vetoriais e euclidianos -----	4
1.2. Matrizes -----	4
1.3. Aplicações lineares -----	6
1.4. Normas de vetores e matrizes -----	7
1.5. Redução de uma matriz à forma diagonal através de equivalência ortogonal -----	13
1.6. Alguns lemas úteis -----	17
1.7. Representação interna dos números nos computadores -----	19
1.7.1. Sistema de representação em vírgula fixa -----	19
1.7.2. Sistema de representação normalizada em vírgula flutuante -----	20
1.7.3. Precisão simples e múltipla -----	21
<u>2. Erros de arredondamento</u> -----	22
2.1. Classificação dos erros -----	22
2.2. Arredondamento: definições e técnicas -----	24
2.3. Arredondamento na aritmética em vírgula fixa -----	26
2.3.1. Determinação de x_R e x_G -----	26
2.3.2. Delimitações para os erros de arredondamento -----	27
2.3.3. Erros de arredondamento nas operações aritméticas -----	27
2.4. Arredondamento na aritmética em vírgula flutuante -----	33
2.4.1. Determinação de x_R e x_G -----	33
2.4.2. Delimitações para os erros de arredondamento -----	34
2.4.3. Erros de arredondamento nas operações aritméticas -----	35
2.5. Conclusões -----	54

	<u>página</u>
3. <u>Análise de erros</u> - - - - -	56
3.1. Sensibilidade - - - - -	56
3.1.1. Problemas mal condicionados - - - - -	56
3.1.2. Números de condição - - - - -	56
3.1.3. Sensibilidade das soluções de sistemas de equações lineares - - - - -	60
3.1.4. Sobre o mal condicionamento de sistemas de equações lineares - - - - -	65
3.2. Análises progressiva e regressiva - - - - -	70
3.3. Evolução histórica - - - - -	77
 4. <u>Análise de erros do método de eliminação de Gauss.</u>	
<u>Aplicações</u> - - - - -	80
4.1. Resolução de sistemas triangulares - - - - -	80
4.2. Eliminação de Gauss - - - - -	85
4.3. Cálculo de determinantes - - - - -	100
4.4. Inversão de matrizes - - - - -	101
4.5. Matrizes especiais - - - - -	103
4.6. Refinamento de solução aproximada - - - - -	107
 <u>Bibliografia</u> - - - - -	115

Introdução

O computador eletrônico constitui atualmente um elemento de grande utilidade na resolução numérica de problemas em todos os ramos do conhecimento. Entretanto, tais máquinas utilizam uma quantidade finita de dígitos para a representação dos números com os quais operam, o que pode ocasionar o aparecimento de erros de arredondamento, quer nos dados, quer nos resultados de operações intermediárias. Como a propagação desses erros não é susceptível de controle, pode-se chegar a resultados finais completamente sem significado em relação ao problema original; daí, ser de suma importância complementar a resolução numérica com uma análise de erros do processo utilizado.

Apesar de sua fundamental importância, a análise de erros ainda não está difundida entre os usuários de computadores, principalmente no Brasil, que costumam considerá-la como assunto por demais especializado. Essa concepção é consequência não só do desenvolvimento relativamente recente de uma teoria bem estruturada a respeito do assunto, mas, também, do estudo irregular da matéria, através da leitura de análises efêmeras segundo orientações diversas.

No desenvolvimento deste trabalho tivemos como principal objetivo o estudo da análise de erros relativamente a processos algébricos. Esse estudo foi realizado em quatro etapas sucessivas, as quais deram origem aos quatro capítulos que constituem o texto em questão.

No primeiro capítulo são apresentados alguns assuntos necessários ao entendimento do texto, assim como explicadas as notações utilizadas no mesmo; também estão incluídas, de forma sucinta e sem demonstrações, algumas definições, propriedades e resultados que são empregados no desenvolvimento dos capítulos subsequentes. As normas de vetores e matrizes

zes foram consideradas em particular, por serem muito úteis tanto no terceiro quanto no quarto capítulos.

No segundo capítulo são estudados os erros de arredondamento nas aritméticas de vírgula fixa e flutuante. O segundo caso recebe um tratamento mais completo, já que, atualmente, os sistemas de representação em vírgula flutuante são os mais utilizados.

No capítulo três são descritas as duas principais formas de análise de erros: progressiva e regressiva. Nossa atenção concentra-se sobre a análise regressiva que, por apresentar inúmeras vantagens sobre a progressiva, é a mais utilizada atualmente. Uma breve evolução histórica é incluída neste capítulo, com a finalidade de situar cronologicamente o desenvolvimento da análise de erros, assim como citar os principais autores responsáveis pelo mesmo.

No quarto capítulo a análise regressiva de erros é aplicada à eliminação de Gauss e, também, a alguns problemas que podem ser resolvidos através desse método.

Durante toda a fase de redação procuramos dar a este trabalho um aspecto didático, visando a obtenção de um texto de utilidade para todos os que se interessam ou venham a se interessar por esse assunto, em nosso país.

As referências bibliográficas foram feitas de acordo com as normas PNB-66, de 1970, da ABNT. Entretanto, a chamada da bibliografia no texto não foi efetuada de acordo com tais normas e sim utilizando uma citação do tipo (a, p. m-n), em que a representa o número que identifica o livro na relação bibliográfica e m e n os números das páginas inicial e final do trecho citado; esses dois últimos números são excluídos quando o assunto de interesse se encontra tratado ao longo de todo o livro em questão.

Embora os resultados contidos neste trabalho também se apliquem a cálculos realizados manualmente,

com ou sem o auxílio de calculadoras de mesa, a sua utilização é mais necessária em cálculos que, por envolverem um número muito grande de operações, exigem o uso de um computador eletrônico digital; por esse motivo, preocupamo-nos exclusivamente com esse último caso. Por simplicidade, utilizamos, em todo o texto, a palavra "computador" em substituição à expressão "computador eletrônico digital".

Apesar de ser a binária a base de operação da maioria dos computadores utilizados para fins científicos, os exemplos são apresentados na base dez, tendo sido efetuados em uma calculadora de mesa. Com esse procedimento buscamos não só facilitar o entendimento dos exemplos, como, também, evitar os erros de arredondamento, que podem ser envolvidos na conversão da base 10 para a base 2, e vice-versa.

1 - Pré-requisitos

1.1 - Espaços vetoriais e euclidianos

No texto é bastante utilizada a estrutura de espaço vetorial real para \mathbb{R}^n . Os vetores e escalares desse espaço são representados, respectivamente, por letras latinas e gregas minúsculas. Um vetor x , de coordenadas x_1, x_2, \dots, x_n , relativamente a alguma base, é indicado por

$$(1.1) \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \text{ou} \quad x^T = (x_1, x_2, \dots, x_n)$$

Em todo o trabalho consideramos, para \mathbb{R}^n , a base canônica, cujos elementos são

$$x_1^T = (1, 0, \dots, 0), \quad x_2^T = (0, 1, \dots, 0), \quad \dots, \quad x_n^T = (0, 0, \dots, 1).$$

A multiplicação escalar utilizada no texto é a usual de \mathbb{R}^n , definida por

$$(1.2) \quad \langle, \rangle : \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R} \\ (x, y) \longrightarrow \langle x, y \rangle = \sum_{i=1}^n x_i \cdot y_i,$$

em que $x^T = (x_1, \dots, x_n)$; $y^T = (y_1, \dots, y_n)$

1.2. Matrizes

São utilizadas no texto, principalmente, as definições de matriz, submatriz, partição de matriz e determinante. As matrizes consideradas são reais, quadradas, de or-

dem n e indicadas por $A = (a_{ij})$; o conjunto dessas matrizes é indicado por $M_n(\mathbb{R})$. O determinante de uma matriz A é indicado por $\det A$. Também são utilizadas as definições e propriedades das operações adição e multiplicação de matrizes e a lei de composição externa multiplicação de um escalar por uma matriz, além das estruturas de espaço vetorial e anel com elemento unidade que essas leis de composição definem em $M_n(\mathbb{R})$. O elemento unidade de $M_n(\mathbb{R})$ é indicado por I . Dada $A \in M_n(\mathbb{R})$, A^T e A^{-1} indicam, respectivamente, a transposta e a inversa de A (é lógico que, no último caso, A é suposta não singular). Sejam $A, B \in M_n(\mathbb{R})$, $A = (a_{ij})$ e $B = (b_{ij})$; as notações $|A|$ e $|A| \leq |B|$ indicam:

$$(1.3) \quad |A| = (|a_{ij}|) \quad , \quad |A| \in M_n(\mathbb{R})$$

$$(1.4) \quad |A| \leq |B| \text{ se e só se } |a_{ij}| \leq |b_{ij}| \quad , \quad i, j = 1, \dots, n,$$

onde $|x|$ representa o valor absoluto de x , $x \in \mathbb{R}$.

Apresentaremos, a seguir, algumas definições que são muito utilizadas nos capítulos subsequentes. Todas as matrizes consideradas são supostas pertencentes a $M_n(\mathbb{R})$.

(1.5) Definição - Uma matriz A é dita densa quando apenas uma pequena porcentagem de seus elementos é diferente de zero. Em caso contrário, A é dita esparsa ou rarefeita.

Como se pode facilmente concluir, a definição (1.5) tem caráter prático, sendo uma questão de convenção o estabelecimento de um limite entre as duas classes de matrizes.

(1.6) Definição - Denominam-se valores singulares de $A = (a_{ij})$ as raízes quadradas não negativas dos autovalores de AA^T .

(1.7) Definição - $A = (a_{ij})$ diz-se ortogonal se $A^T = A^{-1}$.

(1.8) Definição - $A = (a_{ij})$ e $B = (b_{ij})$ são ditas semelhantes

se e somente se existe uma matriz P não singular tal que $A = P^{-1}BP$

(1.9) Definição - $A = (a_{ij})$ e $B = (b_{ij})$ são ditas equivalentes se e somente se existem matrizes P e Q não singulares tais que $A = Q^{-1}BP$.

(1.10) Definição - Se na definição (1.9) as matrizes P e Q são ortogonais, então, A e B dizem-se ortogonalmente equivalentes. Nesse caso, $A = Q^{-1}BP$ pode ser escrita $A = Q^TBP$, em virtude de (1.7)

(1.11) Definição - $A = (a_{ij})$ é dita Hessenberg superior (inferior) se $a_{ij} = 0$ para $i \geq j+2$ ($j \geq i+2$)

(1.12) Definição - $A = (a_{ij})$ é chamada tridiagonal se $a_{ij} = 0$ para $|i-j| > 1$.

(1.13) Definição - $A = (a_{ij})$ é dita diagonalmente dominante se $|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|$, $\forall i$.

1.3. Aplicações lineares

São utilizadas no texto, principalmente, a definição de aplicação linear e a sua representação através de matrizes. As aplicações usadas são do tipo $A: \mathbb{R}^m \rightarrow \mathbb{R}^n$; o conjunto dessas aplicações é representado por $L(\mathbb{R}^n)$. Também são utilizadas as definições e propriedades das operações de adição e multiplicação de aplicações lineares e a lei de composição externa multiplicação de um escalar por uma aplicação linear, além das estruturas de espaço vetorial e anel com elemento unidade que essas leis de composição definem em $L(\mathbb{R}^n)$.

As aplicações lineares consideradas no texto são identificadas com suas respectivas matrizes, relativamente à base canônica de \mathbb{R}^n , sendo as aplicações e respectivas ma

trizes representadas por um mesmo símbolo. Essa identificação é justificada pelo isomorfismo existente entre $L(\mathbb{R}^n)$ e $M_m(\mathbb{R})$, desde que fixada uma base em \mathbb{R}^n , quer ambos sejam considerados como espaços vetoriais, quer como anéis.

1.4. Normas de vetores e matrizes

(1.14) Definição - Seja V um espaço vetorial sobre o corpo K dos números reais ou complexos. Uma norma definida em V é uma função que associa a cada $x \in V$ um número real $\|x\|$, chamado norma de x , que satisfaz a

$$\|x\| > 0, \text{ se } x \neq \theta \text{ e } \|\theta\| = 0, \text{ sendo } \theta \text{ o vetor nulo.}$$

$$\|\lambda x\| = |\lambda| \cdot \|x\|, \quad \forall x \in V, \forall \lambda \in K$$

$$\|x+y\| \leq \|x\| + \|y\|$$

A partir de (1.14) pode ser provada a relação

$$(1.15) \quad \left| \|x\| - \|y\| \right| \leq \|x \pm y\| \leq \|x\| + \|y\|$$

No espaço vetorial \mathbb{R}^n , utilizado neste trabalho, as normas de vetores mais utilizadas são casos particulares da norma de Hölder, que é definida

$$(1.16) \quad \|x\|_p = \left[\sum_{i=1}^n |x_i|^p \right]^{\frac{1}{p}}, \quad p \geq 1,$$

em que x_i é a i -ésima coordenada do vetor $x^T = (x_1, x_2, \dots, x_n)$, $x^T \in \mathbb{R}^n$. Normalmente são utilizadas as normas obtidas pela substituição de p por 1, 2, ∞ em (1.16), sendo $\|x\|_\infty$ interpretada como $\max_{1 \leq i \leq n} |x_i|$, ou seja,

$$(1.17) \quad \|x\|_1 = \sum_{i=1}^n |x_i|$$

$$(1.15) \quad \|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}}$$

$$(1.16) \quad \|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

Referimos-nos à norma obtida pela substituição de p por a ($a \geq 1$) em (1.16) por norma a ou $\| \cdot \|_a$. A norma 2 costuma ser chamada norma euclideana; $\|x\|_2$ representa o comprimento do vetor x , da maneira que é definido usualmente. É importante observar que para a norma ∞ vale

(1.20) Propriedade - Se $x^T = (x_1, x_2, \dots, x_n)$ e se $|x_i| \leq c$, então, $\|x\|_\infty \leq c$.

(1.20) não se verifica para as normas 1 e 2.

Por outro lado, a partir de (1.17) - (1.19) prova-se:

$$(1.21) \quad \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty$$

$$(1.22) \quad \|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty$$

Em virtude da definição de produto escalar dada em (1.2) e de (1.18), temos:

$$(1.23) \quad \|x\|_2^2 = \langle x, x \rangle$$

Consequentemente, a conhecida desigualdade de Cauchy-Schwarz-Bunyakovskii pode ser escrita

$$(1.24) \quad |\langle x, y \rangle| \leq \|x\|_2 \|y\|_2$$

(1.25) Definição - Dada uma sequência de vetores $(x^{(k)}) \in V$ e uma norma $\| \cdot \|_a$ em V , diz-se que a sequência $(x^{(k)})$ converge para $x \in V$ se $\|x^{(k)} - x\|_a \rightarrow 0$.

Como vimos em 1.2, o conjunto $M_n(\mathbb{R})$, das matrizes quadradas de ordem n , possui estrutura de espaço vetorial

torial, relativamente à operações adição de matrizes e à lei de composição multiplicação de um escalar por uma matriz. Portanto, a definição (1.14) vale para $M_n(\mathbb{R})$; basta que se substituam, nos axiomas citados nessa definição, os vetores x e y por matrizes A e B e θ por Θ , sendo Θ a matriz nula.

Podemos associar, a cada norma de vetor $\|\cdot\|_k$, uma norma de matriz definida por

$$(1.26) \quad \|A\|_k = \sup_{x \neq \theta} \frac{\|Ax\|_k}{\|x\|_k}$$

ou, equivalentemente,

$$(1.27) \quad \|A\|_k = \sup_{\|x\|_k=1} \|Ax\|_k$$

Como a norma é uma função contínua, e como o domínio definido por $\|x\|_k=1$ é fechado e limitado, podemos substituir "sup" por "max" em (1.26) e (1.27). Por outro lado, a novidade do isomorfismo citado em 1.3, podemos dizer que (1.27) define $\|A\|_k$ como sendo igual ao comprimento do mais longo vetor do conjunto imagem $\{Ax\}$ da esfera unitária $\{x \text{ tal que } \|x\|_k=1\}$, sob a aplicação $x \rightarrow Ax$.

(1.28) Definição. A norma de matriz definida em (1.26) ou (1.27) é dita subordinada à norma de vetor $\|\cdot\|_k$.

Subordinadas às normas de vetores 1, 2, ∞ , temos:

$$(1.29) \quad \|A\|_1 = \max_j \sum_i |a_{ij}|, \quad i, j = 1, 2, \dots, n$$

$$(1.30) \quad \|A\|_\infty = \max_i \sum_j |a_{ij}|, \quad i, j = 1, 2, \dots, n$$

$$(1.31) \quad \|A\|_2 = \left[\max_i \lambda_i(A^T A) \right]^{\frac{1}{2}}, \text{ onde } \lambda_i(A^T A) \text{ indica os autovalores de matriz } A^T A.$$

A norma definida em (1.31) é, freqüentemente, denominada norma espectral ou de Hilbert.

(1.32) Propriedade - Sejam $A, B \in M_n(\mathbb{R})$, $A = (a_{ij})$ e $B = (b_{ij})$.
Se $|a_{ij}| \leq |b_{ij}|$ ($i, j = 1, \dots, n$), então, $\|A\|_\infty \leq \|B\|_\infty$.

A propriedade (1.32) também é verificada pela norma 1, porém, não é válida para a norma 2.

Facilmente se verifica que, para as normas subordinadas, vale:

$$(1.33) \quad \|I\|_k = 1$$

(1.34) Definição - Sejam $\|\cdot\|_p$ e $\|\cdot\|_k$ normas de vetores e matrizes, respectivamente. Essas normas são ditas consistentes ou compatíveis se

$$\|Ax\|_p \leq \|A\|_k \|x\|_p$$

Uma norma de vetores e a norma de matrizes a ela subordinada são sempre compatíveis. Por outro lado, quando são consideradas normas subordinadas em (1.34), existe sempre um vetor x não nulo para o qual se verifica a igualdade.

Podemos determinar uma norma de matriz, a partir de uma norma de vetor dada, de modo que ambas sejam consistentes, através de (1.26) ou (1.27). Em caso contrário, ou seja, para determinar uma norma de vetor, a partir de uma norma de matriz dada, de modo que ambas sejam consistentes, basta tomar

$$(1.35) \quad \|x\|_p = \|A\|_k, \quad \text{com} \quad A = \begin{pmatrix} x_1 & 0 & \dots & 0 \\ x_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ x_n & 0 & \dots & 0 \end{pmatrix},$$

sendo $x^T = (x_1, \dots, x_n)$.

Com a finalidade de evitar a repetição exagerada das expressões "norma de vetores" e "norma de matrizes", utilizaremos, em sua substituição, $\|x\|$ e $\|A\|$, respectivamente.

Além de $\|A\|_1$, $\|A\|_2$ e $\|A\|_\infty$, existe ainda uma importante norma de matrizes, denominada norma euclidiana ou de Schur e representada por $\|A\|_E$, que é definida por

$$(1.36) \quad \|A\|_E = \left(\sum_{i,j} |a_{ij}|^2 \right)^{\frac{1}{2}} = \left[\text{tr}(A^T A) \right]^{\frac{1}{2}} = \left[\text{tr}(A A^T) \right]^{\frac{1}{2}}$$

Embora compatível com $\|x\|_2$, $\|A\|_E$ não é subordinada a nenhuma norma de vetor, o que pode ser concluído de

$$(1.37) \quad \|I\|_E = n^{\frac{1}{2}},$$

já que para normas subordinadas vale (1.33).

Em trabalhos estritamente matemáticos a norma utilizada em conexão com $\|x\|_2$ é, normalmente, $\|A\|_2$. Entretanto, na prática, geralmente $\|A\|_E$ é utilizada em substituição a $\|A\|_2$, por ser de cálculo mais fácil que esta e por gozar da propriedade

$$(1.38) \quad \| |A| \|_E = \|A\|_E,$$

o que não ocorre com a mesma. (As normas $\|A\|_1$ e $\|A\|_\infty$ também satisfazem (1.38)). As normas $\|A\|_E$ e $\|A\|_2$ são relacionadas por:

$$(1.39) \quad \|A\|_2 \leq \|A\|_E \leq n^{\frac{1}{2}} \|A\|_2$$

$$(1.40) \quad \| |A| \|_2 \leq \| |A| \|_E = \|A\|_E \leq n^{\frac{1}{2}} \|A\|_2$$

Apesar da segunda desigualdade de (1.39), $\|A\|_E$ é, frequentemente, uma boa aproximação para $\|A\|_2$.

Para as normas subordinadas e para $\|A\|_E$, além dos axiomas da definição de norma, vale

$$(1.41) \quad \|AB\|_k \leq \|A\|_k \|B\|_k \quad (k=1, 2, \infty, E)$$

As relações dadas em (1.34) e (1.41) são muito úteis no estudo das aplicações lineares e, em particular, no estudo dos erros de arredondamento contraídos durante a resolução numérica de sistemas de equações lineares.

Teremos, a seguir, um lema que é muito utilizado nos capítulos 3 e 4.

(1.42) Lema. Seja F uma matriz tal que $\|F\|_k < 1$, em que k representa uma norma subordinada ou $k = E$. Então, vale

(1.43) $I + F$ é não singular

$$(1.44) \quad \|(I+F)^{-1}\|_k \leq \frac{1}{1 - \|F\|_k}$$

$$(1.45) \quad \|I - (I+F)^{-1}\|_k \leq \frac{\|F\|_k}{1 - \|F\|_k}$$

Dejamos a demonstração de (1.42). Por absurdo, suponhamos que $I + F$ seja singular. Então, existe um vetor x não nulo, tal que

$$(1.46) \quad (I+F)x = \theta \quad \text{ou} \quad Fx = -x$$

Utilizando normas consistentes em (1.46), temos:

$$\|F\|_k \|x\|_p \geq \|Fx\|_p = \|x\|_p$$

e, portanto, $\|F\|_k \geq 1$, pois $x \neq \theta$ por hipótese. Logo, (1.43) está provado.

Consideremos a identidade

$$(1.47) \quad I = (I+F)^{-1}(I+F),$$

na qual a existência de $(I+F)^{-1}$ é garantida por (1.43)

Seja

$$(1.48) \quad R = (I + F)^{-1} \quad \text{e, portanto,}$$

$$(1.49) \quad I = R + RF$$

Aplicando (1.15), (1.41) e (1.33) a (1.49) temos:

$$(1.50) \quad 1 = \|I\|_k \geq \|R\|_k - \|RF\|_k \geq \|R\|_k - \|R\|_k \|F\|_k$$

e, como $1 - \|F\|_k$ é positivo, por hipótese,

$$(1.51) \quad \|R\|_k \leq \frac{1}{1 - \|F\|_k}$$

(1.51) prova (1.44) para as normas subordinadas.

Essa prova não se aplica a $\|\cdot\|_E$, já que em (1.50) foi utilizada (1.33). Entretanto, como a prova vale para a norma 2, e como por hipótese $\|F\|_k < 1$, então, por (1.39) podemos dizer que (1.44) também vale para $\|\cdot\|_E$, "a fortiori".

Por outro lado, de (1.49) e (1.41) temos

$$\|I - R\|_k = \|RF\|_k \leq \|R\|_k \|F\|_k, \quad \text{ou,}$$

$$\|I - R\|_k \leq \frac{\|F\|_k}{1 - \|F\|_k},$$

o que prova (1.45).

5. Redução de uma matriz à forma diagonal através de equivalência ortogonal.

O teorema apresentado em (1.52) possibilita uma interpretação geométrica muito simples a respeito da representação de uma aplicação linear $A \in L_m(\mathbb{R})$, através de uma matriz $A \in M_m(\mathbb{R})$. Neste trabalho são feitas apenas algumas

considerações a respeito do referido teorema, cuja demonstração pode ser encontrada em (7, p. 9-10)

(1.52) Teorema - Dada $A \in M_n(\mathbb{R})$, existem duas matrizes ortogonais U e V , $U, V \in M_n(\mathbb{R})$, tais que $D = U^T A V$ é uma matriz diagonal. Além disso, é possível determinar U e V de modo que os elementos diagonais de D sejam

$$(1.53) \quad \mu_1 \geq \mu_2 \geq \dots \geq \mu_r > \mu_{r+1} = \dots = \mu_n = 0,$$

onde r é o posto de A . Em particular, se A é não singular,

$$(1.54) \quad \mu_1 \geq \mu_2 \geq \dots \geq \mu_n > 0$$

Os números $\mu_1, \mu_2, \dots, \mu_n$ são os valores singulares de A .

Em muitos livros de Álgebra as matrizes U^T e V são substituídas por matrizes gerais não singulares P^{-1} , Q e todos os μ_i são tomados iguais a 1. Em cálculos realizados em computadores, entretanto, as matrizes ortogonais são mais úteis porque gozam da propriedade $\|Ux\| = \|x\|$, $\forall U$ ortogonal e $\forall x \in \mathbb{R}^n$.

Interpretemos A como a representação de um elemento de $L(\mathbb{R}^n)$, ou seja, $\forall x, x \in \mathbb{R}^n$, $(y = Ax) \in \mathbb{R}^n$. Consideremos novas bases ortogonais em \mathbb{R}^n , de modo que as matrizes de mudança da base canônica para essas bases sejam V e U , ou seja

$$(1.55) \quad x = V x', \quad V \text{ ortogonal}$$

$$(1.56) \quad y = U y', \quad U \text{ ortogonal}$$

As matrizes U e V , definidas em (1.55) e (1.56) são as matrizes citadas em (1.52), pois, de (1.7), (1.55) e (1.56) temos

$$y' = U^T y = U^T A x = U^T A (V x') = (U^T A V) x' = D x'$$

Portanto, em virtude das mudanças de bases, a aplicação, que originalmente era representada por A , obtém uma nova representação que é D . Por outro lado, relativamente às novas bases, essa aplicação possui uma representação bastante simples que, em termos de coordenadas, é dada por

$$(1.57) \quad \begin{cases} y'_1 = \mu_1 x'_1 \\ y'_2 = \mu_2 x'_2 \\ \vdots \\ y'_n = \mu_n x'_n \\ y'_{n+1} = 0 \\ \vdots \\ y'_m = 0 \end{cases}$$

Podemos então afirmar que a aplicação considerada leva o eixo Ox'_1 de \mathbb{R}^n no eixo Oy'_1 , também de \mathbb{R}^n , com um fator de amplificação $\mu_1 > 0$. O mesmo também ocorre com os eixos $Ox'_2, Ox'_3, \dots, Ox'_n$, que são levados em $Oy'_2, Oy'_3, \dots, Oy'_n$, com, respectivamente, fatores de amplificação μ_2, \dots, μ_n . Os eixos Ox'_{n+1}, \dots, Ox'_m são todos levados ao vetor nulo de \mathbb{R}^n .

Facilmente se verifica que os μ_k ($k=1, 2, \dots, n$) são os valores singulares de A . De fato, por (1.7)

$$(1.58) \quad D^T D = (U^T A V)^T (U^T A V) = V^T A^T U U^T A V = V^T (A^T A) V = V^T (A^T A) V$$

Portanto, $D^T D$ e $A^T A$ possuem os mesmos autovalores, já que são semelhantes, como pode ser deduzido de (1.58). Como os autovalores de $D^T D$ são μ_1^2, \dots, μ_n^2 , podemos concluir que os valores singulares de A são μ_1, \dots, μ_n .

Através de (1.52) podemos mostrar que D transforma a esfera unitária $S = \{x' \text{ tal que } \|x'\|_2 = 1\}$

em um hiper elipsóide r -dimensional $E = DS$, de vetores y^i tais que

$$\frac{y_1'^2}{\mu_1^2} + \dots + \frac{y_r'^2}{\mu_r^2} = 1 \quad \text{e} \quad y_{r+1}' = \dots = y_n' = 0$$

Um dos pontos de E mais afastados da origem é o ponto $(\mu_1, 0, \dots, 0)$. Se $r < n$, então E contém a origem θ ; se $r = n$ então E não contém a origem θ e um dos pontos de E mais próximos de θ é $(0, 0, \dots, \mu_n)$. Se $r < n$, então, D e potentes A , são matrizes singulares; se $r = n$, D e A são não singulares e, a partir de (1.57) temos

$$D^{-1} = \begin{pmatrix} \mu_1^{-1} & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & \mu_n^{-1} \end{pmatrix}$$

Portanto, os valores singulares de A^{-1} são $\mu_1^{-1}, \dots, \mu_n^{-1}$

A partir dessas observações que fizemos e de (1.2.6) temos

$$(1.59) \quad \|A\|_2 = \|D\|_2 = \mu_1$$

e, no caso $r = n$,

$$(1.60) \quad \|A^{-1}\|_2 = \|D^{-1}\|_2 = \mu_n^{-1}$$

Resumindo as conclusões mais importantes para uma matriz $A \in M_n(\mathbb{R})$, A não singular, com valores singulares $\mu_1 \geq \dots \geq \mu_n > 0$, podemos afirmar que existe uma linha L_1 em \mathbb{R}^n que é prolongada (ou encurtada) pelo fator μ_1 por A , quando L_1 é levada em uma linha AL_1 de \mathbb{R}^n . Por outro lado, existe uma segunda li-

nha L_m , ortogonal a L_1 , que é prolongada (ou encurtada) do fator μ_n . Além disso, AL_1 e AL_m são ortogonais em R^n . Um círculo unitário, pertencente ao plano de L_1 e L_m , é levado por A numa elipse com semi-eixos μ_1 e μ_n ; essa é a maior distorção que pode ocorrer a qualquer círculo em R^n .

Finalmente, observe-se que o determinante de A satisfaz

$$(1.61) \quad |\det A| = \mu_1 \mu_2 \cdots \mu_n,$$

pois $|\det A| = |\det U^T| |\det D| |\det V| = |\det D|$ já que as matrizes ortogonais têm determinante igual a (± 1) .

1.6. Alguns lemas úteis

Os lemas considerados a seguir auxiliam a apresentar os resultados contidos em análises de erros sob forma mais simples e compacta.

$$(1.62) \quad \text{Lema. Se } 0 \leq \mu < 1, \text{ e se } n = 1, 2, \dots, \text{ então,}$$

$$1 - n\mu \leq (1 - \mu)^n.$$

De fato, seja $f(\mu) = (1 - \mu)^n$. Então, pelo teorema de Taylor

$$f(\mu) = f(0) + \mu f'(0) + \frac{\mu^2}{2} f''(\theta\mu), \quad 0 < \theta < 1$$

$$\text{ou} \quad f(\mu) = 1 - n\mu + R(\mu),$$

$$\text{com} \quad R(\mu) = \frac{\mu^2}{2} n(n-1) (1 - \theta\mu)^{n-2}$$

$$\text{Como } R(\mu) \geq 0, \text{ temos } 1 - n\mu \leq (1 - \mu)^n.$$

(1.63) lema - $1+x \leq e^x$, $\forall x \geq 0$

Desenvolvendo e^x segundo Taylor, temos:

$$e^x = 1 + x + R(x),$$

com $R(x) = \frac{x^2}{2} e^{\theta x}$, $0 < \theta < 1$

Como $R(x)$ é não negativo, então, $1+x \leq e^x$.

(1.64) lema - $e^x \leq 1 + 1,01x$, para $0 \leq x \leq 0,01$

Utilizando o desenvolvimento de e^x apresentada em (1.63) e levando em consideração a hipótese de (1.64) temos:

$$e^x \leq 1 + x + \frac{x^2}{2} e^{0,01} = 1 + x \left(1 + \frac{x}{2} e^{0,01} \right)$$

$$e^x \leq 1 + x \left(1 + \frac{0,01}{2} \cdot 1,01 \right)$$

$$e^x \leq 1 + 1,01x$$

(1.65) lema Se $n = 1, 2, \dots$, e $0 \leq nu \leq 0,01$, então,

$$(1+u)^n \leq 1 + 1,01nu$$

De fato, por (1.63) e (1.64)

$$(1+u)^n \leq (e^u)^n = e^{nu} \leq 1 + 1,01nu$$

(1.66) lema - Se $|\delta_i| \leq u$, $i = 1, \dots, n$, e se $nu \leq 0,01$, então

$$1 - nu \leq \prod_1^n (1 + \delta_i) \leq 1 + 1,01nu$$

Utilizando (1.62) e (1.65) temos

$$1 - nu \leq (1-u)^n \leq \prod_1^n (1 + \delta_i) \leq (1+u)^n \leq 1 + 1,01nu$$

Às vezes é conveniente escrever (1.66) na forma

$$(1.67) \quad \prod_{i=1}^n (1+d_i) = 1 + \theta, \quad |\theta| \leq 1.$$

1.7 - Representação interna dos números nos computadores

Os números são armazenados internamente, nos computadores, através de duas formas principais de representação: vírgula fixa e vírgula flutuante (ou vírgula móvel).

Seja N a base em que opera o particular computador considerado (atualmente são mais comuns computadores para os quais $N=2$, entretanto, também são conhecidos computadores com $N=8, 10$ e 16).

1.7.1 - Sistema de representação em vírgula fixa

Os números representáveis nesse sistema são da forma

$$(1.68) \quad x = \pm 0, d_1 d_2 \dots d_t,$$

em que os d_i ($i=1, \dots, t$) são dígitos que satisfazem a

$$(1.69) \quad 0 \leq d_i \leq N-1,$$

sendo t um número inteiro que ou é fixo para o computador utilizado ($t=35$ para o IBM 7090, $t=98$ para o CDC 1604) ou, em alguns casos, pode ser selecionado pelo usuário (IBM 1620).

De (1.68) e (1.69) pode-se concluir

$$(1.70) \quad -1 < x < 1 \quad \text{ou} \quad |x| < 1$$

O conjunto $F_t = F(N, t)$, dos números representáveis em um certo sistema em vírgula fixa, é um subconjunto finito do conjunto dos números racionais. Na reta real

esse conjunto é representado por pontos equidistantes, pertencentes ao intervalo $(-1, 1)$; a distância entre dois pontos consecutivos desse conjunto é N^{-t} .

1.7.2. Sistema de representação normalizada em vírgula flutuante

Os números representáveis neste sistema são da forma

$$(1.71) \quad x = \pm 0, d_1 d_2 \dots d_t \times N^b,$$

em que os dígitos d_i ($i = 1, \dots, t$) satisfazem a

$$(1.72) \quad \begin{cases} 1 \leq d_1 \leq N-1 \\ 0 \leq d_i \leq N-1 \quad (i = 2, \dots, t) \end{cases}$$

e b , denominado expoente de x , é um inteiro pertencente ao intervalo $-p \leq b \leq P$, sendo t, p, P inteiros que, normalmente, são fixos para um determinado computador. O valor $0, d_1 d_2 \dots d_t$ é denominado mantissa ou parte fracionária de x . A condição $d_1 \neq 0$ caracteriza o sistema como normalizado.

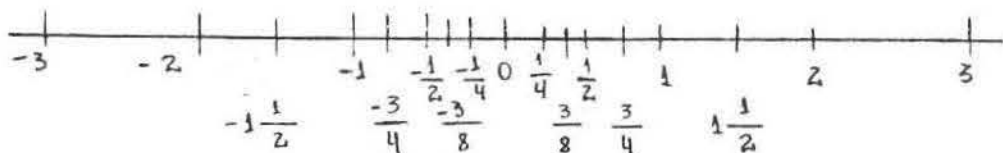
O número zero exige uma representação especial, que varia de computador para computador. As principais existentes são: atribuição do valor zero à mantissa, escolha do valor $(-p)$ para o expoente ou a combinação de ambas. Adotaremos a última, ou seja, representaremos o número zero por

$$(1.73) \quad 0,00\dots0 \times N^{-p}$$

O conjunto $F_L = F(N, t, p, P)$, dos números representáveis em um determinado sistema normalizado em vírgula flutuante, é um subconjunto finito do conjunto dos números racionais. Na reta real esse conjunto é representado por pontos não equidistantes. Esse fato é ilustrado em (1.74), através da representação esquemática de F_L , no caso em que

$$N=2, k=2, p=1 \text{ e } P=2$$

(1.74)



De (1.71) e (1.72) facilmente se conclui que qualquer número x pertencente a F_2 pode ser escrito

$$(1.75) \quad x = \pm a \cdot N^b, \text{ com } N^{-1} \leq a < 1$$

e, portanto,

$$(1.76) \quad N^{b-1} \leq |x| < N^b$$

Muitos autores utilizam a propriedade (1.75) para definir os sistemas normalizados em vírgula flutuante, a exemplo de Wilkinson (26) e Fox (9).

1.7.3. Precisão simples e múltipla

Indiquemos por t ou p o número de dígitos necessários para a representação de um número em vírgula fixa ou o necessário para a representação da mantissa de um número em vírgula flutuante. Diremos, nessas condições, que o computador trabalha com palavras de comprimento t e que a aritmética efetuada com os números assim representáveis é de precisão simples. Às vezes, é interessante trabalhar com precisão maior que uma parte em N^t , utilizando, então, t_1 dígitos, com $t_1 > t$. Se t_1 é múltiplo de t dizemos que o computador trabalha em precisão múltipla; para $t_1 = 2t, t_1 = 3t, \dots$, usam-se os termos específicos precisão dupla, tripla, etc. Nesses casos, as mantissas com $t, 2t, \dots$ dígitos são chamadas, respectivamente, t -digital, $2t$ -digital, ...

2 - Erros de Arredondamento

2.1 - Classificação dos erros

Os valores numéricos obtidos na solução de problemas podem conter erros de dois tipos: inerentes à formulação matemática da questão e contruídos durante a resolução da mesma através de processos numéricos.

Os erros inerentes à formulação matemática do problema podem ser originados ou pelo fato de ser tal formulação apenas uma aproximação de uma situação física ou por imprecisões nos dados físicos. Os primeiros são, em geral, desprezíveis, como acontece, por exemplo, com os efeitos relativísticos nos problemas de mecânica clássica. Os segundos podem ser desprezíveis se originados por imprecisões em constantes físicas, porém, devem ser cuidadosamente examinados quando causados por erros em dados empíricos.

Os erros acumulados durante a solução do problema através de métodos numéricos podem ser de três tipos: erros grosseiros ou enganos, erros de truncamento e erros de arredondamento. Em geral, os enganos têm sua origem em falhas humanas ou no mal funcionamento da máquina. Como exemplos temos, no primeiro caso, a perfuração incorreta de um cartão e, no segundo caso, a execução errônea de uma adição em virtude de falhas no computador. Erros desse tipo são mais comuns em cálculos efetuados manualmente, podendo, também, porém em escala muito menor, ocorrer com o uso de computadores. Nesse caso, a possibilidade de sua ocorrência deve sempre ser analisada, principalmente quando se torna difícil, quando não impossível, verificar se o valor encontrado é a solução procurada. Erros de truncamento

são os que ocorrem devido à resolução de problemas através de métodos de aproximação; geralmente surgem pela substituição de um processo infinito ou infinitesimal por uma aproximação finita. Tais erros aparecem, por exemplo, quando uma integral definida é calculada pela regra de Simpson, quando a soma de uma série infinita é calculada por uma fórmula de aproximação ou quando equações diferenciais são resolvidas por métodos de diferenças.

Erros de arredondamento são os que se cometem quando se utilizam, na solução de um problema, valores numéricos com menor número de dígitos do que aquele que realmente se dispõe; tais aproximações podem ser introduzidas nos dados do problema ou podem ocorrer em resultados intermediários, durante a resolução do mesmo. Como exemplos temos a introdução do valor numérico 3,14 em substituição ao dado π de um problema e a utilização dos números 3,33 e 3,49 como resultados, respectivamente, das operações $1,3 \times 2,41$ e $1,45$. A introdução dos erros de arredondamento é inevitável, quer por limite da capacidade de armazenamento, que ocorre nos computadores e máquinas de calcular em geral, quer por limitação de tempo, quando as operações são efetuadas manualmente.

Os erros que acabamos de conceituar podem ser representados em duas formas principais, absoluta e relativa, definidas por:

$$1) \quad \text{erro absoluto} = \text{valor aproximado} - \text{valor verdadeiro}$$

$$2) \quad \text{erro relativo} = \frac{\text{erro absoluto}}{\text{valor verdadeiro}}$$

Quando julgarmos conveniente, muitos autores consideram o erro relativo definido em relação ao valor aproximado, ou seja,

$$3) \quad \text{erro relativo} = \frac{\text{erro absoluto}}{\text{valor aproximado}}$$

A justificativa para essa atitude é o não conhecimento, de maneira geral, do valor verdadeiro da grandeza considerada.

É importante observar que, freqüentemente, o erro relativo é muito mais informativo do que o absoluto. Assim, por exemplo, sejam $\bar{x} = 0,8976$ e $\bar{y} = 0,0059$ aproximações para grandezas $x = 0,89764$ e $y = 0,00586$. O erro absoluto é igual a $0,4 \times 10^{-4}$ em ambos os casos, enquanto os erros relativos, $\frac{z}{44882}$ e $\frac{z}{293}$ para x e y , respectivamente, diferem bastante.

Neste trabalho são estudados os erros de arredondamento em processos algébricos. Esses erros, quer na forma absoluta, quer na relativa, são representados por letras gregas minúsculas; apesar disso, acreditamos que a notação adotada não causará equívocos porque, em todo o texto, não há possibilidade de confusão entre ambas as formas de erros. As expressões solução exata e resultado exato (soma, produto exatos, etc) indicam a solução de um problema e o resultado de uma operação obtidos sem a incidência de erros. Embora redundantes, essas expressões permitem a distinção entre os valores exatos e os valores fornecidos pelo computador, os últimos geralmente afetados de erros. Os valores fornecidos pelo computador são chamados, neste trabalho, solução e resultado aproximado.

4.2. - Arredondamento: definições e técnicas

Indiquemos por F o conjunto dos números representáveis ou no sistema em vírgula fixa ou no sistema normalizado em vírgula flutuante, ou seja, $F = F_1$ ou $F = F_2$, F_1 e F_2 definidos em 1.7. Seja x , $x \notin F$, o número que queremos aproximar por um elemento de F .

Denomina-se arredondamento ao processo pelo qual se determina um número $\bar{x} \in F$, tal que \bar{x} é uma aproximação de x . O erro cometido pela substituição de \bar{x} por x é chamado erro de arredondamento (absoluto ou relativo, conforme seja considerado de maneira absoluta ou relativa) de x . O valor \bar{x} é dito valor arredondado de x ou representação interna de x .

O arredondamento pode ser obtido através de diferentes técnicas, entre as quais as mais utilizadas são as que são descritas a seguir.

Nas considerações que se seguem, $d(a, b)$ representa a distância entre a e b .

Arredondamento propriamente dito - O valor arredondado de x , usualmente representado por x_R , deve satisfazer às condições:

$$(2.4) \quad \begin{cases} x_R \in F \\ d(x_R, x) \leq d(y, x), \forall y \in F \\ \forall y \in F, y \neq x_R, \text{ se } d(y, x) = d(x_R, x), \text{ então, } |y| < |x_R|. \end{cases}$$

Arredondamento por corte - O valor arredondado de x , normalmente representado por x_c , deve satisfazer às condições:

$$(2.5) \quad \begin{cases} x_c \in F \\ d(x_c, x) \leq d(y, x), \forall y \in F \\ |x_c| \leq |x|. \end{cases}$$

Em muitos livros o arredondamento por corte é denominado truncamento. Preferimos não utilizar esse termo, para evitar equívocos com o erro de truncamento definido em 2.1. Por outro lado, frequentemente o arredondamento propriamente dito é chamado, simplesmente, de arredondamento. Escolhemos a primeira dessas denominações com a finalidade de evitar confusão entre a técnica de obtenção e o processo.

de arredondamento.

2.3. - Arredondamento na aritmética em vírgula fixa

2.3.1 - Determinação de x_R e x_C

Seja $x \in F_i$ o número que queremos aproximar por um elemento de $F_i = F(N, t)$. De acordo com o tipo de arredondamento utilizado pelo computador considerado, x é aproximado por x_R ou x_C .

Determinação de x_R

Para determinar x_R consideramos a representação infinita de $|x| + \frac{1}{2} N^{-t}$ (aquela terminando em zeros, no caso de escolha) e retemos os t primeiros dígitos após a vírgula, desprezando os demais. A aproximação x_R é formada pelo sinal de x e pelos dígitos assim obtidos.

Como $x_R \in F_i$, então, por (1.70), podemos afirmar que só podem ser aproximados por elementos de F_i os números x pertencentes ao intervalo

$$(2.6) \quad R_i = \left(-1 + \frac{1}{2} N^{-t}, \quad 1 - \frac{1}{2} N^{-t} \right)$$

Referir-nos-emos a R_i como o conjunto dos números aproximáveis por elementos de F_i , quando é utilizado o arredondamento propriamente dito.

(2.7.) Exemplos. Seja $F_i = F(10, 4)$

$$x = 0,00157469$$

$$x_R = 0,0016$$

$$y = -0,896512236$$

$$y_R = -0,8965$$

$$z = 0,99999$$

z não pode ser aproximado por um elemento de F_i .

$$t = 0,682351$$

$$t_R = 0,6824$$

$$p = -0,1999782$$

$$p_R = -0,2000$$

Determinação de x_c

Para obter x_c pegamos a expansão infinita de x (aquela terminando em zeros, no caso de escolha) e retemos os t primeiros dígitos após a vírgula, desprezando os demais.

É claro que o conjunto dos números aproximáveis por elementos de F_i , quando é utilizado o arredondamento por corte, difere daquele apresentado em (2.6.). Indicaremos por C_i esse conjunto, que é definido por

$$(2.8) \quad C_i = (-1, +1)$$

(2.9) Exemplos - Consideremos o arredondamento por corte dos números utilizados em (2.7):

$$x_c = 0,0015$$

$$t_c = 0,6823$$

$$y_c = -0,8965$$

$$p_c = -0,1999$$

$$z_c = 0,9999$$

2.3.2. Delimitações para os erros de arredondamento

Como observado em 1.7.1., os elementos de F_i são representados, na reta real, por pontos equidistantes. Em virtude de ser N^{-t} a distância entre dois pontos consecutivos desse conjunto, facilmente se deduz de (2.4) e (2.5) que:

$$(2.10.) \quad |x_R - x| \leq \frac{1}{2} N^{-t}$$

$$(2.11.) \quad |x_c - x| \leq N^{-t}$$

2.3.3. Erros de arredondamento nas operações aritméticas

Um número real x pode ser aproximado, dependendo do computador considerado, por x_R ou x_c . Em virtude das correspondentes delimitações para os erros de arredondamento

sejam diferentes, adotaremos um procedimento que permita, em ambos os casos, a utilização dos resultados contidos neste trabalho.

A representação interna de um número real x , num computador que utilize sistema em vírgula fixa, será indicada por $f_i(x)$. Assim sendo, de acordo com o tipo de arredondamento efetuado no computador considerado, temos:

$$(2.12) \quad f_i(x) = x_R \quad \text{ou} \quad f_i(x) = x_C$$

A expressão $z = f_i(x * y)$ indicará que x , y e z são representações internas de números reais, no computador considerado, e que z é obtido da seguinte maneira: a operação $*$ é efetuada e, quando necessário, o resultado $x * y$ obtido é arredondado, de acordo com o procedimento próprio do computador em questão. Quando $*$ representa a divisão, y é suposto diferente de zero. Portanto, na realidade, o computador executa, em substituição à operação $*$, uma operação aproximada, definida por:

$$(2.13) \quad f_i(x * y) = (x * y)_R \quad \text{ou} \quad f_i(x * y) = (x * y)_C,$$

de acordo com o procedimento de arredondamento do computador em questão. Como observado em 2.1., $f_i(x * y)$ é chamado valor aproximado de $(x * y)$.

O conjunto dos números aproximáveis por elementos de F_i será representado por A_i , ou seja,

$$(2.14) \quad A_i = R_i \quad \text{ou} \quad A_i = C_i,$$

de acordo com o tipo de arredondamento do computador utilizado. (R_i e C_i são definidos em (2.6) e (2.8), respectivamente).

Denominaremos unidade absoluta de arredondamento, e indicaremos por μ_A ,

$$(2.15) \quad \mu_A = \frac{1}{2} N^{-t} \quad \text{ou} \quad \mu_A = N^{-t},$$

conforme o computador considerado utilize arredondamento propriamente dito ou por corte, respectivamente.

Descreveremos a realização das operações aritméticas fundamentais e da multiplicação escalar em um computador que utiliza sistema em vírgula fixa. Para a descrição dessas operações será usado um sistema genérico $F_i = F(N, t)$ e para os exemplos serão utilizados $F_i = F(10, 4)$ e arredondamento propriamente dito.

Adição e subtração

Sejam x e y dois elementos quaisquer de F_i . A determinação de $f_i(x \pm y)$ não é sempre possível porque pode ocorrer $(x \pm y) \notin A_i$. Quando realizáveis, a adição e subtração não envolvem erros de arredondamento e satisfazem:

$$(2.16) \quad f_i(x \pm y) = x \pm y$$

Exemplos

$$x = 0,1542$$

$$y = 0,6321$$

$$f_i(x + y) = 0,7863$$

$$r = 0,8314$$

$$s = 0,2765$$

$f_i(r + s)$ não pode ser determinado no sistema considerado.

Multiplicação

Sejam x e y elementos quaisquer de F_i . O computador determina o produto $x \times y$, que geralmente é um número $2t$ -digital e sempre pertence a A_i , e a seguir o arredonda para t dígitos. Portanto, $f_i(x \times y)$ satisfaz:

$$(2.17) \quad f_i(x \times y) = x y + \epsilon, \quad |\epsilon| \leq \mu_A,$$

onde $x y$ representa o produto exato de x e y .

Exemplo

$$x = 0,1832$$

$$y = 0,4327$$

$$x y = 0,07927064$$

$$f_i(x \times y) = 0,0793$$

$$|\epsilon| = 0,2936 \times 10^{-4}$$

Divisão

Sejam x e y elementos quaisquer de F_i . A divisão $x : y$ somente pode ser efetuada quando $|y| > |x|$; quando isso ocorre, o computador determina o quociente $\frac{x}{y}$, que geralmente possui mais de t dígitos, e o arredonda. Portanto, quando realizável, a divisão em vírgula fixa satisfaz a:

$$(2.18) \quad fi(x:y) = \frac{x}{y} + \varepsilon, \quad |\varepsilon| \leq u_A,$$

onde $\frac{x}{y}$ representa o quociente exato de x e y .

Exemplo

$$\begin{aligned} x &= 0,1234 & y &= 0,8976 \\ \frac{x}{y} &= 0,13747731\dots & fi(x:y) &= 0,1375 & |\varepsilon| &= 0,2228\dots \times 10^{-4} \end{aligned}$$

Multiplicação escalar

A determinação do produto escalar $\langle x, y \rangle$, definido em (1.2), é considerada neste trabalho por ser esse cálculo muito comum nos problemas que envolvem matrizes. O valor $fi(\langle x, y \rangle)$ pode ser obtido através de dois processos diferentes, que originam erros de arredondamento finais distintos. O primeiro consiste na determinação dos produtos $x_i y_i$, arredondamento para t dígitos dos valores obtidos e posterior soma desses valores arredondados. O segundo processo consta da determinação dos produtos $x_i y_i$ em precisão dupla, soma desses valores $2t$ -digitais e posterior arredondamento da soma obtida para t dígitos. O produto escalar fornecido pelo segundo método é usualmente representado por $fi_2(\langle x, y \rangle)$, notação que o diferencia do valor fornecido pelo primeiro método, que é representado por $fi(\langle x, y \rangle)$. Costuma-se dizer que no cálculo de $fi_2(\langle x, y \rangle)$ há acumulação de produtos escalares. Os produtos escalares fornecidos pelos métodos descritos verificam:

$$(2.19) \quad f_i(\langle x, y \rangle) = \sum_{i=1}^n x_i y_i + \varepsilon, \quad |\varepsilon| \leq n \mu_A$$

$$(2.20) \quad f_2(\langle x, y \rangle) = \sum_{i=1}^n x_i y_i + \varepsilon, \quad |\varepsilon| \leq \mu_A$$

Exemplos - Consideremos a determinação da multiplicação escalar:

$$(2.21) \quad \langle x, y \rangle = 0,7865 \times 0,4361 + 0,0053 \times 0,2410 + 0,4126 \times 0,6325$$

Armazenando os produtos parciais na forma 2t-digital e adicionando-os obtemos:

$$(2.22) \quad \langle x, y \rangle = 0,34299265 + 0,00127730 + 0,26096950 = 0,60523945$$

Arredondando para t dígitos o resultado obtido em (2.22) temos:

$$f_2(\langle x, y \rangle) = 0,6052, \quad \text{com } |\varepsilon| = 0,3945 \times 10^{-4}$$

Arredondando separadamente cada um dos produtos parciais e somando-os, obtemos:

$$(2.23) \quad f_i(\langle x, y \rangle) = 0,3430 + 0,0013 + 0,2610 = 0,6053, \\ \text{com } |\varepsilon| = 0,6055 \times 10^{-4}, \quad \text{como pode ser concluído da comparação de (2.23) e (2.22)}$$

Como pode ser observado em (2.19) e (2.20), a determinação de $\langle x, y \rangle$ é efetuada de maneira mais satisfatória em computadores que possibilitam a acumulação de produtos escalares e que, portanto, possuem acumulador de precisão dupla. Nesses computadores a divisão também é efetuada, usualmente, de forma mais precisa do que nos computadores que só possuem acumulador de precisão simples, através do armazenamento do dividendo.

em precisão dupla e do divisor em precisão simples, o que permite a obtenção de um quociente corretamente arredondado; os dividendos t -digitais são estendidos à forma $2t$ -digital pela adição de zeros. A utilização do acumulador de precisão dupla para a determinação de produtos escalares e quocientes é bastante vantajosa no cálculo de expressões do tipo:

$$(2.24) \quad \frac{\sum_{i=1}^n x_i y_i}{z},$$

cujó numerador pode ser determinado exatamente, originando um valor $2t$ -digital que é dividido por z . O arredondamento para t dígitos, desse quociente, fornece:

$$(2.25) \quad \left\{ \begin{array}{l} f_z \left(\frac{\sum_{i=1}^n x_i y_i}{z} \right) = \frac{\sum_{i=1}^n x_i y_i}{z} + \varepsilon, \\ \text{com } |\varepsilon| \leq \mu_A \end{array} \right.,$$

enquanto que em computadores que dispõem apenas de acumulador de precisão simples o cálculo de (2.24) fornece:

$$(2.26) \quad \left\{ \begin{array}{l} f_z \left(\frac{\sum_{i=1}^n x_i y_i}{z} \right) = \frac{\sum_{i=1}^n x_i y_i}{z} + \varepsilon', \\ \text{com } |\varepsilon'| \leq \left(\frac{n}{z} + 1 \right) \mu_A \end{array} \right.,$$

Como pode ser observado em (2.25) e (2.26), a expressão (2.24) é calculada de forma mais precisa em computadores que dispõem de ambos os recursos citados.

2.4. - Arredondamento na aritmética em vírgula flutuante

2.4.1. - Determinação de x_R e x_C

Seja $x \notin F_L$ o número que queremos aproximar por um elemento de $F_L = F(N, t, p, P)$. De acordo com o computador utilizado, x é aproximado por x_R ou x_C .

Determinação de x_R

Se $x = 0$, tomamos $x_R = 0, 0 \dots 0 \times N^{-p}$

Se $x \neq 0$, escolhemos a real e b inteiro tais que

$$(2.27) \quad |x| = a \times N^b, \text{ com } N^{-1} \left(1 - \frac{1}{2} N^{-t}\right) \leq a < \left(1 - \frac{1}{2} N^{-t}\right).$$

Se b não pertence ao intervalo $-p \leq b \leq P$, então, x não pode ser representado no sistema considerado. Em caso contrário, tomamos a expansão infinita de $a + \frac{1}{2} N^{-t}$ (aquela terminando em zeros no caso de escolha) e retemos apenas os t primeiros dígitos após a vírgula, desprezando os demais. A aproximação x_R é constituída pelo sinal de x , pelos t dígitos assim obtidos e pelo expoente b determinado em (2.27).

O conjunto dos números aproximáveis por elementos de F_L , quando é utilizado o arredondamento propriamente dito, é:

$$(2.28) \quad R_L = (-N^P, N^P)$$

(2.29) Exemplos - Seja $F_L = F(10, 4, 18, 19)$

$$p = 0,163173 \times 10^5$$

$$q = 0,000697412$$

$$x = 1308,564 \times 10^7$$

$$y = 9,99999 \dots$$

$$z = 99,99999 \times 10^{18}$$

$$p_R = 0,1632 \times 10^5$$

$$q_R = 0,6974 \times 10^{-3}$$

$$x_R = 0,1309 \times 10^{11}$$

$$y_R = 0,1000 \times 10^2$$

z_R não pode ser representado no sistema em questão.

Determinação de x_c Se $x=0$, $x_c = 0,0 \dots 0 \times N^{-p}$ Se $x \neq 0$, escolhemos a real e b inteiro tais que

$$(2.30) \quad |x| = a \times N^b, \quad \text{com} \quad N^{-1} \leq |a| < 1$$

Se $b \notin [-p, p]$, então, x não pode ser representado no sistema em questão. Em caso contrário, tomamos a expansão infinita de a na base N (aquela terminando em zeros, no caso de escolha) e retemos os t primeiros dígitos após a vírgula, desprezando os demais. A aproximação x_c é formada pelo sinal de x , pelos t dígitos assim obtidos e pelo expoente b determinado em (2.30). O conjunto C_2 , dos números aproximáveis por elementos de F_2 , coincide com R_2 , de finido em (2.28)

(2.31) Exemplos - Consideremos o arredondamento por corte dos números utilizados em (2.29)

$$p_c = 0,1631 \times 10^5$$

$$y_c = 0,9999 \times 10$$

$$q_c = 0,6974 \times 10^{-3}$$

z_c não pode ser representado no sistema em questão.

$$x_c = 0,1308 \times 10^{11}$$

2.4.2-1 Delimitações para os erros de arredondamento

O caso $x=0$ é trivial. Seja $x \neq 0$ e, sem perda de generalidade, suponhamos $x > 0$.

Representemos x na forma

$$(2.32) \quad x = a \times N^b, \quad N^{-1} \leq a < 1.$$

Em virtude de (2.32) podemos afirmar

$$(2.33) \quad N^{b-1} \leq x < N^b$$

No intervalo $[N^{b-1}, N^b]$, definido em (2.33), os elementos de $F_2(N, t, p, P)$ são representados por pontos equidistantes. Em consequência de ser N^{b-t} a distân-

cia entre dois pontos consecutivos desse conjunto, podemos afirmar, por (2.4) e (2.5),

$$(2.34) \quad |x_R - x| \leq \frac{1}{2} N^{b-t}$$

$$(2.35) \quad |x_C - x| \leq N^{b-t}$$

Indicando por \mathcal{E}_R (\mathcal{E}'_R) e \mathcal{E}_C (\mathcal{E}'_C) os erros relativos definidos em relação ao valor verdadeiro (aproximado), cometidos pela substituição de x por x_R e x_C , respectivamente, podemos escrever, por (2.33) - (2.35):

$$(2.36) \quad |\mathcal{E}_R| = \frac{|x_R - x|}{|x|} \leq \frac{\frac{1}{2} N^{b-t}}{N^{b-1}} = \frac{1}{2} N^{1-t}$$

$$(2.37) \quad |\mathcal{E}_C| = \frac{|x_C - x|}{|x|} \leq \frac{N^{b-t}}{N^{b-1}} = N^{1-t}$$

$$(2.38) \quad |\mathcal{E}'_R| = \frac{|x_R - x|}{|x_R|} \leq \frac{\frac{1}{2} N^{b-t}}{N^{b-1}} = \frac{1}{2} N^{1-t}$$

$$(2.39) \quad |\mathcal{E}'_C| = \frac{|x_C - x|}{|x_C|} \leq \frac{N^{b-t}}{N^{b-1}} = N^{1-t}$$

Utilizando as relações (2.36) - (2.39), obtemos:

$$(2.40) \quad x_R = x(1 + \mathcal{E}_R) = \frac{x}{1 + \mathcal{E}'_R}; \quad |\mathcal{E}_R|, |\mathcal{E}'_R| \leq \frac{1}{2} N^{1-t}$$

$$(2.41) \quad x_C = x(1 + \mathcal{E}_C) = \frac{x}{1 + \mathcal{E}'_C}; \quad |\mathcal{E}_C|, |\mathcal{E}'_C| \leq N^{1-t}$$

2.4.3. - Erros de arredondamento nas operações aritméticas

Procedendo de maneira análoga à utilizada em 2.3.3., adotaremos um procedimento que permita o uso dos resultados contidos neste trabalho para computadores que utili-

zam qualquer um dos tipos de arredondamento considerados. Indicaremos:

$$(2.42) \quad fl(x) = x_R \quad \text{ou} \quad fl(x) = x_C$$

$$(2.43) \quad fl(x * y) = (x * y)_R \quad \text{ou} \quad fl(x * y) = (x * y)_C$$

$$(2.44) \quad A_\ell = R_\ell \quad \text{ou} \quad A_\ell = C_\ell$$

de acordo com o tipo de arredondamento do computador em questão.

Em expressões do tipo $fl(x * y * z * w)$, suporemos que as operações sejam efetuadas da esquerda para a direita.

Denominaremos unidade absoluta e relativa de arredondamento, para computadores que utilizam sistema normalizado de representação em vírgula flutuante em que $F_\ell = F(N, t, p, P)$, aos valores μ_A e μ definidos, respectivamente, por

$$(2.45) \quad \mu_A = \frac{1}{2} N^{b-t} \quad \text{ou} \quad \mu_A = N^{b-t}$$

$$(2.46) \quad \mu = \frac{1}{2} N^{1-t} \quad \text{ou} \quad \mu = N^{1-t}$$

conforme o tipo de arredondamento utilizado pelo computador considerado.

Na descrição das operações aritméticas, efetuada para um sistema $F_\ell = F(N, t, p, P)$, os operandos serão representados por $x_1, x_2, \dots, x_i, \dots$, sendo $x_i = a_i N^{b_i}$, e o resultado por $x = a N^b$. Nos exemplos, realizados em um sistema $F_\ell = F(10, 4, -35, 38)$ que utiliza arredondamento propriamente dito, as mantissas são indicadas entre parênteses e os erros de arredondamento escritos na forma $|E| = x \cdot 10^{-3}$, em que x é um número real qualquer.

Adotamos essa forma de representação para tornar imediata a comparação entre o erro obtido e a unidade relativa de arredondamento do sistema considerado, $\frac{1}{2} 10^{-3}$. Chamaremos deslocamento da mantissa à direita (esquerda)

ao deslocamento ocasionado pela multiplicação da mantissa por uma potência negativa (positiva) da base N . Descreveremos a realização das operações aritméticas fundamentais e da multiplicação escalar, em vírgula flutuante, para computadores que não dispõem de acumulador de precisão dupla e para computadores que possuem tal acumulador; nesse último caso, também será considerada a acumulação de somas e produtos escalares.

Adição e subtração em computadores que dispõem de acumulador de precisão dupla.

Consideremos a determinação de $x = x_1 + x_2$.

Seja $b_1 \geq b_2$. O computador efetua, inicialmente,

$$x = x_1 + x_2 = N^{b_1} (a_1 + N^{b_2-b_1} a_2) ,$$

ou seja, desloca a_2 de $(b_1 - b_2)$ posições à direita e, a seguir, efetua a adição em precisão dupla. Quando necessário, a soma obtida é multiplicada por uma potência conveniente de N , através de um deslocamento à esquerda ou à direita, de modo que sua mantissa satisfaça (1.75); nesse caso, o expoente b_1 é ajustado de maneira a concordar com o deslocamento efetuado. Essa última operação é denominada, usualmente, normalização da mantissa. Finalmente, o computador arredonda, para t dígitos, a parte fracionária assim obtida.

O número máximo de deslocamentos à direita, que pode tornar-se necessário, é um, pois, por (1.75):

$$|a_1 + N^{b_2-b_1} a_2| \leq |a_1| + N^{b_2-b_1} |a_2| < 2$$

Entretanto, em virtude de cancelamentos, podem ser necessários até $(t-1)$ deslocamentos à esquerda, como ocorre no exemplo apresentado em (2.53).

Se $b_1 - b_2 > t$, então, x_2 não contribui pa-

na soma aproximada, por não ter influência em seus primeiros t dígitos, e ocorre $fl(x_1 + x_2) = x_1$.

Se $x_1 = 0$, então, $fl(x_1 + x_2) = x_2$. Analogamente, se $x_2 = 0$, $fl(x_1 + x_2) = x_1$. Em ambos os casos o erro de arredondamento é nulo.

Como no cálculo de $fl(x_1 + x_2)$ é efetuada somente o arredondamento, para t dígitos, da soma exata normalizada, verifica-se:

$$(2.47) \quad |fl(x_1 + x_2) - (x_1 + x_2)| \leq \mu_A$$

Essa forma de erro também é utilizada, porém, de maneira usual, na aritmética em vírgula flutuante é utilizado o erro de arredondamento relativo dado por:

$$(2.48) \quad fl(x_1 + x_2) = (x_1 + x_2)(1 + \varepsilon), \quad |\varepsilon| \leq \mu$$

Procedendo com a subtração de forma análoga à utilizada no tratamento da adição, pode ser encontrado o resultado

$$(2.49) \quad fl(x_1 - x_2) = (x_1 - x_2)(1 + \varepsilon), \quad |\varepsilon| \leq \mu$$

Exemplos

$$(2.50) \quad fl(10^{-3}(0,1572) + 10^4(0,9813)) = 10^4(0,9813)$$

$fl(x_1 + x_2) = x_1$ porque $b_1 - b_2 = 7$ e $7 > 4$.
 $|\varepsilon| \approx 0,00002 \times 10^{-3}$

$$(2.51) \quad fl(10^6(0,7513) + 10^3(0,4624)) = fl(10^6(0,75130000 + 0,00046240)) =$$

$$= fl(10^6(0,75176240)) = 10^6(0,7518)$$

$|\varepsilon| \approx 0,05002 \times 10^{-3}$

$$(2.52) \quad fl(10^3(0,3561) + 10^3(0,8225)) = fl(10^3(0,35610000 + 0,82250000)) =$$

$$= fl(10^3(1,17860000)) = fl(10^4(0,11786000)) = 10^4(0,1179)$$

$|\varepsilon| \approx 0,33939 \times 10^{-3}$.

$$(2.53) \quad \begin{aligned} fl(10^{-5}(0,4318) + 10^{-5}(-0,4316)) &= fl(10^{-5}(0,43180000 - 0,43160000)) = \\ &= fl(10^{-5}(0,00020000)) = fl(10^{-8}(0,20000000)) = 10^{-8}(0,2000). \\ |E| &= 0 \end{aligned}$$

$$(2.54) \quad \begin{aligned} fl(10^4(0,1001) + 10^3(-0,9941)) &= fl(10^4(0,10010000 - 0,09941000)) = \\ &= fl(10^4(0,00069000)) = fl(10(0,69000000)) = 10(0,6900). \\ |E| &= 0 \end{aligned}$$

$$(2.55) \quad \begin{aligned} fl(10^3(0,9652) + 10^2(0,5828)) &= fl(10^3(0,96520000 + 0,05828000)) = \\ &= fl(10^3(1,02348000)) = fl(10^4(0,10234800)) = 10^4(0,1023). \\ |E| &\approx 0,46899 \times 10^{-3} \end{aligned}$$

Como pode ser observado através dos exemplos (2.53) e (2.54), em adições efetuadas em computadores que dispõem de acumulador de precisão dupla, a ocorrência de cancelamento dos dígitos localizados imediatamente à direita da vírgula implica na nulidade do erro de arredondamento. Em computadores que não possuem tal acumulador esse fato nem sempre se verifica, como poderá ser notado em (2.96).

Consideraremos, a seguir, adições de mais de duas parcelas. No estudo dessas, assumi como, de maneira geral, de outras operações com mais de dois operandos, frequentemente são encontradas expressões do tipo

$$(2.56) \quad \prod_{i=1}^m (1 + \delta_i) = (1 + \delta_1)(1 + \delta_2) \cdots (1 + \delta_m), \quad |\delta_i| \leq u$$

em que os δ_i representam erros de arredondamento. Procurando simplificar tais expressões, aplicaremos (1.67) as mesmas, obtendo, em seu lugar:

$$(2.57) \quad \prod_{i=1}^m (1 + \delta_i) = 1 + 1,01 m \theta u, \quad |\theta| \leq 1.$$

A hipótese $mu \leq 0,01$, que torna possível a aplicação de (1.67), certamente se verifica nos proble-

mas em que (2.56) será aplicada, já que nesses, de maneira geral, n e μ representam, respectivamente, a ordem de uma matriz e a unidade relativa de arredondamento. Por esse motivo, por simplicidade, sempre que utilizarmos (2.57) não faremos referências a essa hipótese, supondo-a subentendida.

Descreveremos a determinação de

$$(2.58) \quad s_n = fl(x_1 + x_2 + \dots + x_n) ; x_i \in F_L ; i=1, 2, \dots, n.$$

O computador efetua um processo de "adições encaixantes" que, para $n=4$, é representado por

$$fl(x_1 + x_2 + x_3 + x_4) = fl(fl(fl(x_1 + x_2) + x_3) + x_4).$$

No caso genérico, o computador determina, sequencialmente, as quantidades s_1, s_2, \dots, s_n , definidas recursivamente por:

$$(2.59) \quad \begin{cases} s_1 = x_1 \\ s_i = fl(s_{i-1} + x_i) = (s_{i-1} + x_i)(1 + \delta_i^r) ; |\delta_i^r| \leq \mu ; i=2, \dots, n \end{cases}$$

Desenvolvendo as relações (2.59) obtemos:

$$(2.60) \quad s_n = fl(x_1 + x_2 + \dots + x_n) = x_1(1 + \varepsilon_1) + x_2(1 + \varepsilon_2) + \dots + x_n(1 + \varepsilon_n)$$

$$(2.61) \quad \begin{cases} 1 + \varepsilon_1 = (1 + \delta_2^r)(1 + \delta_3^r) \dots (1 + \delta_n^r) \\ 1 + \varepsilon_i = (1 + \delta_i^r)(1 + \delta_{i+1}^r) \dots (1 + \delta_n^r) ; i=2, \dots, n \end{cases}$$

Aplicando (2.57) a (2.61),

$$(2.62) \quad \begin{cases} 1 + \varepsilon_1 = 1 + 1,01(n-1)\theta_1 \mu ; |\theta_1| \leq 1 \\ 1 + \varepsilon_i = 1 + 1,01(n-i+1)\theta_i \mu ; |\theta_i| \leq 1 ; i=2, \dots, n \end{cases}$$

Com a finalidade de evitar delimitações diferentes para ε_1 e ε_i ($i=2, \dots, n$), podemos escrever (2.62) na forma:

$$(2.63) \quad 1 + \varepsilon_i = 1 + 1,01(n-i+1)\theta_i \mu ; |\theta_i| \leq 1 ; i=1, 2, \dots, n,$$

onde, para $n=1$, a delimitação foi aumentada de $1,01(n-1)\theta_1 \mu$

para $1,01 n \theta_1 u$.

Embora o erro relativo originado no cálculo de $fl(x_1+x_2)$, em computadores que possuem acumulador de precisão dupla, seja sempre pequeno, não podemos afirmar que, necessariamente, o mesmo ocorre na determinação de $fl(x_1+x_2+\dots+x_n)$. Para entender esse fato, basta considerar o erro relativo

$$(2.64) \quad |E| = \left| \frac{fl(x_1+\dots+x_n) - (x_1+\dots+x_n)}{(x_1+\dots+x_n)} \right| = \left| \frac{fl(x_1+\dots+x_n)}{(x_1+\dots+x_n)} - 1 \right|$$

Pode ocorrer que a soma aproximada seja nula, sem que a soma exata o seja, ou vice-versa, razão de nem sempre existir uma boa delimitação para (2.64). Assim, por exemplo, quando a soma em vírgula móvel é nula, sem que o mesmo ocorra com a soma exata, obtemos $|E|=1$ e, portanto, um erro de arredondamento que, em valor absoluto, é igual à soma exata.

Apesar das delimitações para os E_i não dependerem da ordem em que a adição se efetue, a delimitação da soma depende de tal ordem. Essa delimitação é menor quando a adição é efetuada segundo a ordem crescente dos valores absolutos das parcelas porque, dessa maneira, a maior delimitação dos E_i é associada ao menor x_i .

Exemplos

$$(2.65) \quad 10^3(0,1418) + 10^1(0,5368) + 10^2(0,9465) + 10^0(0,5812)$$

A soma exata dos números considerados em (2.65) é $10^3(0,2423992)$, cujo arredondamento para 4 dígitos fornece $10^3(0,2424)$, com $|E|=0,0033 \times 10^{-3}$.

Efetuada (2.65) na ordem dada obtemos:

$$s_2 = fl(10^3(0,1418) + 10(0,5368)) = fl(10^3(0,14716800))$$

$$s_2 = 10^3(0,1472) \quad , \quad |E_2| \approx 0,21744 \times 10^{-3}$$

$$s_3 = fl(10^3(0,1472) + 10^2(0,9465)) = fl(10^3(0,24185000))$$

$$s_3 = 10^3(0,2419) \quad ; \quad |\delta_3| = 0,20674 \times 10^{-3}$$

$$s_4 = fl(10^3(0,2419) + 10^0(0,5812)) = fl(10^3(0,24248120))$$

$$s_4 = 10^3(0,2425) \quad ; \quad |\delta_4| \approx 0,077532 \times 10^{-3}$$

do que resulta uma soma aproximada igual a $10^3(0,2425)$ com $|\mathcal{E}| \approx 0,41584 \times 10^{-3}$.

Efetuada (2.65) segundo a ordem crescente dos valores absolutos das parcelas temos:

$$s'_2 = fl(10^0(0,5812) + 10^1(0,5368)) = fl(10^1(0,59492))$$

$$s'_2 = 10^1(0,5949) \quad ; \quad |\delta'_2| \approx 0,033618 \times 10^{-3}$$

$$s'_3 = fl(10^1(0,5949) + 10^2(0,9465)) = fl(10^3(0,100599))$$

$$s'_3 = 10^3(0,1006) \quad ; \quad |\delta'_3| = 0,009940 \times 10^{-3}$$

$$s'_4 = fl(10^3(0,1006) + 10^3(0,1418)) = fl(10^3(0,2424))$$

$$s'_4 = 10^3(0,2424) \quad ; \quad |\delta'_4| = 0$$

A soma aproximada obtida nesse caso é $10^3(0,2424)$, com $|\mathcal{E}| = 0,0033 \times 10^{-3}$.

$$(2.66) \quad 10^{-4}(0,4863) + 10^0(0,1436) + 10^{-4}(0,3985) + 10^{-4}(0,4599)$$

$$\text{Soma exata} = 0,14373447$$

$$\text{arredondamento soma exata} = 0,1437 \quad ; \quad |\mathcal{E}| \approx 0,23982 \times 10^{-3}$$

Efetuada (2.66) na ordem dada obtemos:

$$s_2 = 10^0(0,1436) \quad ; \quad |\delta_2| \approx 0,33853 \times 10^{-3}$$

$$s_3 = 10^0(0,1436) \quad ; \quad |\delta_3| \approx 0,27743 \times 10^{-3}$$

$$s_4 = 10^0(0,1436) \quad ; \quad |\delta_4| \approx 0,32016 \times 10^{-3}$$

e a soma aproximada $10^0(0,1436)$ tem erro relativo $\approx 0,93555 \times 10^{-3}$.

Quando é utilizada a ordem crescente dos valores absolutos das parcelas temos:

$$s'_2 = 10^{-4}(0,8584) \quad ; \quad |\delta'_2| = 0$$

$$s'_3 = 10^{-3}(0,1345) \quad ; \quad |\delta'_3| \approx 0,2231 \times 10^{-3}$$

$$s'_4 = 10^0(0,1437) \quad ; \quad |\delta'_4| \approx 0,24003 \times 10^{-3}$$

de modo que a soma aproximada, 0,1437, tem erro relativo $\approx 0,23982 \times 10^{-3}$.

A escolha da ordem crescente dos valores absolutos das parcelas é particularmente vantajosa quando algumas parcelas têm valor absoluto muito menor que outras, como ocorre em (2.66).

$$(2.67) \quad 10^3(0,1084) + 10^2(-0,9712) + 10(-0,9537) + 10(-0,5014)$$

$$\text{soma exata} = 10(-0,3271)$$

Efetuada (2.67) na ordem dada temos:

$$s_2 = 10^2(0,1128) \quad ; \quad |\delta_2| = 0$$

$$s_3 = 10(0,1743) \quad ; \quad |\delta_3| = 0$$

$$s_4 = 10(-0,3271) \quad ; \quad |\delta_4| = 0 \quad ,$$

de modo que a soma obtida é exata.

Utilizando a ordem crescente dos valores absolutos das parcelas:

$$s'_2 = 10^2(-0,1455) \quad ; \quad |\delta'_2| = 0,068724 \times 10^{-3}$$

$$s'_3 = 10^3(-0,1117) \quad ; \quad |\delta'_3| = 0,268649 \times 10^{-3}$$

$$s'_4 = 10(-0,3300) \quad ; \quad |\delta'_4| = 0 \quad ,$$

e a soma aproximada obtida, $10(-0,3300)$ tem $|E| = 9,86579 \times 10^{-3}$.

Embora forneça sempre a menor delimitação para o erro, a escolha da ordem crescente dos valores absolutos das parcelas não fornece, necessariamente, o menor erro, como pode ser observado em (2.67). É importante, porém, que se observe que esse exemplo é muito especial, tendo sido construído para evidenciar tal fato; o caso geral é dado por (2.65) e (2.66).

Multiplicação em computadores que dispõem de acumulador de precisão dupla.

Dados $x_1, x_2 \in F_t$, o produto aproximado $fl(x_1 x_2)$ é obtido da seguinte maneira: os expoentes b_1 e b_2 são adicionados, originando b , e o produto $a_1 a_2$ é realizado em precisão dupla; quando necessário, esse produto é normalizado, e o expoente b ajustado de maneira conveniente. Finalmente, a parte fracionária é arredondada para t dígitos.

Quando requerida, a normalização se efetua através de um único deslocamento à esquerda, pois, por (1.75)

$$N^{-2} \leq |a_1 a_2| < 1$$

Se x_1 ou x_2 (ou ambos) são nulos, então, é atribuído o valor zero para $fl(x_1 \times x_2)$.

Assim como ocorre na determinação de $fl(x_1 \pm x_2)$, somente um erro de arredondamento é cometido no cálculo de $fl(x_1 \times x_2)$; aquele originado pelo arredondamento, para t dígitos, do produto exato. Portanto:

$$(2.68) \quad fl(x_1 \times x_2) = x_1 x_2 (1 + \epsilon), \quad |\epsilon| \leq u,$$

em que $x_1 x_2$ representa o produto exato de x_1 e x_2 .

Exemplos

$$(2.69) \quad fl(10^{-3}(0,7823) \times 10^5(0,5432)) = fl(10^2(0,42494536)) = 10^2(0,4249) \quad ; \quad |\epsilon| \approx 0,10674 \times 10^{-3}$$

$$(2.70) \quad fl(10^{-2}(0,2351) \times 10^4(0,1634)) = fl(10^{-6}(0,03841534)) = fl(10^{-7}(0,38415340)) = 10^{-7}(0,3842) \quad ; \quad |\epsilon| \approx 0,12131 \times 10^{-3}$$

Consideraremos, a seguir, multiplicações com mais de dois fatores. Descreveremos a determinação de

$$(2.71) \quad p_m = fl(x_1 \times x_2 \times \dots \times x_m) \quad ; \quad x_i \in F_t \quad ; \quad i = 1, 2, \dots, m$$

Assim como ocorre com adições de mais de duas parcelas, o computador realiza um processo de "multiplicações encaixantes" que, para $n=4$ é representado por:

$$fl(x_1 \times x_2 \times \dots \times x_m) = fl(fl(fl(x_1 \times x_2) \times x_3) \times x_4)$$

No caso genérico, o computador determina, sucessivamente, os valores p_1, p_2, \dots, p_m , definidos recursivamente por:

$$(2.72) \quad \begin{cases} p_1 = x_1 \\ p_m = fl(p_{i-1} \times x_i) = p_{i-1} x_i (1 + \delta_i) \quad ; \quad |\delta_i| \leq u \quad ; \quad i = 2, \dots, m \end{cases}$$

Desenvolvendo as relações dadas em (2.72) obtemos:

$$(2.73) \quad p_m = fl(x_1 \times x_2 \times \dots \times x_m) = x_1 x_2 \dots x_m (1 + \varepsilon),$$

$$(2.74) \quad 1 + \varepsilon = (1 + \delta_2)(1 + \delta_3) \dots (1 + \delta_m); \quad |\delta_i| \leq \mu, \quad i = 2, \dots, m$$

Aplicando (2.57) a (2.74), temos:

$$(2.75) \quad 1 + \varepsilon = 1 + 1,01(n-1)\theta\mu, \quad |\theta| \leq 1$$

De maneira geral, o produto aproximado tem erro relativo pequeno e, embora o seu valor dependa da ordem em que as multiplicações se efetuam, a delimitação (2.75) não depende da mesma.

Divisão em computadores que possuem acumulador de precisão dupla.

Enquanto na aritmética em vírgula fixa só podem ser efetuadas divisões $x_1 : x_2$ para as quais $x_1 < x_2$, na aritmética em vírgula flutuante a única divisão que não pode ser realizada é aquela para a qual $x_2 = 0$.

O quociente de x_1 por x_2 é determinado da maneira que se segue. O expoente b , do quociente, é tomado igual à diferença $b_1 - b_2$ e a_1 é colocado nos t dígitos mais significativos do acumulador. Se $|a_1| > |a_2|$, então, o número armazenado no acumulador é deslocado uma posição à direita e b aumentado de uma unidade. A seguir, o número do acumulador é dividido por a_2 e o quociente, que já deve satisfazer a (2.75), é arredondado para t dígitos.

Em virtude de no cálculo de $fl(x_1 : x_2)$ só ocorrer um arredondamento, o do quociente exato para t dígitos:

$$(2.76) \quad fl(x_1 : x_2) = \frac{x_1}{x_2} (1 + \varepsilon); \quad |\varepsilon| \leq \mu; \quad x_2 \neq 0,$$

onde $\frac{x_1}{x_2}$ representa o quociente exato de x_1 por x_2 .

Exemplos

$$(2.77) \quad fl(10^5(0,5632) : 10^{-3}(0,1827)) = fl(10^{-2}(0,56320000 : 0,1827)) =$$

$$= fl(10^{-1}(0,05632000 : 0,1827)) = fl(10^{-1}(0,30826491\dots)) = \\ = 10^{-1}(0,3083) \quad ; \quad |\varepsilon| \approx 0,11381 \times 10^{-3}.$$

$$(2.78) \quad fl(10^6(0,2643) : 10^3(0,9741)) = fl(10^3(0,26430000 : 0,9741)) = \\ = fl(10^3(0,27132737\dots)) = 10^3(0,2713) \\ |\varepsilon| \approx 0,10091 \times 10^{-3}.$$

A igualdade (2.76) pode ser escrita

$$(2.79) \quad fl(x_1 : x_2) = x_1 x_2^{-1} (1 + \varepsilon) ; \quad |\varepsilon| \leq u ; \quad x_2 \neq 0.$$

Em virtude da aplicação de (2.68) a $fl(x_1 x_2^{-1}) = fl(x_1 : x_2)$ também fornecer o resultado (2.79), podemos determinar delimitações para o erro de expressões do tipo

$$(2.80) \quad fl\left(\frac{x_1 x_2 \dots x_m}{y_1 y_2 \dots y_n}\right)$$

através da aplicação de (2.73) - (2.75). Assim procedendo, obtemos

$$(2.81) \quad fl\left(\frac{x_1 x_2 \dots x_m}{y_1 y_2 \dots y_n}\right) = \frac{x_1 x_2 \dots x_m}{y_1 y_2 \dots y_n} (1 + \varepsilon) ,$$

$$(2.82) \quad 1 + \varepsilon = 1 + 1,01(m+n-1)\theta u \quad , \quad |\theta| \leq 1.$$

De maneira geral, o erro relativo cometido no cálculo de (2.80) é pequeno. Embora a aproximação determinada dependa da ordem de realização das multiplicações e divisões, a delimitação (2.82) não depende da mesma. A influência da ordem de realização das operações sobre a aproximação determinada é pequena para (2.71) e (2.80), sendo muito mais acentuada relativamente a (2.58); por esse motivo, tal ordem é, geralmente, considerada apenas em cálculos do último tipo.

Multiplicação escalar em computadores que dispõem de acumulador de precisão dupla.
Consideremos a determinação de

$$(2.83) \quad e_m = fl(x_1 y_1 + x_2 y_2 + \dots + x_m y_m) ; \quad x_i, y_i \in F_\ell ; \quad i = 1, 2, \dots, m$$

A determinação de e_m é feita da seguinte maneira: em primeira etapa são calculados todos os produtos $x_i y_i$ que, em segunda etapa, são somados na ordem em que estão escritos em (2.83). Esse processo é descrito pelas relações:

$$(2.84) \quad \begin{cases} p_i = fl(x_i \times y_i) & , i = 1, \dots, n \\ s_1 = p_1 & \text{e} \quad s_i = fl(s_{i-1} + p_i), i = 2, \dots, n \end{cases}$$

Aplicando (2.68) e (2.48) a (2.84)

$$(2.85) \quad \begin{cases} p_i = x_i y_i (1 + \delta_i) & ; |\delta_i| \leq \mu & ; i = 1, 2, \dots, n \\ s_1 = p_1 & ; \quad s_i = (s_{i-1} + p_i)(1 + \eta_i), |\eta_i| \leq \mu, i = 2, \dots, n \end{cases}$$

Desenvolvendo (2.85) obtemos:

$$(2.86) \quad \begin{aligned} e_m &= fl(x_1 \times y_1 + x_2 \times y_2 + \dots + x_n \times y_n) = \\ &= x_1 y_1 (1 + \varepsilon_1) + x_2 y_2 (1 + \varepsilon_2) + \dots + x_n y_n (1 + \varepsilon_n) \end{aligned}$$

$$(2.87) \quad \begin{cases} 1 + \varepsilon_1 = (1 + \delta_1)(1 + \eta_2) \dots (1 + \eta_n) \\ 1 + \varepsilon_i = (1 + \delta_i)(1 + \eta_i) \dots (1 + \eta_n), i = 2, \dots, n \end{cases}$$

com $|\delta_i|, |\eta_i| \leq \mu, i = 1, 2, \dots, n$.

Aplicando (2.57) a (2.87) podemos, ainda, escrever:

$$(2.88) \quad \begin{cases} 1 + \varepsilon_1 = 1 + 1,01 n \theta_1 \mu & , |\theta_1| \leq 1 \\ 1 + \varepsilon_i = 1 + 1,01 (n - i + 2) \theta_i \mu, & |\theta_i| \leq 1, i = 2, \dots, n \end{cases}$$

ou, para evitar delimitações diferentes para ε_1 e ε_i :

$$(2.89) \quad 1 + \varepsilon_i = 1 + 1,01 (n - i + 2) \theta_i \mu; |\theta_i| \leq 1; i = 1, 2, \dots, n$$

Adição e subtração em computadores que dispõem de acumulador de precisão simples.

Nos computadores que não dispõem de acumulador de precisão dupla, as operações são realizadas de for-

ma que difere um pouco daquela descrita para os computadores que o possuem. Descreveremos apenas a adição e a subtração porque as principais diferenças ocorrem em relação a essas operações; a multiplicação e a divisão não sofrem grandes variações e o produto e quociente aproximados têm sempre erro relativo pequeno.

Na adição e subtração, realizadas através de acumulador de precisão simples, a existência de cancelamento dos dígitos localizados imediatamente à direita da vírgula não implica na nulidade do erro de arredondamento. Por outro lado, a soma e diferença aproximadas não possuem, necessariamente, erro relativo pequeno e as equações (2.48) e (2.49) são substituídas por

$$(2.90) \quad fl(x_1 \pm x_2) = x_1(1 + \varepsilon_1) \pm x_2(1 + \varepsilon_2), \text{ com}$$

$$(2.91) \quad |\varepsilon_1|, |\varepsilon_2| \leq \frac{1+N}{N} u,$$

onde ε_1 e ε_2 são, em geral, diferentes, porém, possuem delimitações da ordem de grandeza de u .

Ilustraremos a realização da adição e subtração, em computadores que utilizam acumulador de precisão simples, através da resolução dos exemplos (2.50)-(2.55) de maneira análoga à efetuada por um computador desse tipo

$$(2.92) \quad \begin{aligned} fl(10^{-3}(0,1572) + 10^4(0,9813)) &= fl(10^4(0,0000 + 0,9813)) = \\ &= fl(10^4(0,9813)) = 10^4(0,9813). \\ |\varepsilon| &\approx 0,00002 \times 10^{-3} \end{aligned}$$

$$(2.93) \quad \begin{aligned} fl(10^6(0,7513) + 10^3(0,4624)) &= fl(10^6(0,7513 + 0,00005)) = \\ &= fl(10^6(0,7518)) = 10^6(0,7518). \\ |\varepsilon| &\approx 0,05002 \times 10^{-3}. \end{aligned}$$

$$\begin{aligned}
 (2.94) \quad & fl(10^3(0,3561) + 10^3(0,8225)) = fl(10^3(0,3561 + 0,8225)) = \\
 & = fl(10^3(1,1786)) = fl(10^4(0,11786)) = 10^4(0,1179) \\
 & |E| \approx 0,33939 \times 10^{-3}
 \end{aligned}$$

$$\begin{aligned}
 (2.95) \quad & fl(10^{-5}(0,4318) + 10^{-5}(-0,4316)) = fl(10^{-5}(0,4318 - 0,4316)) = \\
 & = fl(10^{-5}(0,0002)) = fl(10^{-8}(0,2000)) = 10^{-8}(0,2000) \\
 & |E| = 0
 \end{aligned}$$

$$\begin{aligned}
 (2.96) \quad & fl(10^4(0,1001) + 10^3(-0,9941)) = fl(10^4(0,1001 - 0,0994)) = \\
 & = fl(10^4(0,0007)) = fl(10^4(0,7000)) = 10^4(0,7000) \\
 & |E| = 14,478 \times 10^{-3}
 \end{aligned}$$

$$\begin{aligned}
 (2.97) \quad & fl(10^3(0,9652) + 10^2(0,5828)) = fl(10^3(0,9652 + 0,0583)) = \\
 & = fl(10^3(1,0235)) = fl(10^4(0,10235)) = 10^4(0,1024) \\
 & |E| = 0,50807 \times 10^{-3}
 \end{aligned}$$

Comparamos os resultados fornecidos por ambos os tipos de computadores, com e sem disponibilidade de acumulados de precisão simples.

Os resultados obtidos em (2.92), (2.94) e (2.95) não diferem daqueles fornecidos em, respectivamente, (2.50), (2.52) e (2.53), porque, para a realização desses últimos, não foi utilizada a segunda parte do acumulador de precisão dupla.

Em (2.93) a mantissa 0,4624 foi deslocada três posições para a direita e, então, arredondada. A soma aproximada coincidiu com a aproximação determinada em (2.31)

Em (2.96) a mantissa (-0,9941) foi deslocada uma posição para a direita e, a seguir, arredondada. A soma aproximada difere daquela obtida em (2.54) e tem erro relativo considerável. Esse exemplo, além de mostrar que o erro relativo da aproximação determinada não é, necessariamente, $\leq \epsilon$, ilustra o fato de que, quando são utilizados acumuladores de precisão simples, a existência de

cancelamento dos dígitos colocados à direita da vírgula não implica, obrigatoriamente, na nulidade do erro de arredondamento.

Em (2.97) a mantissa 0,5828 foi deslocada uma posição à direita e arredondada. Esse arredondamento influenciou o arredondamento da soma exata, ocasionando um resultado diferente daquele apresentado em (2.55).

Veremos, a seguir, que (2.90) e (2.91) se verificam não só nos exemplos considerados como, também, de maneira geral.

Em (2.92), (2.94) e (2.95), como a segunda parte do acumulador não foi utilizada para a resolução de (2.50), (2.52) e (2.53), respectivamente, então, (2.90) vale para $\varepsilon_1 = \varepsilon_2 = \varepsilon$ e, como $|\varepsilon| \leq u$, então, (2.91) também vale, pois $\varepsilon_1 = \varepsilon_2 = \varepsilon \leq u \leq \frac{1+N}{N} u$.

Em (2.93) só foi cometido um erro, aquele ocasionado pelo arredondamento da mantissa 0,4624, de modo que vale

$$(2.98) \quad fl(x_1 + x_2) = x_1 + x_2 (1 + \varepsilon), \quad |\varepsilon| \leq u$$

e, portanto, (2.90) e (2.91) se verificam.

Em (2.96) também só foi efetuado um arredondamento, o da mantissa 0,9941, de modo que (2.98) vale também para esse caso.

Em (2.97) ocorreram dois arredondamentos: o da mantissa 0,5828 e o da soma 0,10235. Esses erros de arredondamento satisfazem:

$$(2.99) \quad |\eta_2| \leq N^b \left(\frac{1}{2}\right) N^{-(t+1)}$$

$$(2.100) \quad |\eta_5| \leq N^b \left(\frac{1}{2}\right) N^{-t}$$

em que, como convenção, $x = a \cdot N^b$ representa a soma das parcelas consideradas e, portanto,

$$(2.101) \quad |x_1 + x_2| = N^b |a| \geq N^b \cdot N^{-1} = N^{b-1}$$

Em (2.99) e (2.100) o fator $\frac{1}{2}$ foi colocado entre parênteses porque aparece no caso de utilização de arredondamento propriamente dito e inexistente no caso de arredondamento por corte. Procuramos, dessa maneira, tornar válidas, para ambos os tipos de arredondamentos, as considerações que estamos fazendo relativamente a (2.97).

Por (2.99), (2.100), (2.101) e (2.46) temos:

$$\begin{aligned}
 (2.102) \quad |\eta_2 + \eta_3| &\leq |\eta_2| + |\eta_3| \leq \left(\frac{1}{2}\right) N^{b-t-1} + \left(\frac{1}{2}\right) N^{b-t} \\
 &\leq \frac{1+N}{N} \left(\frac{1}{2}\right) N^{b-t} = \frac{1+N}{N} \left(\frac{1}{2}\right) N^{1-t} N^{b-1} \\
 &\leq \frac{1+N}{N} \left(\frac{1}{2}\right) N^{1-t} |x_1 + x_2| = \frac{1+N}{N} u |x_1 + x_2|
 \end{aligned}$$

A relação (2.102) nos fornece (2.90) e (2.91).

Acumulação de somas e produtos escalares.

Quando consideramos a realização das operações aritméticas através de acumulador de precisão dupla, vimos que para determinar $fl(x_1 + x_2 + x_3)$, por exemplo, o computador efetua $fl(x_1 + x_2)$ temporariamente em precisão dupla e arredonda o valor obtido para precisão simples, antes de adicioná-lo a x_3 . Por outro lado, no cálculo de $fl(x_1 x y_1 + x_2 x y_2 + x_3 x y_3)$, a máquina calcula os produtos parciais em precisão dupla, arredonda-os para precisão simples e, a seguir, os soma da maneira que acabamos de descrever para $fl(x_1 + x_2 + x_3)$. Às vezes, é mais interessante, no caso da multiplicação escalar, não arredondar os produtos parciais e, em ambos os casos, efetuar todas as adições no acumulador de precisão dupla e arredondar somente o resultado final. Essa maneira de operação, embora exija maior tempo de execução, apresenta várias vantagens, principalmente em alguns problemas de álgebra linear. Poucas vezes a acumulação em precisão dupla é efetuada de forma exata, visto que

os expoentes dos diversos termos podem diferir consideravelmente, porém, quando ocorrem, os erros de arredondamento ficam restritos aos últimos dígitos do segundo registro e são, portanto, múltiplos pequenos de $\frac{1}{2} N^{-2t}$ ou N^{-2t} . As operações acumuladas serão indicadas por $fl_{\frac{2}{2}}$.

Adição - Consideremos a determinação de $fl_{\frac{2}{2}}(x_1 + \dots + x_n)$. Cada x_i pode ter mantissa em precisão simples ou dupla. As somas parciais são efetuadas em precisão dupla e a soma final é arredondada para t dígitos. É razoável supor-se a disponibilidade de apenas um acumulador de precisão dupla (não quádrupla), de modo que as delimitações utilizadas nesse caso são as válidas para acumulador de precisão simples, em que t é substituído por $2t$. Procedendo como em (2.58), com $|\delta_i| \leq \frac{1+N}{N} N^{1-2t}$, como em (2.91), obtemos:

$$(2.103) \quad x_1(1+\varepsilon_1) + x_2(1+\varepsilon_2) + \dots + x_n(1+\varepsilon_n), \text{ com}$$

$$(2.104) \quad \begin{cases} 1+\varepsilon_1 = 1 + 1,01(n-1)\theta_1 \frac{1+N}{N} N^{1-2t}, & |\theta_1| \leq 1 \\ 1+\varepsilon_i = 1 + 1,01(n-i+1)\theta_i \frac{1+N}{N} N^{1-2t}; & |\theta_i| \leq 1; i=2, \dots, n \end{cases}$$

ou ainda, para evitar delimitações diferentes para ε_1 e ε_i :

$$(2.105) \quad 1+\varepsilon_i = 1 + 1,01(n-i+1)\theta_i \frac{1+N}{N} N^{1-2t}; \quad |\theta_i| \leq 1; i=1, 2, \dots, n$$

O resultado apresentado em (2.103) é, finalmente, arredondado para t dígitos originando:

$$(2.106) \quad fl_{\frac{2}{2}}(x_1 + x_2 + \dots + x_n) = [x_1(1+\varepsilon_1) + \dots + x_n(1+\varepsilon_n)](1+\varepsilon),$$

$$(2.107) \quad 1+\varepsilon = 1 + \theta_0 u; \quad |\theta_0| \leq 1$$

O fator associado a x_i é $(1+\varepsilon_i)(1+\varepsilon)$ que, em virtude de N^{1-2t} ser, em geral, muito pequeno em relação a N^{1-t} , tem valor comparável a $(1+\varepsilon)$.

Multiplicação escalar. Consideremos a determinação de $fl_z(x_1 \times y_1 + \dots + x_n \times y_n)$, em que os x_i e y_i são números em precisão simples. Procedendo de maneira análoga à utilizada para a acumulação de somas, temos:

$$(2.108) \quad fl_z(x_1 \times y_1 + \dots + x_n \times y_n) = \left[\frac{x_1 y_1 (1 + \epsilon_1)}{z} + \dots + \frac{x_n y_n (1 + \epsilon_n)}{z} \right] (1 + \epsilon)$$

$$(2.109) \quad \begin{cases} 1 + \epsilon_1 = 1 + 1,01(n-1)\theta_1 \frac{1+N}{N} N^{1-2t}, & |\theta_1| \leq 1 \\ 1 + \epsilon_i = 1 + 1,01(n-i+1)\theta_i \frac{1+N}{N} N^{1-2t}; & |\theta_i| \leq 1, \\ & i = 2, \dots, n \end{cases}$$

$$(2.110) \quad 1 + \epsilon = 1 + \theta_0 \mu, \quad |\theta_0| \leq 1$$

Para evitar delimitações diferentes para ϵ_1 e ϵ_i , podemos escrever (2.109) na forma:

$$(2.111) \quad 1 + \epsilon_i = 1 + 1,01(n-i+1)\theta_i \frac{1+N}{N} N^{1-2t}, \quad \text{para } i = 1, 2, \dots, n \text{ e } |\theta_i| \leq 1.$$

Observe-se que em (2.109) e (2.88) há diferença nos fatores que envolvem n ; isso ocorre porque na acumulação de produtos escalares os produtos parciais são efetuados sem erro, visto que o produto de dois números t -digitais pode ser representado exatamente por um número $2t$ -digital.

A operação $fl_z \frac{x_1 \times y_1 + \dots + x_n \times y_n}{z}$, em

que x_i , y_i e z são números em precisão simples, é efetuada de maneira bastante satisfatória. O numerador é acumulado em precisão dupla e esse resultado $2t$ -digital é dividido por z ; finalmente, o quociente obtido é arredondado para t dígitos, fornecendo:

$$(2.112) \quad fl_z \left(\frac{x_1 \times y_1 + \dots + x_n \times y_n}{z} \right) = \left[\frac{x_1 y_1 (1 + \epsilon_1)}{z} + \dots + \frac{x_n y_n (1 + \epsilon_n)}{z} \right] (1 + \epsilon)$$

em que ϵ_i e ϵ satisfazem (2.109), (2.111) e (2.110), respectivamente.

2.5. Conclusões

O arredondamento por corte pode ocasionar erros maiores do que os originados pelo arredondamento propriamente dito, como pode ser observado em (2.10) e (2.11), (2.34) e (2.35), (2.40) e (2.41). Por outro lado, o corte pode introduzir "tendências" em cálculos muito longos, como acontece, por exemplo, na adição de uma grande quantidade de números positivos, em que, em virtude dos erros das somas parciais serem todos de mesmo sinal, pode ser introduzido um erro total considerável após n etapas. A utilização do arredondamento propriamente dito, para esse caso citado, geralmente fornece um erro total menor, após n etapas, já que normalmente ocorre cancelamento de erros parciais de sinais diferentes. Entretanto, apesar do arredondamento propriamente dito apresentar essas vantagens sobre o corte, esse último é bastante utilizado nos computadores que se encontram em uso atualmente, sendo encontrado, por exemplo, no FORTRAN II e no FORTRAN IV dos computadores IBM-7090 e IBM-1130.

Como pode ser observado em (2.6), (2.8) e (2.28), os sistemas em vírgula flutuante possibilitam a representação de uma quantidade de números muito maior do que aquela permitida pelos sistemas em vírgula fixa. Além disso, nos primeiros as operações aritméticas são realizadas de forma mais satisfatória do que nos segundos, a exemplo do que ocorre com a divisão. Principalmente por esses motivos os sistemas em vírgula flutuante são os mais utilizados atualmente.

As operações aritméticas são efetuadas de forma mais satisfatória em computadores que dispõem de acumulador de precisão dupla do que naqueles que não o possuem, como pode ser observado em (2.25) e

(2.26) e em (2.48), (2.49) e (2.90), (2.91).

Em virtude de grande parte dos computadores em uso atualmente utilizar sistema de representação em vírgula flutuante e dispor de acumulador de precisão dupla, utilizaremos, nos capítulos subsequentes, a aritmética de vírgula flutuante e as delimitações de erros válidas para o caso de disponibilidade de tal acumulador. Acreditamos que essa atitude não prejudicará os resultados apresentados neste trabalho, os quais, com pequenas adaptações, também poderão ser aplicados a computadores que trabalham com acumulador de precisão simples e/ou aritmética de vírgula fixa.

Como pode ser concluído de (2.48), (2.49), (2.68) e (2.76), em computadores que utilizam sistema em vírgula flutuante e dispõem de acumulador de precisão dupla verifica-se

$$(2.113) \quad fl(x * y) = (x * y)(1 + \epsilon), \quad |\epsilon| \leq u,$$

em que $*$ representa uma operação aritmética fundamental e y é suposto diferente de zero no caso de $*$ indicar a divisão. A igualdade (2.113) foi obtida a partir de conclusões válidas para o erro relativo definido em relação ao valor verdadeiro; quando esse erro é definido em relação ao valor aproximado, obtém-se, a partir de (2.40) e (2.41), em correspondência a (2.113), :

$$(2.114) \quad fl(x * y) = \frac{x * y}{1 + \epsilon'}, \quad |\epsilon'| \leq u,$$

com $y \neq 0$ no caso em que $*$ representa a divisão.

3- Análises de erros

3.1- Sensibilidade

3.1.1- Problemas mal condicionados

Consideremos o uso de um computador para a resolução numérica de um problema e suponhamos que o mesmo trabalhe com palavras de comprimento t . Não será que os dados sejam definidos exatamente e possam ser representados de maneira exata por t dígitos, o computador inicia a resolução com o que podemos chamar aproximação t -digital do problema verdadeiro. Assim, ao fazer uma análise de erros deve-se, necessariamente, levar em consideração os efeitos ocasionados, nas soluções, por perturbações nos dados. Um problema é dito mal condicionado quando possui soluções muito sensíveis a perturbações em seus dados ou parâmetros, ou seja, quando pequenas perturbações relativas nesses últimos causam, comparativamente, grandes variações nas primeiras. Em caso contrário, ou seja, quando as soluções são pouco sensíveis a perturbações nos parâmetros, o problema é dito bem condicionado.

3.1.2- Números de condição

São valores associados a um problema, de acordo com uma certa regra, que dão informações quanto à sensibilidade das soluções do mesmo, relativamente a perturbações nos dados. Veremos, a seguir, uma formalização desse conceito.

Sejam E e F espaços normados. Seja $x_0 \in E$ e seja $\epsilon > 0$. Seja G o conjunto das aplicações de E em F ,

determinadas por uma certa classe de processos numéricos que relacionam dados e resultados. Seja $g \in G$. Define-se como número de condição de g a um número κ , associado a g de acordo com uma certa regra, que satisfaz a:

$$\|x - x_0\| < \varepsilon \Rightarrow \|g(x) - g(x_0)\| \leq \kappa \|x - x_0\|.$$

Dados $g_1, g_2 \in G$, com números de condição κ_1 e κ_2 , respectivamente, g_1 é dita melhor condicionada que g_2 se $\kappa_1 < \kappa_2$.

Às vezes, também é chamado número de condição o número κ definido pelo limite, se existir,

$$\kappa = \lim_{x \rightarrow x_0} \frac{\|g(x) - g(x_0)\|}{\|x - x_0\|}$$

Quando $E = \mathbb{R}^n$ e $F = \mathbb{R}^m$, também são definidos os números de condição relativos à variação de apenas algumas componentes dos dados ou referentes a apenas algumas componentes dos resultados. Assim, por exemplo, se $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ e se $y = g(x)$, $y \in \mathbb{R}^m, \forall x \in \mathbb{R}^n$, são frequentemente utilizadas como números de condição as quantidades $\kappa_{rs}, \kappa_j, \kappa'_j$ e κ''_j definidas por:

$$\kappa_{rs} = \left| \frac{\partial y_r}{\partial x_s} \right| ; r = 1, \dots, m ; s = 1, \dots, n$$

$$\left(\sum_{j=1}^m \delta y_j^2 \right)^{\frac{1}{2}} \leq \kappa \left(\sum_{i=1}^n \delta x_i^2 \right)^{\frac{1}{2}}$$

$$|\delta y_j| \leq \kappa_j \left(\sum_{i=1}^n \delta x_i^2 \right)^{\frac{1}{2}}$$

$$\left| \frac{\delta y_j}{y_j} \right| \leq \kappa''_j \left[\sum_{i=1}^n \left(\frac{\delta x_i}{x_i} \right)^2 \right]^{\frac{1}{2}}$$

$$\frac{\delta y_j}{\left(\sum y_j^2\right)^{\frac{1}{2}}} \leq \kappa \frac{\left(\sum_{i=1}^n \delta x_i^2\right)^{\frac{1}{2}}}{\left(\sum_{i=1}^n x_i^2\right)^{\frac{1}{2}}}$$

onde ∂ é o símbolo usual de derivada e δ representa uma variação.

Números de condição de matrizes

Seja $A \in M_n(\mathbb{R})$, A não singular. Como diversos problemas relacionados com A são associados a um mesmo número de condição κ , costuma-se atribuir a κ a denominação de "número de condição da matriz A ". Embora os números de condição sejam associados a problemas, esse procedimento facilita bastante o estudo da sensibilidade possibilitando, por exemplo, um estudo geral das propriedades de κ , estudo esse que pode depois ser aplicado quando da consideração dos diversos problemas aos quais κ é associado. Adotaremos essa atitude porque os problemas considerados no capítulo 4 são todos associados a um mesmo número de condição.

(3.1) Definição - Seja $A \in M_n(\mathbb{R})$, A não singular. Denomina-se número de condição de A , e se indica por $\text{cond}_k(A)$, a:

$$\text{cond}_k(A) = \|A\|_k \|A^{-1}\|_k$$

Como pode ser observado em (3.1), o número de condição de uma matriz depende da particular norma adotada. A norma $\|\cdot\|_2$ é bastante utilizada e o correspondente número de condição, denominado número de condição espectral, pode ser expresso por:

$$(3.2) \quad \text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\mu_1}{\mu_n}$$

onde μ_1 e μ_n representam o maior e o menor valor singular de A , respectivamente, como visto em 1.5. Portanto, $\text{cond}_2(A)$ dá a medida da máxima distorção que uma aplicação linear com matriz A produz na esfera unitária.

O número de condições definido em (3.1) goza das propriedades:

$$\text{cond}_k(A) \geq 1$$

$$\text{cond}_k(A) = \text{cond}_k(A^{-1})$$

$$\text{cond}_k(A) = \text{cond}_k(\alpha A),$$

sendo α um escalar e $\|\cdot\|_k$ uma norma subordinada ou a norma euclidiana.

As propriedades citadas são de verificação bastante simples, como pode ser observado a seguir:

$$\|A\|_k \|A^{-1}\|_k \geq \|AA^{-1}\|_k = \|I\|_k = 1, \text{ para normas subordinadas}$$

$$\|A\|_E \|A^{-1}\|_E \geq \|AA^{-1}\|_E = \|I\|_E = n^{\frac{1}{2}} \geq 1$$

$$\text{cond}_k(A) = \|A\|_k \|A^{-1}\|_k = \|A^{-1}\|_k \|A\|_k = \text{cond}_k(A^{-1})$$

$$\text{cond}_k(\alpha A) = \|\alpha A\|_k \|(\alpha A)^{-1}\|_k = |\alpha| \|A\|_k \left| \frac{1}{\alpha} \right| \|A^{-1}\|_k = \text{cond}_k(A)$$

A dificuldade do cálculo efetivo de $\text{cond}_k(A)$ está na determinação de $\|A^{-1}\|_k$. Entretanto, em diversos casos há possibilidade de majorar esse valor e, assim, determinar estimativas para $\text{cond}_k(A)$. Podemos majorar $\|A^{-1}\|_k$, por exemplo, quando temos uma aproximação X de A^{-1} para a qual $F = I - AX$ verifica $\|F\|_k < 1$, já que, nesse caso,

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

$$\|A^{-1}\|_k \|F\|_k \geq \|A^{-1}F\|_k = \|A^{-1} - X\|_k \geq \|A^{-1}\|_k - \|X\|_k$$

$$\|A^{-1}\|_k \leq \frac{\|X\|_k}{1 - \|F\|_k}$$

O valor de F pode ser calculado de forma bastante precisa no caso de disponibilidade de acumulação de produtos escalares.

3.1.3. Sensibilidade das soluções de sistemas de equações lineares

Seja o sistema de equações lineares :

$$(3.3) \quad Ax = b, \text{ com } A \in M_m(\mathbb{R}) \text{ e } A \text{ não singular, onde}$$

x representa o vetor solução e b o vetor dos termos independentes. Como A tem inversa única, então, x é determinado univocamente por

$$(3.4) \quad x = A^{-1}b$$

Estudaremos, a seguir, a sensibilidade dessa solução x relativamente a perturbações nos dados, elementos de A e b . Esse estudo é feito por partes: supõe-se, inicialmente, que A seja conhecida exatamente e que b seja afetada de erros; a seguir, supõe-se que b seja conhecido exatamente e A esteja sujeita a erros. Finalmente, supõe-se que ambos, A e b , sejam sujeitos a erros.

Sensibilidade da solução relativamente a perturbações em b .

Admitamos que, em (3.3), A seja conhecida exatamente e que b seja afetada de erros. Supondo que o vetor solução x se transforme em $x + \delta x$ quando b sofre uma perturbação δb , podemos escrever:

$$(3.5) \quad A(x + \delta x) = b + \delta b$$

e, portanto,

$$(3.6) \quad \delta x = A^{-1} \delta b$$

Sejam $\|\cdot\|_p$ e $\|\cdot\|_k$ normas de vetores e matrizes, respectivamente, consistentes. Adotando-as em (3.6) obtemos:

$$(3.7) \quad \|\delta x\|_p \leq \|A^{-1}\|_k \|\delta b\|_p,$$

sendo possível, para alguns vetores, a obtenção da igualdade.

Utilizando em (3.3) as mesmas normas adotadas

em (3.7) temos:

$$(3.8) \quad \|b\|_p \leq \|A\|_k \|x\|_p$$

Multiplicando (3.7) e (3.8) membro a membro:

$$(3.9) \quad \|\delta x\|_p \|b\|_p \leq \|A\|_k \|A^{-1}\|_k \|x\|_p \|\delta b\|_p$$

Suponhamos $b \neq \theta$ em (3.3). Então, em consequência da não singularidade de A , temos que $x \neq \theta$ e podemos dividir (3.9) por $\|b\|_p \|x\|_p$ obtendo

$$(3.10) \quad \frac{\|\delta x\|_p}{\|x\|_p} \leq \|A\|_k \|A^{-1}\|_k \frac{\|\delta b\|_p}{\|b\|_p}$$

Aplicando (3.1) a (3.10) temos:

$$(3.11) \quad \frac{\|\delta x\|_p}{\|x\|_p} \leq \underset{k}{\text{cond}}(A) \frac{\|\delta b\|_p}{\|b\|_p},$$

que nos permite concluir que o número de condições $\underset{k}{\text{cond}}(A)$ fornece uma indicação da sensibilidade da solução de (3.3) relativamente a variações dos termos conhecidos.

Como, para determinados vetores, pode ser obtida igualdade em (3.6) e (3.8), então, conseqüentemente, isso também pode ocorrer em (3.11); portanto, a majoração fornecida por essa última expressão, para $\frac{\|\delta x\|_p}{\|x\|_p}$, é a melhor possível. Utilizaremos (1.52) para mostrar como é possível, para a norma 2, obter igualdade em (3.11). Para isso, consideremos as mudanças ortogonais de coordenadas citadas em 1.5. Relativamente ao novo sistema de coordenadas, (3.3) será escrito

$$(3.12) \quad D x' = b',$$

ou, ainda,

$$(3.13) \quad \begin{cases} \mu_1 x'_1 = b'_1 \\ \mu_2 x'_2 = b'_2 \\ \vdots \\ \mu_n x'_n = b'_n \end{cases}$$

e a equação $A \delta x = \delta b$, dos erros,

$$(3.14) \quad D \delta x' = \delta b' \quad , \quad \text{ou,}$$

$$(3.15) \quad \begin{cases} \mu_1 \delta x'_1 = \delta b'_1 \\ \mu_2 \delta x'_2 = \delta b'_2 \\ \vdots \\ \mu_n \delta x'_n = \delta b'_n \end{cases}$$

Como A é não singular, então, por (1.59) e (1.60):

$$\|A\|_2 = \|D\|_2 = \mu_1 \quad \text{e} \quad \|A^{-1}\|_2 = \|D^{-1}\|_2 = \mu_n^{-1}.$$

Seja $b' = (\beta, 0, 0, \dots, 0)^T$, onde $b'_1 = \beta$ e $b'_i = 0$ ($i=2, \dots, n$). Então, por (3.13), $x' = (\mu_1^{-1} \beta, 0, 0, \dots, 0)^T$, de modo que é obtida a igualdade em (3.8), para qualquer constante β não nula.

Seja $\delta b' = (0, 0, \dots, 0, \delta)^T$, onde $\delta b'_i = 0$ ($i=1, \dots, n-1$) e $\delta b'_n = \delta$. Então, por (3.15), $\delta x' = (0, 0, \dots, \mu_n^{-1} \delta)^T$ e, para qualquer constante δ não nula, é obtida a igualdade em (3.7).

Portanto, a igualdade em (3.11) pode ser obtida quando b está em uma direção que recebe a maior amplificação por A (e portanto a menor amplificação por A^{-1}) e δb está em uma direção que sofre a menor amplificação por A (e portanto a maior amplificação por A^{-1}). Essas duas direções são, necessariamente, ortogonais, já que $\mu_1 \neq \mu_n$.

Sensibilidade da solução relativamente a perturbações em A.

Suponhamos que em (3.3) b seja conhecido exatamente e A seja afetada de erros. Admitindo que x se transforme em $x + \delta x$, quando A sofre a perturbação δA , temos:

$$(3.16) \quad (A + \delta A)(x + \delta x) = b \quad \text{e, portanto,}$$

$$(3.17) \quad \delta x = -(A + \delta A)^{-1} \delta A x$$

Em virtude da não singularidade de A não implicar, necessariamente, na não singularidade de $A + \delta A$, procuraremos estabelecer condições de existência para (3.17). Seja

$$(3.18) \quad \|A^{-1} \delta A\|_k < 1,$$

onde $\|\cdot\|_p$ é uma norma subordinada ou a norma euclidiana e consideremos a identidade:

$$(3.19) \quad A + \delta A = A(I + A^{-1} \delta A)$$

Então, por (1.43), podemos afirmar que $A + \delta A$ é não singular em todos os casos que (3.18) se verifique. Admitamos, então, no que se segue, que (3.18) seja verificada.

Em consequência de (3.17) e (3.19), temos:

$$(3.20) \quad \delta x = -(I + A^{-1} \delta A)^{-1} A^{-1} \delta A x$$

Aplicando a (3.20) uma norma $\|\cdot\|_p$ compatível com $\|\cdot\|_k$, e utilizando (1.44)

$$(3.21) \quad \|\delta x\|_p \leq \frac{\|A^{-1} \delta A\|_k \|x\|_p}{1 - \|A^{-1} \delta A\|_k} \leq \frac{\|A^{-1}\|_k \|\delta A\|_k \|x\|_p}{1 - \|A^{-1}\|_k \|\delta A\|_k},$$

desde que

$$\|A^{-1}\|_k \|\delta A\|_k < 1.$$

Finalmente, a aplicação de (3.1) a (3.21) fornece:

$$(3.22) \quad \frac{\|\delta x\|_p}{\|x\|_p} \leq \frac{\text{cond}_k(A) \frac{\|\delta A\|_k}{\|A\|_k}}{1 - \text{cond}_k(A) \frac{\|\delta A\|_k}{\|A\|_k}}$$

Quando consideramos o erro relativo de x definido em relação ao valor aproximado $x + \delta x$, obtemos uma expressão correspondente a (3.22) que é mais simples e análoga a (3.11). Vamos determiná-la a seguir. Considerando o valor de δx obtido de (3.16) e (3.4), dado por

$$(3.23) \quad \delta x = \left[(A + \delta A)^{-1} - A^{-1} \right] b,$$

e substituindo B por $A + \delta A$ na identidade

$$(3.24) \quad B^{-1} - A^{-1} = A^{-1} (A - B) B^{-1}, \quad \text{obtemos}$$

$$(3.25) \quad (A + \delta A)^{-1} - A^{-1} = -A^{-1} \delta A (A + \delta A)^{-1}$$

Admitindo a verificação de (3.18), para garantir a não singularidade de $(A + \delta A)$, substituindo (3.25) em (3.23) e utilizando (3.16), obtemos:

$$(3.26) \quad \delta x = -A^{-1} \delta A (A + \delta A)^{-1} b = -A^{-1} \delta A (x + \delta x)$$

Utilizando normas compatíveis em (3.26), temos:

$$\|\delta x\|_p \leq \|A^{-1}\|_k \|\delta A\|_k \|x + \delta x\|_p$$

Finalmente, supondo $x + \delta x \neq \theta$ e utilizando (3.1), obtemos a expressão procurada

$$(3.27) \quad \frac{\|\delta x\|_p}{\|x + \delta x\|_p} \leq \text{cond}_k(A) \frac{\|\delta A\|_k}{\|A\|_k}$$

Sensibilidade da solução relativamente a perturbações em A e b .

Suponhamos que A e b sejam afetados de erros e que o vetor solução x se transforme em $x + \delta x$ quando

A e b sofrem, respectivamente, as perturbações δA e δb . Então,

$$(3.28) \quad (A + \delta A)(x + \delta x) = b + \delta b \quad \text{e, portanto,}$$

$$(3.29) \quad \delta x = (A + \delta A)^{-1} [b + \delta b - (A + \delta A)x] = (A + \delta A)^{-1} (\delta b - \delta A x)$$

Admitindo (3.18) como hipótese e utilizando (3.19):

$$(3.30) \quad \delta x = (I + A^{-1} \delta A)^{-1} (A^{-1} \delta b - A^{-1} \delta A x)$$

Utilizando normas consistentes em (3.30), por (1.44)

$$(3.31) \quad \|\delta x\|_p \leq \frac{\|A^{-1}\|_k \|\delta b\|_k + \|A^{-1}\|_k \|\delta A\|_k \|x\|_p}{1 - \|A^{-1}\|_k \|\delta A\|_k},$$

desde que $\|A^{-1}\|_k \|\delta A\|_k < 1$.

Aplicando a (3.3) as normas adotadas em (3.31):

$$(3.32) \quad \|A\|_k \|x\|_p \geq \|b\|_p$$

Dividindo (3.31) e (3.32) membro a membro

$$(3.33) \quad \frac{\|\delta x\|_p}{\|A\|_k \|x\|_p} \leq \frac{\|A^{-1}\|_k \left(\frac{\|\delta b\|_p}{\|b\|_p} + \frac{\|\delta A\|_k \|x\|_p}{\|b\|_p} \right)}{1 - \|A^{-1}\|_k \|\delta A\|_k}$$

Finalmente, em virtude de (3.1) e (3.4)

$$(3.34) \quad \frac{\|\delta x\|_p}{\|x\|_p} \leq \frac{\text{cond}(A)_k \left(\frac{\|\delta b\|_p}{\|b\|_p} + \text{cond}(A)_k \frac{\|\delta A\|_k}{\|A\|_k} \right)}{1 - \text{cond}(A)_k \frac{\|\delta A\|_k}{\|A\|_k}}$$

3.1.4 - Sobre o mal condicionamento de sistemas de equações lineares

Existem algumas concepções relativas ao mal condicionamento de sistemas de equações lineares que, embora difundidas

na prática, não são verdadeiras. Consideraremos, a seguir, algumas dessas idéias e mostremos, através de exemplos, a não validade das mesmas.

Comumente, na prática, são considerados mal condicionados os sistemas de equações lineares cujas matrizes possuem ordem muito elevada ou têm determinantes pequenos. O exemplo que se segue mostra que esse fato não é verdadeiro. Consideremos um sistema cuja matriz A possui valores singulares $\mu_1 = \mu_2 = \dots = \mu_m = 10^{-30}$. Então, por (3.2), (1.61) e (3.11), temos

$$\text{cond}_2(A) = \mu_1 \mu_m^{-1} = 1$$

$$|\det A| = 10^{-30m}$$

$$\frac{\|\delta x\|_2}{\|x\|_2} = \frac{\|\delta b\|_2}{\|b\|_2}$$

Portanto, $\text{cond}_2(A) = 1$, qualquer que seja m , ordem de A . Por outro lado, embora o determinante de A seja um número bastante pequeno, as perturbações da solução são iguais às do vetor dos termos independentes. Na realidade, a pequena grandeza de $\det A$ só pode ser relacionada ao mal condicionamento do sistema em algumas circunstâncias especiais. Assim, por exemplo, isso pode ser feito para uma matriz A de ordem n fixa desde que, de alguma maneira, multipliquemos A por fatores de escala que façam com que μ_1 permaneça fixo; se μ_1 e n são fixos, então, se $\det A \rightarrow 0$, necessariamente $\mu_m \rightarrow 0$ e $\text{cond}(A) \rightarrow \infty$. Entretanto, essa relação é muito fraca, como pode ser observado para $n = 101$, $\mu_1 = 1$ e $\mu_2 = \mu_3 = \dots = \mu_m = 10^{-1}$. Embora $\det A = \mu_1 (\mu_2)^{m-1} = 10^{-100}$, temos $\text{cond}_2(A) = \mu_1 \mu_2^{-1} = 10$ e, portanto, o sistema que possui A como matriz pode ser considerado bem condicionado.

Também é usual, na prática, considerar mal condicionados os sistemas lineares $Ax = b$ para os quais A^{-1} pos-

sui elementos grandes em comparação aos de A , o que pode originar grandes erros absolutos na solução. Essa concepção não é correta, já que o mal condicionamento do sistema depende de sua solução, como pode ser observado em (3.35). Os diferentes sistemas de equações lineares considerados neste exemplo possuem uma mesma matriz A , cuja inversa A^{-1} tem elementos grandes comparados com os de A , porém, enquanto um deles é bem condicionado, o outro é extremamente mal condicionado. O que ocorre, na realidade, é que se A^{-1} tem elementos grandes, então, sempre é possível determinar um vetor b , dos termos independentes, que origina soluções pequenas e, portanto, erros relativos consideráveis. É óbvio que, se $\|A\|_k = 1$, então, se os elementos de A^{-1} são grandes em comparação aos de A , o sistema associado será mal condicionado, já que $\text{cond}(A) = \|A^{-1}\|_k$.

(3.35) Exemplo - Consideremos o sistema de equações lineares apresentado em (9, p. 75-77)

$$(3.36) \quad \begin{cases} \frac{1}{2}x_1 + \frac{1}{3}x_2 + \frac{1}{4}x_3 + \frac{1}{5}x_4 = b_1 + \varepsilon_1 \\ \frac{1}{3}x_1 + \frac{1}{4}x_2 + \frac{1}{5}x_3 + \frac{1}{6}x_4 = b_2 + \varepsilon_2 \\ \frac{1}{4}x_1 + \frac{1}{5}x_2 + \frac{1}{6}x_3 + \frac{1}{7}x_4 = b_3 + \varepsilon_3 \\ \frac{1}{5}x_1 + \frac{1}{6}x_2 + \frac{1}{7}x_3 + \frac{1}{8}x_4 = b_4 + \varepsilon_4 \end{cases}$$

em que os coeficientes das incógnitas são valores exatos, mas os termos independentes podem ser afetados de erros ε_i ($i=1, \dots, 4$) para os quais é conhecida uma majoração η , ou seja, $|\varepsilon_i| \leq \eta$.

A inversa exata da matriz de (3.36) é

$$(3.37) \quad A^{-1} = \begin{pmatrix} 200 & -1200 & 2100 & -1120 \\ -1200 & 8100 & -15120 & 8400 \\ 2100 & -15120 & 29400 & -16800 \\ -1120 & 8400 & -16800 & 9800 \end{pmatrix}$$

A partir de A^{-1} , podemos determinar como varia a solução de (3.36), em relação aos diferentes ε_i , já que as quantidades δx_i podem ser descritas pelas equações:

$$(3.38) \quad A(x + \delta x) = b + E$$

$$(3.39) \quad \delta x = A^{-1} E,$$

em que $E = (\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4)^T$.

Suponhamos que ocorra a pior combinação possível de grandezas e sinais, de modo que $|\delta x_i|$ seja o máximo possível. Nesse caso,

$$(3.40) \quad \begin{cases} |\delta x_1| = 4620 \eta & |\delta x_2| = 32820 \eta \\ |\delta x_3| = 63420 \eta & |\delta x_4| = 36120 \eta \end{cases}$$

Se tomarmos $\eta = 10^{-4}$, então, mesmo que todas as operações aritméticas necessárias à determinação da solução de (3.36) sejam realizadas exatamente, ainda assim os erros absolutos das componentes do vetor solução poderão ser

$$(3.41) \quad |\delta x_1| \approx 0,46 \quad |\delta x_2| \approx 3,28 \quad |\delta x_3| \approx 6,34 \quad |\delta x_4| \approx 3,61$$

Consideremos, inicialmente, que os b_i originem uma solução que seja da ordem de grandeza de

$$(3.42) \quad x_1 = 4000 \quad x_2 = -32000 \quad x_3 = 63000 \quad x_4 = -36000,$$

de modo que os erros relativos originados pelos ε_i serão, no máximo,

$$(3.43) \quad \begin{cases} \left| \frac{\delta x_1}{x_1} \right| = 0,1150 \times 10^{-3} & \left| \frac{\delta x_2}{x_2} \right| = 0,1025 \times 10^{-3} \\ \left| \frac{\delta x_3}{x_3} \right| = 0,1006 \times 10^{-3} & \left| \frac{\delta x_4}{x_4} \right| = 0,1003 \times 10^{-3}, \end{cases}$$

os quais, para muitos problemas, podem ser considerados pequenos. Uma solução de ordem de grandeza de (3.42) pode ser obtida, por exemplo, com os seguintes termos independentes:

$$(3.44) \quad b_1 = 1 \quad b_2 = -1 \quad b_3 = 1 \quad b_4 = -1$$

Consideremos, a seguir, os termos independentes

$$(3.45) \quad b_1 = b_2 = b_3 = b_4 = 1,$$

que são de ordem de grandeza comparável a (3.44). Nesse caso,

$$x_1 = -20 \quad x_2 = 180 \quad x_3 = -420 \quad x_4 = 280,$$

de modo que, embora os erros absolutos máximos ainda sejam dados em (3.41), os erros relativos são bem maiores:

$$(3.46) \quad \begin{cases} \left| \frac{\delta x_1}{x_1} \right| = 23,00 \times 10^{-3} & \left| \frac{\delta x_2}{x_2} \right| \approx 18,22 \times 10^{-3} \\ \left| \frac{\delta x_3}{x_3} \right| \approx 15,10 \times 10^{-3} & \left| \frac{\delta x_4}{x_4} \right| \approx 12,89 \times 10^{-3} \end{cases}$$

Finalmente, consideremos os termos independentes

$$(3.47) \quad b_1 = 1 \quad b_2 = 0,7227 \quad b_3 = 0,5697 \quad b_4 = 0,4714,$$

que também são de ordem de grandeza comparável a (3.44) e (3.45) e fornecem:

$$x_1 = 1,162 \quad x_2 = -0,234 \quad x_3 = 2,436 \quad x_4 = -0,560.$$

Novamente os erros absolutos da solução são dados por (3.41), porém, os erros relativos são muito maiores que aqueles apresentados em (3.43) e (3.46), sendo dados por:

$$(3.48) \quad \begin{cases} \left| \frac{\delta x_1}{x_1} \right| = 395,87 \times 10^{-3} & \left| \frac{\delta x_2}{x_2} \right| = 1407,1 \times 10^{-3} \\ \left| \frac{\delta x_3}{x_3} \right| = 2602,6 \times 10^{-3} & \left| \frac{\delta x_4}{x_4} \right| = 6446,4 \times 10^{-3} \end{cases}$$

É lógico que, nesse último caso, se $\eta = 10^{-4}$ representar a unidade absoluta de arredondamento do computador no

qual (3.36) deve ser resolvido, então, mesmo que todas as operações aritméticas utilizadas nessa resolução sejam efetuadas exatamente, o que é difícil de ocorrer, ainda assim os resultados obtidos não serão satisfatórios.

Concluindo, podemos afirmar que, embora relacionados com uma mesma matriz, os problemas obtidos pela substituição de (3.44), (3.45) e (3.47) em (3.36) possuem soluções muito diferentes entre si quanto à sensibilidade às perturbações dos termos independentes. O problema resultante da associação de (3.44) e (3.36) é bem condicionado, enquanto aquele obtido de (3.47) e (3.36) é extremamente mal condicionado.

3.2 - Análises progressiva e regressiva

A análise de erros de um processo numérico tem por finalidade determinar uma delimitação para a grandeza do erro do resultado final. Neste trabalho vamos tratar somente dos de arredondamento e, dessa forma, o erro do resultado final será causado pelo efeito propagado dos erros de arredondamento nos dados e nas operações aritméticas realizadas.

As principais formas de análise de erros de processos algébricos são: progressiva ou direta e regressiva ou inversa. Para caracterizá-las precisamos, inicialmente, fazer algumas considerações a respeito de processos numéricos.

Um processo numérico pode ser descrito por uma sequência de equações do tipo:

$$(3.49) \quad x_{i+1} = g_i(x_1, x_2, \dots, x_i), \quad i = 1, 2, \dots, n-1,$$

que representam, em cada estágio i , o cálculo de novos valores em função de dados iniciais e dados previamente calculados. Nas equações (3.49), x_i não representa, necessariamente, um número,

mas pode, também, representar um conjunto de valores; nesse caso, x_i é dito o estado da i -ésima etapa de cálculo. Assim, por exemplo, x_1 representa o estado inicial.

Em particular, no caso de problemas algébricos a obtenção de x_{i+1} a partir dos x_i ($i = 1, 2, \dots, n-1$) envolve somente as operações aritméticas fundamentais.

Quando o processo numérico é executado no computador, são utilizadas, em substituição a x_1, x_2, \dots, x_i , representações internas $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_i$, em virtude das limitações dessa máquina. Por outro lado, as operações aritméticas necessárias à determinação de x_{i+1} são, geralmente, resolvidas de forma aproximada, de modo que a função g_i é substituída por uma aproximação \bar{g}_i . Em virtude dessas aproximações, um resultado intermediário do processo pode ser representado da seguinte maneira:

$$(3.50) \quad \bar{x}_{i+1} = \bar{g}_i(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_i)$$

Análise progressiva ou direta

Neste tipo de análise procuramos, em cada etapa, estimar $|\bar{x}_{i+1} - x_{i+1}|$, a partir de estimativas prévias dos $|\bar{x}_j - x_j|$ ($j = 1, 2, \dots, i$).

A análise progressiva é, em geral, extremamente complicada e laboriosa, em virtude de serem muito confusas as inter-relações entre os diversos erros. Além disso, as regras ordinárias da aritmética não podem ser aplicadas durante a sua realização, já que as operações aproximadas nem sempre satisfazem às mesmas propriedades das operações exatas, a exemplo da não associatividade da adição e multiplicação em vírgula flutuante, como vimos em 2.4.3

Análise regressiva ou inversa

Nesse tipo de análise de erros não nos preocu-

podemos em comparar \bar{x}_{i+1} e x_{i+1} em cada estágio do processo. Em lugar disso, procuramos mostrar que o valor \bar{x}_{i+1} pode ser obtido exatamente em função de valores iniciais perturbados

$$x_1 + \delta_i$$

e procuramos determinar estimativas para os δ_i . Assim procedendo, podemos escrever (3.50) na forma

$$\bar{x}_{i+1} = \phi_i(x_1 + \delta_i),$$

em que ϕ_i é obtida por composição de g_1, g_2, \dots, g_i , e, de maneira geral, os δ_i não são determinados de forma unívoca. No final do processo obtemos:

$$\bar{x}_n = \phi_{n-1}(x_1 + \varepsilon), \quad \varepsilon = \delta_{n-1}$$

De maneira geral, a análise regressiva não determina estimativas para $|\bar{x}_n - x_n|$, preocupação fundamental da maioria das análises de erros. Nesses casos, entretanto, $|\bar{x}_n - x_n|$ pode ser estimada a partir de estimativas de

$$\left| \phi_{n-1}(x_1 + \varepsilon) - \phi_{n-1}(x_1) \right|,$$

através de um estudo da sensibilidade do processo.

As propriedades usuais da aritmética podem ser utilizadas durante a realização da análise inversa, já que são exatas as operações aritméticas utilizadas nessa análise.

Atualmente, a análise regressiva é a mais utilizada, não só por ser de aplicação mais simples, como também por possuir outras vantagens em relação à progressiva, o exemplo da avaliação da eficácia do processo numérico utilizado, como será visto a seguir. Por esse motivo, e por ser desse tipo a análise efetuada no capítulo 4, apresentaremos alguns exemplos que elucidam melhor a sua aplicação e destacam algumas de suas particularidades.

(3.51) Exemplo - Consideremos a determinação do valor numérico do polinômio

$$(3.52) \quad P(x) = b_0 + b_1 x + \dots + b_n x^n,$$

para um valor x dado, através da sequência de operações:

$$(3.53) \quad \begin{cases} v_n = b_n \\ v_i = x v_{i+1} + b_i, \quad i = n-1, \dots, 0 \end{cases},$$

de modo que v_0 é a solução procurada.

Suponhamos que os coeficientes de x não possam ser armazenados exatamente no computador e sejam substituídos por representações internas

$$(3.54) \quad \bar{b}_i = b_i (1 + \delta_i); \quad |\delta_i| \leq \mu, \quad i = 0, 1, \dots, n.$$

Em virtude de (2.113) e (3.54), as expressões (3.53) são substituídas, no computador, por

$$(3.55) \quad \begin{cases} v_n = \bar{b}_n = b_n (1 + \delta_n), \quad |\delta_n| \leq \mu \\ v_i = fl(x v_{i+1} + \bar{b}_i) = x v_{i+1} (1 + \alpha_i) (1 + \beta_i) + b_i (1 + \delta_i) (1 + \beta_i), \\ |\alpha_i|, |\beta_i|, |\delta_i| \leq \mu \end{cases}$$

Desenvolvendo as relações (3.55) obtemos:

$$(3.56) \quad v_0 = b_0 (1 + \varepsilon_0) + b_1 (1 + \varepsilon_1) x + \dots + b_n (1 + \varepsilon_n) x^n,$$

$$(3.57) \quad \begin{cases} 1 + \varepsilon_0 = (1 + \delta_0) (1 + \beta_0) \\ 1 + \varepsilon_1 = (1 + \delta_1) (1 + \beta_1) (1 + \alpha_0) (1 + \beta_0) \\ \dots \\ 1 + \varepsilon_i = (1 + \delta_i) (1 + \beta_i) (1 + \alpha_{i-1}) (1 + \beta_{i-1}) \dots (1 + \alpha_1) (1 + \beta_1) (1 + \alpha_0) (1 + \beta_0) \\ \dots \\ 1 + \varepsilon_n = (1 + \delta_n) (1 + \alpha_{n-1}) (1 + \beta_{n-1}) \dots (1 + \alpha_1) (1 + \beta_1) (1 + \alpha_0) (1 + \beta_0) \end{cases}$$

Aplicando (2.57) a (3.57) obtemos:

$$(3.58) \quad \begin{cases} 1 + \varepsilon_i = 1 + 1,01 (2i+2) \theta_i u ; & |\theta_i| \leq 1, \quad i = 0, 1, \dots, n-1 \\ 1 + \varepsilon_n = 1 + 1,01 (2n+1) \theta_n u ; & |\theta_n| \leq 1. \end{cases}$$

Com a finalidade de evitar delimitações diferentes para ε_i ($i = 0, 1, \dots, n-1$) e ε_n , podemos utilizar:

$$(3.59) \quad 1 + \varepsilon_i = 1 + 1,01 (2i+2) \theta_i u ; \quad |\theta_i| \leq 1, \quad i = 0, 1, \dots, n$$

Podemos interpretar v_0 , dado em (3.56), como o valor numérico exato do polinômio $\bar{P}(x)$ que possui os coeficientes:

$$(3.60) \quad B_i = b_i (1 + \varepsilon_i) = b_i + \delta b_i,$$

com ε_i dado em (3.59) e

$$(3.61) \quad \delta b_i = b_i \varepsilon_i, \quad \text{ou seja,}$$

$$(3.62) \quad \bar{P}(x) = b_0 + \delta b_0 + (b_1 + \delta b_1)x + \dots + (b_n + \delta b_n)x^n$$

E efetuando a análise de sensibilidade do polinômio à variação dos coeficientes, obtemos, por (3.52) e (3.62):

$$\delta P(x) = \bar{P}(x) - P(x) = \delta b_0 + \delta b_1 x + \dots + \delta b_n x^n$$

e, por (3.61),

$$|\delta P(x)| = |\delta b_0| + |\delta b_1| |x| + \dots + |\delta b_n| |x|^n, \quad \text{com}$$

$$|\delta b_i| \leq |b_i| |\varepsilon_i|, \quad \varepsilon_i \text{ dado em (3.59)}$$

e, portanto, encontramos uma majoração para o erro cometido durante o processo de determinação do valor numérico.

(3.63) Exemplo - Em 4.2 veremos que, através da análise regressiva, pode-se mostrar que a solução aproximada \bar{x} , obtida pelo método de eliminação de Gauss com pivoteamento para o sistema linear $Ax = b$ é a solução exata de

$$(3.64) \quad (A + \delta A) \bar{x} = b, \quad \text{em que}$$

$$(3.65) \quad \|\delta A\|_{\infty} \leq 1,01 (n^3 + 3n^2) \rho \|A\|_{\infty} \mu,$$

onde n é a ordem de A , ρ é um certo número que pode ser calculado durante o processo de resolução e μ é a unidade relativa de arredondamento.

Como pode ser observado em (3.64) e (3.65), nesse caso, assim como em (3.51), a análise regressiva não fornece estimativas para $\|\bar{x} - x\|_{\infty}$. Entretanto, essas estimativas podem ser obtidas através do estudo da sensibilidade do problema. Assim, de (3.27) e (3.65) temos:

$$(3.66) \quad \|\bar{x} - x\|_{\infty} \leq \text{cond}_{\infty}(A) \cdot 1,01 (n^3 + 3n^2) \rho \|\bar{x}\|_{\infty} \mu$$

(3.67) Exemplo - Consideremos a solução de um sistema de equações lineares como visto em (3.63).

Se após a aplicação da análise regressiva obtivermos $\|\delta A\| < \epsilon$, com ϵ pré-estabelecido, e se soubermos previamente que os elementos de A estão sujeitos a um erro maior que ϵ , então, não teremos necessidade de efetuar a análise de sensibilidade do processo, pois já saberemos, de antemão, que a aproximação obtida é uma resposta satisfatória para o problema em questão. Portanto, nesse caso, a análise regressiva não representa uma etapa intermediária no estudo dos erros.

(3.68) Exemplo - Suponhamos que um vetor x represente as diversas posições ocupadas por um foguete que esteja realizando uma viagem a Marte, tendo a Terra como ponto de partida. Seja Ax a representação da derivada da quantidade de movimento do foguete nos mesmos instantes. Finalmente, suponhamos que b represente, nos mesmos instantes, a força resultante sobre o foguete, incluindo ambas as forças, as externas e aquelas de pro-

pulsão do próprio foguete. Então, o sistema de equações lineares $Ax = b$ representa o equilíbrio das forças reais e inerciais, de modo que, para qualquer b dado, a trajetória será representada pela solução exata x . Considerar a solução aproximada \bar{x} como solução exata do sistema $A\bar{x} = b + \delta b$ significa afirmar que o foguete pode ser conservado na trajetória \bar{x} através da aplicação de forças adicionais δb ; consequentemente, se soubermos que $\|\delta b\| < \epsilon$, teremos uma majoração para as forças corretivas necessárias. Ora, se ϵ for seguramente menor do que a grandeza das forças corretivas disponíveis para o foguete na forma de impulso de reserva, então, poderemos estar certos de que a trajetória \bar{x} será alcançada e que, portanto, os erros de arredondamento que tenham sido cometidos durante a determinação de \bar{x} serão compensados por uma disponível mudança de b . Portanto, nesse caso, é mais interessante examinar a grandeza de δb do que procurar estimar $\|\bar{x} - x\|$ para um b fixo.

Utilizando os exemplos que acabamos de apresentar, veremos, a seguir, mais algumas características da análise regressiva.

A análise inversa tem a vantagem, relativamente à direta, de permitir uma comparação entre o efeito dos erros de arredondamento cometidos durante a execução do processo numérico e o dos erros inerentes aos dados; fornece, portanto, uma preciosa avaliação da eficácia do processo considerado. Assim, por exemplo, o processo numérico utilizado em (3.63) poderá ser considerado bom se $\frac{\|\delta A\|_\infty}{\|A\|_\infty}$ for pequeno

Ocasionalmente a análise regressiva fornece estimativas "a priori" para o erro de arredondamento,

o exemplo do que ocorre em relação ao cálculo dos autovalores de matrizes tridiagonais pelo método de bissecções, como mostra Wilkinson (24, p. 323-327). De maneira geral, porém, a determinação dessas estimativas é feita "a posteriori" e exige cálculos que, frequentemente, são difíceis de efetuar. Assim, por exemplo, em (3.63) a estimativa de $\|\bar{x} - x\|_\infty$ exige a avaliação de $\text{cond}(A)$, o que exige, pelo menos, uma estimativa de $\|A^{-1}\|_\infty$. Apesar disso, é razoável supor-se que o processo numérico que determina uma estimativa razoável para $\| \delta A \|_\infty$ é melhor do que aquele para o qual isso não é possível.

3.3 - Evolução histórica

O primeiro tratamento rigoroso dos erros de arredondamento nas operações aritméticas foi efetuado por J. von Neumann e H. Goldstine no artigo "Numerical inverting of matrix of high order" publicado no Bulletin of the American Mathematical Society de número 53, em 1947, às páginas 1021-1099. Nesse artigo von Neumann e Goldstine introduziram as desigualdades que devem ser satisfeitas pelos erros de arredondamento nas operações aritméticas fundamentais realizadas em vírgula fixa e estabeleceram um padrão básico para a análise rigorosa de um algoritmo de cálculo. Tanto a notação quanto as idéias gerais contidas nesse trabalho passaram a ser adotadas pelos autores que os sucederam.

Os erros de arredondamento nas operações aritméticas fundamentais realizadas em vírgula flutuante foram estudados principalmente por James H. Wilkinson. Desde 1957, quando introduziu as desigualdades satisfeitas por esses erros, tem o referido autor contribuído significativamente para o desenvolvimento da análise de erros na aritmética de vír-

gula móvel. Em 1960 essas desigualdades, já em forma mais aperfeiçoada, foram rerepresentadas por Wilkinson em um artigo em que foram aplicadas a diferentes tipos de problemas (24). Além de inúmeros trabalhos, publicados em diferentes revistas especializadas do ramo, esse autor publicou dois livros, (26) e (22), em 1963 e 1965, respectivamente, cujo conteúdo é básico a todos que pretendam trabalhar nesse campo da ciência.

A idéia da análise regressiva de erros estava, até certo ponto, implícita nos artigos de von Neumann e Goldstine, porém, somente foi descrita explicitamente por J. W. Givens no artigo "Numerical computation of the characteristic values of a real symmetric matrix", publicado pelo Oak Ridge National Laboratory, do Tennessee, sob a especificação Rep. ORNL-1574, em 1954. Nesse artigo Givens analisou a redução de uma matriz real simétrica à forma tridiagonal por rotações planas e o cálculo dos autovalores de uma matriz tridiagonal através do processo da sequência de Sturm. Ambas as análises foram realizadas na aritmética de vírgula fixa; a primeira delas não trouxe contribuição significativa para os conhecimentos da época, porém, a segunda tomou-se de fundamental importância no desenvolvimento da análise de erros. Entretanto, essa análise, ao ser publicada, não recebeu a devida atenção, o que se deu em consequência de dois fatores principais: o primeiro foi o fato de sua apresentação não ter sido feita através de uma revista especializada; o segundo foi consequência da mesma ter sido precedida por exaustiva análise de erros que só foi compreendida por reduzido número de pessoas. A dificuldade de entendimento dessa análise introdutória resultou do uso da aritmética de vírgula fixa, o que tornou as demonstrações muito trabalhosas, em virtude da necessidade de introdução de vários fatores de escala. Posteriormente, Wilkinson (25, p. 552-554) mostrou que, na aritmética

de vírgula flutuante, essas demonstrações se tornam suscintas e claras.

A análise regressiva de erros constitui um traço característico dos trabalhos de Wilkinson, o qual contribuiu de forma decisiva para o aperfeiçoamento e divulgação dessa forma de análise de erros.

O termo "número de condição" parece ter sido inicialmente empregado por Turing, entretanto, o termo "mal condicionamento" já era de uso comum entre os analistas numéricos há algum tempo antes disso.

Atualmente os erros de arredondamento continuam sendo muito estudados e freqüentemente são publicados trabalhos e realizados congressos visando o seu desenvolvimento.

4. Análise de erros do método de eliminação de Gauss. Aplicações

Neste capítulo é apresentada uma análise regressiva de erros do método de eliminação de Gauss, e são considerados alguns problemas que podem ser resolvidos pela aplicação do mesmo: cálculo de determinantes, inversões de matrizes e refinamento da solução aproximada. Também são feitas algumas considerações a respeito de matrizes que, por suas formas especiais, apresentam simplificações em relação à eliminação de Gauss. Na parte inicial considera-se a resolução de sistemas triangulares, visto que, pela aplicação do método em questão, um sistema geral é transformado em dois triangulares. A norma adotada em todas essas aplicações é a norma ∞ , em virtude de ser de cálculo fácil e por apresentar algumas vantagens sobre as demais, como vimos em (1.20)-(1.22), (1.32), (1.33) e (1.38). O uso da "barra", para indicar o valor aproximado, é dispensado em todo o capítulo, já que a análise regressiva só utiliza os valores calculados, não efetuando comparações entre esses e os verdadeiros. Em alguns poucos pontos em que esses últimos são usados, fizemos observações a respeito e os indicamos de forma a evitar confusões.

4.1. Resolução de sistemas triangulares

Consideremos o sistema de equações lineares

$$(4.1) \quad Lx = b \quad ,$$

em que $L = (l_{ij})$ é uma matriz triangular inferior pertencente a $M_n(\mathbb{R})$. As componentes do vetor solução x são deter-

arranjados na ordem x_1, \dots, x_n , por:

$$(4.2) \quad \begin{cases} x_1 \approx fl \left(\frac{b_1}{l_{11}} \right) \\ x_i \approx fl \left(\frac{-l_{i1} x_1 - l_{i2} x_2 \dots - l_{i,i-1} x_{i-1} + b_i}{l_{ii}} \right), \quad i=2, \dots, n \end{cases}$$

Aplicando (2.113) e (2.114) a (4.2)

$$(4.3) \quad \begin{cases} x_1 = \frac{b_1}{l_{11}(1+\delta_{11})} \\ x_i = \frac{-l_{i1} x_1 (1+\delta_{i1}) - \dots - l_{i,i-1} x_{i-1} (1+\delta_{i,i-1}) x_{i-1} + b_i}{l_{ii}(1+\delta_{ii})(1+\delta_{ii}^2)}, \\ i=2, \dots, n \end{cases}$$

$$(4.4) \quad \begin{cases} |\delta_{ii}|, |\delta'_{ii}| \leq \mu, \quad i=1, \dots, n \\ |\delta_{i1}| \leq (i-1) 1,01 \mu, \quad i=2, \dots, n \\ |\delta_{ij}| \leq (i+j-1) 1,01 \mu, \quad i=2, \dots, n; \quad j=2, \dots, i-1 \end{cases}$$

Em (4.2) foram usadas as relações (2.113) e (2.114), de modo a se obter que os fatores $(1+\delta)$ multipliquem elementos de L .

As equações (4.3) também podem ser escritas

$$(4.5) \quad \begin{cases} l_{11}(1+\delta_{11})x_1 = b_1 \\ l_{i1}(1+\delta_{i1})x_1 + \dots + l_{ij}(1+\delta_{ij})x_{j-1} + l_{ii}(1+\delta_{ii})(1+\delta'_{ii})x_i = b_i, \\ i=2, \dots, n \end{cases}$$

originando

$$(4.6) \quad (L + \delta L)x = b, \quad \text{onde, por (4.4),}$$

$$(4.7) \quad |\delta L| \leq 1,01 \mu \begin{pmatrix} |l_{11}| & & & & & & & & & & \\ & |l_{21}| & 2|l_{22}| & & & & & & & & \\ & 2|l_{31}| & 2|l_{32}| & 2|l_{33}| & & & & & & & \\ & 3|l_{41}| & 3|l_{42}| & 2|l_{43}| & 2|l_{44}| & & & & & & \\ & \vdots & \vdots & \vdots & \vdots & \ddots & & & & & \\ & (n-1)|l_{n1}| & (n-1)|l_{n2}| & (n-2)|l_{n3}| & (n-3)|l_{n4}| & \cdots & 2|l_{nn}| & & & & \end{pmatrix}$$

Conseqüentemente, por (1.38)

$$(4.8) \quad \|\delta L\|_{\infty} \leq \frac{n(n+1)}{2} 1,01 \mu \max_{i,j} |l_{ij}|$$

Por outro lado, de (4.7) também obtemos:

$$|\delta L| \leq 1,01 m \mu \|L\|$$

e, portanto, por (1.38),

$$(4.9) \quad \|\delta L\|_{\infty} \leq 1,01 m \mu \|L\|_{\infty}$$

Reunindo todas as conclusões até agora obtidas para sistemas triangulares inferiores, temos:

(4.10) Teorema - A solução aproximada do sistema triangular inferior $Lx = b$, determinada através de (4.2) e da aritmética de vírgula flutuante com unidade de arredondamento μ , é a solução exata de um sistema $(L + \delta L)x = b$, em que a perturbação δL verifica (4.7) - (4.9).

Para sistemas de equações lineares

$$(4.11) \quad Ux = b,$$

em que $U = (u_{ij})$ é uma matriz triangular superior por

tenente a $M_n(\mathbb{R})$, a determinação das componentes do vetor solução x é feita na ordem x_n, x_{n-1}, \dots, x_1 , por:

$$(4.12) \quad \begin{cases} x_n = fl\left(\frac{b_n}{\mu_{nn}}\right) \\ x_r = fl\left(\frac{-\mu_{r,r+1}x_{r+1} - \mu_{r,r+2}x_{r+2} - \dots - \mu_{rn}x_n + b_r}{\mu_{rr}}\right), \\ r = n-1, n-2, \dots, 1 \end{cases}$$

Procedendo de forma análoga à utilizada para sistemas triangulares superiores obtemos:

(4.13) Teorema - A solução aproximada do sistema triangular superior $Ux = b$, obtida através de (4.12) e do uso da aritmética de régua flutuante com unidade de arredondamento μ é a solução exata de um sistema $(U + \delta U)x = b$, em que a perturbação δU satisfaz (4.14) - (4.16).

$$(4.14) \quad |\delta U| \leq 1,01 \mu \begin{pmatrix} 2|\mu_{11}| & (n-1)|\mu_{12}| & (n-1)|\mu_{13}| & (n-2)|\mu_{14}| & \dots & 3\mu_{1,n-1} & 2|\mu_{1n}| \\ 0 & 2|\mu_{22}| & (n-2)|\mu_{23}| & (n-2)|\mu_{24}| & \dots & 3\mu_{2,n-1} & 2|\mu_{2n}| \\ 0 & 0 & 2|\mu_{33}| & (n-3)|\mu_{34}| & \dots & 3\mu_{3,n-1} & 2|\mu_{3n}| \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 2\mu_{n-1,n-1} & |\mu_{nm}| \\ 0 & 0 & 0 & 0 & \dots & 0 & |\mu_{nn}| \end{pmatrix}$$

$$(4.15) \quad \|\delta U\|_{\infty} \leq \frac{n(n+1)}{2} 1,01 \mu \max_{i,j} |\mu_{ij}|$$

$$|\delta U| \leq 1,01 n \mu |U|$$

$$(4.16) \quad \|\delta U\|_{\infty} \leq 1,01 n \mu \|U\|_{\infty}$$

Poderíamos utilizar acumulação de produtos escalares no cálculo de (4.2) e (4.12), com a finalidade de obter uma solução aproximada mais precisa. Não adotamos essa atitude porque, se ditarmos que o tempo adicional de uso de computador requerido não seria justificável, pois, de maneira geral, os sistemas triangulares considerados possuem coeficientes afetados de erros, em virtude de resultarem da aplicação do método de Gauss a sistemas gerais. Por outro lado, as soluções aproximadas determinadas através de (4.2) e (4.12) são, usualmente, muito precisas, ainda que o sistema considerado seja muito mal condicionado, como veremos a seguir. O uso da acumulação de produtos escalares, para a resolução de sistemas triangulares, é considerado por Wilkinson (26, p. 103-104), (22, p. 247-249).

Seja z a solução exata de (4.1) e seja x uma aproximação de z , obtida pela forma descrita em (4.10). Consideremos o erro relativo $\frac{\|x-z\|}{\|z\|}$, em função do erro relativo $\frac{\|\delta b\|}{\|b\|}$. Por (3.22) e (4.9):

$$\frac{\|x-z\|_{\infty}}{\|z\|_{\infty}} \leq \frac{\text{cond}_{\infty}(L) \frac{1,01 n u \|L\|_{\infty}}{\|L\|_{\infty}}}{1 - \text{cond}_{\infty}(L) \frac{1,01 n u \|L\|_{\infty}}{\|L\|_{\infty}}}$$

$$(4.17) \quad \frac{\|x-z\|_{\infty}}{\|z\|_{\infty}} \leq \frac{1,01 n u \text{cond}_{\infty}(L)}{1 - 1,01 n u \text{cond}_{\infty}(L)}$$

Em virtude de (4.17), podemos afirmar que, se $1,01 n u \text{cond}_{\infty}(L) \ll 1$, o erro relativo da solução é pequeno. Por outro lado, Wilkinson (26, p. 105-107), (22, p. 249-251) afirma que as soluções aproximadas determinadas pelo processo descrito são, freqüentemente, muito mais precisas do que a delimitação (4.17) pode sugerir, mesmo quando L é mal condicionada. Afirma, ainda, que para algumas classes de matrizes o erro relativo não depende de $\text{cond}_{\infty}(L)$.

Consideremos o vetor resíduo $b - Lx = \delta Lx$, que fornece:

$$\|b - Lx\|_{\infty} \leq \|\delta L\|_{\infty} \|x\|_{\infty}$$

Em virtude de (4.8):

$$(4.18) \quad \|b - Lx\|_{\infty} \leq \frac{n(n+1)}{2} 1,01 u \max_{i,j} |l_{ij}| \cdot \|x\|_{\infty}$$

A desigualdade (4.18) fornece uma majoração para o elemento de maior valor absoluto do vetor resíduo, em função do elemento de maior valor absoluto de L e de x , respectivamente. Uma análise dessa majoração permite-nos concluir que a existência de resíduos pequenos não implica, necessariamente, na obtenção de uma solução precisa; se, por exemplo, L for tal que $|l_{ij}| \leq 1$ ($i, j = 1, \dots, n$) e se $\|x\|_{\infty}$ também for aproximadamente 1, então, o resíduo será pequeno quer x seja ou não uma solução precisa.

4.2 - Eliminação de Gauss

Os métodos numéricos utilizados para a resolução de sistemas de equações lineares podem ser divididos em dois grupos, diretos e iterativos. Recebem a denominação de diretos ou exatos os métodos que, com um número finito de operações aritméticas elementares, forneceriam a solução exata de um sistema dado, não fosse a ocorrência de erros de arredondamento durante a sua execução. São chamados iterativos os métodos que possibilitam a resolução de um sistema linear de forma aproximada; a solução é obtida como o limite de sucessivas aproximações, calculadas por algum processo uniforme. Os métodos diretos são recomendáveis para sistemas cuja matriz é densa e os iterativos para sistemas cuja matriz é esparsa ou rarefeita. Con-

sideraremos apenas os primeiros, os quais se baseiam na idéia de eliminação atribuída a Gauss, que será considerada a seguir.

Seja o sistema de equações lineares

$$(4.19) \quad Ax = b,$$

em que A é a matriz dos coeficientes, x o vetor solução e b o vetor dos termos independentes. Suponhamos $A \in M_n(\mathbb{R})$ e A não singular.

São conhecidas diversas variantes do método de eliminação de Gauss, todas equivalentes à decomposição da matriz A em um produto LU de duas matrizes triangulares, L triangular inferior e U triangular superior. Descreveremos aquela em que os elementos de L são fixados e iguais a 1. Obtida essa fatoração LU , denominada decomposição triangular ou decomposição LU , o sistema (4.19) pode ser escrito

$$(4.20) \quad LUx = b, \text{ com}$$

$$(4.21) \quad Ly = b$$

$$(4.22) \quad Ux = y$$

e, portanto, a resolução de (4.19) fica reduzida à resolução dos sistemas triangulares (4.21) e (4.22), o que é simples, como vimos em 4.1.

A decomposição triangular é conseguida através da determinação de uma sequência $A^{(1)} = A, A^{(2)}, \dots, A^{(n)}$ de matrizes tais que os elementos das $(k-1)$ primeiras colunas de $A^{(k)}$, situados abaixo da diagonal principal, são nulos. A matriz $A^{(k+1)}$ é obtida de $A^{(k)}$ pela subtração de um múltiplo da k -ésima linha de cada uma das linhas que lhe estão abaixo (linhas $k+1, k+2, \dots, n$). Cada um desses múltiplos é determinado através de um fator, deno-

minado multiplicador de Gauss, escolhido de forma que seriam nulos os elementos da k -ésima coluna de $A^{(k+1)}$, situados abaixo da diagonal principal, caso não houvesse a ocorrência de erros de arredondamento no processo; de maneira geral, esses elementos não são calculados, mas são tomados iguais a zero por definição. Seja $A^{(k)} = (a_{ij}^{(k)})$, $i, j = 1, \dots, n$. A k -ésima linha de $A^{(k)}$ é chamada k -ésima linha pivô e o elemento $a_{kk}^{(k)}$, k -ésimo pivô. $A^{(1)}, A^{(2)}, \dots, A^{(n)}$ são chamadas primeira, segunda, ..., k -ésima, ..., n -ésima matriz reduzida, respectivamente. Indicaremos por $\bar{A}^{(k)}$ a matriz constituída pelas linhas e colunas $k, (k+1), \dots, n$ de $A^{(k)}$, na mesma ordem em que se encontram em $A^{(k)}$.

Os multiplicadores de Gauss são dados por

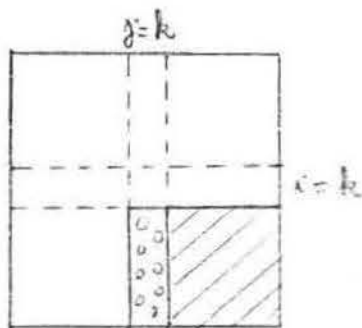
$$(4.23) \quad m_{ik} = fl \left(\frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \right), \quad i \geq k+1, \quad k = 1, 2, \dots, n$$

e os elementos de $A^{(k+1)}$ por

$$(4.24) \quad a_{ij}^{(k+1)} = \begin{cases} 0, & i \geq k+1, j = k \\ fl \left(a_{ij}^{(k)} - m_{ik} \times a_{kj}^{(k)} \right), & i \geq k+1, j \geq k+1 \\ a_{ij}^{(k)}, & \text{nos demais casos} \end{cases}$$

em que $k = 1, 2, \dots, n$.

As três regiões de $A^{(k+1)}$, envolvidas nos três diferentes casos de (4.24), são ilustrados no diagrama abaixo, encontrado em Forsythe (7, p. 99).



Em virtude da ocorrência de erros de arredondamento, em geral não é obtida a decomposição triangular exata de A . Mostraremos, entretanto, que as matrizes L e U assim determinadas constituem a decomposição triangular exata de uma matriz obtida por ligeira perturbação de A , ou seja,

$$LU = A + E,$$

e daremos delimitações para $|E|$.

Consideremos, inicialmente, a obtenção dos elementos $a_{ij}^{(k+1)}$, dados em (4.24). Por (4.23) e (2.113), podemos afirmar que, para os multiplicadores, são encontrados os valores

$$(4.25) \quad m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} (1 + \delta_{ik}), \quad |\delta_{ik}| \leq \mu, \quad \text{ou,}$$

$$(4.26) \quad 0 = a_{ik}^{(k)} - m_{ik} a_{kk}^{(k)} + E_{ik}^{(k)},$$

$$(4.27) \quad E_{ik}^{(k)} = a_{ik}^{(k)} \delta_{ik},$$

Portanto, por (4.26), podemos afirmar que $a_{ik}^{(k+1)}$ ($i \geq k+1$) pode ser tomado exatamente igual a zero desde que $a_{ik}^{(k)}$ seja perturbado de $E_{ik}^{(k)}$, em que $E_{ik}^{(k)}$ é dado em (4.27)

Por outro lado, para $i \geq k+1$ e $j \geq k+1$, (4.24) fornece

$$(4.28) \quad a_{ij}^{(k+1)} = fl(a_{ij}^{(k)} - fl(m_{ik} \times a_{kj}^{(k)}))$$

Aplicando (2.113) e (2.114) a (4.28) obtemos:

$$(4.29) \quad a_{ij}^{(k+1)} = \frac{a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} (1 + \delta_{ij})}{1 + \delta'_{ij}}; \quad |\delta_{ij}|, |\delta'_{ij}| \leq \mu,$$

que fornece

$$(4.30) \quad a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} + E_{ij}^{(k)}, \quad \text{com}$$

$$(4.31) \quad \varepsilon_{ij}^{(k)} = -m_{ik} a_{kj}^{(k)} \delta_{ij}^{(k)} - a_{ij}^{(k+1)} \delta_{ij}^{(k)}$$

Logo, por (4.30) podemos concluir que, para $i \geq k+1$ e $j \geq k+1$, o valor $a_{ij}^{(k+1)}$ seria obtido exatamente na decomposição triangular regular, a partir de elementos $a_{ij}^{(k)}$ perturbados pelos $\varepsilon_{ij}^{(k)}$ dados em (4.31).

Finalmente, os elementos restantes de $A^{(k+1)}$ são obtidos exatamente, pois os $a_{ij}^{(k+1)}$ são iguais aos $a_{ij}^{(k)}$.

Resumindo, $A^{(k+1)}$ pode ser obtida de forma exata na decomposição triangular, a partir da matriz resultante da perturbação de $A^{(k)}$ por

$$(4.32) \quad E = \left(\varepsilon_{ij}^{(k)} \right), \quad \text{com}$$

$$(4.33) \quad \varepsilon_{ij}^{(k)} = \begin{cases} a_{ik}^{(k)} \delta_{ik}^{(k)} & , \text{ para } i \geq k+1, j = k \\ -m_{ik} a_{kj}^{(k)} \delta_{ij}^{(k)} - a_{ij}^{(k+1)} \delta_{ij}^{(k)} & , \text{ para } i \geq k+1, j \geq k+1 \\ 0 & , \text{ nos demais casos} \end{cases}$$

Seja

$$(4.34) \quad L = \begin{pmatrix} \bigcirc & 0 & & & \\ & \bigcirc & 0 & & \\ & & \vdots & & \\ & & & \bigcirc & \\ & & & & m_{k+1,k} \\ & & & & m_{k+2,k} \\ & & & & \vdots \\ & & & & m_{m,k} \end{pmatrix}$$

Então, a k -ésima etapa da decomposição triangular, que produzirá a matriz $A^{(k+1)}$, será completamente descrita pela equação (4.35), incluindo os erros de arredondamento.

$$(4.35) \quad A^{(k+1)} = A^{(k)} - L A^{(k)} + E$$

Somando as equações (4.35) para $k = 1, 2, \dots, n-1$

temos

$$L^{(1)} A^{(1)} + L^{(2)} A^{(2)} + \dots + L^{(n-1)} A^{(n-1)} + L^{(n)} A^{(n)} = A^{(1)} + E^{(1)} + E^{(2)} + \dots + E^{(n-1)}$$

Como a matriz $L^{(k)} A^{(k)}$ depende somente da k -ésima linha de $A^{(k)}$, que é igual à k -ésima linha de $A^{(n)}$, temos:

$$(L^{(1)} + L^{(2)} + \dots + L^{(n-1)} + I) A^{(n)} = A^{(1)} + E^{(1)} + E^{(2)} + \dots + E^{(n-1)}$$

$$(4.36) \quad LU = A + E, \text{ com}$$

$$(4.37) \quad L = L^{(1)} + L^{(2)} + \dots + L^{(n-1)} + I$$

$$(4.38) \quad U = A^{(n)}$$

$$(4.39) \quad A = A^{(1)}$$

$$(4.40) \quad E = E^{(1)} + E^{(2)} + \dots + E^{(n-1)}$$

e que facilmente se verifica que L e U são matrizes triangulares inferior e superior, respectivamente, sendo L diagonal unitária e igual a

$$L = \begin{pmatrix} 1 & & & & 0 \\ m_{21} & 1 & & & \\ m_{31} & m_{32} & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & m_{n3} & \dots & 1 \end{pmatrix}$$

Como pode ser observado em (4.35), as perturbações (de arredondamento) poderão ser grandes se ou os multiplicadores ou os elementos das matrizes reduzidas, ou ambos, forem grandes. Por outro lado, o processo não funcionará se o elemento pivô for muito pequeno. Esse último problema é resolvido pela troca de linhas pivô com outra, de modo que o novo elemento pivô seja diferente de zero (certamente existe um elemento

não nulo na coluna k , caso contrário o determinante de A seria nulo e A singular, contrariando nossa hipótese inicial). A grandeza dos multiplicadores é controlada por uma escolha conveniente do pivô. Na prática, de maneira geral, essa escolha é feita de modo a assegurar

$$(4.42) \quad |m_{ij}^{(k)}| \leq 1, \quad i \geq j+1, \quad j = 1, 2, \dots, n-1$$

e é obtida através de duas técnicas principais. Na primeira delas, denominada pivotamento parcial, é escolhido para k -ésimo pivô o elemento de maior valor absoluto da primeira coluna de $\bar{A}^{(k)}$. Na segunda, denominada pivotamento completo, é eleito k -ésimo pivô o elemento de maior valor absoluto em toda a matriz $\bar{A}^{(k)}$.

O uso de pivotamento, com a finalidade de assegurar (4.42), permite, também, a determinação de uma estimativa para a grandeza dos elementos das matrizes reduzidas. De fato, facilmente pode-se verificar que, quando é utilizado pivotamento parcial, se os elementos de $A^{(1)}$ satisfazem a

$$(4.43) \quad |a_{ij}^{(1)}| \leq a,$$

então, os elementos de $A^{(k)}$ verificam

$$(4.44) \quad |a_{ij}^{(k)}| \leq 2^{(k-1)} a.$$

Por outro lado, quando é utilizado o pivotamento completo, Wilkinson (23, p. 283-284) prova que, se (4.43) se verifica, então, os elementos de $A^{(k)}$ satisfazem a

$$(4.45) \quad |a_{ij}^{(k)}| < a k^{\frac{1}{2}} \left(2^{\frac{1}{2}} \cdot 3^{\frac{1}{2}} \cdot 4^{\frac{1}{3}} \cdots k^{\frac{1}{k-1}} \right)^{\frac{1}{2}}.$$

Para algumas classes de matrizes especiais podem ser obtidas delimitações melhores do que as dadas em (4.44) e (4.45), a exemplo do que ocorre com as matrizes de Hessenberg, como veremos em 4.5.

Relativamente às delimitações para os elementos das matrizes reduzidas, dadas em (4.44) e (4.45), Wilkinson,

(23, p. 283-284) e (26, p. 97), afirma ter conhecimento de algumas matrizes especiais para as quais se verifica a igualdade em (4.44), mas que, na prática, de maneira geral (4.44) é substituída por

$$(4.46) \quad |a_{ij}^{(k)}| \leq \delta a$$

e que, no caso de matrizes extremamente mal condicionadas, geralmente os elementos das sucessivas $A^{(k)}$ decrescem em valor absoluto. Afirma, ainda, desconhecer matrizes que, durante a decomposição triangular com pivotamento completo, originem elementos para os quais

$$(4.47) \quad |a_{ij}^{(k)}| > ka,$$

o que mostra que a majoração (4.45) deve ser exagerada.

De maneira geral, o pivotamento parcial é o único utilizado na prática. Esse fato decorre, principalmente, da grande dificuldade de programação do pivotamento completo para computadores, dificuldade essa não compensada por grande melhoria das delimitações de erros. Por outro lado, o pivotamento parcial conserva a disposição particular dos elementos de algumas matrizes que possuem muitos elementos nulos, enquanto o pivotamento completo a destrói; essa propriedade é muito importante para a resolução de certos problemas, a exemplo do que ocorre relativamente ao cálculo de autovalores e autovetores de matrizes.

Utilizaremos o pivotamento parcial em todo o estudo de erros que efetuaremos a seguir. Ao iniciar a fatoração LU de uma matriz A , com aplicação dessa forma de pivotamento, admitiremos que as linhas de A tenham sido inicialmente permutadas, de modo que nenhuma troca de linhas se torne necessária durante a decomposição triangular. Assim procedendo, poderemos utilizar a descrição do método de Gauss apresentada anteriormente, assim como

as diversas conclusões obtidas para as matrizes reduzidas e para as matrizes L e U . Embora, na prática, a referida permutação inicial não seja conhecida "a priori", a sua determinação não é necessária, já que permutações de linhas não originam erros de arredondamento. Simplificaremos a notação indicando a matriz reordenada e a matriz original pelo mesmo símbolo A . É lógico que o determinante da matriz reordenada é igual ao da matriz original, a menos do fator (± 1) .

Determinaremos, a seguir, uma delimitação para $|E|$, em que E é dada em (4.40), no caso de utilização do pivotamento parcial. Seja

$$(4.48) \quad \rho = \frac{\max_{i,j,k} |a_{ij}^{(k)}|}{\|A\|_{\infty}} \quad ; \quad i, j, k = 1, 2, \dots, n$$

Embora não se tenha uma boa delimitação para ρ "a priori", pode-se calculá-lo de forma relativamente fácil durante a decomposição triangular.

A partir de (4.48) obtemos

$$(4.49) \quad |a_{ij}^{(k)}| \leq \rho \|A\|_{\infty} \quad ; \quad i, j, k = 1, 2, \dots, n$$

que, juntamente com (4.33), (4.29) e (4.42), fornece

$$(4.50) \quad |E_{ij}^{(k)}| \leq \rho \|A\|_{\infty} \begin{cases} u & , \quad i \geq k+1, \quad j = k \\ 2u & , \quad i \geq k+1, \quad j \geq k+1 \\ 0 & , \quad \text{nos demais casos} \end{cases}$$

As relações (4.50) e (4.40) nos dão, finalmente, a majoração para $|E|$, que é apresentada em (4.51). A matriz à direita dessa desigualdade é facilmente reconstruída, já que seus elementos representam o número de operações aritméticas necessárias à determinação das matrizes L e U .

$$(4.51) \quad |E| \leq \rho \|A\|_{\infty} \mu \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ 1 & 2 & 2 & \cdots & 2 & 2 \\ 1 & 3 & 4 & \cdots & 4 & 4 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 3 & 5 & \cdots & 2n-4 & 2n-4 \\ 1 & 3 & 5 & \cdots & 2n-3 & 2n-2 \end{pmatrix}$$

Por (4.51), (4.38) e (4.30)

$$\|E\|_{\infty} \leq \rho \|A\|_{\infty} \mu \left\{ \left[\sum_{j=1}^m (2j-1) \right] - 1 \right\}, \text{ ou,}$$

$$(4.52) \quad \|E\|_{\infty} \leq m^2 \rho \|A\|_{\infty} \mu$$

Concluindo, temos

(4.53) Teorema - Seja $A \in M_n(\mathbb{R})$, A não singular. Consideremos a decomposição LU de A , através da eliminação de Gauss com pivotamento, na aritmética de máquina flutuante, com unidade de arredondamento μ . As matrizes L e U assim determinadas satisfazem a $LU = A + E$, em que E verifica (4.51) e (4.52).

Uma vez efetuada a decomposição triangular de A , a solução aproximada de (4.19) pode ser obtida pela resolução dos sistemas triangulares (4.21) e (4.22). Aplicando (4.30) e (4.13) a esses sistemas, temos, por (4.20), que a aproximação x é a solução exata do sistema

$$(4.54) \quad (L + \delta L)(U + \delta U)x = b,$$

com $\|\delta L\|_{\infty}$ e $\|\delta U\|_{\infty}$ dados em (4.8) e (4.15). Por outro lado, em virtude de (4.36) e (4.54), temos

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

ou ainda

$$(4.55) \quad (A + \delta A) x = b, \text{ com}$$

$$(4.56) \quad \delta A = E + L(\delta U) + (\delta L)U + (\delta L)(\delta U)$$

Como, em consequência de (4.41), (4.42), (4.38) e (4.49)

$$(4.57) \quad |l_{ij}| \leq 1$$

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

$$(4.58) \quad |u_{ij}| \leq \rho \|A\|_{\infty},$$

temos

$$(4.59) \quad \|L\|_{\infty} \leq n$$

$$(4.60) \quad \|U\|_{\infty} \leq m \rho \|A\|_{\infty}$$

e, ainda, utilizando (4.57) e (4.58) em (4.8) e (4.15), respectivamente

$$(4.61) \quad \|\delta L\|_{\infty} \leq \frac{n(n+1)}{2} 1,01 u$$

$$(4.62) \quad \|\delta U\|_{\infty} \leq \frac{n(n+1)}{2} 1,01 \rho \|A\|_{\infty} u$$

Em virtude de $nu \ll 1$ no caso em que esses resultados são utilizados, a partir (4.61) e (4.62) podemos afirmar

$$(4.63) \quad \|\delta L\|_{\infty} \|\delta U\|_{\infty} \leq n^2 \rho \|A\|_{\infty} u$$

Duplicando (4.52), (4.59), (4.61), (4.60), (4.62) e (4.63) e (4.56) obtemos

$$(4.64) \quad \|\delta A\|_{\infty} \leq 1,01 (n^3 + 3n^2) \rho \|A\|_{\infty} u$$

É importante observar que, embora δA dependa de b , pois δU e δL dependem de b , a delimitação de $\|\delta A\|_{\infty}$ não envolve b . Wilkinson (36, p.108) afirma que, na prática, de maneira geral ocorre

$$(4.65) \quad \|\delta A\|_{\infty} \leq 1,01 m \rho \|A\|_{\infty} u$$

Resumindo todos esses resultados temos:

(4.66) Teorema - Seja $A \in M_n(\mathbb{R})$, A não singular. Consideremos a resolução do sistema de equações lineares $Ax = b$, através da eliminação de Gauss com pivotamento, na aritmética de vírgula flutuante, com unidade relativa de arredondamento u . A aproximação assim determinada é a solução exata do sistema $(A + \delta A)x = b$, em que δA , definida em (4.56), satisfaz (4.64).

Em virtude de (4.42) verificar-se para ambos os tipos de pivotamento, parcial e completo, podemos afirmar que (4.53) e (4.66) também valem para essa última forma de pivotamento, embora em nossas considerações só tenhamos mencionado a primeira. É claro, entretanto, que o valor de p depende do tipo de pivotamento considerado e do algoritmo utilizado para sua obtenção.

A seguir faremos algumas considerações a respeito da precisão da aproximação x , determinada da forma descrita em (4.66). Indicando por z a solução exata e utilizando (3.22) e (4.64), temos

$$(4.67) \quad \frac{\|x - z\|_\infty}{\|z\|_\infty} \leq \frac{1,01(n^3 + 3n^2)p \operatorname{cond}_\infty(A)u}{1 - 1,01(n^3 + 3n^2)p \operatorname{cond}_\infty(A)u}$$

e se $1,01(n^3 + 3n^2)p \operatorname{cond}_\infty(A)u \ll 1$, então, o erro absoluto será pequeno em comparação com a solução exata. Entretanto, o erro δA é formado pelos erros cometidos na decomposição triangular de A e por aqueles cometidos na resolução dos sistemas (4.21) e (4.22). Como as soluções aproximadas de sistemas triangulares são, de maneira geral, muito precisas, de modo que as delimitações dadas em (4.61) e (4.62) para $\|\delta L\|_\infty$ e $\|\delta U\|_\infty$, respectivamente, são dificilmente atingidas, podemos utilizar

(4.65) em substituição a (4.64) e concluir que, na prática, usualmente

$$(4.68) \quad \frac{\|x-z\|_{\infty}}{\|z\|_{\infty}} \leq \frac{1,01 n p \operatorname{cond}(A) u}{1 - 1,01 n p \operatorname{cond}(A) u}$$

e que, portanto, o erro absoluto é pequeno em comparação à solução aproximada, desde que $1,01 n p \operatorname{cond}(A) u \ll 1$.

Consideremos o resíduo da solução aproximada, $b - Ax$. Por (4.66)

$$(4.69) \quad \|b - Ax\|_{\infty} \leq 1,01 (n^3 + 3n^2) p \|A\|_{\infty} \|x\|_{\infty} u,$$

sendo que, por (4.65), podemos concluir que, frequentemente, (4.69) é substituído por (4.70)

$$(4.70) \quad \|b - Ax\|_{\infty} \leq 1,01 n p \|A\|_{\infty} \|x\|_{\infty} u$$

Através de (4.69) e (4.70) facilmente se pode concluir que o resíduo pode ser pequeno mesmo que a solução aproximada não seja muito precisa; isso pode ocorrer, por exemplo, quando $\|A\|_{\infty}$ é próxima de 1 e $\|x\|_{\infty}$ é de ordem unitária.

Wilkinson, (23, p. 284-285) e (22, p. 193), insiste em que se "equilibre" a matriz A , antes de aplicar ao sistema (4.19) qualquer método numérico de solução. De maneira grosseira, uma matriz é dita equilibrada quando possui todas as linhas e todas as colunas com aproximadamente o mesmo comprimento, relativamente a alguma norma. A análise de erros do método de Gauss fornece resultados mais efetivos quando aplicada a matrizes equilibradas porque, então, uma pequena perturbação em uma linha (ou coluna) da matriz é da mesma ordem de grandeza daquela de qualquer outra linha (ou coluna). Apesar dessa vantagem, não utilizamos matrizes equi-

libradas neste trabalho porque, infelizmente, o problema de equilibrção ainda não está bem esclarecido. Não existem muitos algoritmos de equilibrção e os poucos conhecidos não fornecem uma solução geral para o problema, pois às vezes valem apenas para algumas normas, enquanto que outras vezes valem apenas para determinadas matrizes especiais. Além disso, não existe uma forma única para se obter o equilíbrio de uma matriz e não se conhece, "a priori" qual das formas possíveis é a melhor; procedimentos de equilibrção diversos podem transformar A em diferentes matrizes equilibrçadas, cujos números de condição podem diferir consideravelmente entre si, como exemplificado em (7, p. 37-46) e (8, p. 169).

Apresentaremos, a seguir, um exemplo bastante simples, através do qual podem ser observados os resultados obtidos em 4.1 e 4.2.

(4.71) Exemplo - Consideremos a resolução do sistema de equações lineares $Ax = b$, com A e b dados em (4.72), através da eliminação de Gauss com pivotamento parcial, na aritmética decimal de vírgula flutuante, com armazenamento de apenas dois dígitos.

$$(4.72) \quad A = \begin{pmatrix} 0,69 & 0,26 & 0,12 \\ 0,15 & 0,45 & 0,27 \\ 0,31 & 0,49 & 0,15 \end{pmatrix} \quad b = \begin{pmatrix} 0,67 \\ 0,24 \\ 0,12 \end{pmatrix}$$

Vejam, inicialmente, a decomposição triangular de A . As matrizes reduzidas e demais matrizes citadas em (4.36) - (4.40) são, nesse caso,

$$(2) \quad A = \begin{pmatrix} 0,69 & 0,26 & 0,12 \\ 0,00 & 0,39 & 0,24 \\ 0,00 & 0,37 & 0,96 \times 10^{-1} \end{pmatrix}$$

$$(3) \quad A = \begin{pmatrix} 0,69 & 0,26 & 0,12 \\ 0,00 & 0,39 & 0,24 \\ 0,00 & 0,00 & -0,13 \end{pmatrix}$$

$$(1) \quad E = \begin{pmatrix} 0,00 & 0,00 & 0,00 \\ 0,18 \times 10^{-2} & -0,28 \times 10^{-2} & -0,36 \times 10^{-2} \\ 0,50 \times 10^{-3} & -0,30 \times 10^{-2} & 0,00 \end{pmatrix}$$

$$(2) \quad E = \begin{pmatrix} 0,00 & 0,00 & 0,00 \\ 0,00 & 0,00 & 0,00 \\ 0,00 & 0,50 \times 10^{-3} & 0,20 \times 10^{-2} \end{pmatrix}$$

$$(1) \quad L = \begin{pmatrix} 0,00 & 0,00 & 0,00 \\ 0,22 & 0,00 & 0,00 \\ 0,45 & 0,00 & 0,00 \end{pmatrix}$$

$$(2) \quad L = \begin{pmatrix} 0,00 & 0,00 & 0,00 \\ 0,00 & 0,00 & 0,00 \\ 0,00 & 0,95 & 0,00 \end{pmatrix}$$

Portanto:

$$(4.73) \quad L = \begin{pmatrix} 0,10 \times 10 & 0,00 & 0,00 \\ 0,22 & 0,10 \times 10 & 0,00 \\ 0,45 & 0,95 & 0,10 \times 10 \end{pmatrix}$$

$$U = \begin{pmatrix} 0,69 & 0,26 & 0,12 \\ 0,00 & 0,39 & 0,24 \\ 0,00 & 0,00 & -0,13 \end{pmatrix}$$

$$E = \begin{pmatrix} 0,00 & 0,00 & 0,00 \\ 0,18 \times 10^{-2} & -0,28 \times 10^{-2} & -0,36 \times 10^{-2} \\ 0,50 \times 10^{-3} & -0,25 \times 10^{-2} & 0,20 \times 10^{-2} \end{pmatrix}$$

$$LU = \begin{pmatrix} 0,69 & 0,26 & 0,12 \\ 0,1518 & 0,4472 & 0,2664 \\ 0,3105 & 0,4875 & 0,1520 \end{pmatrix}$$

Observe-se que $LU = A + E$

Consideremos, a seguir, a resolução dos sistemas triangulares obtidos

$$\left. \begin{array}{l} \text{solução} \\ \text{aproximada} \\ \text{de } Lx = b \end{array} \right\} = \begin{pmatrix} 0,67 \\ 0,90 \times 10^{-1} \\ -0,27 \end{pmatrix}$$

$$\left. \begin{array}{l} \text{solução} \\ \text{aproximada} \\ \text{de } Ux = y \end{array} \right\} = \begin{pmatrix} 0,10 \times 10 \\ -0,11 \times 10 \\ 0,23 \times 10 \end{pmatrix}$$

$$\left. \begin{array}{l} \text{solução exata} \\ \text{de } LUx = b \end{array} \right\} = \begin{pmatrix} 0,1001 \dots \times 10 \\ -0,1038 \dots \times 10 \\ 0,2072 \dots \times 10 \end{pmatrix} \quad \left. \begin{array}{l} \text{solução exata} \\ \text{de } Ax = b \end{array} \right\} = \begin{pmatrix} 0,10 \times 10 \\ -0,10 \times 10 \\ 0,20 \times 10 \end{pmatrix}$$

$$\left. \begin{array}{l} \text{a solução} \\ \text{aproximada} \end{array} \right\} = \begin{pmatrix} 0,6560 \\ 0,2220 \\ 0,8600 \times 10^{-1} \end{pmatrix} \quad \text{vetor resíduos} = \begin{pmatrix} 0,14 \times 10^{-1} \\ 0,18 \times 10^{-1} \\ 0,34 \times 10^{-1} \end{pmatrix}$$

4.3 - Cálculo de determinantes

Seja $A \in M_n(\mathbb{R})$ e consideremos a determinação do valor de $\det A$.

Aplicando à matriz A a eliminação de Gauss com pivotamento parcial, obtemos matrizes triangulares L e U , por (4.53), satisfazem $LU = A + E$, com E satisfazendo (4.51) e (4.52). Como o determinante de LU é igual ao de A , a menos do sinal, que depende do número de permutações efetuadas, temos:

$$|\det(A^{(1)} + E)| = \begin{vmatrix} a_{11}^{(1)} & a_{12}^{(2)} & \dots & a_{1n}^{(n)} \\ & a_{22} & & \\ & & & \\ & & & a_{nn} \end{vmatrix}$$

Aplicando (2.73) e (2.75) a (4.74), obtemos:

$$|\det(A^{(1)} + E)| \approx \left| \prod_{i=1}^n a_{ii}^{(i)} \right| = \begin{vmatrix} a_{11}^{(1)} & a_{12}^{(2)} & \dots & a_{1n}^{(n)} \\ & a_{22} & & \\ & & & \\ & & & a_{nn} \end{vmatrix} (1 + \epsilon),$$

$$1 + \epsilon = 1 + 1,01^{(n-1)} \theta, \quad |\theta| \leq 1$$

O fator $(1 + \epsilon)$ tem pouca influência no resultado de modo que, a menos desse fator, o valor calculado pelo computador é igual ao determinante exato de A .

4.4. Inversão de matrizes

Dada a matriz não singular $A \in M_n(\mathbb{R})$, podemos determinar $A^{-1} \approx X$ através da resolução de n sistemas de equações lineares

$$(4.77) \quad Ax_i = e_i, \quad i = 1, 2, \dots, n,$$

em que x_i e e_i são, respectivamente, a i -ésima coluna de X e I . Aplicando (4.66) a (4.77), podemos afirmar que cada aproximação x_i , determinada através da eliminação de Gauss com pivoteamento parcial, satisfaz a

$$(4.78) \quad (A + \delta_i A) x_i = e_i,$$

em que, como o indicam os índices em (4.78), as perturbações são diferentes para cada coluna da matriz I . Entretanto, a delimitação geral apresentada em (4.64) vale para todos os $\delta_i A$, de modo que a inversa aproximada X satisfaz a

$$(4.79) \quad (A + \delta A) X = I,$$

com $\|\delta A\|_\infty$ dado em (4.64)

Consideremos a matriz resíduo

$$(4.80) \quad F = I - AX = \delta A X$$

Em virtude de (4.64)

$$(4.81) \quad \|F\|_\infty \leq 1,01 (n^3 + 3n^2) \rho \|A\|_\infty \|X\|_\infty,$$

sendo que, ao contrário do que ocorre para a solução de sistemas lineares, em que soluções aproximadas que possuem resíduo pequeno não são, necessariamente, precisas, na determinação de inversas de matrizes a obtenção de resíduos pequenos implica na determinação de uma inversa precisa, ou seja, (4.81) fornece uma estimativa de precisão da aproximação determinada. Isso decorre de

(4.80), pois

$$\bar{A}^{-1} - X = \bar{A}^{-1} F \quad e, \text{ portanto,}$$

$$(4.82) \quad \frac{\|\bar{A}^{-1} - X\|_k}{\|\bar{A}^{-1}\|_k} \leq \|F\|_k,$$

qualquer que seja a norma k considerada.

Como (4.81) foi obtida a partir de (4.80), em que X é inversa de A à direita, precisamos considerar o caso em que X é inversa de A à esquerda. Em consequência de (4.80) temos:

$$\|I - XA\|_\infty \leq \|\bar{A}^{-1}\|_\infty \|F\|_\infty \|A\|_\infty, \text{ ou,}$$

$$(4.83) \quad \|I - XA\|_\infty \leq \text{cond}_\infty(A) \|F\|_\infty,$$

ou seja, o resíduo de X , como inversa à esquerda, pode ser maior que o de X , como inversa à direita, do fator $\text{cond}_\infty(A)$. Entretanto, Wilkinson (26, p. 110-111) afirma que, na prática, o método de Gauss com pivotamento parcial fornece resíduos $I - AX$ e $I - XA$ de mesma ordem de grandeza, mesmo quando A é mal condicionada.

Daremos, a seguir, um exemplo de cálculo de inversa.

(4.84) Exemplo - Consideremos a inversão da matriz A dada em (4.72). Utilizando, para isso, as matrizes L e U dadas em (4.73), obtemos:

$$\left. \begin{array}{l} \text{inversa} \\ \text{aproximada} \\ \text{de } LU \end{array} \right\} = \begin{pmatrix} 0,18 \times 10 & -0,48 & -0,55 \\ -0,17 \times 10 & -0,21 \times 10 & 0,49 \times 10 \\ 0,18 \times 10 & 0,73 \times 10 & -0,77 \times 10 \end{pmatrix},$$

que pode ser comparada com

$$\left. \begin{array}{l} \text{inversa exata de } LU \\ \text{(arredondada para} \\ \text{2 dígitos)} \end{array} \right\} = \begin{pmatrix} 0,18 \times 10 & -0,54 & -0,45 \\ -0,17 \times 10 & -0,19 \times 10 & 0,47 \times 10 \\ 0,19 \times 10 & 0,73 \times 10 & -0,77 \times 10 \end{pmatrix}$$

$$\left. \begin{array}{l} \text{inversa exata de } A \\ \text{(arredondada para} \\ \text{2 dígitos)} \end{array} \right\} = \begin{pmatrix} 0,18 \times 10 & -0,54 & -0,44 \\ -0,17 \times 10 & -0,18 \times 10 & 0,46 \times 10 \\ 0,18 \times 10 & 0,70 \times 10 & -0,74 \times 10 \end{pmatrix}$$

Por outro lado,

$$AX = \begin{pmatrix} 0,96 & 0 & 0 \\ -0,10 \times 10^{-1} & 0,10 \times 10 & 0 \\ 0 & -0,10 & 0,10 \times 10 \end{pmatrix} \quad XA = \begin{pmatrix} 0,96 & -0,20 \times 10^{-1} & 0,10 \times 10^{-1} \\ -0,20 \times 10^{-1} & 0,10 \times 10 & -0,30 \times 10^{-1} \\ -0,10 & 0 & 0,10 \times 10 \end{pmatrix}$$

Observe-se que AX e XA são da mesma ordem de grandeza.

4.5- Matrizes especiais

A eliminação de Gauss, até agora considerada em relação a matrizes gerais, apresenta simplificações quando aplicada a algumas classes de matrizes especiais. Assim, por exemplo, para o cálculo de determinantes de matrizes de Hessenberg a eliminação de Gauss pode ser realizada sem a utilização de qualquer forma de pivotamento porque, apesar da possibilidade de ocorrência de multiplicadores e elementos das matrizes reduzidas, com ordem de grandeza elevada, ainda assim é sempre obtida uma boa aproximação. Para verificar esse fato, consideremos uma matriz superior de Hessenberg e determinemos os elementos de $A^{(k)}$, que são os utilizados no cálculo de $\det A$.

Como pode ser concluído pela análise de (1.11), a k -ésima etapa é iniciada com uma matriz

reduzida $A^{(k)}$, cujas linhas $(k+1)$ a n são constituídas por elementos de $A^{(k)}$. Por outro lado, a primeira coluna de $\bar{A}^{(k)}$, formada pelos elementos que serão zerados, apresenta apenas dois elementos não nulos, um dos quais, aquele que será anulado, pertence a uma das linhas ainda não modificadas. Como exemplo, apresentaremos $A^{(3)}$ para o caso $n=6$

$$A^{(3)} = \left(\begin{array}{cc|cc|cc} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & a_{14}^{(1)} & a_{15}^{(1)} & a_{16}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & a_{24}^{(2)} & a_{25}^{(2)} & a_{26}^{(2)} \\ \hline 0 & 0 & a_{33}^{(3)} & a_{34}^{(3)} & a_{35}^{(3)} & a_{36}^{(3)} \\ 0 & 0 & a_{43}^{(1)} & a_{44}^{(1)} & a_{45}^{(1)} & a_{46}^{(1)} \\ 0 & 0 & 0 & a_{54}^{(1)} & a_{55}^{(1)} & a_{56}^{(1)} \\ 0 & 0 & 0 & 0 & a_{65}^{(1)} & a_{66}^{(1)} \end{array} \right)$$

A matriz separada por linhas tracejadas, no canto inferior direito de $A^{(3)}$ e $\bar{A}^{(3)}$.

Portanto, existe apenas um multiplicador de Gauss por coluna, definido por

$$(4.85) \quad m_{i+1,i} = fl \left(\frac{a_{i+1,i}^{(1)}}{a_{ii}^{(i)}} \right) = \frac{a_{i+1,i}^{(1)}}{a_{ii}^{(i)}} (1 + \epsilon_{i+1,i}),$$

$$|\epsilon_{i+1,i}| \leq \mu,$$

sendo que no numerador de (4.85) utilizamos $a_{i+1,i}^{(1)}$, em lugar de $a_{i+1,i}^{(i)}$, porque a linha $(i+1)$ ainda não foi modificada.

Em virtude da existência de um único multiplicador por coluna, podemos simplificar a notação, adotando a forma de representação utilizada por Wilkinson (24, p. 329-331), na qual os multiplicadores

$m_{i+1,i}$ são denotados por m_i e qualquer quantidade δ , não necessariamente a mesma em cada etapa, que satisfaça $|\delta| \leq \epsilon$, é representada por ϵ . Assim procedendo, podemos escrever as equações que originam as aproximações para os $a_{i,i}^{(i)}$ na forma:

$$(4.86) \quad \begin{cases} a_{2i}^{(2)} = \left[a_{2i}^{(1)} - m_1 a_{1,i}^{(1)} (1+\epsilon) \right] (1+\epsilon) \\ a_{3i}^{(3)} = \left[a_{3i}^{(2)} - m_2 a_{2i}^{(2)} (1+\epsilon) \right] (1+\epsilon) \\ \vdots \\ a_{i,i}^{(i)} = \left[a_{i,i}^{(i-1)} - m_{i-1} a_{i-1,i}^{(i-1)} (1+\epsilon) \right] (1+\epsilon) \end{cases},$$

para as quais utilizamos (2.113).

Eliminando de (4.86) todas as quantidades $a_{j,i}^{(j)}$, para as quais $j \neq 1$ ou $j \neq i$, obtemos:

$$(4.87) \quad a_{i,i}^{(i)} = a_{i,i}^{(1)} (1+\epsilon) - m_{i-1} (1+\epsilon) a_{i-1,i}^{(1)} + m_{i-1} m_{i-2} (1+\epsilon) a_{i-2,i}^{(1)} - \dots + (-1)^{i-2} m_{i-1} m_{i-2} \dots m_2 (1+\epsilon) a_{2i}^{(1)} + (-1)^{i-1} m_{i-1} m_{i-2} \dots m_1 (1+\epsilon) a_{1i}^{(1)}$$

A partir de (4.85) e (4.87), podemos concluir que os elementos diagonais de $A^{(n)}$ são exatamente iguais aos elementos diagonais que teriam sido obtidos se tivéssemos realizado, de maneira exata, uma eliminação de Gauss com uma matriz A^2 , obtida por perturbação de A , cujos elementos da i -ésima coluna fossem:

$$(4.88) \quad \begin{cases} a'_{1i} = a_{1i}^{(1)} (1+\epsilon)^{2i-2} \\ a'_{ji} = a_{ji}^{(j)} (1+\epsilon)^{2(i-j)+1} & , j = 2, 3, \dots, i \\ a'_{i+1,i} = a_{i+1,i}^{(i)} (1+\epsilon) \end{cases}$$

As modificações são tão pequenas, que praticamente nenhum prejuízo decorreu da não utilização de pivota-
mento. Observe-se, ainda, que a matriz A' também ori-
ginária, exatamente, os m_i obtidos na prática.

É importante que se observe que, se o cál-
culo do determinante de uma matriz Hessenberg superior é
efetuado de maneira satisfatória sem uso de pivota-
mento, o mesmo não ocorre quanto à resolução de sistemas de equa-
ções lineares cuja matriz é desse tipo. Wilkinson
(24, p. 330-331) ilustra esse fato através de um interes-
sante exemplo.

Quando o pivota-
mento parcial é utilizado
na aplicação da eliminação de Gauss a matrizes Hessen-
berg superiores, então, se $|a_{ij}^{(1)}| \leq a$, ocorre
 $\max_{i,j} |a_{ij}^{(k)}| \leq ka$, uma delimitação melhor do que a
apresentada em (4.44) para matrizes gerais. Isso decor-
re do fato de $\bar{A}^{(k)}$ também ser Hessenberg superior.

As matrizes tridiagonais, definidas em (1.12),
são matrizes Hessenberg tanto superiores quanto inferio-
res. Portanto, os resultados que vimos para matrizes
Hessenberg superiores também valem para as tridiago-
nais. Para essas últimas matrizes, a eliminação de
Gauss com pivota-
mento parcial fornece $|a_{ij}^{(k)}| \leq 2a$,
se $|a_{ij}^{(1)}| \leq a$, como mostra Wilkinson (23, p. 288-289).

Outra classe de matrizes que apresenta
simplificações para a eliminação de Gauss é a das
diagonalmente dominantes, definidas em (1.13). Para
essas matrizes os multiplicadores de Gauss satisfa-
zem $|m_{ij}| < 1$, mesmo sem a utilização de qualquer
forma de pivota-
mento. Por outro lado, se $|a_{ij}^{(1)}| \leq a$,
então, $\max_{i,j} |a_{ij}^{(k)}| \leq 2 \max_{i,j} |a_{ij}^{(1)}|$ ($i, j = 1, \dots, n$), como

mostra Wilkinson (23, p. 288-289).

4.6. Refinamento da solução aproximada

Vimos, em 4.1 e 4.2, que, ao obter a solução aproximada de um sistema de equações lineares pelo método de eliminação de Gauss, não podemos, de imediato, afirmar nada a respeito da precisão da mesma. O refinamento iterativo fornece esse tipo de informação, dando orientações quanto ao condicionamento do sistema original, e possibilita uma melhoria, um refinamento, das aproximações utilizáveis. Exponemos apenas as características fundamentais desse processo, cuja descrição detalhada pode ser encontrada em Forsythe (7, p. 49-54), e, a seguir, faremos considerações quanto à sua convergência.

Consideremos o sistema (4.19) e indiquemos por x_1 a solução aproximada obtida através de (4.66). Para $m = 1, 2, \dots$, a m -ésima iteração envolve três etapas:

(4.89) determinação do vetor resíduo, r_m , definido por

$$r_m = b - Ax_m$$

(4.90) cálculo do vetor correção, c_m , através da resolução de

$$Ac_m = r_m$$

(4.91) obtenção da aproximação refinada x_{m+1} , dada por

$$x_{m+1} = x_m + c_m$$

Quando o processo é executado em um computador, os ciclos de refinamento podem ser interrompidos pela verificação de um dos critérios:

(4.92) $\frac{\|c_m\|_\infty}{\|x_1\|_\infty} < \varepsilon$, $\varepsilon > 0$ pré-estabelecido

(4.93) o número m , de iterações, excede um certo valor

N pré-especificado.

Em algumas circunstâncias (4.92) é substituído por

$$(4.94) \quad \left\| r_k - r_{k+1} \right\|_{\infty} < \delta, \quad \delta > 0 \text{ pré-estabelecido.}$$

No caso de interrupção do processo em virtude da verificação de (4.92) ou (4.94), sem que (4.93) tenha ocorrido, obtém-se um refinamento da aproximação. Quando a interrupção decorre da verificação de (4.93), então, a aproximação inicial não é utilizável.

Caso não houvesse ocorrência de erros de arredondamento, o processo de refinamento convergiria em uma iteração, pois teríamos, por (4.89)-(4.91):

$$A x_2 = A(x_1 + c_1) = A x_1 + A c_1 = A x_1 + r_1 = b.$$

A única etapa crítica é (4.89), pois b e $A x_m$ são, no caso da obtenção de x_m pela eliminação de Gauss com pivotamento, muito próximos. Por esse motivo, na prática, o cálculo dos resíduos é efetuado com precisão maior do que a utilizada nas demais etapas. Em geral, os resíduos são calculados em precisão dupla e as correções e aproximações refinadas em precisão simples. É lógico que, nesse caso, a precisão que se pode alcançar para a solução aproximada, através de refinamentos sucessivos, é limitada pelo número de dígitos significativos utilizados em precisão simples, entretanto, paramente é necessário utilizar precisão múltipla em (4.90) e (4.91); isso só ocorre quando o sistema linear é muito mal condicionado.

Na etapa (4.90) pode ser utilizado qualquer método de resolução apropriado, desde que o mesmo forneça uma precisão razoável, porém, o que torna prático o uso do refinamento é justamente a utilização de uma mes-

ma matriz na etapa (4.90) de todas as iterações. O método de eliminação de Gauss possibilita essa simplificação, visto que a decomposição LU de A , determinada durante o cálculo de x_1 , pode ser utilizada para esse fim. Com esse procedimento conseguimos obter um bom refinamento de x_1 , através de pequeno trabalho adicional, aquele requerido para a determinação da referida aproximação. (aproximadamente 25% em tempo de computador).

Consideremos, então, a convergência do processo que acabamos de descrever, com utilização da eliminação de Gauss com pivotamento parcial no cálculo de x_1 e aproveitamento das matrizes L e U , determinadas para esse fim, em (4.90). Uma análise completa da convergência do processo de refinamento na aritmética de vírgula flutuante é tão laboriosa algebricamente, que importantes aspectos do problema podem ficar mascarados pelos inúmeros detalhes algébricos. Por esse motivo, neste trabalho é descrita uma versão idealizada, apresentada em Forsythe (7, p. 109-113), na qual se supõe que a ocorrência dos erros de arredondamento se dá apenas no cálculo das correções a partir dos resíduos, ou seja, na etapa (4.90). Essa hipótese é razoável no caso de ser utilizada a acumulação de produtos escalares em (4.89) o que, como vimos anteriormente, é normalmente empregado na prática. Uma análise completa da convergência em questão, sem o uso dessa hipótese simplificadora, é efetuada por Moler (17).

A seguir, apresentaremos a versão adotada. A aplicação de (4.66) a (4.90) fornece:

$$(4.95) \quad (A + \delta A) c_m = r_m,$$

em que δA , que verifica (4.56) e (4.64), depende de x_m , porém a delimitação de $\|\delta A\|_\infty$ não depende do mesmo.

Para evidenciar a dependência em m e simplificar as expressões que serão obtidas no decorrer da análise, podemos escrever (4.95) na forma

$$(4.96) \quad A(I + F_m) x_m = r_m, \text{ com}$$

$$A F_m = \delta A \quad \text{e, portanto,}$$

$$(4.97) \quad F_m = A^{-1} \delta A$$

As matrizes F_m , definidas acima, determinarão, de maneira completa, o comportamento da versão idealizada em questão. Isso será melhor entendido através do resultado que se segue.

(4.98) Teorema. Sejam x_m ($m=1,2,\dots$) as aproximações da solução do sistema (4.19), obtidas por sucessivas iterações segundo a forma descrita. Seja x^* a solução exata de (4.19), ou seja, $x^* = A^{-1}b$. Se ocorrer:

$$(4.99) \quad \|F_m\|_{\infty} \leq \theta < \frac{1}{2}, \quad \forall m,$$

então a sequência de aproximações $x_1, x_2, \dots, x_m, \dots$ converge para a solução exata de (4.19), ou seja, em virtude de (3.25), $\|x_m - x^*\|_{\infty} \rightarrow 0$, com $m \rightarrow \infty$.

Provejamos a demonstração de (4.98). A partir de (4.93) e (4.96) temos

$$A(I + F_m) x_m = b - A x_m$$

e, como A é não singular por hipótese,

$$(I + F_m) x_m = A^{-1}b - x_m = x^* - x_m$$

Por outro lado, por (4.91),

$$(4.100) \quad (I + F_m) x_{m+1} = F_m x_m + x^*$$

Subtraindo $(I + F_m)x^*$ de ambos os membros de (4.100), temos

$$(4.101) \quad (I + F_m)(x_{m+1} - x^*) = F_m(x_m - x^*).$$

Em virtude de (4.99) e de (1.43), podemos afirmar que $(I + F_m)$ é não singular e escrever (4.101) na forma

$$(4.102) \quad x_{m+1} - x^* = (I + F_m)^{-1} F_m(x_m - x^*)$$

A equação (4.102) é muito importante para o estudo da convergência do processo de refinamento, por evidenciar que o erro da $(m+1)$ -ésima iteração é igual ao produto de uma certa matriz (desconhecida) pelo erro da m -ésima iteração.

Considerando a norma $\|\cdot\|_\infty$ para (4.102) e aplicando (1.44) e (4.99) à expressão resultante, obtemos

$$\|x_{m+1} - x^*\|_\infty \leq \frac{1}{1 - \|F_m\|_\infty} \|F_m\|_\infty \|x_m - x^*\|_\infty$$

$$(4.103) \quad \|x_{m+1} - x^*\|_\infty \leq \frac{1}{1 - \sigma} \sigma \|x_m - x^*\|_\infty$$

Seja

$$(4.104) \quad \tau = \frac{\sigma}{1 - \sigma}$$

Então, (4.103) pode ser escrita

$$(4.105) \quad \|x_{m+1} - x^*\|_\infty \leq \tau \|x_m - x^*\|_\infty,$$

que fornece

$$(4.106) \quad \begin{aligned} \|x_m - x^*\|_\infty &\leq \tau \|x_{m-1} - x^*\|_\infty \leq \tau^2 \|x_{m-2} - x^*\|_\infty \leq \\ &\leq \tau^3 \|x_{m-3} - x^*\|_\infty \leq \dots \leq \tau^{m-1} \|x_1 - x^*\|_\infty \end{aligned}$$

Fazendo $x_0 = \theta$, ou seja, considerando $m=0$ em (4.91), o que corresponde a tomar x_1 igual à aproximação de (4.19), determinada pela eliminação de Gauss com pivotamento parcial, obtemos, de (4.106),

$$(4.107) \quad \|x_1 - x^*\|_{\infty} \leq \tau \|x_0 - x\|_{\infty} = \tau \|x^*\|_{\infty}$$

Portanto, (4.106) e (4.107) nos dão:

$$(4.108) \quad \|x_m - x^*\|_{\infty} \leq \tau^m \|x^*\|_{\infty}$$

Como, em virtude de (4.99) e (4.104), $\tau < 1$, então, quando $m \rightarrow \infty$, $\tau^m \rightarrow 0$, e, por (4.108), $\|x_m - x^*\|_{\infty} \rightarrow 0$, o que prova (4.98).

O teorema (4.98) afirma que é possível tornar $\|x_m - x^*\|_{\infty}$ arbitrariamente pequeno, desde que efetuemos um número conveniente de iterações. Facilmente, porém, se pode concluir que isso é impossível, visto que as componentes de x_m devem ter precisão simples, enquanto o mesmo não ocorre, de maneira geral, com as de x^* . Através de alguma reflexão sobre o problema, pode-se, porém, entender que essa dificuldade decorre, principalmente, da suposição de que (4.91) se verifica exatamente, o que não ocorre na prática, quando, nessa etapa, cometemos um erro de arredondamento em cada componente. Uma análise mais detalhada da convergência do processo de refinamento, com a consideração dos erros nas três etapas (4.89) - (4.91), revela que o erro, em vez de tender a zero, tornar-se pequeno.

Consideremos o refinamento da aproximação determinada em (4.71). Os cálculos necessários para a realização das etapas (4.89) foram realizados com acumulação de produto escalar.

(4.109) Exemplo: Consideremos o sistema $Ax = b$, para o qual

$$(4.110) \quad A = \begin{pmatrix} 0,6900 & 0,2600 & 0,1200 \\ 0,1500 & 0,4500 & 0,2700 \\ 0,3100 & 0,4900 & 0,1500 \end{pmatrix} \quad b = \begin{pmatrix} 0,6700 \\ 0,2400 \\ 0,1200 \end{pmatrix}$$

e determinemos o refinamento da aproximação

$$x_1 = \begin{pmatrix} 0,10 \times 10 \\ -0,11 \times 10 \\ 0,21 \times 10 \end{pmatrix}, \text{ obtida em (4.71).}$$

A eliminação de Gauss com pivotamento parcial fornece a seguinte decomposição triangular de A

$$(4.111) \quad L = \begin{pmatrix} 0,1000 \times 10 & & \\ 0,2174 & 0,1000 \times 10 & \\ 0,4493 & 0,9484 & 0,1000 \times 10 \end{pmatrix}$$

$$(4.112) \quad U = \begin{pmatrix} 0,6900 & 0,2600 & 0,1200 \\ 0,0000 & 0,3935 & 0,2439 \\ 0,0000 & 0,0000 & -0,1352 \end{pmatrix}$$

Utilizando como x_1 a aproximação determinada em (4.71)

$$(4.113) \quad x_1 = \begin{pmatrix} 0,10 \times 10 \\ -0,11 \times 10 \\ 0,21 \times 10 \end{pmatrix}$$

e utilizando para a realização das diversas etapas (4.90) as matrizes L e U dadas, respectivamente, em (4.111) e (4.112) temos:

$$r_1 = \begin{pmatrix} 0,14 \times 10^{-1} \\ 0,18 \times 10^{-1} \\ 0,34 \times 10^{-1} \end{pmatrix} \quad c_1 = \begin{pmatrix} 0,00 \\ 0,10 \\ -0,10 \end{pmatrix} \quad x_2 = \begin{pmatrix} 0,10 \times 10 \\ -0,10 \times 10 \\ 2,00 \times 10 \end{pmatrix}$$

$$(4.114) \quad \tilde{\mu}_2 = \begin{pmatrix} 0,00 \\ 0,00 \\ 0,00 \end{pmatrix}$$

Como pode ser observado através de (4.114),
o processo convergiu rapidamente.

- _____ // _____

BIBLIOGRAFIA

1. ALBASINY, Ernest L. Error in digital solution of linear problems. In: ADVANCED SEMINAR CONDUCTED BY THE MATHEMATICS RESEARCH CENTER, University of Wisconsin, oct. 5-7, 1964. Error in digital computation, edited by Louis B. Pall. New York, John Wiley /c1965/ v.1, p.131 - 184.
2. BARROS, Ivan de Queiroz. Elementos de análise numérica. Pocos de Caldas /IMPA/ 1969. 139p. Trab. apres. no 7º Colóquio Brasileiro de Matemática.
3. _____. Algebra linear. In: _____. Métodos numéricos. /Campos, UEC. Instituto de Matemática, Estatística e Ciência de Computação/ 1970. v.1, p.1-121. Mineografado.
4. BARTLE, Robert G. The elements of real analysis. New York, John Wiley /c1964/ p.58-69.
5. FADDEEVA, V.N. Computational methods of linear algebra. New York, Dover /c1959/ p.1-145.
6. FORSYTHE, Alexandra I.; KEELMAN, Thomas A.; ORCANICK, Elliot I.; STENBERG, Warren. Computer science: a first course. New York, John Wiley - /c1969/ p.3-232.
7. FORSYTHE, George E. & MOLEK, Cleve B. Computer solution of linear algebraic systems. Englewood Cliffs, Prentice-Hall /c1967/ p.2-131.
8. FOX, L. An introduction to numerical linear algebra. Oxford, Clarendon Press /c1964/ p.1-214.
9. FOX, L. & MAYERS, D.F. Computing methods for scientists and engineers. Oxford, Clarendon Press, 1968. p.1-121.
10. FRÖBERG, Carl-Erik. Introduction to numerical analysis. Reading, Addison-Wesley /c1965/ cap.1, p.1-14; cap.3-6, 44-125.
11. GANTMACHER, F.R. The theory of matrices. New York, Chelsea /c1959/ v.1.
12. GEL'FAND, I.M. Lectures on linear algebra. 2.ed. New York, Interscience /c1961/ p.1-163. (Lectures on Linear Algebra, 9).
13. HENRICI, Peter. Elements of numerical analysis. New York, John Wiley /c1964/ p.291-321.
14. HILDEBRAND, F.B. Introduction to numerical analysis. New York, McGraw-Hill, 1956, cap.1, p.1-34; cap.10, p.424-485.
15. LIPSCHUTZ, Seymour. Algebra linear. Rio de Janeiro, McGraw-Hill - /c1968/ 403p. (Coleção Schaum).
16. McCracken, Daniel D. & DORN, William S. Numerical methods and FORTRAN programming with applications in engineering and science. New York, John Wiley /c1964/ cap.2, p.43-67; cap.8, p.226-283.

17. MOLER, Cleve B. Iterative refinement in floating point. J. ACM, New York, 14(2): 316-321, apr. 1967.
18. MONTEIRO, L.H. Jacz. Algebra linear. 5.ed. São Paulo, Nobel /1969/ v.1.
19. _____. _____. São Paulo, Nobel /1970/ v.2.
20. NATIONAL PHYSICAL LABORATORY. Modern computing methods. 2.ed. London, Her Majesty's Stationery Office, 1961. cap.1-5, p.1-52. (Notes on Applied Science, 16).
21. RALSTON, Anthony. A first course in numerical analysis. New York, McGraw-Hill /c1965/ cap.1, p.1-22; cap.9, p.394-463.
22. WILKINSON, J.H. The algebraic eigenvalue problem. Oxford, Clarendon Press /c1965/ p.1-264.
23. _____. Error analysis of direct methods of matrix inversion. J. ACM, New York, 8(3): 281-330, jul. 1961.
24. _____. Error analysis of floating point computation. Numerische Mat., Berlin, 2(5): 319-340, okt. 1960.
25. _____. Modern error analysis. SIAM Rev., Philadelphia, 13(4): 548-568, oct. 1971.
26. _____. Bounding errors in algebraic processes. Englewood Cliffs, Prentice-Hall /c1963/ 161p.