

Lecture Notes for Inf-Mat 4350, 2012

Tom Lyche

August 17, 2012

Contents

Preface	ix
0 Preliminaries	1
0.1 Notation	1
0.2 Vector Spaces and Subspaces	4
0.2.1 Linear Independence and Bases	7
0.2.2 Subspaces	8
0.2.3 The vector spaces \mathbb{R}^n and \mathbb{C}^n	11
0.3 Vector Norms	12
0.3.1 Convergence of Vectors	14
0.3.2 Convergence of Series of Vectors	17
0.4 Inner Products	17
0.4.1 Orthogonality	20
0.4.2 Orthogonal Projection, and the Column Space Decomposition	22
0.5 Linear Systems	24
0.5.1 The Inverse Matrix	26
0.6 Determinants	28
0.7 Eigenpairs	32
0.8 Algorithms and Numerical Stability	35
I Direct Methods for Linear Systems	37
1 Gaussian Elimination	39
1.1 Block Multiplication	40
1.2 Triangular matrices	43
1.2.1 Algorithms for Triangular Systems	45
1.3 Naive Gaussian Elimination and LU factorization	48
1.3.1 Operation count	51

1.4	Gaussian Elimination with Row Interchanges	54
1.4.1	Pivoting	54
1.4.2	Permutation matrices	55
1.4.3	The PLU-Factorization	55
1.5	An Algorithm for Finding the PLU-Factorization	57
1.5.1	Pivot strategies	58
1.6	Review Questions	60
2	Examples of Linear Systems	61
2.1	The Second Derivative Matrix	61
2.2	LU Factorization of a Tridiagonal System	62
2.2.1	Diagonal Dominance	64
2.3	Cubic Spline Interpolation	68
2.3.1	The Runge Phenomenon	68
2.3.2	Piecewise Linear and Cubic Spline Interpolation .	69
2.3.3	Give me a Moment	72
2.4	Review Questions	78
3	LU Factorizations	81
3.1	The LU Factorization	81
3.2	Block LU Factorization	86
3.3	The Symmetric LU Factorization	88
3.4	Positive Definite- and Positive Semidefinite Matrices	90
3.4.1	Definition and Examples	90
3.4.2	Principal Submatrices	92
3.4.3	The Symmetric Positive Definite Case	93
3.5	The Cholesky Factorization	95
3.6	The Symmetric Positive Semidefinite Case	97
3.6.1	An Algorithm for SemiCholesky Factorization of a Banded Matrix	100
3.7	Review Questions	101
4	The Kronecker Product	103
4.1	Test Matrices	103
4.1.1	The 2D Poisson Problem	103
4.1.2	The Test Matrices	107
4.2	The Kronecker Product	108
4.3	Properties of the 1D and 2D Test Matrices	112
4.4	Review Questions	116
5	Fast Direct Solution of a Large Linear System	117
5.1	Algorithms for a Banded Positive Definite System	117

5.1.1	Cholesky Factorization	118
5.1.2	Block LU Factorization of a Block Tridiagonal Matrix	118
5.1.3	Other Methods	119
5.2	A Fast Poisson Solver based on Diagonalization	119
5.3	A Fast Poisson Solver based on the Discrete Sine and Fourier Transforms	121
5.3.1	The Discrete Sine Transform (DST)	121
5.3.2	The Discrete Fourier Transform (DFT)	121
5.3.3	The Fast Fourier Transform (FFT)	123
5.3.4	A Poisson Solver based on the FFT	126
5.4	Review Questions	128
II	Some Matrix Theory	131
6	Matrix Reduction by Similarity Transformations	133
6.1	Orthonormal, Unitary, and Similar Matrices	133
6.1.1	Orthonormal and Unitary Matrices	133
6.1.2	Similarity Transformations	135
6.2	Linear Independence of Eigenvectors	137
6.2.1	Algebraic and Geometric Multiplicity of Eigenvalues	138
6.3	Normal Matrices	142
6.3.1	The Schur Decomposition	142
6.3.2	Matrices with Orthonormal Eigenvectors	144
6.4	Hermitian Matrices	146
6.4.1	The Rayleigh Quotient	146
6.4.2	Minmax Theorems	147
6.4.3	The Hoffman-Wielandt Theorem	149
6.5	Left Eigenvectors	149
6.6	The Jordan Form and the Minimal Polynomial	151
6.6.1	The Minimal Polynomial	154
6.7	Proof of the Real Schur Form	156
6.8	Conclusions	157
6.9	Review Questions	158
7	The Singular Value Decomposition	159
7.1	Singular Values and Singular Vectors	159
7.1.1	SVD and SVF	160
7.1.2	Examples	163
7.2	SVD and the Four Fundamental Subspaces	166
7.3	A Geometric Interpretation	168
7.4	Determining the Rank of a Matrix Numerically	169

7.4.1	The Frobenius Norm	170
7.4.2	Low Rank Approximation	171
7.5	The Minmax Theorem for Singular Values and the Hoffman-Wielandt Theorem	172
7.5.1	Proof of the Hoffman-Wielandt theorem for singular values	173
7.6	Review Questions	174
8	Matrix Norms	177
8.1	Matrix Norms	177
8.1.1	Consistent and Subordinate Matrix Norms	178
8.1.2	Operator Norms	179
8.1.3	The Operator p -Norms	181
8.1.4	Unitary Invariant Matrix Norms	183
8.1.5	Absolute and Monotone Norms	185
8.2	The Condition Number with Respect to Inversion	185
8.3	Proof that the p -Norms are Norms	191
8.4	Review Questions	195
III	Iterative Methods for Large Linear Systems	197
9	The Classical Iterative Methods	199
9.1	Classical Iterative Methods; Component Form	200
9.1.1	The Discrete Poisson System	201
9.1.2	Matrix Formulations of the Classical Methods	204
9.1.3	The Splitting Matrices for the Classical Methods	205
9.2	Convergence and Spectral Radius	206
9.2.1	Neumann Series	209
9.3	Convergence of Fixed-point Iteration	210
9.3.1	Stopping the Iteration	212
9.3.2	Richardson's Method (R method)	213
9.4	Convergence of the Classical Methods for the Discrete Poisson Matrix	214
9.4.1	Number of Iterations	216
9.5	Convergence Analysis for SOR	217
9.6	The Optimal SOR Parameter ω	219
9.7	Review Questions	222
10	The Conjugate Gradient Method	223
10.1	Quadratic Minimization	223
10.2	Steepest Descent	226

10.2.1	Convergence Analysis for Steepest Descent	227
10.3	The Conjugate Gradient Method	230
10.3.1	The Best Approximation Property	232
10.4	The Conjugate Gradient Algorithm	234
10.4.1	Numerical Example	234
10.4.2	Implementation Issues	235
10.4.3	The Spectral Condition Numbers	236
10.5	Proof of Convergence	238
10.5.1	Chebyshev Polynomials	240
10.5.2	Monotonicity of the error	243
10.6	Preconditioning	244
10.7	Preconditioning Example	247
10.7.1	A Banded Matrix	248
10.7.2	Applying Preconditioning	250
10.8	Review Questions	252
IV	Orthonormal Transformations and Least Squares	253
11	Orthonormal and Unitary Transformations	255
11.1	The Householder Transformation	255
11.2	Householder Triangulation	259
11.2.1	Solving Linear Systems using Unitary Transformations	260
11.2.2	The number of Arithmetic Operations	261
11.3	The QR Decomposition and QR Factorization	262
11.3.1	QR and Gram-Schmidt	264
11.4	Givens Rotations	266
11.5	Review Questions	268
12	Least Squares	269
12.1	Existence, Uniqueness, and Characterization	269
12.2	Some Curve Fitting examples	272
12.3	The Least Squares Problem and the Singular Value Decomposition	276
12.3.1	Orthogonal Projections	276
12.3.2	The Generalized Inverse	277
12.4	Numerical Solution	280
12.4.1	Normal Equations	280
12.4.2	QR Factorization	281
12.4.3	Singular Value Factorization	283
12.5	Perturbation Theory for Least Squares	283

12.5.1	Perturbing the Right Hand Side	284
12.5.2	Perturbing the Matrix	286
12.6	Perturbation Theory for Singular Values	287
12.7	Review Questions	288
V	Eigenvalues and Eigenvectors	289
13	Numerical Eigenvalue Problems	291
13.1	Eigenpars	291
13.2	Perturbation of Eigenvalues	292
13.2.1	Gerschgorin's Theorem	294
13.3	Unitary Similarity Transformation of a Matrix into Upper Hessenberg Form	296
13.4	Computing a Selected Eigenvalue of a Symmetric Matrix	299
13.4.1	The Inertia Theorem	301
13.4.2	Approximating λ_m	303
13.5	Perturbation Proofs	304
13.6	Review Questions	306
14	The QR Algorithm	307
14.1	The Power Method	307
14.1.1	The Inverse Power Method	310
14.2	The basic QR Algorithm	312
14.2.1	The Relation to the Power Method	314
14.2.2	Invariance of the Hessenberg Form	315
14.2.3	Deflation	315
14.3	The Shifted QR Algorithms	316
14.4	A Convergence Theorem	317
14.5	Review Questions	318
VI	Appendix	319
A	Determinants	321
A.1	Permutations	321
A.2	Basic Properties of Determinants	323
A.3	The Adjoint Matrix and Cofactor Expansion	327
A.4	Computing Determinants	330
A.5	Some Useful Determinant Formulas	330
B	Computer Arithmetic	333
B.1	Absolute and Relative Errors	333

B.2	Floating Point Numbers	334
B.3	Rounding and Arithmetic Operations	337
B.3.1	Rounding	337
B.3.2	Arithmetic Operations	338
B.4	Backward Rounding-Error Analysis	338
B.4.1	Computing a Sum	338
B.4.2	Computing an Inner Product	341
B.4.3	Computing a Matrix Product	341
C	Differentiation of Vector Functions	343
Bibliography		347
Index		370

Preface

These lecture notes contains the text for a course in matrix analysis and numerical linear algebra given at the beginning graduate level at the University of Oslo. Most of the chapters correspond approximately to one week of lectures.

Oslo, 20 August, 2012

Tom Lyche

Chapter 0

Preliminaries

In this book we consider methods for solving square linear systems, solving over-determined linear systems by least squares, and computing eigenvalues and eigenvectors of matrices. These problems can be large and we will show that much can be gained by taking advantage of their special structures. To achieve this goal, we need a good understanding of linear algebra and matrix theory.

In this introductory chapter we give a compact introduction to linear algebra with emphasis on \mathbb{R}^n and \mathbb{C}^n . For a more elementary introduction, see for example the book [18].

We start by introducing the notation used.

0.1 Notation

The following sets will be used throughout this book.

1. The sets of natural numbers, integers, rational numbers, real numbers, and complex numbers are denoted by $\mathbb{N}, \mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$, respectively.
2. We use the “colon equal” symbol $v := e$ to indicate that the symbol v is defined by the expression e .
3. \mathbb{R}^n is the set of n -tuples of real numbers which we will represent as column vectors. Thus $\mathbf{x} \in \mathbb{R}^n$ means

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix},$$

where $x_i \in \mathbb{R}$ for $i = 1, \dots, n$. Row vectors are normally identified using the transpose operation. Thus if $\mathbf{x} \in \mathbb{R}^n$ then \mathbf{x} is a column vector and \mathbf{x}^T is a row vector.

4. Addition and scalar multiplication are denoted and defined by

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix}, \quad a\mathbf{x} = \begin{bmatrix} ax_1 \\ \vdots \\ ax_n \end{bmatrix}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad a \in \mathbb{R}.$$

5. $\mathbb{R}^{m \times n}$ is the set of matrices¹ \mathbf{A} with real elements. The integers m and n are the number of rows and columns in the tableau

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

The element in the i th row and j th column of \mathbf{A} will be denoted by $a_{i,j}$, a_{ij} , $\mathbf{A}(i, j)$ or $(\mathbf{A})_{i,j}$. We use the notations

$$\mathbf{a}_{:j} = \begin{bmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{bmatrix}, \quad \mathbf{a}_{i:}^T = [a_{i1}, a_{i2}, \dots, a_{in}], \quad \mathbf{A} = [\mathbf{a}_{:1}, \mathbf{a}_{:2}, \dots, \mathbf{a}_{:n}] = \begin{bmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \vdots \\ \mathbf{a}_{m:}^T \end{bmatrix}$$

for the columns $\mathbf{a}_{:j}$ and rows $\mathbf{a}_{i:}^T$ of \mathbf{A} . We often drop the colon and write \mathbf{a}_j and \mathbf{a}_i^T when no confusion can arise. If $m = 1$ then \mathbf{A} is a row vector, if $n = 1$ then \mathbf{A} is a column vector, while if $m = n$ then \mathbf{A} is a square matrix. In this text we will denote matrices by boldface capital letters $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$ and vectors most often by boldface lower case letters $\mathbf{x}, \mathbf{y}, \mathbf{z}, \dots$.

6. The imaginary unit $\sqrt{-1}$ is denoted by i . The complex conjugate and the modulus of a complex number z is denoted by \bar{z} and $|z|$, respectively. Thus if $z = x + iy = re^{i\phi} = r(\cos \phi + i \sin \phi)$, with $x, y \in \mathbb{R}$, is a complex number then $\bar{z} := x - iy = re^{-i\phi} = \cos \phi - i \sin \phi$ and $|z| := \sqrt{x^2 + y^2} = r$. $\text{Re}(z) := x$ and $\text{Im}(z) := y$ denote the real and imaginary part of the complex number z .
7. For matrices and vectors with complex elements we use the notation $\mathbf{A} \in \mathbb{C}^{m \times n}$ and $\mathbf{x} \in \mathbb{C}^n$. We define complex row vectors using either the transpose \mathbf{x}^T or the conjugate transpose operation $\mathbf{x}^* := \bar{\mathbf{x}}^T = [\bar{x}_1, \dots, \bar{x}_n]$.

¹The word matrix to denote a rectangular array of numbers, was first used by Sylvester in 1850

8. For $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ and $a \in \mathbb{C}$ the operations of vector addition and scalar multiplication is defined by component operations as in the real case (cf. 4.).
9. The arithmetic operations on rectangular matrices are
 - **matrix addition** $\mathbf{C} = \mathbf{A} + \mathbf{B}$ if $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are matrices of the same size, i.e., with the same number of rows and columns, and $c_{ij} = a_{ij} + b_{ij}$ for all i, j .
 - **multiplication by a scalar** $\mathbf{C} = \alpha \mathbf{A}$, where $c_{ij} = \alpha a_{ij}$ for all i, j .
 - **matrix multiplication** $\mathbf{C} = \mathbf{A}\mathbf{B}$, $\mathbf{C} = \mathbf{A} \cdot \mathbf{B}$ or $\mathbf{C} = \mathbf{A} * \mathbf{B}$, where $\mathbf{A} \in \mathbb{C}^{m \times p}$, $\mathbf{B} \in \mathbb{C}^{p \times n}$, $\mathbf{C} \in \mathbb{C}^{m \times n}$, and $c_{ij} = \sum_{k=1}^p a_{ik}b_{kj}$ for $i = 1, \dots, m$, $j = 1, \dots, n$.
 - **element-by-element matrix operations** $\mathbf{C} = \mathbf{A} \times \mathbf{B}$, $\mathbf{D} = \mathbf{A}/\mathbf{B}$, and $\mathbf{E} = \mathbf{A} \wedge r$ where all matrices are of the same dimension and $c_{ij} = a_{ij}b_{ij}$, $d_{ij} = a_{ij}/b_{ij}$ and $e_{ij} = a_{ij}^r$ for all i, j and suitable r . The element-by-element product $\mathbf{C} = \mathbf{A} \times \mathbf{B}$ is known as the **Schur product** and also the **Hadamard product**.
10. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ or $\mathbf{A} \in \mathbb{C}^{m \times n}$. The **transpose** \mathbf{A}^T and **conjugate transpose** \mathbf{A}^* are $n \times m$ matrices with elements $a_{ij}^T = a_{ji}$ and $a_{ij}^* = \bar{a}_{ji}$, respectively. If \mathbf{B} is an n, p matrix then $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$ and $(\mathbf{AB})^* = \mathbf{B}^* \mathbf{A}^*$.
11. The **unit vectors** in \mathbb{R}^n and \mathbb{C}^n are denoted by

$$\mathbf{e}_1 := \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 := \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{e}_3 := \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad \mathbf{e}_n := \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix},$$

while $\mathbf{I}_n = \mathbf{I} := [\delta_{ij}]_{i,j=1}^n$, where

$$\delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

is the **identity matrix** of order n . Both the columns and the transpose of the rows of \mathbf{I} are the unit vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$.

12. Some matrices with many zeros have names indicating their “shape”. Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ or $\mathbf{A} \in \mathbb{C}^{n \times n}$. Then \mathbf{A} is

- **diagonal** if $a_{ij} = 0$ for $i \neq j$.
- **upper triangular** or **right triangular** if $a_{ij} = 0$ for $i > j$.
- **lower triangular** or **left triangular** if $a_{ij} = 0$ for $i < j$.

- **upper Hessenberg** if $a_{ij} = 0$ for $i > j + 1$.
- **lower Hessenberg** if $a_{ij} = 0$ for $i < j + 1$.
- **tridiagonal** if $a_{ij} = 0$ for $|i - j| > 1$.
- **d -banded** if $a_{ij} = 0$ for $|i - j| > d$.
- **block upper triangular** if there is an integer k such that $a_{ij} = 0$ for $i = k + 1, \dots, n$ and $j = 1, \dots, k$.
- **block lower triangular** if \mathbf{A}^T is block upper triangular.

13. We use the following notations for diagonal- and tridiagonal $n \times n$ matrices

$$\text{diag}(d_i) = \text{diag}(d_1, \dots, d_n) := \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} = \begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & d_n \end{bmatrix},$$

$$\mathbf{B} = \text{tridiag}(a_i, d_i, c_i) = \text{tridiag}(\mathbf{a}, \mathbf{d}, \mathbf{c}) := \begin{bmatrix} d_1 & c_1 & & & \\ a_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-2} & d_{n-1} & c_{n-1} \\ & & & a_{n-1} & d_n \end{bmatrix}.$$

Here $b_{ii} = d_i$ for $i = 1, \dots, n$, $b_{i+1,i} = a_i$, $b_{i,i+1} = c_i$ for $i = 1, \dots, n-1$, and $b_{ij} = 0$ otherwise.

14. Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ and $1 \leq i_1 < i_2 < \dots < i_r \leq m$, $1 \leq j_1 < j_2 < \dots < j_c \leq n$. The matrix $\mathbf{A}(\mathbf{i}, \mathbf{j}) \in \mathbb{C}^{r \times c}$ is the submatrix of \mathbf{A} consisting of rows $\mathbf{i} := [i_1, \dots, i_r]$ and columns $\mathbf{j} := [j_1, \dots, j_c]$

$$\mathbf{A}(\mathbf{i}, \mathbf{j}) := \mathbf{A} \left(\begin{array}{cccc} i_1 & i_2 & \cdots & i_r \\ j_1 & j_2 & \cdots & j_c \end{array} \right) = \begin{bmatrix} a_{i_1,j_1} & a_{i_1,j_2} & \cdots & a_{i_1,j_c} \\ a_{i_2,j_1} & a_{i_2,j_2} & \cdots & a_{i_2,j_c} \\ \vdots & \vdots & \ddots & \vdots \\ a_{i_r,j_1} & a_{i_r,j_2} & \cdots & a_{i_r,j_c} \end{bmatrix}.$$

For the special case of consecutive rows and columns we use the notation

$$\mathbf{A}(r_1 : r_2, c_1 : c_2) := \begin{bmatrix} a_{r_1,c_1} & a_{r_1,c_1+1} & \cdots & a_{r_1,c_2} \\ a_{r_1+1,c_1} & a_{r_1+1,c_1+1} & \cdots & a_{r_1+1,c_2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{r_2,c_1} & a_{r_2,c_1+1} & \cdots & a_{r_2,c_2} \end{bmatrix}.$$

0.2 Vector Spaces and Subspaces

Many mathematical systems have analogous properties to vectors in \mathbb{R}^2 or \mathbb{R}^3 .

Definition 0.1 (Real vector space)

A **real vector space** is a nonempty set \mathcal{V} , whose objects are called **vectors**, together with two operations $+ : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$ and $\cdot : \mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V}$, called **addition** and **scalar multiplication**, satisfying the following axioms for all vectors $\mathbf{u}, \mathbf{v}, \mathbf{w}$ in \mathcal{V} and scalars c, d in \mathbb{R} .

(V1) The sum $\mathbf{u} + \mathbf{v}$ is in \mathcal{V} ,

(V2) $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$,

(V3) $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$,

(V4) There is a **zero vector** $\mathbf{0}$ such that $\mathbf{u} + \mathbf{0} = \mathbf{u}$,

(V5) For each \mathbf{u} in \mathcal{V} there is a vector $-\mathbf{u}$ in \mathcal{V} such that $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$,

(S1) The scalar multiple $c \cdot \mathbf{u}$ is in \mathcal{V} ,

(S2) $c \cdot (\mathbf{u} + \mathbf{v}) = c \cdot \mathbf{u} + c \cdot \mathbf{v}$,

(S3) $(c + d) \cdot \mathbf{u} = c \cdot \mathbf{u} + d \cdot \mathbf{u}$,

(S4) $c \cdot (d \cdot \mathbf{u}) = (cd) \cdot \mathbf{u}$,

(S5) $1 \cdot \mathbf{u} = \mathbf{u}$.

The scalar multiplication symbol \cdot is often omitted, writing $c\mathbf{v}$ instead of $c \cdot \mathbf{v}$. We define $\mathbf{u} - \mathbf{v} := \mathbf{u} + (-\mathbf{v})$. We call \mathcal{V} a **complex vector space** if the scalars consist of all complex numbers \mathbb{C} . In this book a vector space is either real or complex.

From the axioms it follows that

1. The zero vector is unique.
2. For each $\mathbf{u} \in \mathcal{V}$ the **negative** $-\mathbf{u}$ of \mathbf{u} is unique.
3. $0\mathbf{u} = \mathbf{0}$, $c\mathbf{0} = \mathbf{0}$, and $-\mathbf{u} = (-1)\mathbf{u}$.

Here are some examples

1. The space \mathbb{R}^n , where $n \in \mathbb{N}$, is a real vector space.
2. Similarly, \mathbb{C}^n is a complex vector space.
3. Let \mathcal{D} be a subset of \mathbb{R} and $d \in \mathbb{N}$. The set \mathcal{V} of all functions $\mathbf{f}, \mathbf{g} : \mathcal{D} \rightarrow \mathbb{R}^d$ is a real vector space with

$$(\mathbf{f} + \mathbf{g})(t) := \mathbf{f}(t) + \mathbf{g}(t), \quad (c\mathbf{f})(t) := c\mathbf{f}(t), \quad t \in \mathcal{D}, \quad c \in \mathbb{R}.$$

Two functions \mathbf{f}, \mathbf{g} in \mathcal{V} are equal if $\mathbf{f}(t) = \mathbf{g}(t)$ for all $t \in \mathcal{D}$. The zero element is the **zero function** given by $\mathbf{f}(t) = \mathbf{0}$ for all $t \in \mathcal{D}$ and the negative of \mathbf{f} is given by $-\mathbf{f} = (-1)\mathbf{f}$. In the following we will use boldface letters for functions only if $d > 1$.

4. For $n \geq 0$ the space Π_n of polynomials of degree at most n consists of all polynomials $p : \mathbb{R} \rightarrow \mathbb{R}$, $p : \mathbb{R} \rightarrow \mathbb{C}$, or $p : \mathbb{C} \rightarrow \mathbb{C}$ of the form

$$p(t) = a_0 + a_1 t + a_2 t^2 + \cdots + a_n t^n, \quad (2)$$

where the coefficients a_0, \dots, a_n are real or complex numbers. p is called the **zero polynomial** if all coefficients are zero. All other polynomials are said to be **nontrivial**. The **degree** of a nontrivial polynomial p given by (2) is the smallest integer $0 \leq k \leq n$ such that $p(t) = a_0 + \cdots + a_k t^k$ with $a_k \neq 0$. The degree of the zero polynomial is not defined. Π_n is a vector space if we define addition and scalar multiplication as for functions.

Definition 0.2 (Linear combination)

For $n \geq 1$ let $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be a set of vectors in a vector space \mathcal{V} and let c_1, \dots, c_n be scalars.

1. The sum $c_1 \mathbf{x}_1 + \cdots + c_n \mathbf{x}_n$ is called a **linear combination** of $\mathbf{x}_1, \dots, \mathbf{x}_n$.
2. The linear combination is **nontrivial** if $c_j \mathbf{x}_j \neq \mathbf{0}$ for at least one j .
3. The set of all linear combinations of elements in \mathcal{X} is denoted $\text{span}(\mathcal{X})$.
4. A vector space is **finite dimensional** if it has a finite spanning set; i.e., there exists $n \in \mathbb{N}$ and $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ in \mathcal{V} such that $\mathcal{V} = \text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\})$.

Example 0.3 (Linear combinations)

1. Any $\mathbf{x} = [x_1, \dots, x_m]^T$ in \mathbb{C}^m can be written as a linear combination of the unit vectors as $\mathbf{x} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \cdots + x_m \mathbf{e}_m$. Thus, $\mathbb{C}^m = \text{span}(\{\mathbf{e}_1, \dots, \mathbf{e}_m\})$ and \mathbb{C}^m is finite dimensional. Similarly \mathbb{R}^m is finite dimensional.
2. Let $\Pi = \cup_n \Pi_n$ be the space of all polynomials. Π is a vector space that is not finite dimensional. For suppose Π is finite dimensional. Then $\Pi = \text{span}(\{p_1, \dots, p_m\})$ for some polynomials p_1, \dots, p_m . Let d be an integer such that the degree of p_j is less than d for $j = 1, \dots, m$. A polynomial of degree d cannot be written as a linear combination of p_1, \dots, p_m , a contradiction.

0.2.1 Linear Independence and Bases

Definition 0.4 (Linear independence)

A set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ of nonzero vectors in a vector space is **linearly dependent** if $\mathbf{0}$ can be written as a nontrivial linear combination of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. Otherwise \mathcal{X} is **linearly independent**.

A set of vectors $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is linearly independent if and only if

$$c_1\mathbf{x}_1 + \cdots + c_n\mathbf{x}_n = \mathbf{0} \quad \Rightarrow \quad c_1 = \cdots = c_n = 0. \quad (3)$$

Suppose $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is linearly independent. Then

1. If $\mathbf{x} \in \text{span}(\mathcal{X})$ then the scalars c_1, \dots, c_n in the representation $\mathbf{x} = c_1\mathbf{x}_1 + \cdots + c_n\mathbf{x}_n$ are unique.
2. Any nontrivial linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is nonzero,

Lemma 0.5 (Linear independence and span)

Suppose $\mathbf{v}_1, \dots, \mathbf{v}_n$ span a vector space \mathcal{V} and that $\mathbf{w}_1, \dots, \mathbf{w}_k$ are linearly independent vectors in \mathcal{V} . Then $k \leq n$.

Proof. Suppose $k > n$. Write \mathbf{w}_1 as a linear combination of elements from the set $\mathcal{X}_0 := \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$, say $\mathbf{w}_1 = c_1\mathbf{v}_1 + \cdots + c_n\mathbf{v}_n$. Since $\mathbf{w}_1 \neq \mathbf{0}$ not all the c 's are equal to zero. Pick a nonzero c , say c_{i_1} . Then \mathbf{v}_{i_1} can be expressed as a linear combination of \mathbf{w}_1 and the remaining \mathbf{v} 's. So the set $\mathcal{X}_1 := \{\mathbf{w}_1, \mathbf{v}_1, \dots, \mathbf{v}_{i_1-1}, \mathbf{v}_{i_1+1}, \dots, \mathbf{v}_n\}$ must also be a spanning set for \mathcal{V} . We repeat this for \mathbf{w}_2 and \mathcal{X}_1 . In the linear combination $\mathbf{w}_2 = d_{i_1}\mathbf{w}_1 + \sum_{j \neq i_1} d_j\mathbf{v}_j$, we must have $d_{i_2} \neq 0$ for some i_2 with $i_2 \neq i_1$. For otherwise $\mathbf{w}_2 = d_1\mathbf{w}_1$ contradicting the linear independence of the \mathbf{w} 's. So the set \mathcal{X}_2 consisting of the \mathbf{v} 's with \mathbf{v}_{i_1} replaced by \mathbf{w}_1 and \mathbf{v}_{i_2} replaced by \mathbf{w}_2 is again a spanning set for \mathcal{V} . Repeating this process $n - 2$ more times we obtain a spanning set \mathcal{X}_n where $\mathbf{v}_1, \dots, \mathbf{v}_n$ have been replaced by $\mathbf{w}_1, \dots, \mathbf{w}_n$. Since $k > n$ we can then write \mathbf{w}_k as a linear combination of $\mathbf{w}_1, \dots, \mathbf{w}_n$ contradicting the linear independence of the \mathbf{w} 's. We conclude that $k \leq n$. \square

Definition 0.6 (basis)

A finite set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ in a vector space \mathcal{V} is a **basis** for \mathcal{V} if

1. $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\} = \mathcal{V}$.
2. $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is linearly independent.

Theorem 0.7 (Basis subset of a spanning set)

Suppose \mathcal{V} is a vector space and that $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is a spanning set for \mathcal{V} . Then we can find a subset $\{\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}\}$ that forms a basis for \mathcal{V} .

Proof. If $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is linearly dependent we can express one of the \mathbf{v} 's as a nontrivial linear combination of the remaining \mathbf{v} 's and drop that \mathbf{v} from the spanning set. Continue this process until the remaining \mathbf{v} 's are linearly independent. They still span the vector space and therefore form a basis. \square

Corollary 0.8 (Existence of a basis)

A vector space is finite dimensional if and only if it has a basis.

Proof. Let $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ be a finite dimensional vector space. By Theorem 0.7, \mathcal{V} has a basis. Conversely, if $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is a basis then it is by definition a finite spanning set. \square

Theorem 0.9 (Dimension of a vector space)

Every basis for a vector space \mathcal{V} has the same number of elements. This number is called the **dimension** of the vector space and denoted $\dim \mathcal{V}$.

Proof. Suppose $\mathcal{X} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ and $\mathcal{Y} = \{\mathbf{w}_1, \dots, \mathbf{w}_k\}$ are two bases for \mathcal{V} . By Lemma 0.5 we have $k \leq n$. Using the same Lemma with \mathcal{X} and \mathcal{Y} switched we obtain $n \leq k$. We conclude that $n = k$. \square

The set of unit vectors $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ form a basis for both \mathbb{R}^n and \mathbb{C}^n .

Theorem 0.10 (Enlarging vectors to a basis)

Every linearly independent set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ in a finite dimensional vector space \mathcal{V} can be enlarged to a basis for \mathcal{V} .

Proof. If $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ does not span \mathcal{V} we can enlarge the set by one vector \mathbf{v}_{k+1} which cannot be expressed as a linear combination of $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$. The enlarged set is also linearly independent. Continue this process. Since the space is finite dimensional it must stop after a finite number of steps. \square

0.2.2 Subspaces

Definition 0.11 (Subspace)

A nonempty subset \mathcal{S} of a real or complex vector space \mathcal{V} is called a **subspace** of \mathcal{V} if

(V1) The sum $\mathbf{u} + \mathbf{v}$ is in \mathcal{S} for any $\mathbf{u}, \mathbf{v} \in \mathcal{S}$.

(S1) The scalar multiple $c\mathbf{u}$ is in \mathcal{S} for any scalar c and any $\mathbf{u} \in \mathcal{S}$.

Using the operations in \mathcal{V} , any subspace \mathcal{S} of \mathcal{V} is a vector space, i. e., all 10 axioms $V1 - V5$ and $S1 - S5$ are satisfied for \mathcal{S} . In particular, \mathcal{S} must contain the zero element in \mathcal{V} . This follows since the operations of vector addition and scalar multiplication are inherited from \mathcal{V} .

Example 0.12 (Examples of subspaces)

1. $\{\mathbf{0}\}$, where $\mathbf{0}$ is the zero vector is a subspace, the **trivial subspace**. The dimension of the trivial subspace is defined to be zero. All other subspaces are **nontrivial**.
2. \mathcal{V} is a subspace of itself.
3. $\text{span}(\mathcal{X})$ is a subspace of \mathcal{V} for any $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{V}$. Indeed, it is easy to see that **(V1)** and **(S1)** hold.
4. The **sum** of two subspaces \mathcal{R} and \mathcal{S} of a vector space \mathcal{V} is defined by

$$\mathcal{R} + \mathcal{S} := \{\mathbf{r} + \mathbf{s} : \mathbf{r} \in \mathcal{R} \text{ and } \mathbf{s} \in \mathcal{S}\}. \quad (4)$$

Clearly **(V1)** and **(S1)** hold and it is a subspace of \mathcal{V} .

5. The **intersection** of two subspaces \mathcal{R} and \mathcal{S} of a vector space \mathcal{V} is defined by

$$\mathcal{R} \cap \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ and } \mathbf{x} \in \mathcal{S}\}. \quad (5)$$

It is a subspace of \mathcal{V} .

6. The **union** of two subspaces \mathcal{R} and \mathcal{S} of a vector space \mathcal{V} is defined by

$$\mathcal{R} \cup \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ or } \mathbf{x} \in \mathcal{S}\}. \quad (6)$$

In general it is not a subspace of \mathcal{V} .

7. A sum of two subspaces \mathcal{R} and \mathcal{S} of a vector space \mathcal{V} is called a **direct sum** and denoted $\mathcal{R} \oplus \mathcal{S}$ if $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$. The subspaces \mathcal{R} and \mathcal{S} are called **complementary** in the subspace $\mathcal{R} \oplus \mathcal{S}$.

Theorem 0.13 (Dimension formula for sums of subspaces)

Let \mathcal{R} and \mathcal{S} be two finite dimensional subspaces of a vector space \mathcal{V} . Then

$$\dim(\mathcal{R} + \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S}) - \dim(\mathcal{R} \cap \mathcal{S}). \quad (7)$$

In particular, for a direct sum

$$\dim(\mathcal{R} \oplus \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S}). \quad (8)$$

Proof. Let $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ be a basis for $\mathcal{R} \cap \mathcal{S}$, where $\{\mathbf{u}_1, \dots, \mathbf{u}_p\} = \emptyset$, the empty set, in the case $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$. We use Theorem 0.10 to extend $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ to a basis $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q\}$ for \mathcal{R} and a basis $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{s}_1, \dots, \mathbf{s}_t\}$ for \mathcal{S} . Every $\mathbf{x} \in \mathcal{R} + \mathcal{S}$ can be written as a linear combination of $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q, \mathbf{s}_1, \dots, \mathbf{s}_t\}$ so these vectors span $\mathcal{R} + \mathcal{S}$. We show that they are linearly independent and hence a basis. Suppose $\mathbf{u} + \mathbf{r} + \mathbf{s} = \mathbf{0}$, where $\mathbf{u} := \sum_{j=1}^p \alpha_j \mathbf{u}_j$, $\mathbf{r} := \sum_{j=1}^q \rho_j \mathbf{r}_j$, and $\mathbf{s} := \sum_{j=1}^t \sigma_j \mathbf{s}_j$. Now $\mathbf{r} = -(\mathbf{u} + \mathbf{s})$ belongs to both \mathcal{R} and to \mathcal{S} and hence $\mathbf{r} \in \mathcal{R} \cap \mathcal{S}$. Therefore \mathbf{r} can be written as a linear combination of $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ say $\mathbf{r} := \sum_{j=1}^p \beta_j \mathbf{u}_j$. But then $\mathbf{0} = \sum_{j=1}^p \beta_j \mathbf{u}_j - \sum_{j=1}^q \rho_j \mathbf{r}_j$ and since $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q\}$ is linearly independent we must have $\beta_1 = \dots = \beta_p = \rho_1 = \dots = \rho_q = 0$ and hence $\mathbf{r} = \mathbf{0}$. We then have $\mathbf{u} + \mathbf{s} = \mathbf{0}$ and by linear independence of $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{s}_1, \dots, \mathbf{s}_t\}$ we obtain $\alpha_1 = \dots = \alpha_p = \sigma_1 = \dots = \sigma_t = 0$. We have shown that the vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q, \mathbf{s}_1, \dots, \mathbf{s}_t\}$ constitute a basis for $\mathcal{R} + \mathcal{S}$. But then

$$\dim(\mathcal{R} + \mathcal{S}) = p + q + t = (p + q) + (p + t) - p = \dim(\mathcal{R}) + \dim(\mathcal{S}) - \dim(\mathcal{R} \cap \mathcal{S})$$

and (7) follows. (7) implies (8) since $\dim\{\mathbf{0}\} = 0$. \square

Theorem 0.14 (Direct sum decomposition)

Let \mathcal{R} and \mathcal{S} be two subspaces of a vector space \mathcal{V} and assume $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$. Every $\mathbf{x} \in \mathcal{R} \oplus \mathcal{S}$ can be decomposed uniquely in the form $\mathbf{x} = \mathbf{r} + \mathbf{s}$, where $\mathbf{r} \in \mathcal{R}$ and $\mathbf{s} \in \mathcal{S}$. If $\{\mathbf{r}_1, \dots, \mathbf{r}_k\}$ is a basis for \mathcal{R} and $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ is a basis for \mathcal{S} then $\{\mathbf{r}_1, \dots, \mathbf{r}_k, \mathbf{s}_1, \dots, \mathbf{s}_n\}$ is a basis for $\mathcal{R} \oplus \mathcal{S}$.

Proof. To show uniqueness, suppose we could write $\mathbf{x} = \mathbf{r}_1 + \mathbf{s}_1 = \mathbf{r}_2 + \mathbf{s}_2$ for $\mathbf{r}_1, \mathbf{r}_2 \in \mathcal{R}$ and $\mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}$. Then $\mathbf{r}_1 - \mathbf{r}_2 = \mathbf{s}_2 - \mathbf{s}_1$ and it follows that $\mathbf{r}_1 - \mathbf{r}_2$ and $\mathbf{s}_2 - \mathbf{s}_1$ belong to both \mathcal{R} and \mathcal{S} and hence to $\mathcal{R} \cap \mathcal{S}$. But then $\mathbf{r}_1 - \mathbf{r}_2 = \mathbf{s}_2 - \mathbf{s}_1 = \mathbf{0}$ so $\mathbf{r}_1 = \mathbf{r}_2$ and $\mathbf{s}_2 = \mathbf{s}_1$. Thus uniqueness follows. Suppose $\{\mathbf{r}_1, \dots, \mathbf{r}_k\}$ is a basis for \mathcal{R} and $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ is a basis for \mathcal{S} . Since $\dim(\mathcal{R} + \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S})$ the vectors $\{\mathbf{r}_1, \dots, \mathbf{r}_k, \mathbf{s}_1, \dots, \mathbf{s}_n\}$ span $\mathcal{R} + \mathcal{S}$. To show linear independence suppose $\sum_{j=1}^k \rho_j \mathbf{r}_j + \sum_{j=1}^n \sigma_j \mathbf{s}_j = \mathbf{0}$. The first sum belongs to \mathcal{R} and the second to \mathcal{S} and the sum is a decomposition of $\mathbf{0}$. By uniqueness of the decomposition both sums must be zero. But then $\rho_1 = \dots = \rho_k = \sigma_1 = \dots = \sigma_n = 0$ and linear independence follows. \square

It is convenient to introduce a matrix transforming a basis in a subspace into a basis for the space itself.

Lemma 0.15 (Change of basis matrix)

Suppose \mathcal{S} is a subspace of a finite dimensional vector space \mathcal{V} and let $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ be a basis for \mathcal{S} and $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ a basis for \mathcal{V} . Then each \mathbf{s}_j can be expressed

as a linear combination of $\mathbf{v}_1, \dots, \mathbf{v}_m$, say

$$\mathbf{s}_j = \sum_{i=1}^m a_{ij} \mathbf{v}_i \text{ for } j = 1, \dots, n. \quad (9)$$

If $\mathbf{x} \in \mathcal{S}$ then $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j = \sum_{i=1}^m b_i \mathbf{v}_i$ for some coefficients $\mathbf{b} := [b_1, \dots, b_m]^T$, $\mathbf{c} := [c_1, \dots, c_n]^T$. Moreover $\mathbf{b} = \mathbf{A}\mathbf{c}$, where $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{m \times n}$ is given by (9). The matrix \mathbf{A} has linearly independent columns.

Proof. (9) holds for some a_{ij} since $\mathbf{s}_j \in \mathcal{V}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ spans \mathcal{V} . Since $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ is a basis for \mathcal{S} and $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ a basis for \mathcal{V} , every $\mathbf{x} \in \mathcal{S}$ can be written $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j = \sum_{i=1}^m b_i \mathbf{v}_i$ for some scalars (c_j) and (b_i) . But then

$$\sum_{i=1}^m b_i \mathbf{v}_i = \mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j \stackrel{(9)}{=} \sum_{j=1}^n c_j \left(\sum_{i=1}^m a_{ij} \mathbf{v}_i \right) = \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij} c_j \right) \mathbf{v}_i.$$

Since $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ is linearly independent it follows that $b_i = \sum_{j=1}^n a_{ij} c_j$ for $i = 1, \dots, m$ or $\mathbf{b} = \mathbf{A}\mathbf{c}$. Finally, to show that \mathbf{A} has linearly independent columns suppose $\mathbf{b} := \mathbf{A}\mathbf{c} = \mathbf{0}$ for some $\mathbf{c} = [c_1, \dots, c_n]^T$. Define $\mathbf{x} := \sum_{j=1}^n c_j \mathbf{s}_j$. Then $\mathbf{x} = \sum_{i=1}^m b_i \mathbf{v}_i$ and since $\mathbf{b} = \mathbf{0}$ we have $\mathbf{x} = \mathbf{0}$. But since $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ is linearly independent it follows that $\mathbf{c} = \mathbf{0}$. \square

The matrix \mathbf{A} in Lemma 0.15 is called a **change of basis matrix**.

0.2.3 The vector spaces \mathbb{R}^n and \mathbb{C}^n

When $\mathcal{V} = \mathbb{R}^m$ we can think of n vectors in \mathbb{R}^n , say $\mathbf{x}_1, \dots, \mathbf{x}_n$, as a set $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ or as the columns of a matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$. A linear combination can then be written as a matrix times vector $\mathbf{X}\mathbf{c}$, where $\mathbf{c} = [c_1, \dots, c_n]^T$ is the vector of scalars. Thus

$$\text{span}(\mathcal{X}) = \text{span}(\mathbf{X}) = \{\mathbf{X}\mathbf{c} : \mathbf{c} \in \mathbb{R}^n\}.$$

Of course the same holds for \mathbb{C}^m .

In \mathbb{R}^m and \mathbb{C}^m each of the following statements is equivalent to linear independence of \mathcal{X} .

- (i) $\mathbf{X}\mathbf{c} = \mathbf{0} \Rightarrow \mathbf{c} = \mathbf{0}$,
- (ii) \mathbf{X} has linearly independent columns,

Definition 0.16 (Column space and null space)

Associated with a matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ are the following subspaces

1. The subspace $\text{span}(\mathbf{X})$ is called the **column space** of \mathbf{X} . It is the smallest subspace containing $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.
2. $\text{span}(\mathbf{X}^T)$ is called the **row space** of \mathbf{X} . It is generated by the rows of \mathbf{X} written as column vectors.
3. The subspace $\ker(\mathbf{X}) := \{\mathbf{y} \in \mathbb{R}^n : \mathbf{X}\mathbf{y} = \mathbf{0}\}$ is called the **null space** or **kernel space** of \mathbf{X} .

Note that the subspace $\ker(\mathbf{X})$ is nontrivial if and only if \mathcal{X} is linearly dependent.

0.3 Vector Norms

To measure the size of a vector we use norms.

Definition 0.17 (Vector norm)

A **(vector) norm** in a real (resp. complex) vector space \mathcal{V} is a function $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ that satisfies for all \mathbf{x}, \mathbf{y} in \mathcal{V} and all a in \mathbb{R} (resp. \mathbb{C})

1. $\|\mathbf{x}\| \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$. (positivity)
2. $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$. (homogeneity)
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$. (subadditivity)

The triple $(\mathcal{V}, \mathbb{R}, \|\cdot\|)$ (resp. $(\mathcal{V}, \mathbb{C}, \|\cdot\|)$) is called a **normed vector space** and the inequality 3. is called the **triangle inequality**.

In this book we will use the following family of vector norms on $\mathcal{V} = \mathbb{C}^n$ and $\mathcal{V} = \mathbb{R}^n$.

Definition 0.18 (Vector p-norms)

We define for $p \geq 1$ and $\mathbf{x} \in \mathbb{R}^n$ or $\mathbf{x} \in \mathbb{C}^n$ the **p -norms** by

$$\|\mathbf{x}\|_p := \left(\sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad (10)$$

$$\|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|. \quad (11)$$

The most important cases are $p = 1, 2, \infty$:

1. $\|\mathbf{x}\|_1 = \sum_{j=1}^n |x_j|$, (the **one-norm** or l_1 -norm)

2. $\|\mathbf{x}\|_2 = \sqrt{\sum_{j=1}^n |x_j|^2}$, (the two-norm, l_2 -norm, or Euclidian norm)
3. $\|\mathbf{x}\|_\infty = \max_{1 \leq j \leq n} |x_j|$, (the infinity-norm, l_∞ -norm, or max norm)

Some remarks are in order.

1. In Section 8.3, we show that the p -norms are vector norms for $1 \leq p \leq \infty$.
2. The triangle inequality $\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$ is called **Minkowski's inequality**.
3. To prove it one first establishes **Hölder's inequality**

$$\sum_{j=1}^n |x_j y_j| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1, \quad \mathbf{x}, \mathbf{y} \in \mathbb{C}^n. \quad (12)$$

The relation $\frac{1}{p} + \frac{1}{q} = 1$ means that if $p = 1$ then $q = \infty$ and if $p = 2$ then $q = 2$.

4. The infinity norm is related to the other p -norms by

$$\lim_{p \rightarrow \infty} \|\mathbf{x}\|_p = \|\mathbf{x}\|_\infty \text{ for all } \mathbf{x} \in \mathbb{C}^n. \quad (13)$$

5. The equation (13) clearly holds for $\mathbf{x} = \mathbf{0}$. For $\mathbf{x} \neq \mathbf{0}$ we write

$$\|\mathbf{x}\|_p := \|\mathbf{x}\|_\infty \left(\sum_{j=1}^n \left(\frac{|x_j|}{\|\mathbf{x}\|_\infty} \right)^p \right)^{1/p}.$$

Now each term in the sum is not greater than one and at least one term is equal to one, and we obtain

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_p \leq n^{1/p} \|\mathbf{x}\|_\infty, \quad p \geq 1. \quad (14)$$

Since $\lim_{p \rightarrow \infty} n^{1/p} = 1$ for any $n \in \mathbb{N}$ we see that (13) follows.

We return now to the general case.

Definition 0.19 (Equivalent norms)

We say that two norms $\|\cdot\|$ and $\|\cdot\|'$ on \mathcal{V} are **equivalent** if there are positive constants m and M such that for all vectors $\mathbf{x} \in \mathcal{V}$ we have

$$m\|\mathbf{x}\|' \leq \|\mathbf{x}\| \leq M\|\mathbf{x}\|'. \quad (15)$$

By (14) the p - and ∞ -norms are equivalent for any $p \geq 1$. This result is generalized in the following theorem.

Theorem 0.20 (Basic properties of vector norms)

The following holds for a normed vector space $(\mathcal{V}, \mathbb{C}, \|\cdot\|)$.

1. $\|\mathbf{x} - \mathbf{y}\| \geq |\|\mathbf{x}\| - \|\mathbf{y}\||$, for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ (*inverse triangle inequality*).
2. The vector norm is a continuous function $\mathcal{V} \rightarrow \mathbb{R}$.
3. All vector norms on \mathcal{V} are equivalent provided \mathcal{V} is finite dimensional.

Proof.

1. Since $\|\mathbf{x}\| = \|\mathbf{x} - \mathbf{y} + \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{y}\|$ we obtain $\|\mathbf{x} - \mathbf{y}\| \geq \|\mathbf{x}\| - \|\mathbf{y}\|$. By symmetry $\|\mathbf{x} - \mathbf{y}\| = \|\mathbf{y} - \mathbf{x}\| \geq \|\mathbf{y}\| - \|\mathbf{x}\|$ and we obtain the inverse triangle inequality.
2. This follows from the inverse triangle inequality.
3. The following proof can be skipped by those who do not have the necessary background in advanced calculus. Define the $\|\cdot\|'$ unit sphere

$$\mathcal{S} := \{\mathbf{y} \in \mathcal{V} : \|\mathbf{y}\|' = 1\}.$$

The set \mathcal{S} is a closed and bounded set and the function $f : \mathcal{S} \rightarrow \mathbb{R}$ given by $f(\mathbf{y}) = \|\mathbf{y}\|$ is continuous by what we just showed. Therefore f attains its minimum and maximum value on \mathcal{S} . Thus, there are positive constants m and M such that

$$m \leq \|\mathbf{y}\| \leq M, \quad \mathbf{y} \in \mathcal{S}. \tag{16}$$

For any $\mathbf{x} \in \mathcal{V}$ one has $\mathbf{y} := \mathbf{x}/\|\mathbf{x}\|' \in \mathcal{S}$, and (15) follows if we apply (16) to these \mathbf{y} .

□

0.3.1 Convergence of Vectors

Consider an infinite sequence $\{\mathbf{x}_k\} = \mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$ of vectors in \mathbb{R}^n . This sequence converges to zero if and only if each component sequence $\mathbf{x}_k(j)$ converges to zero for $j = 1, \dots, n$. In terms of the natural basis we have $\mathbf{x}_k = \sum_{j=1}^n \mathbf{x}_k(j) \mathbf{e}_j$ and another way of stating convergence to zero is that in terms of the basis $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ for \mathbb{R}^n each coefficient $\mathbf{x}_k(j)$ of \mathbf{x}_k converges to zero.

In this section \mathcal{V} will be a general finite dimensional vector space over \mathbb{R} or \mathbb{C} .

Definition 0.21 (Convergence of vectors)

Let $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ be a basis for a finite dimensional vector space \mathcal{V} and let $\{\mathbf{x}_k\}$ be an infinite sequence of vectors in \mathcal{V} with **basis coefficients** $\{\mathbf{c}_k\}$, i.e. $\mathbf{x}_k = \sum_{j=1}^n c_{kj} \mathbf{v}_j$ for each k . We say that $\{\mathbf{x}_k\}$ converges to zero, or have the limit zero, if $\lim_{k \rightarrow \infty} c_{kj} = 0$ for $j = 1, \dots, n$. We say that $\{\mathbf{x}_k\}$ converge to the limit \mathbf{x} in \mathcal{V} if $\mathbf{x}_k - \mathbf{x}$ converges to zero. We write this as $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}$ or $\mathbf{x}_k \rightarrow \mathbf{x}$ (as $k \rightarrow \infty$).

This definition is actually independent of the basis chosen. If $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ is another basis for \mathcal{V} and $\mathbf{x}_k = \sum_{j=1}^n b_{kj} \mathbf{w}_j$ for each k then from Lemma 0.15 $\mathbf{b}_k = \mathbf{A}\mathbf{c}_k$ for some nonsingular matrix \mathbf{A} independent of k . Hence $\mathbf{c}_k \rightarrow \mathbf{0}$ if and only if $\mathbf{b}_k \rightarrow \mathbf{0}$. If $\{a_k\}$ and $\{b_k\}$ are sequences of scalars and $\{\mathbf{x}_k\}$ and $\{\mathbf{y}_k\}$ are sequences of vectors such that $\{a_k\} \rightarrow a$, $\{b_k\} \rightarrow b$, $\{\mathbf{x}_k\} \rightarrow \mathbf{x}$, and $\{\mathbf{y}_k\} \rightarrow \mathbf{y}$ then $\{a_k \mathbf{x}_k + b_k \mathbf{y}_k\} \rightarrow a\mathbf{x} + b\mathbf{y}$. This shows that scalar multiplication and vector addition are continuous functions with respect to this notion of limit.

Corresponding to a basis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$, we define the **coefficient norm**

$$\|\mathbf{x}\|_c := \max_{1 \leq j \leq n} |c_j| \text{ where } \mathbf{x} = \sum_{j=1}^n c_j \mathbf{v}_j.$$

We leave as an exercise to show that this is a norm on \mathcal{V} . Recall that any two norms on \mathcal{V} are equivalent. This implies that for any other norm $\|\cdot\|$ on \mathcal{V} there are positive constants α, β such that any $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{v}_j$ satisfy

$$\|\mathbf{x}\| \leq \alpha \max_{1 \leq j \leq n} |c_j| \text{ and } |c_j| \leq \beta \|\mathbf{x}\| \text{ for } j = 1, \dots, n. \quad (17)$$

The notion of limit can be stated in terms of convergence in norm.

Theorem 0.22 (Norm convergence)

In a finite dimensional vector space \mathcal{V} with norm $\|\cdot\|$ we have $\mathbf{x}_k \rightarrow \mathbf{x}$ if and only if $\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}\| = 0$.

Proof. Suppose $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is a basis for the vector space and assume $\mathbf{x}_k, \mathbf{x} \in \mathcal{V}$. Then $\mathbf{x}_k - \mathbf{x} = \sum_{j=1}^n c_{kj} \mathbf{v}_j$ for some scalars c_{kj} . By (17) we see that

$$\frac{1}{\beta} \max_{k,j} |c_{kj}| \leq \|\mathbf{x}_k - \mathbf{x}\| \leq \alpha \max_{k,j} |c_{kj}|$$

and hence $\|\mathbf{x}_k - \mathbf{x}\| \rightarrow 0 \Leftrightarrow \lim_k c_{kj} \rightarrow 0$ for each $j \Leftrightarrow \mathbf{x}_k \rightarrow \mathbf{x}$. \square

Since all vector norms are equivalent we have convergence in any norm we can define on a finite dimensional vector space.

Definition 0.23 (Cauchy sequence)

Let \mathcal{V} be a finite dimensional vector space with norm $\| \cdot \|$ and let $\{\mathbf{x}_k\}$ in \mathcal{V} be an infinite sequence.

1. $\{\mathbf{x}_k\}$ is a **Cauchy sequence** if $\lim_{k,l \rightarrow \infty} (\mathbf{x}_k - \mathbf{x}_l) = 0$ or equivalently $\lim_{k,l \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}_l\| = 0$. More precisely, for each $\epsilon > 0$ there is an integer $N \in \mathbb{N}$ such that for each $k, l \geq N$ we have $\|\mathbf{x}_k - \mathbf{x}_l\| \leq \epsilon$.
2. The normed vector space is said to be **complete** if every Cauchy sequence converges to a point in the space.
3. $\{\mathbf{x}_k\}$ is called **bounded** if there is a positive number M such that $\|\mathbf{x}_k\| \leq M$ for all k .
4. $\{\mathbf{x}_{n_k}\}$ is said to be a **subsequence** of $\{\mathbf{x}_k\}_{k \geq 0}$ if $0 \leq n_0 < n_1 < n_2 \dots$.

Theorem 0.24 (Complete vector space)

In a finite dimensional vector space \mathcal{V} the following hold:

1. A sequence in \mathcal{V} is convergent if and only if it is a Cauchy sequence.
2. \mathcal{V} is complete.
3. Every bounded sequence in \mathcal{V} has a convergent subsequence.

Proof.

1. Suppose $\mathbf{x}_k \rightarrow \mathbf{x}$. By the triangle inequality $\|\mathbf{x}_k - \mathbf{x}_l\| \leq \|\mathbf{x}_k - \mathbf{x}\| + \|\mathbf{x}_l - \mathbf{x}\|$ and hence $\|\mathbf{x}_k - \mathbf{x}_l\| \rightarrow 0$. Conversely, let $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ be a basis for \mathcal{V} and $\{\mathbf{x}_k\}$ a Cauchy sequence with $\mathbf{x}_k = \sum_{j=1}^n c_{kj} \mathbf{v}_j$ for each k . Then $\mathbf{x}_k - \mathbf{x}_l = \sum_{j=1}^n (c_{kj} - c_{lj}) \mathbf{v}_j$ and since $\lim_{k,l \rightarrow \infty} (\mathbf{x}_k - \mathbf{x}_l) = 0$ we have by definition of convergence $\lim_{k,l \rightarrow \infty} (c_{kj} - c_{lj}) = 0$ for $j = 1, \dots, n$. Thus for each j we have a Cauchy-sequence $\{c_{kj}\} \in \mathbb{C}$ and since \mathbb{C} is complete $\{c_{kj}\}$ converges to some $c_j \in \mathbb{C}$. But then $\mathbf{x}_k \rightarrow \mathbf{x} := \sum_{j=1}^n c_j \mathbf{v}_j \in \mathcal{V}$.
2. \mathcal{V} is complete since we just showed that every Cauchy sequence converges to a point in the space.
3. Let $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ be a basis for \mathcal{V} and $\{\mathbf{x}_k\}$ be a bounded sequence with $\mathbf{x}_k = \sum_{j=1}^n c_{kj} \mathbf{v}_j$ for each k . By (17) each coefficient sequence $\{c_{kj}\}_k$ is a bounded sequence of complex numbers and therefore, by a well known property of complex numbers, has a convergent subsequence. In particular the sequence of \mathbf{v}_1 coefficients $\{c_{k1}\}$ has a convergent subsequence $c_{k_i,1}$. For the second component the sequence $\{c_{k_i,2}\}$ has a convergent subsequence, say $c_{l_i,2}$. Continuing with $j = 3, \dots, n$ we obtain integers $0 \leq m_0 < m_1 < \dots$ such that $\{c_{m_i,j}\}$ is a convergent subsequence of c_{kj} for $j = 1, \dots, n$. But then $\{\mathbf{x}_{m_i}\}$ is a convergent subsequence of $\{\mathbf{x}_k\}$.

□

0.3.2 Convergence of Series of Vectors

Consider now an infinite series $\sum_{m=0}^{\infty} \mathbf{y}_m$ of vectors in a finite dimensional vector space \mathcal{V} . We say that the series converges if the sequence of partial sums $\{\mathbf{x}_k\}$ given by $\mathbf{x}_k = \sum_{m=0}^k \mathbf{y}_m$ converges. A sufficient condition for convergence is that $\sum_{m=0}^{\infty} \|\mathbf{y}_m\|$ converges for some vector norm. We say that the series converges **absolutely** if this is the case. Note that $\|\sum_{m=0}^{\infty} \mathbf{y}_m\| \leq \sum_{m=0}^{\infty} \|\mathbf{y}_m\|$, and absolute convergence in one norm implies absolute convergence in any norm by equivalence of norms. In an absolute convergent series we may change the order of the terms without changing the value of the sum.

Exercise 0.25 (Linear combinations of convergent sequences)

Show that if $\{a_k\} \rightarrow a$, $\{b_k\} \rightarrow b$, $\{\mathbf{x}_k\} \rightarrow \mathbf{x}$, and $\{\mathbf{y}_k\} \rightarrow \mathbf{y}$ then $\{a_k \mathbf{x}_k + b_k \mathbf{y}_k\} \rightarrow a\mathbf{x} + b\mathbf{y}$.

Exercise 0.26 (Coefficient norm)

Show that $\|\cdot\|_c$ is a norm.

0.4 Inner Products

An **inner product** or **scalar product** in a vector space is a function mapping pairs of vectors into a scalar. We consider first the real case.

Definition 0.27 (Real inner product)

An **inner product** in a real vector space \mathcal{V} is a function $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ satisfying for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$ and all $a, b \in \mathbb{R}$ the following conditions:

1. $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$. (positivity)
2. $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ (symmetry)
3. $\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$. (linearity)

The pair $(\mathcal{V}, \langle \cdot, \cdot \rangle)$ is called a **real inner product space**. The function

$$\|\cdot\| : \mathcal{V} \longrightarrow \mathbb{R}, \quad \mathbf{x} \longmapsto \|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \tag{18}$$

is called the **inner product norm**.

The **standard inner product** in $\mathcal{V} = \mathbb{R}^n$ is given by $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{y}$. It is clearly an inner product in \mathbb{R}^n . The corresponding inner product norm is the Euclidian norm $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2$.

Consider next inner products in a complex vector space. Property 2. in the definition of a real inner product is altered from symmetry to skew symmetry.

Definition 0.28 (Complex inner product)

An **inner product** in a complex vector space \mathcal{V} is a function $\mathcal{V} \times \mathcal{V} \rightarrow \mathbb{C}$ satisfying for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$ and all $a, b \in \mathbb{C}$ the following conditions:

1. $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$. (positivity)
2. $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$ (skew symmetry)
3. $\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$. (linearity)

The pair $(\mathcal{V}, \langle \cdot, \cdot \rangle)$ is called a **complex inner product space**. The function

$$\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto \|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \quad (19)$$

is called the **inner product norm**.

Note the complex conjugate in 2. We find

$$\langle \mathbf{x}, a\mathbf{y} + b\mathbf{z} \rangle = \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle + \bar{b}\langle \mathbf{x}, \mathbf{z} \rangle, \quad \langle a\mathbf{x}, a\mathbf{y} \rangle = |a|^2 \langle \mathbf{x}, \mathbf{y} \rangle. \quad (20)$$

The **standard inner product** in \mathbb{C}^n is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^* \mathbf{x} = \mathbf{x}^T \bar{\mathbf{y}} = \sum_{j=1}^n x_j \bar{y}_j.$$

It is clearly an inner product in \mathbb{C}^n . The corresponding inner product norm is the Euclidian norm $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^* \mathbf{x}}$.

The following inequality holds for any inner product.

Theorem 0.29 (Cauchy-Schwarz inequality)

For any \mathbf{x}, \mathbf{y} in a real or complex inner product space

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|, \quad (21)$$

with equality if and only if \mathbf{x} and \mathbf{y} are linearly dependent.

Proof. If $\mathbf{y} = \mathbf{0}$ then $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, 0\mathbf{y} \rangle = 0$, $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ and $\|\mathbf{y}\| = 0$. Thus the inequality holds with equality, and \mathbf{x} and \mathbf{y} are linearly dependent. So assume $\mathbf{y} \neq \mathbf{0}$. Define

$$z := \mathbf{x} - a\mathbf{y}, \quad a := \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{y} \rangle}.$$

Then $\langle \mathbf{z}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle - a\langle \mathbf{y}, \mathbf{y} \rangle = 0$ so that by 2. and (20)

$$\langle a\mathbf{y}, \mathbf{z} \rangle + \langle \mathbf{z}, a\mathbf{y} \rangle = a\overline{\langle \mathbf{z}, \mathbf{y} \rangle} + \bar{a}\langle \mathbf{z}, \mathbf{y} \rangle = 0. \quad (22)$$

But then

$$\begin{aligned} \|\mathbf{x}\|^2 &= \langle \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{z} + a\mathbf{y}, \mathbf{z} + a\mathbf{y} \rangle \\ &\stackrel{(22)}{=} \langle \mathbf{z}, \mathbf{z} \rangle + \langle a\mathbf{y}, a\mathbf{y} \rangle \stackrel{(20)}{=} \|\mathbf{z}\|^2 + |a|^2\|\mathbf{y}\|^2 \\ &\geq |a|^2\|\mathbf{y}\|^2 = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2}. \end{aligned}$$

Multiplying by $\|\mathbf{y}\|^2$ gives (21). We have equality if and only if $\mathbf{z} = \mathbf{0}$, which means that \mathbf{x} and \mathbf{y} are linearly dependent. \square

Theorem 0.30 (Inner product norm)

The inner product norm is a vector norm.

Proof. For all \mathbf{x}, \mathbf{y} in an inner product space and all a in \mathbb{C} we need to show

1. $\|\mathbf{x}\| \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$. (positivity)
2. $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$. (homogeneity)
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$. (subadditivity)

The first statement is an immediate consequence of positivity, while the second one follows from (20). Expanding $\|\mathbf{x} + a\mathbf{y}\|^2 = \langle \mathbf{x} + a\mathbf{y}, \mathbf{x} + a\mathbf{y} \rangle$ using (20) we obtain

$$\|\mathbf{x} + a\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + a\langle \mathbf{y}, \mathbf{x} \rangle + \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle + |a|^2\|\mathbf{y}\|^2, \quad a \in \mathbb{C}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{V}. \quad (23)$$

Now (23) with $a = 1$ and the Cauchy-Schwarz inequality implies

$$\|\mathbf{x} + \mathbf{y}\|^2 \leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 = (\|\mathbf{x}\| + \|\mathbf{y}\|)^2.$$

Taking square roots completes the proof. \square

We also have

Theorem 0.31 (Parallelogram Identity)

For all \mathbf{x}, \mathbf{y} in a real or complex inner product space

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2.$$

Proof. We set $a = \pm 1$ in (23) and add the two equations. \square

The inner product of two vectors can be written as a sum of inner product norms. In the real case for any $\mathbf{x}, \mathbf{y} \in \mathcal{V}$

$$4\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2, \quad (24)$$

while in the complex case it follows from (23) that

$$4\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2 + i\|\mathbf{x} + i\mathbf{y}\|^2 - i\|\mathbf{x} - i\mathbf{y}\|^2, \quad (25)$$

where $i = \sqrt{-1}$. See Exercise 0.33.

In the real case the Cauchy-Schwarz inequality implies that $-1 \leq \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} \leq 1$ for nonzero \mathbf{x} and \mathbf{y} , so there is a unique angle θ in $[0, \pi]$ such that

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}. \quad (26)$$

This defines the **angle** between vectors in a real inner product space.

Exercise 0.32 (The $\mathbf{A}^T \mathbf{A}$ inner product)

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ has linearly independent columns. Show that $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{y}$ defines an inner product on \mathbb{R}^n .

Exercise 0.33 (Complex inner product as sums of norms)

Show (25).

Exercise 0.34 (Angle between vectors in complex case)

Show that in the complex case there is a unique angle θ in $[0, \pi/2]$ such that

$$\cos \theta = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}. \quad (27)$$

0.4.1 Orthogonality

Definition 0.35 (Orthogonality)

Two vectors \mathbf{x}, \mathbf{y} in a real or complex inner product space are **orthogonal** or **perpendicular**, denoted as $\mathbf{x} \perp \mathbf{y}$, if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. The vectors are **orthonormal** if in addition $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$.

Theorem 0.36 (Pythagoras)

For a real or complex inner product space

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2, \quad \text{if } \mathbf{x} \perp \mathbf{y}. \quad (28)$$

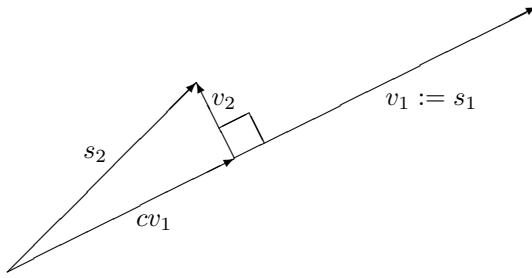


Figure 1: The construction of \mathbf{v}_1 and \mathbf{v}_2 in Gram-Schmidt. The constant c is given by $c := \langle \mathbf{s}_2, \mathbf{v}_1 \rangle / \langle \mathbf{v}_1, \mathbf{v}_1 \rangle$.

Proof. We set $a = 1$ in (23) and use the orthogonality. \square

Definition 0.37 (Orthogonal- and Orthonormal Bases)

A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ in a subspace \mathcal{S} of a real or complex inner product space is an **orthogonal basis** for \mathcal{S} if it is a basis for \mathcal{S} and $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$ for $i \neq j$. It is an **orthonormal basis** for \mathcal{S} if it is a basis for \mathcal{S} and $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}$ for all i, j .

A basis for a subspace of an inner product space can be turned into an orthogonal- or orthonormal basis for the subspace by the following construction.

Theorem 0.38 (Gram-Schmidt)

Let $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$ be a basis for a real or complex inner product space $(\mathcal{S}, \langle \cdot, \cdot \rangle)$. Define

$$\mathbf{v}_1 := \mathbf{s}_1, \quad \mathbf{v}_j := \mathbf{s}_j - \sum_{i=1}^{j-1} \frac{\langle \mathbf{s}_j, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \mathbf{v}_i, \quad j = 2, \dots, k. \quad (29)$$

Then $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is an orthogonal basis for \mathcal{S} and the normalized vectors

$$\{\mathbf{u}_1, \dots, \mathbf{u}_k\} := \left\{ \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}, \dots, \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|} \right\}$$

form an orthonormal basis for \mathcal{S} .

Proof. To show that $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is an orthogonal basis for \mathcal{S} we use induction on k . Let $S_j := \text{span}\{\mathbf{s}_1, \dots, \mathbf{s}_j\}$ for $j = 1, \dots, k$. Clearly $\mathbf{v}_1 = \mathbf{s}_1$ is an orthogonal basis for S_1 . Suppose for some $j \geq 2$ that $\mathbf{v}_1, \dots, \mathbf{v}_{j-1}$ is an orthogonal basis for

S_{j-1} and let \mathbf{v}_j be given by (29) as a linear combination of \mathbf{s}_j and $\mathbf{v}_1, \dots, \mathbf{v}_{j-1}$. Replacing each of these \mathbf{v}_i by a linear combination of $\mathbf{s}_1, \dots, \mathbf{s}_{j-1}$ we obtain $\mathbf{v}_j = \sum_{i=1}^j a_i \mathbf{s}_i$ for some a_0, \dots, a_j with $a_j = 1$. Since $\mathbf{s}_1, \dots, \mathbf{s}_j$ are linearly independent and $a_j \neq 0$ we deduce that $\mathbf{v}_j \neq 0$. By the induction hypothesis

$$\langle \mathbf{v}_j, \mathbf{v}_l \rangle = \langle \mathbf{s}_j, \mathbf{v}_l \rangle - \sum_{i=1}^{j-1} \frac{\langle \mathbf{s}_j, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \langle \mathbf{v}_i, \mathbf{v}_l \rangle = \langle \mathbf{s}_j, \mathbf{v}_l \rangle - \frac{\langle \mathbf{s}_j, \mathbf{v}_l \rangle}{\langle \mathbf{v}_l, \mathbf{v}_l \rangle} \langle \mathbf{v}_l, \mathbf{v}_l \rangle = 0$$

for $l = 1, \dots, j-1$. Thus $\mathbf{v}_1, \dots, \mathbf{v}_j$ is an orthogonal basis for S_j .

If $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is an orthogonal basis for \mathcal{S} then clearly $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ is an orthonormal basis for \mathcal{S} . \square

Sometimes we want to extend an orthogonal basis for a subspace to an orthogonal basis for a larger space.

Theorem 0.39 (Orthogonal Extension of basis)

Suppose $\mathcal{S} \subset \mathcal{T}$ are finite dimensional subspaces of a vector space \mathcal{V} . An orthogonal basis for \mathcal{S} can always be extended to an orthogonal basis for \mathcal{T} .

Proof. Suppose $\dim \mathcal{S} := k < \dim \mathcal{T} = n$. Using Theorem 0.10 we first extend an orthogonal basis $\mathbf{s}_1, \dots, \mathbf{s}_k$ for \mathcal{S} to a basis $\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{s}_{k+1}, \dots, \mathbf{s}_n$ for \mathcal{T} and then apply the Gram-Schmidt process to this basis obtaining an orthogonal basis $\mathbf{v}_1, \dots, \mathbf{v}_n$ for \mathcal{T} . This is an extension of the basis for \mathcal{S} since $\mathbf{v}_i = \mathbf{s}_i$ for $i = 1, \dots, k$. We show this by induction. Clearly $\mathbf{v}_1 = \mathbf{s}_1$. Suppose for some $2 \leq r < k$ that $\mathbf{v}_j = \mathbf{s}_j$ for $j = 1, \dots, r-1$. Consider (29) for $j = r$. Since $\langle \mathbf{s}_r, \mathbf{v}_i \rangle = \langle \mathbf{s}_r, \mathbf{s}_i \rangle = 0$ for $i < r$ we obtain $\mathbf{v}_r = \mathbf{s}_r$. \square

Letting $\mathcal{S} = \text{span}(\mathbf{s}_1, \dots, \mathbf{s}_k)$ and \mathcal{T} be \mathbb{R}^n or \mathbb{C}^n we obtain

Corollary 0.40 (Extending orthogonal vectors to a basis)

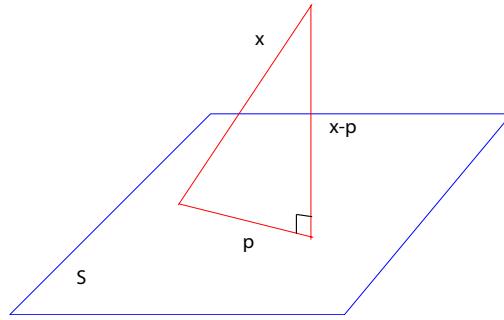
For $1 \leq k < n$ a set $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$ of nonzero orthogonal vectors in \mathbb{R}^n or \mathbb{C}^n can be extended to an orthogonal basis for the whole space.

0.4.2 Orthogonal Projection, and the Column Space Decomposition

Theorem 0.41 (Orthogonal Projection)

Let \mathcal{S} be a subspace of a finite dimensional real or complex inner product space $(\mathcal{V}, \langle \cdot, \cdot \rangle)$. To each $\mathbf{x} \in \mathcal{V}$ there is a unique vector $\mathbf{p} \in \mathcal{S}$, called the **orthogonal projection of \mathbf{x} into \mathcal{S}** , such that

$$\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0, \quad \text{for all } \mathbf{s} \in \mathcal{S}. \quad (30)$$

Figure 2: The orthogonal projection of \mathbf{x} into \mathcal{S} .

If $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is an orthogonal basis for \mathcal{S} then

$$\mathbf{p} = \sum_{i=1}^k \frac{\langle \mathbf{x}, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \mathbf{v}_i. \quad (31)$$

Proof. Define \mathbf{p} by (31). Then

$$\langle \mathbf{p}, \mathbf{v}_j \rangle = \sum_{i=1}^k \frac{\langle \mathbf{x}, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \langle \mathbf{v}_i, \mathbf{v}_j \rangle = \frac{\langle \mathbf{x}, \mathbf{v}_j \rangle}{\langle \mathbf{v}_j, \mathbf{v}_j \rangle} \langle \mathbf{v}_j, \mathbf{v}_j \rangle = \langle \mathbf{x}, \mathbf{v}_j \rangle$$

so that by linearity $\langle \mathbf{x} - \mathbf{p}, \mathbf{v}_j \rangle = 0$ for $j = 1, \dots, k$. But then $\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0$ for all $\mathbf{s} \in \mathcal{S}$. Indeed, if $\mathbf{s} = \sum_{j=1}^n a_j \mathbf{v}_j$ then $\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = \sum_j \bar{a}_j \langle \mathbf{x} - \mathbf{p}, \mathbf{v}_j \rangle = 0$. This shows existence of a \mathbf{p} satisfying (30). For uniqueness suppose $\mathbf{p}_1, \mathbf{p}_2 \in \mathcal{S}$ and $\langle \mathbf{x} - \mathbf{p}_1, \mathbf{s} \rangle = \langle \mathbf{x} - \mathbf{p}_2, \mathbf{s} \rangle = 0$ for all $\mathbf{s} \in \mathcal{S}$. Then $\langle \mathbf{x} - \mathbf{p}_1, \mathbf{s} \rangle - \langle \mathbf{x} - \mathbf{p}_2, \mathbf{s} \rangle = \langle \mathbf{p}_2 - \mathbf{p}_1, \mathbf{s} \rangle = 0$ for all $\mathbf{s} \in \mathcal{S}$ and in particular $\langle \mathbf{p}_2 - \mathbf{p}_1, \mathbf{p}_2 - \mathbf{p}_1 \rangle = 0$ which implies that $\mathbf{p}_2 - \mathbf{p}_1 = \mathbf{0}$ or $\mathbf{p}_1 = \mathbf{p}_2$. Now (30) holds if p is given by (31), and by uniqueness this is the only p . \square

Definition 0.42 (Orthogonal sum)

Let \mathcal{S} and \mathcal{T} be subspaces of a vector space \mathcal{V} and $\langle \cdot, \cdot \rangle$ an inner product on \mathcal{V} . The sum $\mathcal{B} := \mathcal{S} + \mathcal{T}$ is called an **orthogonal sum** if $\langle \mathbf{s}, \mathbf{t} \rangle = 0$ for each $\mathbf{s} \in \mathcal{S}$ and $\mathbf{t} \in \mathcal{T}$. We often write $\mathcal{S} \overset{\perp}{\oplus} \mathcal{T}$ instead of $\mathcal{S} + \mathcal{T}$ to indicate an orthogonal sum.

An orthogonal sum is a direct sum. For if $\mathbf{b} \in \mathcal{S} \cap \mathcal{T}$ then \mathbf{b} is orthogonal to itself, $\langle \mathbf{b}, \mathbf{b} \rangle = 0$, which implies that $\mathbf{b} = 0$. Thus, every $\mathbf{b} \in \mathcal{S} \overset{\perp}{\oplus} \mathcal{T}$ can be written

uniquely as $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$, where $\mathbf{b}_1 \in \mathcal{S}$ and $\mathbf{b}_2 \in \mathcal{T}$. The vectors \mathbf{b}_1 and \mathbf{b}_2 are the orthogonal projections of \mathbf{b} into \mathcal{S} and \mathcal{T} . Indeed, for any $\mathbf{s} \in \mathcal{S}$ and $\mathbf{t} \in \mathcal{T}$, $\langle \mathbf{b} - \mathbf{b}_1, \mathbf{s} \rangle = \langle \mathbf{b}_2, \mathbf{s} \rangle = 0$ and $\langle \mathbf{b} - \mathbf{b}_2, \mathbf{t} \rangle = \langle \mathbf{b}_1, \mathbf{t} \rangle = 0$.

The following orthogonal sum plays a major role when studying least squares methods.

Theorem 0.43 (Column space decomposition)

For any $m, n \in \mathbb{N}$ and $\mathbf{A} \in \mathbb{C}^{m,n}$ we have

$$\mathbb{C}^m = \text{span}(\mathbf{A}) \overset{\perp}{\oplus} \ker(\mathbf{A}^*), \quad (32)$$

where $\langle \mathbf{s}, \mathbf{t} \rangle := \mathbf{t}^* \mathbf{s}$ is the usual inner product on \mathbb{C}^m .

Proof. If $\mathbf{s} \in \mathcal{S} := \text{span}(\mathbf{A})$ and $\mathbf{t} \in \mathcal{T} := \ker(\mathbf{A}^*)$ then $\mathbf{s} = \mathbf{Ax}$ for some $\mathbf{x} \in \mathbb{C}^n$ and $\mathbf{A}^* \mathbf{t} = \mathbf{0}$. But then $\mathbf{t}^* \mathbf{s} = \mathbf{t}^* \mathbf{As} = (\mathbf{A}^* \mathbf{t})^* \mathbf{s} = \mathbf{0}$ and $\mathcal{S} + \mathcal{T}$ is an orthogonal sum. To show that $\mathbb{C}^m = \mathcal{S} \overset{\perp}{\oplus} \mathcal{T}$ we consider any $\mathbf{b} \in \mathbb{C}^m$. Let \mathbf{b}_1 be the orthogonal projection of \mathbf{b} into \mathcal{S} and define $\mathbf{b}_2 := \mathbf{b} - \mathbf{b}_1$. Then $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$ and we need to show that $\mathbf{b}_2 \in \mathcal{T}$ or $\mathbf{A}^* \mathbf{b}_2 = \mathbf{0}$. It follows from Theorem 0.41 that $(\mathbf{b} - \mathbf{b}_1)^* \mathbf{s} = \mathbf{b}_2^* \mathbf{s} = 0$ for any $\mathbf{s} \in \mathcal{S}$. Moreover, $\mathbf{s} = \mathbf{Ax}$ for some $\mathbf{x} \in \mathbb{C}^n$. But then $\mathbf{b}_2^* \mathbf{s} = \mathbf{b}_2^* \mathbf{Ax} = (\mathbf{A}^* \mathbf{b}_2)^* \mathbf{x} = 0$ for all $\mathbf{x} \in \mathbb{C}^n$ and $\mathbf{A}^* \mathbf{b}_2 = \mathbf{0}$. \square

0.5 Linear Systems

Consider a linear system

$$\begin{aligned} a_{11}x_1 &+ a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 &+ a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots &\quad \vdots \quad \vdots \quad \vdots \\ a_{m1}x_1 &+ a_{m2}x_2 + \cdots + a_{mn}x_n = b_m \end{aligned}$$

of m equations in n unknowns. Here for all i, j , the coefficients a_{ij} , the unknowns x_j , and the components of the right hand sides b_i , are real or complex numbers. The system can be written as a vector equation

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \cdots + x_n \mathbf{a}_n = \mathbf{b},$$

where $\mathbf{a}_j = [a_{1j}, \dots, a_{mj}]^T \in \mathbb{C}^m$ for $j = 1, \dots, n$ and $\mathbf{b} = [b_1, \dots, b_m]^T \in \mathbb{C}^m$. It can also be written as a matrix equation

$$\mathbf{Ax} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} = \mathbf{b}.$$

The system is **homogeneous** if $\mathbf{b} = \mathbf{0}$ and it is said to be **underdetermined**, **square**, or **overdetermined** if $m < n$, $m = n$, or $m > n$, respectively.

A linear system may have a unique solution, infinitely many solutions, or no solution. To discuss this we first consider the real case, and a homogeneous underdetermined system.

Lemma 0.44 (Underdetermined system)

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n$. Then there is a nonzero $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \mathbf{0}$.

Proof. Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n$. The n columns of \mathbf{A} span a subspace of \mathbb{R}^m . Since \mathbb{R}^m has dimension m the dimension of this subspace is at most m . By Lemma 0.5 the columns of \mathbf{A} must be linearly dependent. It follows that there is a nonzero $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \mathbf{0}$. \square

Consider now a square linear system. The following definition is essential.

Definition 0.45 (Real nonsingular matrix)

A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is said to be **nonsingular** if the only real solution of the homogeneous system $\mathbf{Ax} = \mathbf{0}$ is $\mathbf{x} = \mathbf{0}$. The matrix is **singular** if it is not nonsingular.

Theorem 0.46 (Linear systems; existence and uniqueness)

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$. The linear system $\mathbf{Ax} = \mathbf{b}$ has a unique solution $\mathbf{x} \in \mathbb{R}^n$ for any $\mathbf{b} \in \mathbb{R}^n$ if and only if the matrix \mathbf{A} is nonsingular.

Proof. Suppose \mathbf{A} is nonsingular. We define $\mathbf{B} = [\mathbf{A} \ \mathbf{b}] \in \mathbb{R}^{n \times (n+1)}$ by adding a column to \mathbf{A} . By Lemma 0.44 there is a nonzero $\mathbf{z} \in \mathbb{R}^{n+1}$ such that $\mathbf{Bz} = \mathbf{0}$. If we write $\mathbf{z} = \begin{bmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{bmatrix}$ where $\tilde{\mathbf{z}} = [z_1, \dots, z_n]^T \in \mathbb{R}^n$ and $z_{n+1} \in \mathbb{R}$, then

$$\mathbf{Bz} = [\mathbf{A} \ \mathbf{b}] \begin{bmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{bmatrix} = \mathbf{A}\tilde{\mathbf{z}} + z_{n+1}\mathbf{b} = \mathbf{0}.$$

We cannot have $z_{n+1} = 0$ for then $\mathbf{A}\tilde{\mathbf{z}} = \mathbf{0}$ for a nonzero $\tilde{\mathbf{z}}$, contradicting the nonsingularity of \mathbf{A} . Define $\mathbf{x} := -\tilde{\mathbf{z}}/z_{n+1}$. Then

$$\mathbf{Ax} = -\mathbf{A}\left(\frac{\tilde{\mathbf{z}}}{z_{n+1}}\right) = -\frac{1}{z_{n+1}}\mathbf{A}\tilde{\mathbf{z}} = -\frac{1}{z_{n+1}}(-z_{n+1}\mathbf{b}) = \mathbf{b},$$

so \mathbf{x} is a solution.

Suppose $\mathbf{Ax} = \mathbf{b}$ and $\mathbf{Ay} = \mathbf{b}$ for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Then $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{0}$ and since \mathbf{A} is nonsingular we conclude that $\mathbf{x} - \mathbf{y} = \mathbf{0}$ or $\mathbf{x} = \mathbf{y}$. Thus the solution is unique.

Conversely, if $\mathbf{A}\mathbf{x} = \mathbf{b}$ has a unique solution for any $\mathbf{b} \in \mathbb{R}^n$ then $\mathbf{A}\mathbf{x} = \mathbf{0}$ has a unique solution which must be $\mathbf{x} = \mathbf{0}$. Thus \mathbf{A} is nonsingular. \square

For the complex case we have

Lemma 0.47 (Complex underdetermined system)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m < n$. Then there is a nonzero $\mathbf{x} \in \mathbb{C}^n$ such that $\mathbf{A}\mathbf{x} = \mathbf{0}$.

Definition 0.48 (Complex nonsingular matrix)

A square matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is said to be **nonsingular** if the only complex solution of the homogeneous system $\mathbf{A}\mathbf{x} = \mathbf{0}$ is $\mathbf{x} = \mathbf{0}$. The matrix is **singular** if it is not nonsingular.

Theorem 0.49 (Complex linear system; existence and uniqueness)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$. The linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ has a unique solution $\mathbf{x} \in \mathbb{C}^n$ for any $\mathbf{b} \in \mathbb{C}^n$ if and only if the matrix \mathbf{A} is nonsingular.

0.5.1 The Inverse Matrix

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a square matrix. A matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ is called a **right inverse** of \mathbf{A} if $\mathbf{AB} = \mathbf{I}$. A matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ is said to be a **left inverse** of \mathbf{A} if $\mathbf{CA} = \mathbf{I}$. We say that \mathbf{A} is **invertible** if it has both a left- and a right inverse. If \mathbf{A} has a right inverse \mathbf{B} and a left inverse \mathbf{C} then

$$\mathbf{C} = \mathbf{CI} = \mathbf{C}(\mathbf{AB}) = (\mathbf{CA})\mathbf{B} = \mathbf{IB} = \mathbf{B}$$

and this common inverse is called the **inverse** of \mathbf{A} and denoted by \mathbf{A}^{-1} . Thus the inverse satisfies $\mathbf{A}^{-1}\mathbf{A} = \mathbf{AA}^{-1} = \mathbf{I}$.

We want to characterize the class of invertible matrices and start with a lemma.

Theorem 0.50 (Product of nonsingular matrices)

If $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times n}$ with $\mathbf{AB} = \mathbf{C}$ then \mathbf{C} is nonsingular if and only if both \mathbf{A} and \mathbf{B} are nonsingular. In particular, if $\mathbf{AB} = \mathbf{I}$ or $\mathbf{BA} = \mathbf{I}$ then \mathbf{A} is nonsingular and $\mathbf{A}^{-1} = \mathbf{B}$.

Proof. Suppose both \mathbf{A} and \mathbf{B} are nonsingular and let $\mathbf{Cx} = \mathbf{0}$. Then $\mathbf{ABx} = \mathbf{0}$ and since \mathbf{A} is nonsingular we see that $\mathbf{Bx} = \mathbf{0}$. Since \mathbf{B} is nonsingular we have $\mathbf{x} = \mathbf{0}$. We conclude that \mathbf{C} is nonsingular.

For the converse suppose first that \mathbf{B} is singular and let $\mathbf{x} \in \mathbb{R}^n$ be a nonzero vector so that $\mathbf{Bx} = \mathbf{0}$. But then $\mathbf{Cx} = (\mathbf{AB})\mathbf{x} = \mathbf{A}(\mathbf{Bx}) = \mathbf{A}\mathbf{0} = \mathbf{0}$ so \mathbf{C} is

singular. Finally suppose \mathbf{B} is nonsingular, but \mathbf{A} is singular. Let $\tilde{\mathbf{x}}$ be a nonzero vector such that $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{0}$. By Theorem 0.46 there is a vector \mathbf{x} such that $\mathbf{Bx} = \tilde{\mathbf{x}}$ and \mathbf{x} is nonzero since $\tilde{\mathbf{x}}$ is nonzero. But then $\mathbf{Cx} = (\mathbf{AB})\mathbf{x} = \mathbf{A}(\mathbf{Bx}) = \mathbf{A}\tilde{\mathbf{x}} = \mathbf{0}$ for a nonzero vector \mathbf{x} and \mathbf{C} is singular. \square

Theorem 0.51 (When is a square matrix invertible?)

A square matrix is invertible if and only if it is nonsingular.

Proof. Suppose first \mathbf{A} is a nonsingular matrix. By Theorem 0.46 each of the linear systems $\mathbf{Ab}_i = \mathbf{e}_i$ has a unique solution \mathbf{b}_i for $i = 1, \dots, n$. Let $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$. Then $\mathbf{AB} = [\mathbf{Ab}_1, \dots, \mathbf{Ab}_n] = [\mathbf{e}_1, \dots, \mathbf{e}_n] = \mathbf{I}$ so that \mathbf{A} has a right inverse \mathbf{B} . By Lemma 0.50 \mathbf{B} is nonsingular since \mathbf{I} is nonsingular and $\mathbf{AB} = \mathbf{I}$. Since \mathbf{B} is nonsingular we can use what we have shown for \mathbf{A} to conclude that \mathbf{B} has a right inverse \mathbf{C} , i.e. $\mathbf{BC} = \mathbf{I}$. But then $\mathbf{AB} = \mathbf{BC} = \mathbf{I}$ so \mathbf{B} has both a right inverse and a left inverse which must be equal so $\mathbf{A} = \mathbf{C}$. Since $\mathbf{BC} = \mathbf{I}$ we have $\mathbf{BA} = \mathbf{I}$, so \mathbf{B} is also a left inverse of \mathbf{A} and \mathbf{A} is invertible.

Conversely, if \mathbf{A} is invertible then it has a right inverse \mathbf{B} . Since $\mathbf{AB} = \mathbf{I}$ and \mathbf{I} is nonsingular, we again use Lemma 0.50 to conclude that \mathbf{A} is nonsingular. \square

If $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$ and $\mathbf{AB} = \mathbf{I}$ then

1. since \mathbf{I} is nonsingular it follows from Lemma 0.50 that both \mathbf{A} and \mathbf{B} are nonsingular.
2. Since \mathbf{A} and \mathbf{B} are nonsingular Theorem 0.51 implies that \mathbf{A} and \mathbf{B} are inverses of each other.

Thus to verify that some matrix \mathbf{B} is an inverse of another matrix \mathbf{A} it is enough to show that \mathbf{B} is either a left inverse or a right inverse of \mathbf{A} . This calculation also proves that \mathbf{A} is nonsingular. We use this observation to give simple proofs of the following results.

Corollary 0.52 (Basic properties of the inverse matrix)

Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ are nonsingular and c is a nonzero constant.

1. \mathbf{A}^{-1} is nonsingular and $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$.
2. $\mathbf{C} = \mathbf{AB}$ is nonsingular and $\mathbf{C}^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$.
3. \mathbf{A}^T is nonsingular and $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T =: \mathbf{A}^{-T}$.
4. $c\mathbf{A}$ is nonsingular and $(c\mathbf{A})^{-1} = \frac{1}{c}\mathbf{A}^{-1}$.

Proof.

1. Since $\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$ the matrix \mathbf{A} is a right inverse of \mathbf{A}^{-1} . Thus \mathbf{A}^{-1} is nonsingular and $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$.
2. We note that $(\mathbf{B}^{-1}\mathbf{A}^{-1})(\mathbf{AB}) = \mathbf{B}^{-1}(\mathbf{A}^{-1}\mathbf{A})\mathbf{B} = \mathbf{B}^{-1}\mathbf{B} = \mathbf{I}$. Thus \mathbf{AB} is invertible with the indicated inverse since it has a left inverse.
3. Now $\mathbf{I} = \mathbf{I}^T = (\mathbf{A}^{-1}\mathbf{A})^T = \mathbf{A}^T(\mathbf{A}^{-1})^T$ showing that $(\mathbf{A}^{-1})^T$ is a right inverse of \mathbf{A}^T .
4. The matrix $\frac{1}{c}\mathbf{A}^{-1}$ is a one sided inverse of $c\mathbf{A}$.

□

Exercise 0.53 (The inverse of a general 2×2 matrix)*Show that*

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \alpha \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}, \quad \alpha = \frac{1}{ad - bc},$$

*for any a, b, c, d such that $ad - bc \neq 0$.***Exercise 0.54 (The inverse of a 2×2 matrix)***Find the inverse of*

$$A = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}.$$

Exercise 0.55 (Sherman-Morrison formula)*Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$, and $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times m}$ for some $n, m \in \mathbb{N}$. If $(\mathbf{I} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{B})^{-1}$ exists then*

$$(\mathbf{A} + \mathbf{BC}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C}^T \mathbf{A}^{-1}.$$

0.6 Determinants

Determinants, denoted by $\det(\cdot)$ or $|\cdot|$,² are useful for studying eigenvalues. Recall that if \mathbf{A}, \mathbf{B} are square matrices of order n with real or complex elements, then (see Appendix A for proofs)

1. $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$.
2. If \mathbf{A} is triangular then $\det(\mathbf{A}) = a_{11}a_{22} \cdots a_{nn}$. In particular, $\det(\mathbf{I}) = 1$.
3. $\det(\mathbf{A}^T) = \det(\mathbf{A})$, and $\det(\mathbf{A}^*) = \overline{\det(\mathbf{A})}$, (complex conjugate).
4. $\det(a\mathbf{A}) = a^n \det(\mathbf{A})$, for $a \in \mathbb{C}$.

²This notation is due to Cayley 1841. Gauss introduced determinants in 1801

5. \mathbf{A} is singular if and only if $\det(\mathbf{A}) = 0$.
6. If $\mathbf{A} = \begin{bmatrix} \mathbf{C} & \mathbf{D} \\ \mathbf{0} & \mathbf{E} \end{bmatrix}$ for some square matrices \mathbf{C}, \mathbf{E} then $\det(\mathbf{A}) = \det(\mathbf{C}) \det(\mathbf{E})$.
7. **Cramer's rule** Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular and $\mathbf{b} \in \mathbb{C}^n$. Let $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ be the unique solution of $\mathbf{Ax} = \mathbf{b}$. Then

$$x_j = \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n,$$

where $\mathbf{A}_j(\mathbf{b})$ denote the matrix obtained from \mathbf{A} by replacing the j th column of \mathbf{A} by \mathbf{b} .

8. **Adjoint.** Let $\mathbf{A}_{i,j}$ denote the submatrix of \mathbf{A} obtained by deleting the i th row and j th column of \mathbf{A} . For $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $1 \leq i, j \leq n$ the determinant $\det(\mathbf{A}_{ij})$ is called the **cofactor** of a_{ij} . The matrix $\text{adj}(\mathbf{A}) \in \mathbb{C}^{n \times n}$ with elements $\text{adj}(\mathbf{A})_{i,j} = (-1)^{i+j} \det(\mathbf{A}_{j,i})$ is called the **adjoint** of \mathbf{A} .
9. **Adjoint formula for the inverse.** If $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular then

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A}).$$

10. **Cofactor expansion.** For any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } i = 1, \dots, n, \quad (33)$$

$$\det(\mathbf{A}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } j = 1, \dots, n. \quad (34)$$

To compute the value of a determinant it is often convenient to use row- or column operations to introduce zeros in a row or column of \mathbf{A} and then use one of the cofactor expansions.

Exercise 0.56 (Cramer's rule; special case)

Solve the following system by Cramers rule:

$$\left[\begin{array}{cc} 1 & 2 \\ 2 & 1 \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \end{array} \right] = \left[\begin{array}{c} 3 \\ 6 \end{array} \right]$$

Exercise 0.57 (Adjoint matrix; special case)

Show that if

$$\mathbf{A} = \begin{bmatrix} 2 & -6 & 3 \\ 3 & -2 & -6 \\ 6 & 3 & 2 \end{bmatrix},$$

then

$$\text{adj}(\mathbf{A}) = \begin{bmatrix} 14 & 21 & 42 \\ -42 & -14 & 21 \\ 21 & -42 & 14 \end{bmatrix}.$$

Moreover,

$$\text{adj}(\mathbf{A})\mathbf{A} = \begin{bmatrix} 343 & 0 & 0 \\ 0 & 343 & 0 \\ 0 & 0 & 343 \end{bmatrix} = \det(\mathbf{A})\mathbf{I}.$$

Example 0.58 (Determinant equation for a straight line)

The equation for a straight line through two points (x_1, y_1) and (x_2, y_2) in the plane can be written as the equation

$$\det(\mathbf{A}) := \begin{vmatrix} 1 & x & y \\ 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \end{vmatrix} = 0$$

involving a determinant of order 3. We can compute this determinant using row operations of type 3. Subtracting row 2 from row 3 and then row 1 from row 2 we obtain

$$\begin{vmatrix} 1 & x & y \\ 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \end{vmatrix} = \begin{vmatrix} 1 & x & y \\ 0 & x_1 - x & y_1 - y \\ 0 & x_2 - x_1 & y_2 - y_1 \end{vmatrix} = (x_1 - x)(y_2 - y_1) - (y_1 - y)(x_2 - x_1).$$

Rearranging the equation $\det(\mathbf{A}) = 0$ we obtain

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1}(x - x_1)$$

which is the slope form of the equation of a straight line.

Exercise 0.59 (Determinant equation for a plane)

Show that

$$\begin{vmatrix} x & y & z & 1 \\ x_1 & y_1 & z_1 & 1 \\ x_2 & y_2 & z_2 & 1 \\ x_3 & y_3 & z_3 & 1 \end{vmatrix} = 0.$$

is the equation for a plane through three points (x_1, y_1, z_1) , (x_2, y_2, z_2) and (x_3, y_3, z_3) is

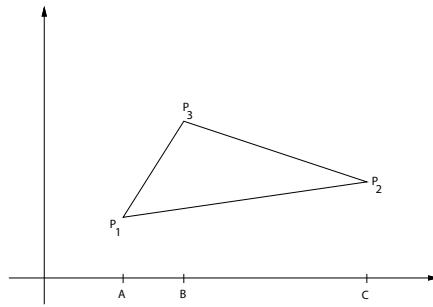


Figure 3: The triangle T defined by the three points P_1 , P_2 and P_3 .

Exercise 0.60 (Signed area of a triangle)

Let $P_i = (x_i, y_i)$, $i = 1, 2, 3$, be three points in the plane defining a triangle T . Show that the area of T is³

$$A(T) = \frac{1}{2} \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{vmatrix}.$$

The area is positive if we traverse the vertices in counterclockwise order.

Exercise 0.61 (Vandermonde matrix)

Show that

$$\begin{vmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{n-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^{n-1} \end{vmatrix} = \prod_{i>j} (x_i - x_j),$$

where $\prod_{i>j} (x_i - x_j) = \prod_{i=2}^n (x_i - x_1)(x_i - x_2) \cdots (x_i - x_{i-1})$. This determinant is called the Vandermonde determinant.⁴

Exercise 0.62 (Cauchy determinant (1842))

Let $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n]^T$, $\boldsymbol{\beta} = [\beta_1, \dots, \beta_n]^T$ be in \mathbb{R}^n .

- a) Consider the matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with elements $a_{i,j} = 1/(\alpha_i + \beta_j)$, $i, j = 1, 2, \dots, n$. Show that

$$\det(\mathbf{A}) = Pg(\boldsymbol{\alpha})g(\boldsymbol{\beta})$$

³Hint: $A(T) = A(ABP_3P_1) + A(P_3BCP_2) - A(P_1ACP_2)$, c.f. Figure 3

⁴Hint: subtract x_n^k times column k from column $k+1$ for $k = n-1, n-2, \dots, 1$.

where $P = \prod_{i=1}^n \prod_{j=1}^n a_{ij}$, and for $\gamma = [\gamma_1, \dots, \gamma_n]^T$

$$g(\gamma) = \prod_{i=2}^n (\gamma_i - \gamma_1)(\gamma_i - \gamma_2) \cdots (\gamma_i - \gamma_{i-1})$$

Hint: Multiply the i th row of \mathbf{A} by $\prod_{j=1}^n (\alpha_i + \beta_j)$ for $i = 1, 2, \dots, n$. Call the resulting matrix \mathbf{C} . Each element of \mathbf{C} is a product of $n-1$ factors $\alpha_r + \beta_s$. Hence $\det(\mathbf{C})$ is a sum of terms where each term contain precisely $n(n-1)$ factors $\alpha_r + \beta_s$. Thus $\det(\mathbf{C}) = q(\alpha, \beta)$ where q is a polynomial of degree at most $n(n-1)$ in α_i and β_j . Since $\det(\mathbf{A})$ and therefore $\det(\mathbf{C})$ vanishes if $\alpha_i = \alpha_j$ for some $i \neq j$ or $\beta_r = \beta_s$ for some $r \neq s$, we have that $q(\alpha, \beta)$ must be divisible by each factor in $g(\alpha)$ and $g(\beta)$. Since $g(\alpha)$ and $g(\beta)$ is a polynomial of degree $n(n-1)$, we have

$$q(\alpha, \beta) = kg(\alpha)g(\beta)$$

for some constant k independent of α and β . Show that $k = 1$ by choosing $\beta_i + \alpha_i = 0$, $i = 1, 2, \dots, n$.

- b) Notice that the cofactor of any element in the above matrix \mathbf{A} is the determinant of a matrix of similar form. Use the cofactor and determinant of \mathbf{A} to represent the elements of $\mathbf{A}^{-1} = (b_{j,k})$. Answer:

$$b_{j,k} = (\alpha_k + \beta_j)A_k(-\beta_j)B_j(-\alpha_k),$$

where

$$A_k(x) = \prod_{s \neq k} \left(\frac{\alpha_s - x}{\alpha_s - \alpha_k} \right), \quad B_k(x) = \prod_{s \neq k} \left(\frac{\beta_s - x}{\beta_s - \beta_k} \right).$$

Exercise 0.63 (Inverse of the Hilbert matrix)

Let $\mathbf{H}_n = (h_{i,j})$ be the $n \times n$ matrix with elements $h_{i,j} = 1/(i+j-1)$. Use Exercise 0.62 to show that the elements $t_{i,j}^n$ in $\mathbf{T}_n = \mathbf{H}_n^{-1}$ are given by

$$t_{i,j}^n = \frac{f(i)f(j)}{i+j-1},$$

where

$$f(i+1) = \left(\frac{i^2 - n^2}{i^2} \right) f(i), \quad i = 1, 2, \dots, \quad f(1) = -n.$$

0.7 Eigenpairs

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a square matrix, $\lambda \in \mathbb{C}$ and $\mathbf{x} \in \mathbb{C}^n$. We say that (λ, \mathbf{x}) is an **eigenpair** for \mathbf{A} if $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ and \mathbf{x} is nonzero. The scalar λ is called an **eigenvalue** and \mathbf{x} is said to be an **eigenvector**.⁵ The set of eigenvalues is called

⁵The word “eigen” is derived from German and means “own”

the **spectrum** of \mathbf{A} and is denoted by $\sigma(\mathbf{A})$. For example, $\sigma(\mathbf{I}) = \{1, \dots, 1\} = \{1\}$.

Lemma 0.64 (Characteristic equation)

For any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have $\lambda \in \sigma(\mathbf{A}) \iff \det(\mathbf{A} - \lambda\mathbf{I}) = 0$.

Proof. Suppose (λ, \mathbf{x}) is an eigenpair for \mathbf{A} . The equation $\mathbf{Ax} = \lambda\mathbf{x}$ can be written $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$. Since \mathbf{x} is nonzero the matrix $\mathbf{A} - \lambda\mathbf{I}$ must be singular with a zero determinant. Conversely, if $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$ then $\mathbf{A} - \lambda\mathbf{I}$ is singular and $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$ for some nonzero $\mathbf{x} \in \mathbb{C}^n$. Thus $\mathbf{Ax} = \lambda\mathbf{x}$ and (λ, \mathbf{x}) is an eigenpair for \mathbf{A} . \square

The expression $\det(\mathbf{A} - \lambda\mathbf{I})$ is a polynomial of exact degree n in λ . For $n = 3$ we have

$$\det(\mathbf{A} - \lambda\mathbf{I}) = \begin{vmatrix} a_{11} - \lambda & a_{12} & a_{13} \\ a_{21} & a_{22} - \lambda & a_{23} \\ a_{31} & a_{32} & a_{33} - \lambda \end{vmatrix}.$$

Expanding this determinant by the first column we find

$$\begin{aligned} \det(\mathbf{A} - \lambda\mathbf{I}) &= (a_{11} - \lambda) \begin{vmatrix} a_{22} - \lambda & a_{23} \\ a_{32} & a_{33} - \lambda \end{vmatrix} - a_{21} \begin{vmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} - \lambda \end{vmatrix} \\ &\quad + a_{31} \begin{vmatrix} a_{12} & a_{13} \\ a_{22} - \lambda & a_{23} \end{vmatrix} = (a_{11} - \lambda)(a_{22} - \lambda)(a_{33} - \lambda) + r(\lambda) \end{aligned}$$

for some polynomial r of degree at most one. In general

$$\det(\mathbf{A} - \lambda\mathbf{I}) = (a_{11} - \lambda)(a_{22} - \lambda) \cdots (a_{nn} - \lambda) + r(\lambda), \quad (35)$$

where each term in $r(\lambda)$ has at most $n - 2$ factors containing λ . It follows that r is a polynomial of degree at most $n - 2$, $\pi_{\mathbf{A}}$ is a polynomial of exact degree n , and the eigenvalues are the roots of this polynomial.

We observe that $\det(\mathbf{A} - \lambda\mathbf{I}) = (-1)^n \det(\lambda\mathbf{I} - \mathbf{A})$ so $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$ if and only if $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$.

Definition 0.65 (Characteristic polynomial of a matrix)

The function $\pi_{\mathbf{A}}: \mathbb{C} \rightarrow \mathbb{C}$ given by $\pi_{\mathbf{A}}(\lambda) = \det(\mathbf{A} - \lambda\mathbf{I})$ is called the **characteristic polynomial** of \mathbf{A} . The equation $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$ is called the **characteristic equation** of \mathbf{A} .

By the fundamental theorem of algebra an $n \times n$ matrix has, counting multiplicities, precisely n eigenvalues $\lambda_1, \dots, \lambda_n$ some of which might be complex even if \mathbf{A} is real. The complex eigenpairs of a real matrix occur in complex conjugate pairs. Indeed, taking the complex conjugate on both sides of the equation $\mathbf{Ax} = \lambda\mathbf{x}$ with \mathbf{A} real gives $\mathbf{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$.

The following results will be useful.

Theorem 0.66 (Derived eigenpairs)

Suppose (λ, \mathbf{x}) is an eigenpair for $\mathbf{A} \in \mathbb{C}^{n \times n}$. Then

1. If \mathbf{A} is nonsingular then $(\lambda^{-1}, \mathbf{x})$ is an eigenpair for \mathbf{A}^{-1} .
2. (λ^k, \mathbf{x}) is an eigenpair for \mathbf{A}^k for $k \in \mathbb{N}$.
3. If p given by $p(t) = a_0 + a_1t + a_2t^2 + \cdots + a_kt^k$ is a polynomial, then $(p(\lambda), \mathbf{x})$ is an eigenpair for the matrix $p(\mathbf{A}) := a_0\mathbf{I} + a_1\mathbf{A} + a_2\mathbf{A}^2 + \cdots + a_k\mathbf{A}^k$.
4. λ is an eigenvalue for \mathbf{A}^T , in fact $\pi_{\mathbf{A}^T} = \pi_{\mathbf{A}}$.
5. $\bar{\lambda}$ is an eigenvalue for \mathbf{A}^* , in fact $\pi_{\mathbf{A}^*}(\bar{\lambda}) = \overline{\pi_{\mathbf{A}}(\lambda)}$ for all $\lambda \in \mathbb{C}$.
6. If $\mathbf{A} = [\begin{smallmatrix} \mathbf{B} & \mathbf{C} \\ 0 & \mathbf{D} \end{smallmatrix}]$ is block triangular then $\pi_{\mathbf{A}} = \pi_{\mathbf{B}} \cdot \pi_{\mathbf{D}}$.

Proof.

1. $\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \implies \mathbf{A}^{-1}\mathbf{A}\mathbf{x} = \lambda^{-1}\mathbf{x} \implies \mathbf{x} = \lambda^{-1}\mathbf{x}$.
2. We use induction on k . The case $k = 1$ is trivial, and if $\mathbf{A}^{k-1}\mathbf{x} = \lambda^{k-1}\mathbf{x}$ then $\mathbf{A}^k\mathbf{x} = \mathbf{A}\mathbf{A}^{k-1}\mathbf{x} = \lambda^{k-1}\mathbf{A}\mathbf{x} = \lambda^k\mathbf{x}$.
3. $p(\mathbf{A})\mathbf{x} = \sum_{j=0}^k a_j \mathbf{A}^j \mathbf{x} \stackrel{2.}{=} \sum_{j=0}^k a_j \lambda^j \mathbf{x} = p(\lambda)\mathbf{x}$.
4. Using Property 3. of determinants we find for any $\lambda \in \mathbb{C}$

$$\pi_{\mathbf{A}^T}(\lambda) = \det(\mathbf{A}^T - \lambda\mathbf{I}) = \det((\mathbf{A} - \lambda\mathbf{I})^T) = \det(\mathbf{A} - \lambda\mathbf{I}) = \pi_{\mathbf{A}}(\lambda).$$

Thus \mathbf{A}^T and \mathbf{A} have the same characteristic polynomial and hence the same eigenvalues.

5. We have $\pi_{\mathbf{A}^*}(\bar{\lambda}) \stackrel{4.}{=} \pi_{\mathbf{A}}(\bar{\lambda}) = \det(\mathbf{A} - \bar{\lambda}\mathbf{I}) = \overline{\det(\mathbf{A} - \lambda\mathbf{I})} = \overline{\pi_{\mathbf{A}}(\lambda)}$. Thus $\pi_{\mathbf{A}}(\lambda) = 0 \Leftrightarrow \pi_{\mathbf{A}^*}(\bar{\lambda}) = 0$ and the result follows.
6. By Property 6. of determinants

$$\pi_{\mathbf{A}}(\lambda) = \begin{vmatrix} \mathbf{B} - \lambda\mathbf{I} & \mathbf{C} \\ 0 & \mathbf{D} - \lambda\mathbf{I} \end{vmatrix} = \det(\mathbf{B} - \lambda\mathbf{I}) \det(\mathbf{D} - \lambda\mathbf{I}) = \pi_{\mathbf{B}}(\lambda) \cdot \pi_{\mathbf{D}}(\lambda).$$

□

In general it is not easy to find all eigenvalues of a matrix. One notable exception is a triangular matrix. By Property 2. of determinants we obtain

Theorem 0.67 (Eigenvalues of a triangular matrix)

The eigenvalues of a triangular matrix are given by its diagonal elements.

There are two useful relations between the elements of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and its eigenvalues $\lambda_1, \dots, \lambda_n$.

Theorem 0.68 (Sums and products of eigenvalues; trace)

For any $\mathbf{A} \in \mathbb{C}^{n \times n}$

$$\text{trace}(\mathbf{A}) = \lambda_1 + \lambda_2 + \cdots + \lambda_n, \quad \det(\mathbf{A}) = \lambda_1 \lambda_2 \cdots \lambda_n, \quad (36)$$

where the **trace** of $\mathbf{A} \in \mathbb{C}^{n \times n}$ is the sum of its diagonal elements

$$\text{trace}(\mathbf{A}) := a_{11} + a_{22} + \cdots + a_{nn}. \quad (37)$$

Proof. We compare two different expansion of $\pi_{\mathbf{A}}$. On the one hand from (35) we find

$$\pi_{\mathbf{A}}(\lambda) = (-1)^n \lambda^n + c_{n-1} \lambda^{n-1} + \cdots + c_0,$$

where $c_{n-1} = (-1)^{n-1} \text{trace}(\mathbf{A})$ and $c_0 = \pi_{\mathbf{A}}(0) = \det(\mathbf{A})$. On the other hand

$$\pi_{\mathbf{A}}(\lambda) = (\lambda_1 - \lambda) \cdots (\lambda_n - \lambda) = (-1)^n \lambda^n + d_{n-1} \lambda^{n-1} + \cdots + d_0,$$

where $d_{n-1} = (-1)^{n-1}(\lambda_1 + \cdots + \lambda_n)$ and $d_0 = \lambda_1 \cdots \lambda_n$. Since $c_j = d_j$ for all j we obtain (36). \square

For a 2×2 matrix the characteristic equation takes the convenient form

$$\lambda^2 - \text{trace}(\mathbf{A})\lambda + \det(\mathbf{A}) = 0. \quad (38)$$

Thus, if $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$ then $\text{trace}(\mathbf{A}) = 4$, $\det(\mathbf{A}) = 3$ so that $\pi_{\mathbf{A}}(\lambda) = \lambda^2 - 4\lambda + 3$.

Using Property 6. of determinants we have an additional characterization of a singular matrix.

Theorem 0.69 (Zero eigenvalue)

The matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is singular if and only if zero is an eigenvalue.

Proof. Zero is an eigenvalue if and only if $\pi_{\mathbf{A}}(0) = \det(\mathbf{A}) = 0$ which happens if and only if \mathbf{A} is singular. \square

0.8 Algorithms and Numerical Stability

In this text we consider mathematical problems (i. e., linear algebra problems) and many detailed numerical algorithms to solve them. Many of these algorithms have built inn Matlab equivalents that are more efficient. Some reasons for giving our own algorithms and asking students to give their own are

1. Increased understanding.
2. Easier to discuss strength and weaknesses like complexity, programming issues, and stability.

Complexity is discussed briefly in Section 1.3.1. As for programming issues we often vectorize the algorithms leading to shorter and more efficient programs. Stability is important both for the mathematical problems and for the numerical algorithms. Stability can be studied in terms of perturbation theory leading to condition numbers, see Chapters 8, 12, 13. We will often use phrases like “the algorithm is numerically stable” or “the algorithm is not numerically stable” without saying precisely what we mean by this. Loosely speaking, an algorithm is numerically stable if the solution, computed in floating point arithmetic, is the exact solution of a slightly perturbed problem. To determine upper bounds for these perturbations is the topic of backward error analysis. We give a rather limited introduction to floating point arithmetic and backward error analysis in Appendix B, but in the text we will not discuss this. This does not mean that numerical stability is not an important issue. In fact, numerical stability is crucial for a good algorithm. For thorough treatments of numerical stability issues we refer to the books [11] and [22, 23].

A list of freely available software for solving linear algebra problems can be found at

<http://www.netlib.org/utk/people/JackDongarra/la-sw.html>

Part I

Direct Methods for Linear Systems

Chapter 1

Gaussian Elimination

Gaussian elimination⁶ is the classical method for solving n linear equations in n unknowns. In component form the system is

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2, \\ \vdots &\quad \vdots \quad \vdots \quad \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b_n, \end{aligned}$$

and in matrix form

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \mathbf{b}.$$

The elements of \mathbf{A} and \mathbf{b} can be either real or complex numbers.

We recall (see Theorem 0.46) that the square system $\mathbf{A}\mathbf{x} = \mathbf{b}$ has a unique solution for all right hand sides \mathbf{b} if and only if \mathbf{A} is nonsingular, i. e., the homogeneous system $\mathbf{A}\mathbf{x} = \mathbf{0}$ only has the solution $\mathbf{x} = \mathbf{0}$. We recall (cf. Theorem 0.51) that a square matrix has an inverse if and only if \mathbf{A} is nonsingular, and the solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be written $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$, where \mathbf{A}^{-1} is the inverse of \mathbf{A} .

In Gaussian elimination with row interchanges we compute a factorization of the coefficient matrix \mathbf{A} known as a PLU factorization. In this chapter we discuss some theoretical and algorithmic aspects of Gaussian elimination.

⁶The method was known long before Gauss used it in 1809. It was further developed by Doolittle in 1881, see [5].

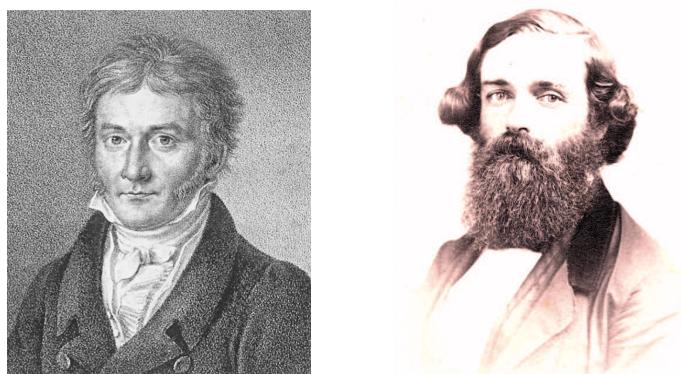


Figure 1.1: Carl Friedrich Gauss, 1777-1855, Lithograph by Siegfried Detlev Bendixen, 1828 (left), Myrick Hascall Doolittle, 1830-1911.

We start with an introduction to block multiplication and triangular matrices.

1.1 Block Multiplication

A rectangular matrix \mathbf{A} can be partitioned into submatrices by drawing horizontal lines between selected rows and vertical lines between selected columns. For example, the matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

can be partitioned as

$$(i) \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \quad (ii) \begin{bmatrix} \mathbf{a}_{\cdot 1}, \mathbf{a}_{\cdot 2}, \mathbf{a}_{\cdot 3} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix},$$

$$(iii) \begin{bmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \mathbf{a}_{3:}^T \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \quad (iv) \begin{bmatrix} \mathbf{A}_{11}, \mathbf{A}_{12} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}.$$

In (i) the matrix \mathbf{A} is divided into four submatrices

$$\mathbf{A}_{11} = [1], \quad \mathbf{A}_{12} = [2, 3], \quad \mathbf{A}_{21} = \begin{bmatrix} 4 \\ 7 \end{bmatrix}, \quad \text{and} \quad \mathbf{A}_{22} = \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix},$$

while in (ii) and (iii) \mathbf{A} has been partitioned into columns and rows, respectively. The submatrices in a partition are often referred to as **blocks** and a partitioned matrix is sometimes called a **block matrix**.

In the following we assume that $\mathbf{A} \in \mathbb{C}^{m \times p}$ and $\mathbf{B} \in \mathbb{C}^{p \times n}$. Here are some rules and observations for block multiplication.

1. If $\mathbf{B} = [\mathbf{b}_{:1}, \dots, \mathbf{b}_{:n}]$ is partitioned into columns then the partition of the product \mathbf{AB} into columns is

$$\mathbf{AB} = [\mathbf{Ab}_{:1}, \mathbf{Ab}_{:2}, \dots, \mathbf{Ab}_{:n}].$$

In particular, if \mathbf{I} is the identity matrix of order p then

$$\mathbf{A} = \mathbf{AI} = \mathbf{A} [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p] = [\mathbf{Ae}_1, \mathbf{Ae}_2, \dots, \mathbf{Ae}_p]$$

and we see that column j of \mathbf{A} can be written \mathbf{Ae}_j for $j = 1, \dots, p$.

2. Similarly, if \mathbf{A} is partitioned into rows then

$$\mathbf{AB} = \begin{bmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \vdots \\ \mathbf{a}_{m:}^T \end{bmatrix} \mathbf{B} = \begin{bmatrix} \mathbf{a}_{1:}^T \mathbf{B} \\ \mathbf{a}_{2:}^T \mathbf{B} \\ \vdots \\ \mathbf{a}_{m:}^T \mathbf{B} \end{bmatrix},$$

and taking $\mathbf{A} = \mathbf{I}$ it follows that row i of \mathbf{B} can be written $\mathbf{e}_i^T \mathbf{B}$ for $i = 1, \dots, m$.

3. It is often useful to write the matrix-vector product \mathbf{Ax} as a linear combination of the columns of \mathbf{A}

$$\mathbf{Ax} = x_1 \mathbf{a}_{:1} + x_2 \mathbf{a}_{:2} + \cdots + x_p \mathbf{a}_{:p}.$$

4. If $\mathbf{B} = [\mathbf{B}_1, \mathbf{B}_2]$, where $\mathbf{B}_1 \in \mathbb{C}^{p \times r}$ and $\mathbf{B}_2 \in \mathbb{C}^{p \times (n-r)}$ then

$$\mathbf{A} [\mathbf{B}_1, \mathbf{B}_2] = [\mathbf{AB}_1, \mathbf{AB}_2].$$

This follows from Rule 1. by an appropriate grouping of columns.

5. If $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix}$, where $\mathbf{A}_1 \in \mathbb{C}^{k \times p}$ and $\mathbf{A}_2 \in \mathbb{C}^{(m-k) \times p}$ then

$$\begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} \mathbf{B} = \begin{bmatrix} \mathbf{A}_1 \mathbf{B} \\ \mathbf{A}_2 \mathbf{B} \end{bmatrix}.$$

This follows from Rule 2. by a grouping of rows.

6. If $\mathbf{A} = [\mathbf{A}_1, \mathbf{A}_2]$ and $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}$, where $\mathbf{A}_1 \in \mathbb{C}^{m \times s}$, $\mathbf{A}_2 \in \mathbb{C}^{m \times (p-s)}$, $\mathbf{B}_1 \in \mathbb{C}^{s \times n}$ and $\mathbf{B}_2 \in \mathbb{C}^{(p-s) \times n}$ then

$$[\mathbf{A}_1, \mathbf{A}_2] \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix} = [\mathbf{A}_1 \mathbf{B}_1 + \mathbf{A}_2 \mathbf{B}_2].$$

Indeed, $(\mathbf{AB})_{ij} = \sum_{k=1}^p a_{ik} b_{kj} = \sum_{k=1}^s a_{ik} b_{kj} + \sum_{k=s+1}^p a_{ik} b_{kj} = (\mathbf{A}_1 \mathbf{B}_1)_{ij} + (\mathbf{A}_2 \mathbf{B}_2)_{ij} = (\mathbf{A}_1 \mathbf{B}_1 + \mathbf{A}_2 \mathbf{B}_2)_{ij}.$

7. If $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$ and $\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}$ then

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11} \mathbf{B}_{11} + \mathbf{A}_{12} \mathbf{B}_{21} & \mathbf{A}_{11} \mathbf{B}_{12} + \mathbf{A}_{12} \mathbf{B}_{22} \\ \mathbf{A}_{21} \mathbf{B}_{11} + \mathbf{A}_{22} \mathbf{B}_{21} & \mathbf{A}_{21} \mathbf{B}_{12} + \mathbf{A}_{22} \mathbf{B}_{22} \end{bmatrix},$$

provided the vertical partition in \mathbf{A} matches the horizontal one in \mathbf{B} , i.e. the number of columns in \mathbf{A}_{11} and \mathbf{A}_{21} equals the number of rows in \mathbf{B}_{11} and \mathbf{B}_{12} and the number of columns in \mathbf{A} equals the number of rows in \mathbf{B} . To show this we use Rule 4. to obtain

$$\mathbf{AB} = \left[\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} \\ \mathbf{B}_{21} \end{bmatrix}, \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{12} \\ \mathbf{B}_{22} \end{bmatrix} \right].$$

We complete the proof using Rules 5. and 6.

8. Consider finally the general case. If all the matrix products $\mathbf{A}_{ik} \mathbf{B}_{kj}$ in

$$\mathbf{C}_{ij} = \sum_{k=1}^s \mathbf{A}_{ik} \mathbf{B}_{kj}, \quad i = 1, \dots, p, \quad j = 1, \dots, q$$

are well defined then

$$\begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1s} \\ \vdots & & \vdots \\ \mathbf{A}_{p1} & \cdots & \mathbf{A}_{ps} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \cdots & \mathbf{B}_{1q} \\ \vdots & & \vdots \\ \mathbf{B}_{s1} & \cdots & \mathbf{B}_{sq} \end{bmatrix} = \begin{bmatrix} \mathbf{C}_{11} & \cdots & \mathbf{C}_{1q} \\ \vdots & & \vdots \\ \mathbf{C}_{p1} & \cdots & \mathbf{C}_{pq} \end{bmatrix}.$$

The requirements are that

- the number of columns in \mathbf{A} is equal to the number of rows in \mathbf{B} .
- the position of the vertical partition lines in \mathbf{A} has to match the position of the horizontal partition lines in \mathbf{B} . The horizontal lines in \mathbf{A} and the vertical lines in \mathbf{B} can be anywhere.

Exercise 1.1 (Matrix element as a quadratic form)

For any matrix \mathbf{A} show that $a_{ij} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_j$ for all i, j .

Exercise 1.2 (Outer product expansion of a matrix)

For any matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ show that $\mathbf{A} = \sum_{i=1}^m \sum_{j=1}^n a_{ij} \mathbf{e}_i \mathbf{e}_j^T$.

Exercise 1.3 (The product $\mathbf{A}^T \mathbf{A}$)

Let $\mathbf{B} = \mathbf{A}^T \mathbf{A}$. Explain why this product is defined for any matrix \mathbf{A} . Show that $b_{ij} = \mathbf{a}_{:,i}^T \mathbf{a}_{:,j}$ for all i, j .

Exercise 1.4 (Outer product expansion)

For $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{p \times n}$ show that

$$\mathbf{AB}^T = \mathbf{a}_{:,1} \mathbf{b}_{:,1}^T + \mathbf{a}_{:,2} \mathbf{b}_{:,2}^T + \cdots + \mathbf{a}_{:,n} \mathbf{b}_{:,n}^T.$$

This is called the **outer product expansion** of the columns of \mathbf{A} and \mathbf{B} .

Exercise 1.5 (System with many right hand sides; compact form)

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m \times p}$, and $\mathbf{X} \in \mathbb{R}^{n \times p}$. Show that

$$\mathbf{AX} = \mathbf{B} \iff \mathbf{Ax}_{:,j} = \mathbf{b}_{:,j}, \quad j = 1, \dots, p.$$

Exercise 1.6 (Block multiplication example)

Suppose $\mathbf{A} = [\mathbf{A}_1, \mathbf{A}_2]$ and $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{0} \end{bmatrix}$. When is $\mathbf{AB} = \mathbf{A}_1 \mathbf{B}_1$?

Exercise 1.7 (Another block multiplication example)

Suppose $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times n}$ are given in block form by

$$\mathbf{A} := \begin{bmatrix} \lambda & \mathbf{a}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix}, \quad \mathbf{B} := \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B}_1 \end{bmatrix}, \quad \mathbf{C} := \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{C}_1 \end{bmatrix},$$

where $\mathbf{A}_1, \mathbf{B}_1, \mathbf{C}_1 \in \mathbb{R}^{(n-1) \times (n-1)}$. Show that

$$\mathbf{CAB} = \begin{bmatrix} \lambda & \mathbf{a}^T \mathbf{B}_1 \\ \mathbf{0} & \mathbf{C}_1 \mathbf{A}_1 \mathbf{B}_1 \end{bmatrix}.$$

1.2 Triangular matrices

Recall that a matrix \mathbf{U} is upper triangular if $u_{ij} = 0$ for $i > j$, and a matrix \mathbf{L} is lower triangular if $l_{ij} = 0$ for $i < j$. If \mathbf{U} is upper triangular then \mathbf{U}^T is lower triangular.

We need some basic facts about triangular matrices and we start with

Lemma 1.8 (Inverse of a block triangular matrix)

Suppose

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ 0 & \mathbf{A}_{22} \end{bmatrix}$$

where \mathbf{A} , \mathbf{A}_{11} and \mathbf{A}_{22} are square matrices. Then \mathbf{A} is nonsingular if and only if both \mathbf{A}_{11} and \mathbf{A}_{22} are nonsingular. In that case

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ 0 & \mathbf{A}_{22}^{-1} \end{bmatrix}. \quad (1.1)$$

Proof. If \mathbf{A}_{11} and \mathbf{A}_{22} are nonsingular then

$$\begin{bmatrix} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ 0 & \mathbf{A}_{22}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ 0 & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \mathbf{I} \end{bmatrix} = \mathbf{I}$$

and \mathbf{A} is nonsingular with the indicated inverse. Conversely, let \mathbf{B} be the inverse of the nonsingular matrix \mathbf{A} . We partition \mathbf{B} conformally with \mathbf{A} and have

$$\mathbf{B}\mathbf{A} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ 0 & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \mathbf{I} \end{bmatrix} = \mathbf{I}$$

Using block-multiplication we find

$$\mathbf{B}_{11}\mathbf{A}_{11} = \mathbf{I}, \quad \mathbf{B}_{21}\mathbf{A}_{11} = \mathbf{0}, \quad \mathbf{B}_{21}\mathbf{A}_{12} + \mathbf{B}_{22}\mathbf{A}_{22} = \mathbf{I}.$$

The first equation implies that \mathbf{A}_{11} is nonsingular, this in turn implies that $\mathbf{B}_{21} = \mathbf{0}$ in the second equation, and then the third equation simplifies to $\mathbf{B}_{22}\mathbf{A}_{22} = \mathbf{I}$. We conclude that also \mathbf{A}_{22} is nonsingular. \square

Consider now a triangular matrix.

Lemma 1.9 (Inverse of a triangular matrix)

An upper (lower) triangular matrix $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{n \times n}$ is nonsingular if and only if the diagonal elements a_{ii} , $i = 1, \dots, n$ are nonzero. In that case the inverse is upper (lower) triangular with diagonal elements a_{ii}^{-1} , $i = 1, \dots, n$.

Proof. We use induction on n . The result holds for $n = 1$. The 1-by-1 matrix $\mathbf{A} = [a_{11}]$ is nonsingular if and only if $a_{11} \neq 0$ and in that case $\mathbf{A}^{-1} = [a_{11}^{-1}]$. Suppose the result holds for $n = k$ and let $\mathbf{A} \in \mathbb{C}^{(k+1) \times (k+1)}$ be upper triangular. We partition \mathbf{A} in the form

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_k & \mathbf{a}_k \\ 0 & a_{k+1,k+1} \end{bmatrix}$$

and note that $\mathbf{A}_k \in \mathbb{C}^{k \times k}$ is upper triangular. By Lemma 1.8 \mathbf{A} is nonsingular if and only if \mathbf{A}_k and $(a_{k+1,k+1})$ are nonsingular and in that case

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_k^{-1} & -\mathbf{A}_k^{-1}\mathbf{a}_k a_{k+1,k+1}^{-1} \\ 0 & a_{k+1,k+1}^{-1} \end{bmatrix}.$$

By the induction hypothesis \mathbf{A}_k is nonsingular if and only if the diagonal elements a_{11}, \dots, a_{kk} of \mathbf{A}_k are nonzero and in that case \mathbf{A}_k^{-1} is upper triangular with diagonal elements a_{ii}^{-1} , $i = 1, \dots, k$. The result for \mathbf{A} follows. \square

Lemma 1.10 (Product of triangular matrices)

The product $\mathbf{C} = \mathbf{AB} = (c_{ij})$ of two upper (lower) triangular matrices $\mathbf{A} = (a_{ij})$ and $\mathbf{B} = (b_{ij})$ is upper (lower) triangular with diagonal elements $c_{ii} = a_{ii}b_{ii}$ for all i .

Proof. Exercise. \square

A matrix is **unit triangular** if it is triangular with 1's on the diagonal.

Lemma 1.11 (Unit triangular matrices)

For a unit upper (lower) triangular matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$:

1. \mathbf{A} is nonsingular and the inverse is unit upper(lower) triangular.
2. The product of two unit upper (lower) triangular matrices is unit upper (lower) triangular.

Proof. 1. follows from Lemma 1.9, while Lemma 1.10 implies 2. \square

1.2.1 Algorithms for Triangular Systems

A nonsingular triangular linear system $\mathbf{Ax} = \mathbf{b}$ is easy to solve. By Lemma 1.9 \mathbf{A} has nonzero diagonal elements. Consider first the lower triangular case. For $n = 3$ the system is

$$\begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

From the first equation we find $x_1 = b_1/a_{11}$. Solving the second equation for x_2 we obtain $x_2 = (b_2 - a_{21}x_1)/a_{22}$. Finally the third equation gives $x_3 = (b_3 - a_{31}x_1 - a_{32}x_2)/a_{33}$. This process is known as forward substitution. In general

$$x_k = \left(b_k - \sum_{j=1}^{k-1} a_{k,j}x_j \right) / a_{kk}, \quad k = 1, 2, \dots, n. \quad (1.2)$$

$$\begin{bmatrix} a_{11} & 0 & 0 & 0 & 0 \\ a_{21} & a_{22} & 0 & 0 & 0 \\ 0 & a_{32} & a_{33} & 0 & 0 \\ 0 & 0 & a_{43} & a_{44} & 0 \\ 0 & 0 & 0 & a_{54} & a_{55} \end{bmatrix}, \quad \begin{bmatrix} a_{11} & 0 & 0 & 0 & 0 \\ a_{21} & a_{22} & 0 & 0 & 0 \\ a_{31} & a_{32} & a_{33} & 0 & 0 \\ 0 & a_{42} & a_{43} & a_{44} & 0 \\ 0 & 0 & a_{53} & a_{54} & a_{55} \end{bmatrix}$$

Figure 1.2: Lower triangular 5×5 band matrices: $d = 1$ (left) and $d = 2$ right.

When \mathbf{A} is a lower triangular band matrix the number of arithmetic operations necessary to find \mathbf{x} can be reduced. Suppose \mathbf{A} is a lower triangular d -banded, so that $a_{k,j} = 0$ for $j \notin \{l_k, l_k + 1, \dots, k\}$ for $k = 1, 2, \dots, n$, and where $l_k := \max(1, k-d)$, see Figure 1.2. For a lower triangular d -band matrix the calculation in (1.2) can be simplified as follows

$$x_k = \left(b_k - \sum_{j=l_k}^{k-1} a_{k,j} x_j \right) / a_{kk}, \quad k = 1, 2, \dots, n. \quad (1.3)$$

Note that (1.3) reduces to (1.2) if $d = n$. Letting $A(k, l_k : k-1) * x(l_k : k-1)$ denote the sum $\sum_{j=l_k}^{k-1} a_{kj} x_j$ we arrive at the following algorithm.

Algorithm 1.12 (forwardsolve)

Given a nonsingular lower triangular d -banded matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\mathbf{b} \in \mathbb{C}^n$. An $\mathbf{x} \in \mathbb{C}^n$ is computed so that $\mathbf{Ax} = \mathbf{b}$.

```

1 function x=rforwardsolve(A,b,d)
2 n=length(b); x=b;
3 x(1)=b(1)/A(1,1);
4 for k=2:n
5     lk=max(1,k-d);
6     x(k)=(b(k)-A(k,lk:k-1)*x(lk:k-1))/A(k,k);
7 end
```

A system $\mathbf{Ax} = \mathbf{b}$, where \mathbf{A} is upper triangular must be solved by back substitution or 'bottom-up'. We first find x_n from the last equation and then move upwards for the remaining unknowns. For an upper triangular d -banded matrix this leads to the following algorithm.

Algorithm 1.13 (backsolve)

Given a nonsingular upper triangular d -banded matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\mathbf{b} \in \mathbb{C}^n$. An $\mathbf{x} \in \mathbb{C}^n$ is computed so that $\mathbf{Ax} = \mathbf{b}$.

```

1 function x=rbacksolve(A,b,d)
2 n=length(b); x=b;
3 x(n)=b(n)/A(n,n);
4 for k=n-1:-1:1
5     uk=min(n,k+d);
6     x(k)=(b(k)-A(k,k+1:uk)*x(k+1:uk))/A(k,k);
7 end

```

Exercise 1.14 (Column oriented backsolve)

The intial "r" in the names of Algorithms 1.12,1.13 signals that these algorithms are row oriented. For each k we take the inner product of a part of a row with the already computed unknowns. In this exercise we develop column oriented vectorized versions of forward and backward substitution. Consider the system $\mathbf{Ax} = \mathbf{b}$, where $\mathbf{A} \in \mathbb{C}^{n \times n}$ is lower triangular. Suppose after $k - 1$ steps of the algorithm we have a reduced system in the form

$$\begin{bmatrix} a_{k,k} & 0 & \cdots & 0 \\ a_{k+1,k} & a_{k+1,k+1} & \cdots & 0 \\ \vdots & \ddots & \vdots & \\ a_{n,k} & \cdots & a_{n \times n} \end{bmatrix} \begin{bmatrix} x_k \\ x_{k+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_k \\ b_{k+1} \\ \vdots \\ b_n \end{bmatrix}.$$

This system is of order $n - k + 1$. The unknowns are x_k, \dots, x_n .

a) We see that $x_k = b_k/a_{k,k}$ and eliminating x_k from the remaining equations we obtain a system of order $n - k$ with unknowns x_{k+1}, \dots, x_n

$$\begin{bmatrix} a_{k+1,k+1} & 0 & \cdots & 0 \\ a_{k+2,k+1} & a_{k+2,k+2} & \cdots & 0 \\ \vdots & \ddots & \vdots & \\ a_{n,k+1} & \cdots & a_{n \times n} \end{bmatrix} \begin{bmatrix} x_{k+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_{k+1} \\ \vdots \\ b_n \end{bmatrix} - x_k \begin{bmatrix} a_{k+1,k} \\ \vdots \\ a_{n,k} \end{bmatrix}.$$

Thus at the k th step, $k = 1, 2, \dots, n$ we set $x_k = b_k/A(k,k)$ and update b as follows:

$$b(k+1:n) = b(k+1:n) - x(k) * A(k+1:n,k).$$

This leads to the following algorithm.

Algorithm 1.15 (Forward Solve (column oriented))

Given a nonsingular lower triangular d -banded matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\mathbf{b} \in \mathbb{C}^n$. An $\mathbf{x} \in \mathbb{C}^n$ is computed so that $\mathbf{Ax} = \mathbf{b}$.

```

1 function x=cforwardsolve(A,b,d)
2 x=b; n=length(b);
3 for k=1:n-1
4   x(k)=b(k)/A(k,k); uk=min(n,k+d);
5   b(k+1:uk)=b(k+1:uk)-A(k+1:uk,k)*x(k);
6 end
7 x(n)=b(n)/A(n,n);
8 end

```

b) Suppose now $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular, upper triangular, d -banded, and $\mathbf{b} \in \mathbb{C}^n$. Justify the following column oriented vectorized algorithms for solving $\mathbf{Ax} = \mathbf{b}$.

Algorithm 1.16 (Backsolve (column oriented))

Given a nonsingular upper triangular d -banded matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\mathbf{b} \in \mathbb{C}^n$. An $\mathbf{x} \in \mathbb{C}^n$ is computed so that $\mathbf{Ax} = \mathbf{b}$.

```

1 function x=cbacksolve(A,b,d)
2 x=b; n=length(b);
3 for k=n:-1:2
4   x(k)=b(k)/A(k,k); lk=max(1,k-d);
5   b(lk:k-1)=b(lk:k-1)-A(lk:k-1,k)*x(k);
6 end
7 x(1)=b(1)/A(1,1);
8 end

```

Exercise 1.17 (Computing the inverse of a triangular matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a nonsingular triangular matrix with inverse $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$. The k th column \mathbf{b}_k of \mathbf{B} is the solution of the linear systems $\mathbf{Ab}_k = \mathbf{e}_k$. Write this system as a 2×2 triangular block system and explain why we can find \mathbf{b}_k by solving the linear systems

$$\mathbf{A}(k:n, k:n)\mathbf{b}_k(k:n) = \mathbf{I}(k:n, k), \quad k = 1, \dots, n \quad \text{lower triangular}, \quad (1.4)$$

$$\mathbf{A}(1:k, 1:k)\mathbf{b}_k(1:k) = \mathbf{I}(1:k, k), \quad k = n, n-1, \dots, 1, \quad \text{upper triangular} \quad (1.5)$$

Is it possible to store the interesting part of \mathbf{b}_k in \mathbf{A} as soon as it is computed?

1.3 Naive Gaussian Elimination and LU factorization

In this section we describe Gaussian elimination without row interchanges, known also as **naive Gaussian elimination**. The process is illustrated in Figure 1.3. We

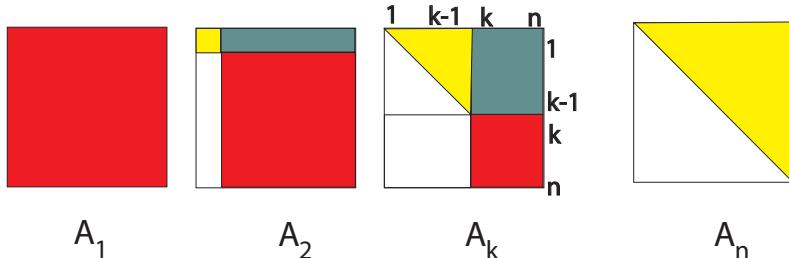


Figure 1.3: Gaussian elimination

start with a linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ and generate a sequence of equivalent systems $\mathbf{A}_k\mathbf{x} = \mathbf{b}_k$ for $k = 1, \dots, n$, where $\mathbf{A}_1 = \mathbf{A}$, $\mathbf{b}_1 = \mathbf{b}$, and \mathbf{A}_k has zeros under the diagonal in its first $k - 1$ columns. Thus \mathbf{A}_n is upper triangular and the system $\mathbf{A}_n\mathbf{x} = \mathbf{b}_n$ can be solved using one of Algorithms 1.13 or 1.16.

The matrix \mathbf{A}_k has the form

$$\mathbf{A}_k = \left[\begin{array}{cccc|cccccc} a_{1,1}^1 & \cdots & a_{1,k-1}^1 & a_{1,k}^1 & \cdots & a_{1,j}^1 & \cdots & a_{1,n}^1 \\ \ddots & & \vdots & \vdots & & \vdots & & \vdots \\ & a_{k-1,k-1}^{k-1} & a_{k-1,k}^{k-1} & \cdots & a_{k-1,j}^{k-1} & \cdots & a_{k-1,n}^{k-1} \\ \hline & a_{k,k}^k & \cdots & a_{k,j}^k & \cdots & a_{k,n}^k \\ & \vdots & & \vdots & & \vdots \\ & a_{i,k}^k & \cdots & a_{i,j}^k & \cdots & a_{i,n}^k \\ & \vdots & & \vdots & & \vdots \\ & a_{n,k}^k & \cdots & a_{n,j}^k & \cdots & a_{n,n}^k \end{array} \right] \quad (1.6)$$

$$= \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \mathbf{D}_k \end{bmatrix}.$$

If $a_{k,k}^k \neq 0$ the process transforming \mathbf{A}_k into \mathbf{A}_{k+1} for $k = 1, \dots, n - 1$ can be described as follows.

$$\begin{aligned} &\text{for } i = k + 1 : n \\ &\quad l_{ik}^k = a_{ik}^k / a_{kk}^k \\ &\quad \text{for } j = k : n \\ &\quad \quad a_{ij}^{k+1} = a_{ij}^k - l_{ik}^k a_{kj}^k \end{aligned} \quad (1.7)$$

For $j = k$ it follows from (1.7) that $a_{ik}^{k+1} = a_{ik}^k - \frac{a_{ik}^k}{a_{kk}^k} a_{kk}^k = 0$ for $i = k+1, \dots, n$. Thus \mathbf{A}_{k+1} will have zeros under the diagonal in its first k columns and the elimination is carried one step further. The numbers l_{ik}^k in (1.7) are called **multipliers**.

Since we get division by zero in (1.7) if $a_{kk}^k = 0$ for some $k \leq n-1$ it is important to know when this can happen.

Theorem 1.18 (When is naive Gaussian elimination possible?)

We have $a_{kk}^k \neq 0$ for $k = 1, \dots, n-1$ if and only if the leading principal submatrices

$$\mathbf{A}_{[k]} := \begin{bmatrix} a_{11} & \dots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \dots & a_{kk} \end{bmatrix}$$

of \mathbf{A} are nonsingular for $k = 1, \dots, n-1$.

Proof. Let \mathbf{B}_k be the upper left $k-1$ corner of \mathbf{A}_k given by (1.6). Observe that the elements of the matrix \mathbf{B}_k is computed from \mathbf{A} by using only elements in $\mathbf{A}_{[k-1]}$. Since the determinant of a matrix does not change under the operation of subtracting a multiple of one row from another row the determinant of $\mathbf{A}_{[k]}$ equals the determinant of \mathbf{B}_{k+1} , that is given by the product of its diagonal elements

$$\det(\mathbf{A}_{[k]}) = a_{11}^1 a_{22}^2 \cdots a_{kk}^k, \quad k = 1, \dots, n. \quad (1.8)$$

But then $a_{11}^1 \cdots a_{kk}^k \neq 0$ for $k = 1, \dots, n-1$ if and only if $\det(\mathbf{A}_{[k]}) \neq 0$ for $k = 1, \dots, n-1$, or equivalently $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n-1$. \square

Naive Gaussian elimination is a way to compute a factorization of the coefficient matrix known as an LU factorization.

Theorem 1.19 (Gauss=LU)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ and that $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n-1$. Then naive Gaussian elimination results in an LU factorization of $\mathbf{A} \in \mathbb{C}^{n \times n}$, i.e., $\mathbf{A} = \mathbf{L}\mathbf{U}$, where

$$\mathbf{L} = \begin{bmatrix} 1 & & & \\ l_{21}^1 & 1 & & \\ \vdots & & \ddots & \\ l_{n1}^1 & l_{n2}^2 & \cdots & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} a_{11}^1 & \cdots & a_{1n}^1 \\ & \ddots & \vdots \\ & & a_{nn}^n \end{bmatrix}, \quad (1.9)$$

and where the l_{ij}^j and a_{ij}^i are given by (1.7).

Proof. We use (1.7). For $i \leq j$

$$(\mathbf{LU})_{ij} = \sum_{k=1}^n l_{ik}^k u_{kj} = \sum_{k=1}^{i-1} l_{ik}^k a_{kj}^k + a_{ij}^i = \sum_{k=1}^{i-1} (a_{ij}^k - a_{ij}^{k+1}) + a_{ij}^i = a_{ij}^1 = a_{ij},$$

while for $i > j$

$$(\mathbf{LU})_{ij} = \sum_{k=1}^n l_{ik}^k u_{kj} = \sum_{k=1}^{j-1} l_{ik}^k a_{kj}^k + l_{ij}^j a_{jj}^j = \sum_{k=1}^{j-1} (a_{ij}^k - a_{ij}^{k+1}) + a_{ij}^j = a_{ij}.$$

□

We note that

1. Theorem 1.18 gives a sufficient condition for the existence of an LU factorization of a matrix. In Chapter 3 we show that this condition is necessary in order to have a **unique** LU factorization.
2. Theorem 1.18 holds even if \mathbf{A} is singular. Since \mathbf{L} is nonsingular the matrix \mathbf{U} is then singular, and since $a_{kk}^k \neq 0$ for $k = 1, \dots, n-1$ we must have $a_{nn}^n = 0$ when \mathbf{A} is singular.
3. To verify the nonsingularity of the leading principal submatrices can be difficult in practice. We show in Chapter 2, 3 that this condition holds for a class of diagonally dominant matrices and for positive definite matrices.

1.3.1 Operation count

It is useful to have a number which indicates the amount of work an algorithm requires. In this book we measure this by estimating the total number of arithmetic operations. We count both additions, subtractions, multiplications and divisions, but not work on indices. As an example it is shown below that the calculation of the LU factorization in Theorem 1.19 requires exactly

$$N_{LU} := \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n$$

arithmetic operations for a full matrix. We are only interested in this number when n is large and for such n the term $\frac{2}{3}n^3$ dominates. We therefore regularly ignore lower order terms and use **number of arithmetic operations** both for the exact count and for the highest order term. We also say more loosely that the number of arithmetic operations is $O(n^3)$. We will use the number of operations counted in one of these ways as a measure of the complexity of an algorithm.

We will compare the number of arithmetic operations of many algorithms with the number of arithmetic operations of naive Gaussian elimination and define for $n \in \mathbb{N}$ the number G_n as follows.

Definition 1.20 ($G_n := \frac{2}{3}n^3$)

Consider now finding the total number of arithmetic operations, N_{LU} for LU factorization using Gaussian elimination. Let M, D, A, S be the number of multiplications, divisions, additions, and subtractions. We do an exact count using (1.7) where we only let j run from $k+1$ to n . Then

- $M = \sum_{k=1}^{n-1} (n-k)^2 = \sum_{m=1}^{n-1} m^2 = \frac{1}{3}n(n-1)(n-\frac{1}{2})$
- $D = \sum_{m=1}^{n-1} m = \frac{1}{2}n(n-1), \quad S = M, \quad A = 0.$

Adding these numbers we obtain the number N_{LU} given above.

There is a quick way to arrive at the estimate $2n^3/3$. We only consider the arithmetic operations contributing to the leading term (the inner loops). Then we replace sums by integrals letting the summation indices be continuous variables and adjust limits of integration in an insightful way to simplify the calculation. In the Gaussian elimination case the contribution to the leading term only comes from M and S and we find

$$M + S = 2 \sum_{k=1}^{n-1} (n-k)^2 \approx 2 \int_1^{n-1} (n-k)^2 dk \approx 2 \int_0^n (n-k)^2 dk = \frac{2}{3}n^3.$$

This is the correct leading term and we obtain the number G_n .

Consider next forward and backward substitutions on a full matrix. The algorithms 1.12 and 1.13 with $d = n$ each require exactly $N_S := n^2$ arithmetic operations. Indeed they need $\sum_{k=1}^n (k-1) = n(n-1)/2$ multiplications, the same number of subtractions and n divisions, a total of n^2 arithmetic operations. Comparing G_n and N_S we see that LU factorization is an $O(n^3)$ process while the solution stage only require $O(n^2)$ arithmetic operations.

In many implementations the computing time T_A for an algorithm A applied to a large problem is proportional to N_A the number of arithmetic operations.⁷ If this is true then we typically have $T_A = \alpha N_A$, where α is in the range 10^{-12} to 10^{-9} on a modern computer. To see what this means for Gaussian elimination let us assume that the computing time for the LU factorization is $T_{LU} = 10^{-9}n^3$ and the computing time for the forward and backward substitution is $T_S = 3 \times 10^{-9}n^2$ corresponding to $\alpha = 3 \times 10^{-9}/2$.

This leads to dramatic differences in computing time as illustrated in the following table:

⁷It should be noted that in scientific computing of today the number of arithmetic operations is not necessarily the best estimate for how fast an algorithm can finish.

n	T_{LU}	T_S
10^3	1s	0.003s
10^4	17min.	0.3s
10^6	32 years	51min

To further illustrate the difference between n^3 and n^2 for large n suppose we want to solve m systems $\mathbf{A}_j \mathbf{x}_j = \mathbf{b}_j$ for $j = 1, \dots, m$, where $\mathbf{A}_j \in \mathbb{R}^{n \times n}$ and $\mathbf{b}_j \in \mathbb{R}^n$. We need $m(\frac{2}{3}n^3 + 2n^2)$ arithmetic operations for this. Thus if $n = 10^4$ and $m = 100$ the table gives a computing time of approximately 1700min. Suppose now $\mathbf{A}_j = \mathbf{A}$, i.e. we have the same coefficient matrix in all systems. We can then write the m systems more compactly as $\mathbf{AX} = \mathbf{B}$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{B} \in \mathbb{R}^{n,m}$ and the matrix $\mathbf{X} \in \mathbb{R}^{n,m}$ is the unknown. To solve $\mathbf{AX} = \mathbf{B}$ we first compute the LU factorization of \mathbf{A} and then apply forward and backward substitution to the columns of \mathbf{B} . If $n = 10^4$ the computing time for this would be 17min for the LU factorization and 30s for the solution phase.

Exercise 1.21 (Gaussian elimination example)

Show that the singular matrix $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ satisfies the condition in Theorem 1.19 and use naive Gaussian elimination to find its LU factorization.

Exercise 1.22 (Finite sums of integers)

Use induction on m , or some other method, to show that

$$1 + 2 + \dots + m = \frac{1}{2}m(m+1), \quad (1.10)$$

$$1^2 + 2^2 + \dots + m^2 = \frac{1}{3}m(m+\frac{1}{2})(m+1), \quad (1.11)$$

$$1 + 3 + 5 + \dots + 2m - 1 = m^2, \quad (1.12)$$

$$1 * 2 + 2 * 3 + 3 * 4 + \dots + (m-1)m = \frac{1}{3}(m-1)m(m+1). \quad (1.13)$$

Exercise 1.23 (Operations)

To solve an upper triangular linear system by back substitution takes n^2 arithmetic operations. Show that the number of arithmetic operations in (1.5) is $\frac{1}{3}n(n+\frac{1}{2})(n+1) \approx \frac{1}{2}G_n$.

Exercise 1.24 (Multiplying triangular matrices)

Show that the matrix multiplication \mathbf{AB} can be done in $\frac{1}{3}n(2n^2 + 1) \approx G_n$ arithmetic operations when $\mathbf{A} \in \mathbb{R}^{n \times n}$ is lower triangular and $\mathbf{B} \in \mathbb{R}^{n \times n}$ is upper triangular. What about \mathbf{BA} ?

Exercise 1.25 (Matrix formulation of Gaussian elimination)

Alternatively, we can describe the transformation $\mathbf{A}_k \rightarrow \mathbf{A}_{k+1}$ as a multiplication of \mathbf{A}_k by a matrix known as an **elementary lower triangular matrix**.

Definition 1.26 (Elementary lower triangular matrix)

For $1 \leq k \leq n-1$ and $\mathbf{l}_k = [l_{k+1,k}, \dots, l_{n,k}]^T \in \mathbb{R}^{n-k}$ we define the matrix $\mathbf{M}_k \in \mathbb{R}^{n \times n}$ by

$$\mathbf{M}_k := \mathbf{I} - \begin{bmatrix} \mathbf{0} \\ \mathbf{l}_k \end{bmatrix} \mathbf{e}_k^T = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & \ddots & \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & -l_{k+1,k} & 1 & \cdots & 0 \\ \vdots & & & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -l_{n,k} & 0 & \cdots & 1 \end{bmatrix}, \quad (1.14)$$

where $\mathbf{0}$ is the zero vector in \mathbb{R}^k . We call \mathbf{M}_k an **elementary lower triangular matrix**.

We have

$$\mathbf{A}_{k+1} = \mathbf{M}_k \mathbf{A}_k, \text{ for } k = 1, \dots, n-1, \quad (1.15)$$

where in (1.14) $l_{ik} = l_{ik}^k$ is given by (1.7) for $i = k+1, \dots, n$.

- (a) Show (1.15).
- (b) Use (1.15) to show that $\mathbf{A}_n = \mathbf{M}\mathbf{A}$, where \mathbf{M} is unit lower triangular.
- (c) What is \mathbf{M}^{-1} ?

1.4 Gaussian Elimination with Row Interchanges

Theorem 1.19 shows that under certain conditions naive Gaussian elimination is equivalent to an LU factorization of the coefficient matrix. These conditions show that naive Gaussian elimination cannot be used on many nonsingular linear systems. A simple example is $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. We show here that any nonsingular linear system can be solved by Gaussian elimination if we incorporate row interchanges.

1.4.1 Pivoting

Interchanging two rows (and/or two columns) during Gaussian elimination is known as **pivoting**. The element which is moved to the diagonal position (k, k)

is called the **pivot element** or **pivot** for short, and the row containing the pivot is called the **pivot row**. Gaussian elimination with row pivoting can be described as follows.

1. Choose $r_k \geq k$ so that $a_{r_k,k}^k \neq 0$.
2. Interchange rows r_k and k of \mathbf{A}_k .
3. Eliminate by computing l_{ik}^k and a_{ij}^{k+1} using (1.7).

Row interchanges can be described in terms of permutation matrices.

1.4.2 Permutation matrices

Definition 1.27 Let the components of $\mathbf{p} = [k_1, \dots, k_n]^T$ be a permutation of the components of $[1, 2, \dots, n]^T$. We call $\mathbf{P} := \mathbf{I}(:, \mathbf{p}) = [\mathbf{e}_{k_1}, \mathbf{e}_{k_2}, \dots, \mathbf{e}_{k_n}] \in \mathbb{R}^{n \times n}$ a **permutation matrix**. When discussing Gaussian elimination a permutation \mathbf{p} is sometimes called a **pivot vector**.

Since $\mathbf{P}^T = \mathbf{I}(\mathbf{p}, :)$ it follows that $(\mathbf{P}^T \mathbf{P})_{i,j} = \mathbf{e}_{k_i}^T \mathbf{e}_{k_j} = \delta_{ij}$. Thus $\mathbf{P}^T \mathbf{P} = \mathbf{I}$, the inverse of \mathbf{P} is equal to its transpose, and $\mathbf{P} \mathbf{P}^T = \mathbf{I}$ as well. If \mathbf{p} and \mathbf{P} are as in Definition 1.27 and $\mathbf{A} \in \mathbb{C}^{n \times n}$ then

$$\mathbf{AP} = \mathbf{A}(:, \mathbf{p}), \quad \mathbf{P}^T \mathbf{A} = \mathbf{A}(\mathbf{p}, :). \quad (1.16)$$

Thus, post-multiplying a matrix \mathbf{A} by a permutation matrix results in a permutation of the columns, while pre-multiplying by the transpose of a permutation matrix gives a permutation of the rows.

We will use a particularly simple permutation matrix.

Definition 1.28 (Interchange matrix)

We define a **(j,k)-Interchange Matrix** \mathbf{I}_{jk} by interchanging column j and k of the identity matrix.

Since $\mathbf{I}_{jk} = \mathbf{I}_{kj}$, and we obtain the identity by applying \mathbf{I}_{jk} twice, we see that $\mathbf{I}_{jk}^2 = \mathbf{I}$ and an interchange matrix is symmetric and equal to its own inverse. Pre-multiplying a matrix by an interchange matrix interchanges two rows of the matrix, while post-multiplication interchanges two columns.

1.4.3 The PLU-Factorization

Naive Gaussian elimination can be described as a factorization of the coefficient matrix. With no row interchanges we obtain an LU factorization. Consider now

Gaussian elimination with row pivoting. We can keep track of the row interchanges using **pivot vectors** \mathbf{p}^k . We define

$$\mathbf{p} := \mathbf{p}^{n-1}, \text{ where } \mathbf{p}^0 := [1, 2, \dots, n]^T, \text{ and } \mathbf{p}^k := \mathbf{I}_{r_k,k} \mathbf{p}^{k-1} \text{ for } k = 1, \dots, n-1. \quad (1.17)$$

We obtain \mathbf{p}^k from \mathbf{p}^{k-1} by interchanging the elements r_k and k in \mathbf{p}^{k-1} . In particular the first $k-1$ components in \mathbf{p}^{k-1} and \mathbf{p}^k are the same.

In an algorithm, instead of interchanging the rows of \mathbf{A} during elimination, we can keep track of the ordering of the rows using the pivot vectors \mathbf{p}^k . The incorporation of row interchanges in Gaussian elimination (1.7) can be described as follows:

$$\begin{aligned} \mathbf{p}^0 &= [p_1^0, \dots, p_n^0]^T = [1, \dots, n]^T; \\ \text{for } k &= 1 : n-1 \\ \text{choose } r_k &\geq k \text{ so that } a_{p_{r_k}^{k-1}, k}^k \neq 0. \\ \mathbf{p}^k &= \mathbf{I}_{r_k, k} \mathbf{p}^{k-1} = [p_1^k, \dots, p_n^k]^T \\ \text{for } i &= k+1 : n \\ a_{p_i^k, k}^k &= a_{p_i^k, k}^k / a_{p_k^k, k}^k \\ \text{for } j &= k : n \\ a_{p_i^k, j}^{k+1} &= a_{p_i^k, j}^k - a_{p_i^k, k}^k a_{p_k^k, j}^k \end{aligned} \quad (1.18)$$

This leads to the following factorization:

Theorem 1.29 (PLU theorem)

Gaussian elimination with row pivoting on a nonsingular matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ leads to a factorization $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{U}$, where \mathbf{P} is a permutation matrix, \mathbf{L} is lower triangular with ones on the diagonal, and \mathbf{U} is upper triangular. More explicitly, $\mathbf{P} = \mathbf{I}(:, \mathbf{p})$, where $\mathbf{p} = \mathbf{I}_{r_{n-1}, n-1} \cdots \mathbf{I}_{r_1, 1} [1, \dots, n]^T$, and

$$\mathbf{L} = \begin{bmatrix} 1 & & & \\ a_{p_2, 1}^1 & 1 & & \\ \vdots & & \ddots & \\ a_{p_n, 1}^1 & a_{p_n, 2}^2 & \cdots & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} a_{p_1, 1}^1 & \cdots & a_{p_1, n}^1 \\ & \ddots & \vdots \\ & & a_{p_n \times n}^n \end{bmatrix}, \quad (1.19)$$

Proof. Consider (1.18). Since a row is not interchanged after it becomes a pivot row we obtain $p_k := p_k^{n-1} = \cdots = p_k^{k+1} = p_k^k$ which implies that $a_{p_k^k, k}^k = a_{p_k^k, k}^k$. Since later row exchanges only affects rows below the pivot row we can compute the quantities in the i and j loops in (1.18) in any order. Thus we can use the

final pivot vector and replace these loops by

$$\begin{aligned}
 & \text{for } i = k + 1 : n \\
 & \quad a_{p_i, k}^k = a_{p_i, k}^k / a_{p_k, k}^k \\
 & \quad \text{for } j = k : n \\
 & \quad \quad a_{p_i, j}^{k+1} = a_{p_i, j}^k - a_{p_i, k}^k a_{p_k, j}^k.
 \end{aligned} \tag{1.20}$$

But this means that if we knew the row pivots in advance then we can carry out naive Gaussian elimination on the matrix $\mathbf{P}^T \mathbf{A}$, where $\mathbf{P}^T = \mathbf{I}_{r_{n-1}, n-1} \cdots \mathbf{I}_{r_1, 1}$. But then the result follows from Theorem 1.19. \square

The factorization $\mathbf{A} = \mathbf{PLU}$, where \mathbf{P} is a permutation matrix, \mathbf{L} is lower triangular, and \mathbf{U} is upper triangular, is called a **PLU factorization**. It can also be written $\mathbf{P}^T \mathbf{A} = \mathbf{LU}$. Thus, if \mathbf{A} is nonsingular then there exists a permutation of the rows of \mathbf{A} so that the matrix with the rows permuted has an LU factorization.

1.5 An Algorithm for Finding the PLU-Factorization

Using pivot vectors we can compute the PLU factorization of \mathbf{A} without physically interchanging the elements a_{ij}^k . We only consider a full matrix. PLU factorization of a band matrix is more complicated to describe. Indeed, a row interchange will increase the bandwidth.

As is clear from (1.19) we can store the elements of \mathbf{L} and \mathbf{U} in \mathbf{A} and work with $\mathbf{A}(\mathbf{p}^k, :)$. We can rewrite (1.18) using outer product notation and dropping the superscript on the $a_{i,j}$'s. We have

$$\begin{bmatrix} a_{p_{k+1}, k+1}^k & \cdots & a_{p_{k+1}, n}^k \\ \vdots & & \vdots \\ a_{p_n^k, k+1}^k & \cdots & a_{p_n^k, n}^k \end{bmatrix} = \begin{bmatrix} a_{p_{k+1}, k+1}^k & \cdots & a_{p_{k+1}, n}^k \\ \vdots & & \vdots \\ a_{p_n^k, k+1}^k & \cdots & a_{p_n^k, n}^k \end{bmatrix} - \begin{bmatrix} a_{p_{k+1}, k}^k \\ \vdots \\ a_{p_n^k, k}^k \end{bmatrix} [a_{p_k^k, k+1}^k \cdots a_{p_k^k, n}^k].$$

The result is a matrix of order $n - k$. At the end the elements of \mathbf{L} and \mathbf{U} will be located under and above the diagonal.

Once we have a PLU factorization of \mathbf{A} the system $\mathbf{Ax} = \mathbf{b}$ is solved easily in three steps. Since $\mathbf{PLUx} = \mathbf{b}$ we have $\mathbf{Pz} = \mathbf{b}$, $\mathbf{Ly} = \mathbf{z}$, and $\mathbf{Ux} = \mathbf{y}$. Using the output $[\mathbf{p}, \mathbf{L}, \mathbf{U}]$ of Algorithm 1.30 the solution can be found from Algorithms 1.12 and 1.13 in two steps.

1. $\mathbf{y} = \text{rforwardsolve}(\mathbf{L}, \mathbf{b}(\mathbf{p}), n);$
2. $\mathbf{x} = \text{rbacksolve}(\mathbf{U}, \mathbf{y}, n);$

In the following algorithm we use a version of row pivoting called **partial pivoting**. We choose r_k so that

$$a_{r_k,k}^k := \max\{|a_{i,k}^k| : k \leq i \leq n\}$$

with r_k the smallest such index in case of a tie.

Algorithm 1.30 (PLU factorization)

Given a nonsingular $\mathbf{A} \in \mathbb{C}^{n \times n}$. This algorithm computes a PLU factorization of \mathbf{A} using Gaussian elimination with partial pivoting. The permutation matrix \mathbf{P} can be recovered from the pivot vector \mathbf{p} as $\mathbf{P} = \mathbf{I}(:, \mathbf{p})$.

```

1 function [p,L,R] = plufactor(A)
2 n = length(A);
3 p = 1:n;
4 for k=1:n-1
5     [maxv, r] = max(abs(A(p(k:n),k)));
6     p([k r+k-1]) = p([r+k-1 k]);
7     ps=p(k+1:n);
8     A(ps,k) = A(ps,k)/A(p(k),k);
9     A(ps,k+1:n) = A(ps,k+1:n)-A(ps,k)*A(p(k),k+1:n);
10 end
11 L = eye(n,n) + tril(A(p,:),-1);
12 R = triu(A(p,:));

```

Exercise 1.31 (Using PLU of \mathbf{A} to solve $\mathbf{A}^T \mathbf{x} = \mathbf{b}$)

Suppose we know the PLU factors $\mathbf{P}, \mathbf{L}, \mathbf{U}$ in a PLU factorization $\mathbf{A} = \mathbf{PLU}$ of $\mathbf{A} \in \mathbb{C}^{n \times n}$. Explain how we can solve the system $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ economically.

Exercise 1.32 (Using PLU to compute the determinant)

Suppose we know the PLU factors $\mathbf{P}, \mathbf{L}, \mathbf{U}$ in a PLU factorization $\mathbf{A} = \mathbf{PLU}$ of $\mathbf{A} \in \mathbb{C}^{n \times n}$. Explain how we can use this to compute the determinant of \mathbf{A} .

Exercise 1.33 (Using PLU to compute the inverse)

Suppose the factors $\mathbf{P}, \mathbf{L}, \mathbf{U}$ in a PLU factorization of $\mathbf{A} \in \mathbb{R}^{n \times n}$ are known. Use Exercises 1.23, 1.24 to show that it takes approximately $2G_n$ arithmetic operations to compute $\mathbf{A}^{-1} = \mathbf{U}^{-1} \mathbf{L}^{-1} \mathbf{P}^T$. Here we have not counted the final multiplication with \mathbf{P}^T which amounts to n row interchanges.

1.5.1 Pivot strategies

In Algorithm 1.30 we used a version of row pivoting called partial pivoting

$$|a_{r_k,k}^k| := \max\{|a_{i,k}^k| : k \leq i \leq n\}$$

with r_k the smallest such index in case of a tie. The following example illustrating that small pivots should be avoided.

Example 1.34 (Row pivoting)

Applying Gaussian elimination without row interchanges to the linear system

$$\begin{aligned} 10^{-4}x_1 + 2x_2 &= 4 \\ x_1 + x_2 &= 3 \end{aligned}$$

we obtain the upper triangular system

$$\begin{aligned} 10^{-4}x_1 + 2x_2 &= 4 \\ (1 - 2 \times 10^4)x_2 &= 3 - 4 \times 10^4 \end{aligned}$$

The exact solution is

$$x_2 = \frac{-39997}{-19999} \approx 2, \quad x_1 = \frac{4 - 2x_2}{10^{-4}} = \frac{20000}{19999} \approx 1.$$

Suppose we round the result of each arithmetic operation to three digits. The solutions $\text{fl}(x_1)$ and $\text{fl}(x_2)$ computed in this way is

$$\text{fl}(x_2) = 2, \quad \text{fl}(x_1) = 0.$$

The computed value 0 of x_1 is completely wrong. Suppose instead we apply Gaussian elimination to the same system, but where we have interchanged the equations. The system is

$$\begin{aligned} x_1 + x_2 &= 3 \\ 10^{-4}x_1 + 2x_2 &= 4 \end{aligned}$$

and we obtain the upper triangular system

$$\begin{aligned} x_1 + x_2 &= 3 \\ (2 - 10^{-4})x_2 &= 4 - 3 \times 10^{-4} \end{aligned}$$

Now the solution is computed as follows

$$x_2 = \frac{3.9997}{1.9999} \approx 2, \quad x_1 = 3 - x_2 \approx 1.$$

In this case rounding each calculation to three digits produces $\text{fl}(x_1) = 1$ and $\text{fl}(x_2) = 2$ which is quite satisfactory since it is the exact solution rounded to three digits.

Related to partial pivoting is **scaled partial pivoting**. Here r_k is the smallest index such that

$$\frac{|a_{r_k,k}^k|}{s_k} := \max\left\{\frac{|a_{i,k}^k|}{s_k} : k \leq i \leq n\right\}, \quad s_k := \max_{1 \leq j \leq n} |a_{kj}|.$$

This can sometimes give more accurate results if the coefficient matrix have coefficients of wildly different sizes. Note that the scaling factors s_k are computed using the initial matrix.

It also is possible to interchange both rows and columns. The choice

$$a_{r_k, s_k}^k := \max\{|a_{i,j}^k| : k \leq i, j \leq n\}$$

with r_k, s_k the smallest such indices in case of a tie, is known as **complete pivoting**. Complete pivoting is known to be more numerically stable than partial pivoting, but requires a lot of search and is seldom used in practice.

1.6 Review Questions

1.6.1 When is a triangular matrix nonsingular?

1.6.2 Approximately how many arithmetic operations are needed for

- the multiplication of two square matrices?
- Gaussian elimination on a matrix?
- the solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$, when \mathbf{A} is triangular?

1.6.3 What is the general condition for Gaussian elimination without row interchanges to be well defined?

1.6.4 What is a PLU factorization? When does it exist?

1.6.5 What is complete pivoting?

Chapter 2

Examples of Linear Systems

Gaussian elimination with row interchanges can in principle be used to solve any nonsingular linear system. Solving a dense system of order n requires $O(n^3)$ arithmetic operations and we saw in Section 1.3.1 that solving large linear systems can require more time than we are willing to spend. Row interchanges are another issue in Gaussian elimination. For example, if we interchange two rows in a tridiagonal matrix then the tridiagonal structure is lost in general. In this chapter we present two problems leading to tridiagonal linear systems. We show that row interchanges are not necessary for the two problems we consider. and derive an algorithm that only requires $O(n)$ arithmetic operations.

2.1 The Second Derivative Matrix

Consider the simple **two point boundary value problem**

$$-u''(x) = f(x), \quad x \in [0, 1], \quad u(0) = 0, \quad u(1) = 0, \quad (2.1)$$

where f is a given continuous function on $[0, 1]$. This problem is also known as the **one-dimensional (1D) Poisson problem**. In principle it is easy to solve (2.1) exactly. We just integrate f twice and determine the two integration constants so that the homogeneous boundary conditions $u(0) = u(1) = 0$ are satisfied. For example, if $f(x) = 1$ then $u(x) = x(x - 1)/2$ is the solution. However, many functions f cannot be integrated exactly, and in such cases a numerical method can be used.

Problem (2.1) can be solved approximately using the **finite difference method**. We first derive a difference approximation to the second derivative.

If g is a function differentiable at x then

$$g'(x) = \lim_{h \rightarrow 0} \frac{g(x + \frac{h}{2}) - g(x - \frac{h}{2})}{h}$$

and applying this to a function u that is twice differentiable at x

$$\begin{aligned} u''(x) &= \lim_{h \rightarrow 0} \frac{u'(x + \frac{h}{2}) - u'(x - \frac{h}{2})}{h} = \lim_{h \rightarrow 0} \frac{\frac{u(x+h)-u(x)}{h} - \frac{u(x)-u(x-h)}{h}}{h} \\ &= \lim_{h \rightarrow 0} \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}. \end{aligned}$$

To define the points where this difference approximation is used we choose a positive integer m , let $h := 1/(m+1)$ be the discretization parameter, and replace the interval $[0, 1]$ by grid points $x_j := jh$ for $j = 0, 1, \dots, m+1$. We then obtain approximations v_j to the exact solution $u(x_j)$ for $j = 1, \dots, m$ by replacing the differential equation by the difference equation

$$\frac{-v_{j-1} + 2v_j - v_{j+1}}{h^2} = f(jh), \quad j = 1, \dots, m, \quad v_0 = v_{m+1} = 0.$$

Moving the h^2 factor to the right hand side this can be written as an $m \times m$ linear system

$$\mathbf{T}\mathbf{v} = \begin{bmatrix} 2 & -1 & 0 & & & \\ -1 & 2 & -1 & & & \\ 0 & \ddots & \ddots & \ddots & & \\ & & & & 0 & \\ & & & & -1 & 2 & -1 \\ & & & & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_{m-1} \\ v_m \end{bmatrix} = h^2 \begin{bmatrix} f(h) \\ f(2h) \\ \vdots \\ f((m-1)h) \\ f(mh) \end{bmatrix} =: \mathbf{b}. \quad (2.2)$$

The matrix \mathbf{T} is called the **second derivative matrix** and will occur frequently in this book. It is our first example of a tridiagonal matrix, $\mathbf{T} = \text{tridiag}(a_i, d_i, c_i) \in \mathbb{R}^{m \times m}$, where in this case $a_i = c_i = -1$ and $d_i = 2$ for all i .

2.2 LU Factorization of a Tridiagonal System

We saw in Theorem 1.19 that if naive Gaussian elimination can be used, it leads to an LU factorization $\mathbf{A} = \mathbf{LU}$ of the coefficient matrix, i.e., where \mathbf{L} is lower triangular with ones on the diagonal and \mathbf{U} is upper triangular. Consider now the special case of a nonsingular tridiagonal linear system $\mathbf{Ax} = \mathbf{b}$, where $\mathbf{A} = \text{tridiag}(a_i, d_i, c_i) \in \mathbb{C}^{n \times n}$. To solve $\mathbf{Ax} = \mathbf{b}$ we can either use Gaussian elimination adapted to the special structure, or as we do here, compute the LU factorization

directly. Let us try to construct triangular matrices \mathbf{L} and \mathbf{U} such that the product $\mathbf{A} = \mathbf{LU}$ has the form

$$\begin{bmatrix} d_1 & c_1 & & \\ a_1 & d_2 & c_2 & \\ & \ddots & \ddots & \ddots \\ & & a_{n-2} & d_{n-1} & c_{n-1} \\ & & & a_{n-1} & d_n \end{bmatrix} = \begin{bmatrix} 1 & & & \\ l_1 & 1 & & \\ & \ddots & \ddots & \\ & & l_{n-1} & 1 \end{bmatrix} \begin{bmatrix} u_1 & c_1 & & \\ & \ddots & \ddots & \\ & & u_{n-1} & c_{n-1} \\ & & & u_n \end{bmatrix}. \quad (2.3)$$

If \mathbf{L} and \mathbf{U} can be determined we can find \mathbf{x} by solving two simpler systems $\mathbf{Ly} = \mathbf{b}$ and $\mathbf{Ux} = \mathbf{y}$.

To find \mathbf{L} and \mathbf{U} we note that \mathbf{L} and \mathbf{U} are bidiagonal, \mathbf{L} has ones on the diagonal, and that we have the same c_i elements on the super-diagonals of \mathbf{A} and \mathbf{U} . By equating elements in (2.3) we find

$$d_1 = u_1, \quad a_k = l_k u_k, \quad d_{k+1} = u_{k+1} + l_k c_k, \quad k = 1, 2, \dots, n-1.$$

Solving for l_k and u_k leads to

$$u_1 = d_1, \quad l_k = \frac{a_k}{u_k}, \quad u_{k+1} = d_{k+1} - l_k c_k, \quad k = 1, 2, \dots, n-1. \quad (2.4)$$

If u_1, u_2, \dots, u_{n-1} are nonzero then (2.4) is well defined. If in addition $u_n \neq 0$ then we can solve $\mathbf{Ly} = \mathbf{b}$ and $\mathbf{Ux} = \mathbf{y}$ for \mathbf{y} and \mathbf{x} .

$$\begin{aligned} y_1 &= b_1, & y_k &= b_k - l_{k-1} y_{k-1}, & k &= 2, 3, \dots, n, \\ x_n &= y_n/u_n, & x_k &= (y_k - c_k x_{k+1})/u_k, & k &= n-1, \dots, 2, 1. \end{aligned} \quad (2.5)$$

We formulate this as two algorithms. Since division by zero can occur, the algorithms will not work in general. We give one condition for success in Theorem 2.5 below.

Algorithm 2.1 (trifactor)

Vectors $\mathbf{l} \in \mathbb{C}^{n-1}$, $\mathbf{u} \in \mathbb{C}^n$ are computed from $\mathbf{a}, \mathbf{c} \in \mathbb{C}^{n-1}$, $\mathbf{d} \in \mathbb{C}^n$. This implements the LU factorization of a tridiagonal matrix.

```

1 function [l,u]=trifactor(a,d,c)
2 u=d; l=a;
3 for k=1:length(a)
4 l(k)=a(k)/u(k);
5 u(k+1)=d(k+1)-l(k)*c(k);
6 end

```

Algorithm 2.2 (trisolve)

The solution x of the tridiagonal system $LUX = b$ is found from (2.5). Here $\mathbf{l}, \mathbf{c} \in \mathbb{C}^{n-1}$, $\mathbf{u}, \mathbf{b} \in \mathbb{C}^n$. The vectors \mathbf{l}, \mathbf{u} can be output from `trifactor`.

```

1 function x = trisolve (l,u,c,b)
2 x=b;
3 n= size(b,1);
4 for k =2:n
5 x(k,:)=b(k,:)-l(k-1)*x(k-1,:);
6 end
7 x(n,:)=x(n,:)/u(n);
8 for k=n-1:-1:1
9 x(k,:)=(x(k,:)-c(k)*x(k+1,:))/ u(k);
10 end
```

The number of arithmetic operations to compute the LU factorization of a tridiagonal matrix using Algorithm 2.1 is only $3n - 3$, while the number of arithmetic operations for Algorithm 2.2 is $5n - 4$. This means that the number of arithmetic operations (the complexity) to solve a tridiagonal system is $O(n)$, or more precisely $8n - 7$, and this number only grows linearly with n , while Gaussian elimination on a full $n \times n$ system is an $O(n^3)$ process.

2.2.1 Diagonal Dominance

We show that Algorithms 2.1, 2.2 are well defined for a class of tridiagonal linear systems. Moreover, these linear systems have unique solutions.

Definition 2.3 (Diagonal dominance)

The matrix $A = [a_{ij}] \in \mathbb{C}^{n \times n}$ is weakly diagonally dominant if

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \quad i = 1, \dots, n. \quad (2.6)$$

It is strictly diagonally dominant if strict inequality holds for $i = 1, \dots, n$.

The following holds for strictly diagonally dominant matrices.

Theorem 2.4 (Strict diagonal dominance)

A strictly diagonally dominant matrix is nonsingular. Moreover, the solution \mathbf{x} of $\mathbf{Ax} = \mathbf{b}$ is bounded as follows:

$$\max_{1 \leq i \leq n} |x_i| \leq \max_{1 \leq i \leq n} \left(\frac{|b_i|}{\sigma_i} \right), \quad \text{where } \sigma_i := |a_{ii}| - \sum_{j \neq i} |a_{ij}|. \quad (2.7)$$

Proof. We first show that the bound (2.7) holds for any solution \mathbf{x} . Choose k so that $|x_k| = \max_i |x_i|$. Then

$$|b_k| = |a_{kk}x_k + \sum_{j \neq k} a_{kj}x_j| \geq |a_{kk}| |x_k| - \sum_{j \neq k} |a_{kj}| |x_j| \geq |x_k| (|a_{kk}| - \sum_{j \neq k} |a_{kj}|),$$

and this implies $\max_{1 \leq i \leq n} |x_i| = |x_k| \leq \frac{|b_k|}{\sigma_k} \leq \max_{1 \leq i \leq n} \left(\frac{|b_i|}{\sigma_i} \right)$. For nonsingularity, if $\mathbf{Ax} = \mathbf{0}$, then $\max_{1 \leq i \leq n} |x_i| \leq 0$ by (2.7), and so $\mathbf{x} = \mathbf{0}$. \square

To answer the question of nonsingularity for weakly diagonally dominant matrices is not so easy. For example, the zero matrix is weakly diagonally dominant. Consider the 3 matrices

$$\mathbf{A}_1 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_3 = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}.$$

They are all weakly diagonally dominant, but \mathbf{A}_1 and \mathbf{A}_2 are singular, while \mathbf{A}_3 is nonsingular. Indeed, for \mathbf{A}_1 column two is the sum of columns one and three, \mathbf{A}_2 has a zero row, and $\det(\mathbf{A}_3) = 4 \neq 0$.

In the literature diagonal dominance is therefore most often defined by including some additional sufficient condition(s). We prove the following result for tridiagonal matrices.

Theorem 2.5 (Weak diagonal dominance)

Suppose $\mathbf{A} = \text{tridiag}(a_i, d_i, c_i) \in \mathbb{C}^{n \times n}$ is tridiagonal and weakly diagonally dominant. If $|d_1| > |c_1|$ and $a_i \neq 0$ for $i = 1, \dots, n-2$, then \mathbf{A} has a unique LU factorization (2.3). If in addition $d_n \neq 0$, then \mathbf{A} is nonsingular.

Proof. The matrix \mathbf{A} has an LU factorization if the u_k 's in (2.4) are nonzero for $k = 1, \dots, n-1$. For this it is sufficient to show by induction that $|u_k| > |c_k|$ for $k = 1, \dots, n-1$. By assumption $|u_1| = |d_1| > |c_1|$. Suppose $|u_k| > |c_k|$ for some $1 \leq k \leq n-2$. Then $|c_k|/|u_k| < 1$ and by (2.4) and since $a_k \neq 0$

$$|u_{k+1}| = |d_{k+1} - l_k c_k| = |d_{k+1} - \frac{a_k c_k}{u_k}| \geq |d_{k+1}| - \frac{|a_k||c_k|}{|u_k|} > |d_{k+1}| - |a_k|. \quad (2.8)$$

This also holds for $k = n-1$ if $a_{n-1} \neq 0$. By weak diagonal dominance $|u_{k+1}| > |c_{k+1}|$ and it follows by induction that an LU factorization exists. It is unique since any LU factorization must satisfy (2.4). For nonsingularity we need to show that $u_n \neq 0$. For then by Lemma 1.9, both \mathbf{L} and \mathbf{U} are nonsingular, and this is equivalent to $\mathbf{A} = \mathbf{LU}$ being nonsingular. If $a_{n-1} \neq 0$ then by (2.4) $|u_n| > |d_n| - |a_{n-1}| \geq 0$ by weak diagonal dominance, while if $a_{n-1} = 0$ then again by (2.8) $|u_n| \geq |d_n| > 0$. \square

Consider now the special system $\mathbf{T}\mathbf{v} = \mathbf{b}$ given by (2.2). The matrix \mathbf{T} is weakly diagonally dominant and satisfies the additional conditions in Theorem 2.5. Thus it is nonsingular and we can solve the system in $O(n)$ arithmetic operations using Algorithms 2.1,2.2.

We could use the explicit inverse of \mathbf{T} , given in Exercise 2.7, to compute the solution of $\mathbf{T}\mathbf{v} = \mathbf{b}$ as $\mathbf{v} = \mathbf{T}^{-1}\mathbf{b}$. However this is not a good idea. In fact, all elements in \mathbf{T}^{-1} are nonzero and the calculation of $\mathbf{T}^{-1}\mathbf{b}$ requires $O(n^2)$ operations.

Exercise 2.6 (LU factorization of 2. derivative matrix)

Show that $\mathbf{T} = \mathbf{L}\mathbf{U}$, where

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ -\frac{1}{2} & 1 & \ddots & & \vdots \\ 0 & -\frac{2}{3} & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & -\frac{m-1}{m} & 1 \end{bmatrix}, \mathbf{U} = \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ 0 & \frac{3}{2} & -1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \frac{m}{m-1} & -1 \\ 0 & \cdots & \cdots & 0 & \frac{m+1}{m} \end{bmatrix} \quad (2.9)$$

is the LU factorization of \mathbf{T} .

Exercise 2.7 (Inverse of 2. derivative matrix)

Let $\mathbf{S} \in \mathbb{R}^{m \times m}$ have elements s_{ij} given by

$$s_{i,j} = s_{j,i} = \frac{1}{m+1} j(m+1-i), \quad 1 \leq j \leq i \leq m. \quad (2.10)$$

Show that $\mathbf{S}\mathbf{T} = \mathbf{I}$ and conclude that $\mathbf{T}^{-1} = \mathbf{S}$.

Exercise 2.8 (Central difference approximation of 2. derivative)

Consider

$$\delta^2 f(x) := \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}, \quad h > 0, \quad f : [x-h, x+h] \rightarrow \mathbb{R}.$$

1. Show that if $f \in C^2[x-h, x+h]$ then for some η_2

$$\delta^2 f(x) = f''(\eta_2), \quad x-h < \eta_2 < x+h.$$

2. If $f \in C^4[x-h, x+h]$ then for some η_4

$$\delta^2 f(x) = f''(x) + \frac{h^2}{12} f^{(4)}(\eta_4), \quad x-h < \eta_4 < x+h.$$

$\delta^2 f(x)$ is known as the **central difference approximation** to the second derivative at x .

Exercise 2.9 (Two point boundary value problem)

We consider a finite difference method for the two point boundary value problem

$$\begin{aligned} -u''(x) + r(x)u'(x) + q(x)u(x) &= f(x), \text{ for } x \in [a, b], \\ u(a) &= g_0, \quad u(b) = g_1. \end{aligned} \quad (2.11)$$

We assume that the given functions f, q and r are continuous on $[a, b]$ and that $q(x) \geq 0$ for $x \in [a, b]$. It can then be shown that (2.11) has a unique solution u .

To solve (2.11) numerically we choose $m \in \mathbb{N}$, $h = (b-a)/(m+1)$, $x_j = a+jh$ for $j = 0, 1, \dots, m+1$ and solve the difference equation

$$\frac{-v_{j-1} + 2v_j - v_{j+1}}{h^2} + r(x_j) \frac{v_{j+1} - v_{j-1}}{2h} + q(x_j)v_j = f(x_j), \quad j = 1, \dots, m, \quad (2.12)$$

with $v_0 = g_0$ and $v_{m+1} = g_1$.

- (a) Show that (2.12) leads to a tridiagonal linear system $\mathbf{Av} = \mathbf{b}$, where $\mathbf{A} = \text{tridiag}(a_j, d_j, c_j) \in \mathbb{R}^{m \times m}$ has elements

$$a_j = -1 - \frac{h}{2}r(x_j), \quad c_j = -1 + \frac{h}{2}r(x_j), \quad d_j = 2 + h^2q(x_j),$$

and

$$b_j = \begin{cases} h^2f(x_1) - a_1g_0, & \text{if } j = 1, \\ h^2f(x_j), & \text{if } 2 \leq j \leq m-1, \\ h^2f(x_m) - c_mg_1, & \text{if } j = m. \end{cases}$$

- (b) Show that the linear system satisfies the conditions in Theorem 2.5 if the spacing h is so small that $\frac{h}{2}|r(x)| < 1$ for all $x \in [a, b]$.
(c) Propose a method to find v_1, \dots, v_m .

Exercise 2.10 (Two point boundary value problem; computation)

- (a) Consider the problem (2.11) with $r = 0$, $f = q = 1$ and boundary conditions $u(0) = 1$, $u(1) = 0$. The exact solution is $u(x) = 1 - \sinh x / \sinh 1$. Write a computer program to solve (2.12) for $h = 0.1, 0.05, 0.025, 0.0125$, and compute the "error" $\max_{1 \leq j \leq m} |u(x_j) - v_j|$ for each h .
- (b) Make a combined plot of the solution u and the computed points v_j , $j = 0, \dots, m+1$ for $h = 0.1$.
- (c) One can show that the error is proportional to h^p for some integer p . Estimate p based on the error for $h = 0.1, 0.05, 0.025, 0.0125$.

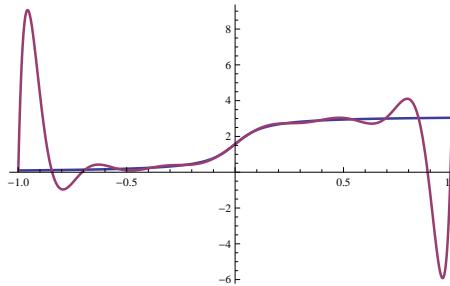


Figure 2.1: The polynomial of degree 13 interpolating $f(x) = \arctan(10x) + \pi/2$ on $[-1, 1]$. See text

2.3 Cubic Spline Interpolation

We next consider an **interpolation problem** leading to a tridiagonal linear system.

Given $n+1 \geq 2$ interpolation sites $\mathbf{x} = [x_1, \dots, x_{n+1}]^T$ with $a := x_1 < \dots < x_{n+1} =: b$ and y values $\mathbf{y} = [y_1, \dots, y_{n+1}]^T$. We seek a function $g : [a, b] \rightarrow \mathbb{R}$ such that

$$g(x_i) = y_i, \text{ for } i = 1, \dots, n+1. \quad (2.13)$$

2.3.1 The Runge Phenomenon

Since there are $n+1$ interpolation conditions in (2.13) a natural choice for a function g is a polynomial of degree at most n . As shown in most books on numerical methods such a g is uniquely defined and there are good algorithms for computing it. Evidently, when $n=1$, g is the straight line

$$g(x) = y_1 + \frac{y_2 - y_1}{x_2 - x_1}(x - x_1), \quad (2.14)$$

known as the **linear interpolation polynomial**.

Polynomial interpolation is an important technique which often gives good results, but it can have problems when n is large and the sites are not carefully chosen. As an example, the polynomial g of degree $n \leq 13$ interpolating the function f given by $f(x) = \arctan(10x) + \pi/2$, $x \in [-1, 1]$ at the points $x_i = -1 + 2(i-1)/n$, $i = 1, \dots, n+1$ is shown in Figure 2.1. The interpolant has large oscillations near the ends of the range. This is called the **Runge phenomenon** and is partly due to the fact that the x 's are uniformly spaced. Using uniform sites with larger n will only make the oscillations bigger.

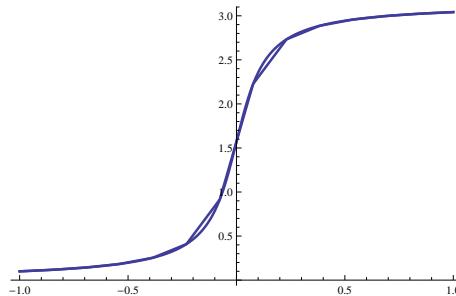


Figure 2.2: The piecewise linear polynomial interpolating $f(x) = \arctan(10x) + \pi/2$ at $n = 14$ uniform points on $[-1, 1]$.

2.3.2 Piecewise Linear and Cubic Spline Interpolation

To avoid oscillations like the one in Figure 2.1 piecewise linear interpolation can be used. An example is shown in Figure 2.2. The interpolant g approximates the original function quite well, and for some applications, like plotting, the linear interpolant using many points is good enough. Note that g is a piecewise polynomial of the form

$$g(x) := \begin{cases} p_1(x), & \text{if } x_1 \leq x < x_2, \\ p_2(x), & \text{if } x_2 \leq x < x_3, \\ \vdots & \\ p_{n-1}(x), & \text{if } x_{n-1} \leq x < x_n, \\ p_n(x), & \text{if } x_n \leq x \leq x_{n+1}, \end{cases} \quad (2.15)$$

where each p_i is a polynomial of degree ≤ 1 . In particular, p_1 is given in (2.14) and the other polynomials p_i are given by similar expressions.

The piecewise linear interpolant is continuous, but the first derivative has jumps at the interior interpolation sites. We can obtain a smoother approximation by letting g be a piecewise polynomial of higher degree. Degree 3 (cubic) is often used.

A piecewise cubic polynomial g of the form (2.15) is called a **(C^2) cubic spline** if each p_i is a cubic polynomial and $g \in C^2[a, b]$, i.e., g is continuous and has a continuous first and second derivative on $[a, b]$. The sites x_2, x_3, \dots, x_n where the third derivative of a C^2 cubic spline can have jumps are called **knots**. The knots and data sites can be different, but in this book the knots will be a subset of the data sites.

We are looking for a C^2 cubic spline that satisfies

1. $p_{i-1}^{(j)}(x_i) = p_i^{(j)}(x_i), \quad j = 0, 1, 2, \quad i = 2, \dots, n, \quad (C^2 \text{ conditions}),$

2. $g(x_i) = y_i, \quad i = 1, 2, \dots, n+1$ (interpolation conditions).
3. Two boundary conditions.

The two extra conditions are needed for uniqueness. Indeed counting requirements we have $3(n - 1)$ C^2 conditions, $n + 1$ interpolation conditions, and two boundary conditions, adding up to $4n$. Since a cubic polynomial has four coefficients, this number is equal to the number of coefficients of the n polynomials p_1, \dots, p_n .

There are many possible choices for boundary conditions. One standard choice is to use the **first derivative boundary conditions**

$$g'(a) = s_1, \quad g'(b) = s_{n+1}, \quad (2.16)$$

where s_1 and s_{n+1} are given end slopes.

Example 2.11 (A cubic spline interpolant)

Show that g given by

$$g(x) := \begin{cases} p_1(x) = -x^2 + 2x^3, & \text{if } 0 \leq x < 1, \\ p_2(x) = 1 + 4(x-1) + 5(x-1)^2 + 6(x-1)^3, & \text{if } 1 \leq x \leq 2, \end{cases} \quad (2.17)$$

is a cubic C^2 spline interpolating the data

$$\mathbf{x} := [0, 1, 2]^T, \quad \mathbf{y} := [0, 1, 16]^T, \quad \mathbf{s} := [0, 32]^T.$$

Discussion: Clearly p_1 and p_2 are cubic polynomials. Also the interpolation conditions are satisfied: $g(0) = p_1(0) = 0$, $g(1) = p_2(1) = 1$, $g(2) = p_2(2) = 16$, $g'(0) = p_1'(0) = 0$, $g'(2) = p_2'(2) = 32$. Moreover, $g \in C^2$ since $p_1(1) = p_2(1)$, $p_1'(1) = p_2'(1)$, $p_1''(1) = p_2''(1)$. The data are sampled from the function $f : [0, 2] \rightarrow \mathbb{R}$ given by the rule $f(x) = x^4$. A plot of f and g is shown in Figure 2.3. It is hard to distinguish the two curves.

It can be shown that there is a unique C^2 cubic spline satisfying

$$g(x_i) = y_i, \quad i = 1, 2, \dots, n+1, \quad g'(x_1) = s_1, \quad g'(x_{n+1}) = s_{n+1}, \quad (2.18)$$

see Exercise 2.21. The following theorem shows that this function has an amazing property. It is the interpolant with the smallest second derivative.

Theorem 2.12 (Cubic spline; minimal 2. derivative)

Suppose g is a cubic C^2 spline. Then

$$\int_a^b (g''(x))^2 dx \leq \int_a^b (h''(x))^2 dx$$

for all $h \in C^2[a, b]$ such that $h(x_i) = g(x_i)$, $i = 1, \dots, n+1$, $h'(a) = g'(a)$, and $h'(b) = g'(b)$.

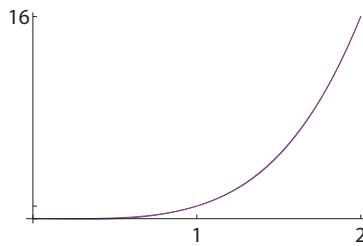


Figure 2.3: A cubic spline with one knot interpolating $f(x) = x^4$ on $[0, 2]$.

For the proof see Exercise 2.22

The name spline is inherited from a “physical analogue”, an elastic ruler that is used to draw smooth curves. Heavy weights, called **ducks**, are used to force the ruler to pass through, or near given locations. (Cf. Figure 2.4). The ruler will take a shape that minimizes its potential energy. Since the potential energy is proportional to the integral of the square of the curvature, and the curvature can be approximated by the second derivative it follows from Theorem 2.12 that the mathematical spline g approximately models the physical spline.



Figure 2.4: A physical spline with ducks.

If the cubic spline g in Theorem 2.12 interpolates a smooth function f ,

$(f \in C^4)$. Then it can be shown that

$$|g(x) - f(x)| \leq \frac{5}{384} h^4 \max_{a \leq t \leq b} |f^{(4)}(t)|, \quad h = \max_{1 \leq i \leq n} x_{i+1} - x_i.$$

Moreover the constant $\frac{5}{384}$ is best possible ([8]). We say that g approximates f with 4th order accuracy, i.e., $f - g = O(h^4)$.

2.3.3 Give me a Moment

To find the spline interpolant we need to solve a tridiagonal linear system. The unknowns can be the first or second derivatives at the knots. We consider here using second derivatives which are sometimes called **moments**.

We will consider alternative boundary conditions that also give 4th order accuracy. These conditions can be written

$$p_1 = p_2, \quad p_n = p_{n-1}, \quad (\text{not-a-knot conditions}). \quad (2.19)$$

The name of these boundary conditions refer to the fact that the first and last knot, i.e., x_2 and x_n , are not knots. The name **free boundary** is also used. An analogue of Theorem 2.12 does not hold for these boundary conditions, but we do not have to specify derivatives at the ends of the range.

The following theorem details the construction. For simplicity, we consider only equidistant sites.

Theorem 2.13 (Cubic spline with not-a-knot boundary conditions)

Given $a < b$, $h = (b - a)/n$ with $n \geq 4$, $x_i = a + (i - 1)h$, and numbers y_i for $i = 1, \dots, n + 1$. There is a unique cubic C^2 spline g of the form (2.15) such that

$$g(x_i) = y_i, \quad i = 1, \dots, n + 1, \quad p_1 = p_2, \quad p_n = p_{n-1}. \quad (2.20)$$

If we represent each p_i in the **shifted power form**

$$p_i(x) = \sum_{j=1}^4 c_{i,j}(x - x_i)^{j-1}, \quad i = 1, \dots, n, \quad (2.21)$$

then for $i = 1, 2, \dots, n$,

$$c_{i1} = y_i, \quad c_{i2} = \frac{y_{i+1} - y_i}{h} - \frac{h}{3}\mu_i - \frac{h}{6}\mu_{i+1}, \quad c_{i,3} = \frac{\mu_i}{2}, \quad c_{i,4} = \frac{\mu_{i+1} - \mu_i}{6h}, \quad (2.22)$$

where $\mu_1 = 2\mu_2 - \mu_3$, $\mu_{n+1} = 2\mu_n - \mu_{n-1}$, $[\mu_2, \dots, \mu_n]^T$ is the unique solution of

the linear system

$$\begin{bmatrix} 6 & 0 & & \\ 1 & 4 & 1 & \\ & \ddots & \ddots & \ddots \\ & & 1 & 4 & 1 \\ & & & 0 & 6 \end{bmatrix} \begin{bmatrix} \mu_2 \\ \mu_3 \\ \vdots \\ \mu_{n-1} \\ \mu_n \end{bmatrix} = \frac{6}{h^2} \begin{bmatrix} \delta^2 y_2 \\ \delta^2 y_3 \\ \vdots \\ \delta^2 y_{n-1} \\ \delta^2 y_n \end{bmatrix}, \quad (2.23)$$

and $\delta^2 y_i := y_{i+1} - 2y_i + y_{i-1}$, $i = 2, \dots, n$.

Proof. Since the linear system (2.23) is strictly diagonally dominant it has a unique solution and the μ_i 's are well defined. Let $p_i(x) := \sum_{j=1}^4 c_{i,j}(x - x_i)^{j-1}$, $i = 1, \dots, n$ have coefficients given by (2.22). It is not too hard to show that

$$p_i(x_i) = y_i, \quad p_i(x_{i+1}) = y_{i+1}, \quad p_i''(x_i) = \mu_i, \quad p_i''(x_{i+1}) = \mu_{i+1}, \quad i = 1, \dots, n. \quad (2.24)$$

But then $g(x_i) = y_i$ for $i = 1, \dots, n+1$. Since $p_{i-1}(x_i) = p_i(x_i) = y_i$ and $p_{i-1}''(x_i) = p_i''(x_i) = \mu_i$ for $i = 2, \dots, n$ it follows that $g \in C^2$ if and only if $p_{i-1}'(x_i) = p_i'(x_i)$ for $i = 2, \dots, n$. By (2.21)

$$\begin{aligned} p_{i-1}'(x_i) &= c_{i-1,2} + 2hc_{i-1,3} + 3h^2c_{i-1,4} \\ &= \frac{y_i - y_{i-1}}{h} - \frac{h}{3}\mu_{i-1} - \frac{h}{6}\mu_i + h\mu_{i-1} + \frac{h}{2}\mu_i - \frac{h}{2}\mu_{i-1} \\ &= \frac{y_i - y_{i-1}}{h} + \frac{h}{6}\mu_{i-1} + \frac{h}{3}\mu_i \\ p_i'(x_i) &= c_{i,2} = \frac{y_{i+1} - y_i}{h} - \frac{h}{3}\mu_i - \frac{h}{6}\mu_{i+1}. \end{aligned} \quad (2.25)$$

But then $p_{i-1}'(x_i) = p_i'(x_i)$ if and only if $\frac{h}{6}\mu_{i-1} + \frac{2h}{3}\mu_i + \frac{h}{6}\mu_{i+1} = \frac{1}{h}(y_{i+1} - 2y_i + y_{i-1})$ leading to the middle equations in (2.23)

$$\mu_{i-1} + 4\mu_i + \mu_{i+1} = \frac{6}{h^2}\delta^2 y_i, \quad i = 2, \dots, n. \quad (2.26)$$

For the first and last equation we need to consider the boundary conditions. Now $p_1 = p_2$ implies that $p_1'''(x_2) = p_2'''(x_2)$ and therefore it follows from (2.22) that $6\frac{\mu_2 - \mu_1}{6h} = 6\frac{\mu_3 - \mu_2}{6h}$ or $\mu_1 = 2\mu_2 - \mu_3$. Inserting this in $\mu_1 + 4\mu_2 + \mu_3 = \frac{6}{h^2}\delta^2 y_2$ we obtain $6\mu_2 = \frac{6}{h^2}\delta^2 y_2$ the first equation in (2.23). The formula for μ_{n+1} and the last equation follow in a similar manner.

To show uniqueness suppose g_1 and g_2 are two cubic C^2 splines interpolating the same data \mathbf{x}, \mathbf{y} . Then $g := g_1 - g_2$ is a cubic C^2 spline interpolating zero y values. By strict diagonal dominance the solution $[\mu_2, \dots, \mu_n]^T$ of (2.23) is zero

and then μ_1 and μ_{n+1} are also zero. It then follows from (2.22) that all coefficients $c_{i,j}$ are zero. We conclude that $g = 0$ and $g_1 = g_2$. \square

Example 2.14 (Not-a-knot)

Find the cubic C^2 spline interpolating the data $(0, 0, 0, \frac{1}{6}, \frac{2}{3}, \frac{1}{6}, 0, 0, 0)$ on the interval $[a, b] = [-2, 6]$. The tridiagonal linear system is

$$\begin{bmatrix} 6 & 0 \\ 1 & 4 & 1 \\ & 1 & 4 & 1 \\ & & 1 & 4 & 1 \\ & & & 1 & 4 & 1 \\ & & & & 1 & 4 & 1 \\ & & & & & 0 & 6 \end{bmatrix} \begin{bmatrix} \mu_2 \\ \mu_3 \\ \mu_4 \\ \mu_5 \\ \mu_6 \\ \mu_7 \\ \mu_8 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 2 \\ -6 \\ 2 \\ 1 \\ 0 \end{bmatrix},$$

with solution $\mu_4 = \mu_6 = 1$, $\mu_5 = -2$, and $\mu_j = 0$ otherwise. Using (2.22) (cf. Exercise 2.16) we find $p_1 = p_2 = p_7 = p_8 = 0$ and

$$g(x) := \begin{cases} p_3(x) = \frac{1}{6}x^3, & \text{if } 0 \leq x < 1, \\ p_4(x) = \frac{1}{6} + \frac{1}{2}(x-1) + \frac{1}{2}(x-1)^2 - \frac{1}{2}(x-1)^3, & \text{if } 1 \leq x < 2, \\ p_5(x) = \frac{2}{3} - (x-2)^2 + \frac{1}{2}(x-2)^3, & \text{if } 2 \leq x < 3, \\ p_6(x) = \frac{1}{6} - \frac{1}{2}(x-3) + \frac{1}{2}(x-3)^2 - \frac{1}{6}(x-3)^3 = \frac{1}{6}(4-x)^3, & \text{if } 3 \leq x < 4, \end{cases} \quad (2.27)$$

A plot of this spline is shown in Figure 2.5. It is known as a C^2 cubic B-spline.

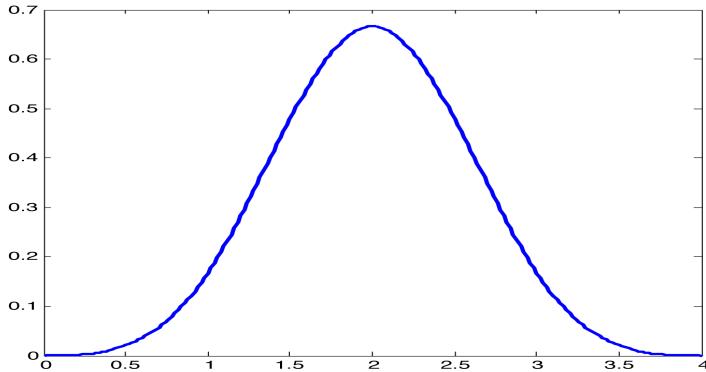


Figure 2.5: A cubic B-spline.

To plot a piecewise polynomial g in the form (2.15) we need to compute y values $q_j = g(r_j)$ at a number of sites $\mathbf{r} = [r_1, \dots, r_m] \in \mathbb{R}^m$ for some reasonably

large integer m . To determine $g(r_j)$ for some j we need to find an integer i_j so that $g(r_j) = p_{i_j}(r_j)$.

Given $k \in \mathbb{N}$, $\mathbf{t} = [t_1, \dots, t_k]$ and a real number x . The problem is to compute an integer i so that $i = 1$ if $x < t_2$, $i = k$ if $x \geq t_k$, and $t_i \leq x < t_{i+1}$ otherwise. If \mathbf{x} is a vector then a vector \mathbf{i} should be computed, such that the j th component of \mathbf{i} gives the location of the j th component of \mathbf{x} . The following Matlab function determines $\mathbf{i} = [i_1, \dots, i_m]$. It uses the built in Matlab functions `length`, `min`, `sort`, `find`.

Algorithm 2.15 (findsubintervals)

```

1 function i = findsu binterv al s1(t,x)
2 k= length(t); m= length(x);
3 if k<2
4     i= ones(m,1);
5 else
6     t(1)= min(x(1),t(1))-1;
7     [~,j]= sort([t(:)',x(:)'']);
8     i=(find(j>k)-(1:m))';
9 end
```

The not-a-knot cubic spline interpolating the function used to illustrate the Runge phenomenon is shown in Figure 2.6, (cf. Exercise 2.16). The spline approximate the function quite well.

Exercise 2.16 (Arctan example)

Given an interval $[a, b]$ and a vector $\mathbf{y} \in \mathbb{R}^{n+1}$ containing at least 5 components. Let g be the interpolating spline in Theorem 2.13 such that $g(x_i) = y_i$, where $x_i = a + (i - 1)\frac{b-a}{n}$ for $i = 1, 2, \dots, n + 1$. The vector $\mathbf{x} = [x_1, \dots, x_n]$ and the shifted power basis coefficients $\mathbf{c} \in \mathbb{R}^{4n}$ in (2.21) are returned in the following algorithm. We use algorithms 2.1 and 2.2 to solve the tridiagonal linear system. Use the following algorithm to compute the $c_{i,j}$ in Examples 2.14 and Figure 2.6.

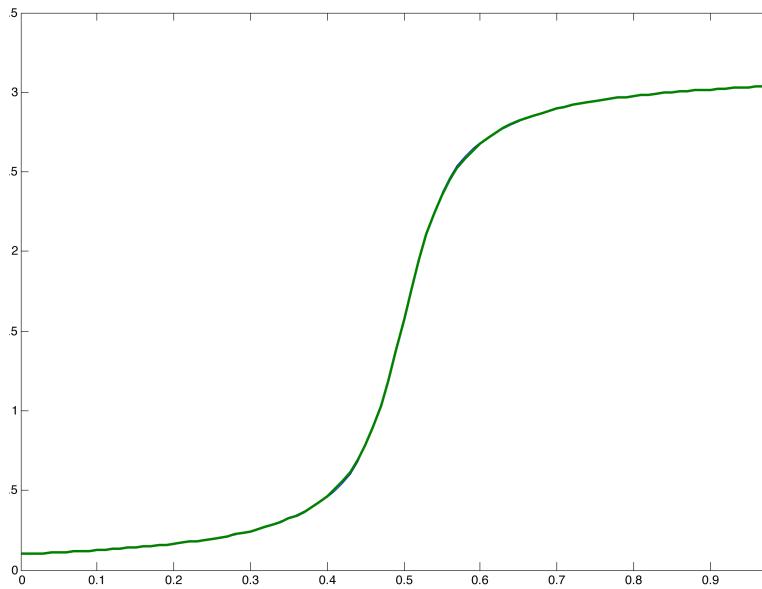


Figure 2.6: The cubic spline interpolating $f(x) = \arctan(10x) + \pi/2$ at 14 equidistant sites on $[-1, 1]$. The exact function is also shown.

Algorithm 2.17 (splineint)

```

1 function [x,C]=splineint(a,b,y)
2 y=y(:); n=length(y)-1;
3 h=(b-a)/n; x=a:h:b-h;
4 a1=[ones(n-3,1);0]; c1=[0;ones(n-3,1)];
5 d1=[6;4*ones(n-3,1);6];
6 [l,u]= trifactor ( a1,d1,c1 );
7 b1=6/h^2*(y(3:n+1)-2*y(2:n)+y(1:n-1));
8 mu= trisolve ( l,u,c1,b1 );
9 mu=[2*mu(1)-mu(2);mu;2*mu(n-1)-mu(n-2)];
10 delta=(y(2:n+1)-y(1:n))/h;
11 C=zeros(4*n,1);
12 C(1:4:n-3)=y(1:n);
13 C(2:4:n-2)=delta-h*mu(1:n)/3-h*mu(2:n+1)/6;
14 C(3:4:n-1)=mu(1:n)/2;
15 C(4:4:n)=(mu(2:n+1)-mu(1:n))/(6*h);

```

Exercise 2.18 (Splineevaluation)

Use the following algorithm to make the plots in Figures 2.14 and 2.6.

Algorithm 2.19 (splineeval) Given the output \mathbf{x}, \mathbf{C} of Algorithm 2.17 defining a cubic spline g , and a vector \mathbf{X} . The vector $\mathbf{G} = g(\mathbf{X})$ is computed.

```

1 function [X,G]=splineval(x,C,X)
2 m=length(X);
3 i=findsubintervals(x,X);
4 G=zeros(m,1);
5 for j=1:m
6     k=i(j);
7     t=X(j)-xs(k);
8     G(j)=[1,t,t^2,t^3]*C(4*k-3:4*k);
9 end

```

Exercise 2.20 (Bounding the moments)

Let $\mathbf{b} \in \mathbb{R}^{n-1}$ be the right-hand-side of the linear system (2.23). Show that the solution of this system is bounded as follows⁸

$$\max_{2 \leq j \leq n} |\mu_j| \leq \frac{1}{2} \max_{1 \leq j \leq n-1} |b_j|.$$

Exercise 2.21 (Moment equations for 1. derivative boundary conditions)

Suppose in Theorem 2.13 we replace the not-a-knot boundary conditions with the first derivative boundary conditions (2.16). Show that the linear system (2.23) now becomes

$$\begin{bmatrix} 2 & 1 & & \\ 1 & 4 & 1 & \\ \ddots & \ddots & \ddots & \\ & 1 & 4 & 1 \\ & & 1 & 2 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \\ \mu_{n+1} \end{bmatrix} = \frac{6}{h^2} \begin{bmatrix} y_2 - y_1 - hs_1 \\ \delta^2 y_2 \\ \delta^2 y_3 \\ \vdots \\ \delta^2 y_{n-1} \\ \delta^2 y_n \\ hs_{n+1} - y_{n+1} + y_n \end{bmatrix}, \quad (2.28)$$

where $\delta^2 y_i := y_{i+1} - 2y_i + y_{i-1}$, $i = 2, \dots, n$. Hint: Use (2.25) to compute $g'(x_1)$ and $g'(x_{n+1})$,

Show that there is a unique cubic C^2 spline g of the form (2.15) such that (2.18) holds.

Exercise 2.22 (Proof of minimal 2. derivative property)

Study the following proof of Theorem 2.12.

⁸Hint, use Theorem 2.4

Proof. Let h be any interpolant as in the theorem. We first show the orthogonality condition

$$\int_a^b g'' e'' = 0, \quad e := h - g. \quad (2.29)$$

Using the piecewise polynomial nature of g and integration by parts

$$\int_a^b g'' e'' = \sum_{i=1}^n \int_{x_i}^{x_{i+1}} p_i'' e'' = \sum_{i=1}^n [p_i'' e']_{x_i}^{x_{i+1}} - \sum_{i=1}^n \int_{x_i}^{x_{i+1}} p_i''' e'.$$

Both terms on the right are zero. First, since $e'(a) = e'(b) = 0$

$$\sum_{i=1}^n [p_i'' e']_{x_i}^{x_{i+1}} = \sum_{i=1}^n (p_i'' e')(x_{i+1}) - (p_i'' e')(x_i) = (p_n'' e')(b) - (p_1'' e')(a) = 0.$$

Second, since p_i''' is equal to a constant and $e(x_i) = 0$, for $i = 1, \dots, n + 1$

$$\sum_{i=1}^n \int_{x_i}^{x_{i+1}} p_i''' e' = \sum_{i=1}^n p_i''' \int_{x_i}^{x_{i+1}} e' = \sum_{i=1}^n p_i''' (e(x_{i+1}) - e(x_i)) = 0.$$

Writing $h = g + e$ and using (2.29)

$$\begin{aligned} \int_a^b (h'')^2 &= \int_a^b (g'' + e'')^2 \\ &= \int_a^b (g'')^2 + \int_a^b (e'')^2 + 2 \int_a^b g'' e'' \\ &= \int_a^b (g'')^2 + \int_a^b (e'')^2 \geq \int_a^b (g'')^2 \end{aligned}$$

and the proof is complete. \square

2.4 Review Questions

2.4.1 Define the second derivative matrix \mathbf{T} . Why is it nonsingular?

2.4.2 Is a weakly diagonally dominant matrix nonsingular?

2.4.3 Why do we not use the explicit inverse of \mathbf{T} to solve the linear system $\mathbf{T}\mathbf{x} = \mathbf{b}$?

2.4.4 What is the Runge phenomenon?

2.4.5 Where does the name "spline" come from?

- 2.4.6** The values of the second derivatives at the knots of an interpolating cubic spline can be found in $O(h^p)$ arithmetic operations. What is the smallest value of p ?
- 2.4.7** What is the approximation order of an interpolating spline with first derivative or not-a-knot boundary conditions?
- 2.4.8** Show that a strictly diagonally dominant matrix is nonsingular.
- 2.4.9** Does a tridiagonal matrix always have an LU factorization?

Chapter 3

LU Factorizations

Numerical methods for solving systems of linear equations are often based on writing a matrix as a product of simpler matrices. Such a **factorization** is useful if the corresponding matrix problem for each of the factors is simple to solve and extra numerical stability issues are not introduced. Examples of factorizations were encountered in Chapter 1 and 2. It was shown that naive Gaussian elimination was equivalent to an **LU factorization** of the coefficient matrix and that Gaussian elimination with row interchanges resulted in a **PLU factorization**. In Chapter 2 we saw how an LU factorization can be used to solve certain tridiagonal systems efficiently. Other factorizations based on unitary matrices will be considered later in this book.

In this chapter we consider the general theory of LU factorizations. We consider some related factorizations called block LU, symmetric LU or LDLT, Cholesky, and semiCholesky. The latter can be used for symmetric positive semidefinite matrices, and we give an introduction to positive definite and positive semidefinite matrices.

3.1 The LU Factorization

We say that $\mathbf{A} = \mathbf{L}\mathbf{U}$ is an **LU factorization** of $\mathbf{A} \in \mathbb{C}^{n \times n}$ if $\mathbf{L} \in \mathbb{C}^{n \times n}$ is lower triangular (**left triangular**) and $\mathbf{U} \in \mathbb{C}^{n \times n}$ is upper triangular (**right triangular**).

Consider finding \mathbf{L} and \mathbf{U} . Equating the i, j element in \mathbf{A} and the product $\mathbf{L}\mathbf{U}$, and noting that $l_{i,j} = 0$ for $j > i$ and $u_{i,j} = 0$ for $i > j$, we obtain an



Figure 3.1: Henry Jensen, 1915-1974 (left), Prescott Durand Crout, 1907-1984.

equation

$$a_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj}, \quad i, j = 1, 2, \dots, n \quad (3.1)$$

involving the unknown elements in \mathbf{L} and \mathbf{U} . . But this is a system of n^2 equations in $n^2 + n$ unknowns. One way to reduce the number of unknowns is to require that one of the triangular matrices should be **unit triangular**, i. e., have ones on the diagonal. Other scalings of the diagonals are also possible, see Section 3.5. Choosing \mathbf{U} to be unit triangular is sometimes known as a **Crout factorization**.

For our discussion we will assume that \mathbf{L} is unit triangular, i. e., it has ones on the diagonal. Three things can happen. An LU factorization exists and is unique, it exists, but it is not unique, or it does not exist. The following 2×2 example illustrates this.

Example 3.1 (LU of 2×2 matrix)

Let $a, b, c, d \in \mathbb{C}$. An LU factorization of $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ must satisfy the equations

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_1 & 1 \end{bmatrix} \begin{bmatrix} u_1 & u_3 \\ 0 & u_2 \end{bmatrix} = \begin{bmatrix} u_1 & u_3 \\ l_1 u_1 & l_1 u_3 + u_2 \end{bmatrix}$$

for the unknowns l_1 in \mathbf{L} and u_1, u_2, u_3 in \mathbf{U} . The equations are

$$u_1 = a, \quad u_3 = b, \quad l_1 a = c, \quad u_2 = d - l_1 b.$$

These equations do not always have a solution. Indeed, the main problem is the nonlinear equation $l_1 a = c$. There are three cases

1. $a \neq 0$: The matrix has a unique LU factorization with $l_1 = c/a$.

2. $a = 0, c \neq 0$: No LU factorization exists.

3. $a = c = 0$: The LU factorization exists, but it is not unique. Any value for l_1 can be used.

Of the four matrices

$$\mathbf{A}_1 := \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_3 := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_4 := \begin{bmatrix} 0 & 1 \\ 0 & 2 \end{bmatrix}.$$

\mathbf{A}_1 has a unique LU factorization, \mathbf{A}_2 has no LU factorization, \mathbf{A}_3 has a unique LU factorization even if it is singular, and \mathbf{A}_4 has an LU factorization, but it is not unique.

Example 3.2 (LU of 3×3 matrices)

The matrix

$$\mathbf{A} := \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 1 & 3 \end{bmatrix}$$

has an LU factorization $\mathbf{A} = \mathbf{LU}$, with

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & y & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 2-y \end{bmatrix}.$$

It is not unique since $\mathbf{A} = \mathbf{LU}$ for any $y \in \mathbb{C}$.

It is possible to characterize matrices with a unique LU factorization. Recall that naive Gaussian elimination is possible if and only if all the leading principal submatrices of order less than n are nonsingular, and in that case naive Gaussian elimination leads to an LU factorization of the coefficient matrix. We will now show that this condition is both necessary and sufficient for the existence of a unique LU factorization. We first recall the definition of (leading) principal submatrices.

Definition 3.3 (Principal submatrix)

For $k = 1, \dots, n$ the matrices $\mathbf{A}_{[k]} \in \mathbb{C}^{k \times k}$ given by

$$\mathbf{A}_{[k]} := \mathbf{A}(1:k, 1:k) = \begin{bmatrix} a_{11} & \cdots & a_{k1} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{bmatrix}$$

are called the **leading principal submatrices** of $\mathbf{A} \in \mathbb{C}^{n \times n}$. More generally, a matrix $\mathbf{B} \in \mathbb{C}^{k \times k}$ is called a **principal submatrix** of \mathbf{A} if $\mathbf{B} = \mathbf{A}(\mathbf{r}, \mathbf{r})$, where $\mathbf{r} = [r_1, \dots, r_k]$ for some $1 \leq r_1 < \dots < r_k \leq n$. Thus,

$$b_{i,j} = a_{r_i, r_j}, \quad i, j = 1, \dots, k.$$

The determinant of a (leading) principal submatrix is called a **(leading) principal minor**.

A principal submatrix is leading if $r_j = j$ for $j = 1, \dots, k$. Also a principal submatrix is special in that it uses the same rows and columns of \mathbf{A} . For $k = 1$ The only principal submatrices of order $k = 1$ are the diagonal elements of \mathbf{A} .

Example 3.4 (Principal submatrices)

The principal submatrices of $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ are

$$[1], [5], [9], [\frac{1}{4} \frac{2}{5}], [\frac{1}{7} \frac{3}{9}], [\frac{5}{8} \frac{6}{9}], \mathbf{A}.$$

The leading principal submatrices are

$$[1], [\frac{1}{4} \frac{2}{5}], \mathbf{A}.$$

In preparation for the main theorem about LU factorization we prove a simple lemma.

Lemma 3.5 (LU of leading principal sub matrices)

Suppose $\mathbf{A} = \mathbf{L}\mathbf{U}$ is an LU factorization of $\mathbf{A} \in \mathbb{C}^{n \times n}$. For $k = 1, \dots, n$ let $\mathbf{A}_{[k]}, \mathbf{L}_{[k]}, \mathbf{U}_{[k]}$ be the leading principal submatrices of $\mathbf{A}, \mathbf{L}, \mathbf{U}$, respectively. Then $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$ is an LU factorization of $\mathbf{A}_{[k]}$ for $k = 1, \dots, n$.

Proof. For $k = 1, \dots, n - 1$ we partition $\mathbf{A} = \mathbf{L}\mathbf{U}$ as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[k]} & \mathbf{B}_k \\ \mathbf{C}_k & \mathbf{F}_k \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{[k]} & \mathbf{0} \\ \mathbf{M}_k & \mathbf{N}_k \end{bmatrix} \begin{bmatrix} \mathbf{U}_{[k]} & \mathbf{S}_k \\ \mathbf{0} & \mathbf{T}_k \end{bmatrix} = \mathbf{L}\mathbf{U}, \quad (3.2)$$

where $\mathbf{F}_k, \mathbf{N}_k, \mathbf{T}_k \in \mathbb{C}^{n-k, n-k}$. Using block multiplication we find $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$. Since $\mathbf{L}_{[k]}$ is unit lower triangular and $\mathbf{U}_{[k]}$ is upper triangular this is an LU factorization of $\mathbf{A}_{[k]}$. \square

The following theorem give a necessary and sufficient condition for existence of a unique LU factorization.

Theorem 3.6 (LU Theorem)

A square matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ has a unique LU factorization if and only if the leading principal submatrices $\mathbf{A}_{[k]}$ of \mathbf{A} are nonsingular for $k = 1, \dots, n - 1$.

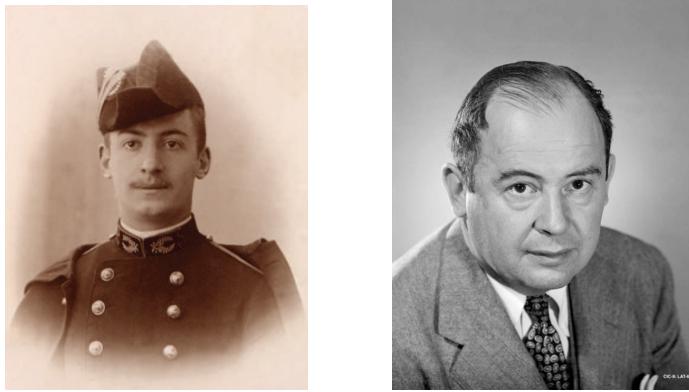


Figure 3.2: André-Louis Cholesky, 1875–1918 (left), John von Neumann, 1903–1957.

Proof. Suppose $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n - 1$. We use induction on n to show that \mathbf{A} has a unique LU factorization. The result is clearly true for $n = 1$, since the unique LU factorization of a 1-by-1 matrix is $[a_{11}] = [1][a_{11}]$. Suppose that $\mathbf{A}_{[n-1]}$ has a unique LU factorization $\mathbf{A}_{[n-1]} = \mathbf{L}_{n-1}\mathbf{U}_{n-1}$, and that $\mathbf{A}_{[1]}, \dots, \mathbf{A}_{[n-1]}$ are nonsingular. By block multiplication

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[n-1]} & \mathbf{b} \\ \mathbf{c}^T & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{n-1} & \mathbf{0} \\ \mathbf{m}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{U}_{n-1} & \mathbf{s} \\ 0 & t \end{bmatrix} = \mathbf{LU}, \quad (3.3)$$

if and only if $\mathbf{m}, \mathbf{s} \in \mathbb{C}^{n-1}$ and $t \in \mathbb{C}$ satisfy $\mathbf{b} = \mathbf{L}_{[n-1]}\mathbf{s}$, $\mathbf{c}^T = \mathbf{m}^T\mathbf{U}_{[n-1]}$, and $a_{nn} = \mathbf{m}^T\mathbf{s} + t$. Since $\mathbf{A}_{[n-1]}$ is nonsingular it follows that \mathbf{L}_{n-1} and \mathbf{U}_{n-1} are nonsingular and therefore \mathbf{m} , \mathbf{s} , and t are uniquely given. Thus (3.3) gives a unique LU factorization of \mathbf{A} .

Conversely, suppose \mathbf{A} has a unique LU factorization $\mathbf{A} = \mathbf{LU}$. By Lemma 3.5 $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$ is an LU factorization of $\mathbf{A}_{[k]}$ for $k = 1, \dots, n - 1$. Suppose $\mathbf{A}_{[k]}$ is singular for some $k \leq n - 1$. We will show that this leads to a contradiction. Let k be the smallest integer so that $\mathbf{A}_{[k]}$ is singular. Since $\mathbf{A}_{[j]}$ is nonsingular for $j \leq k - 1$ it follows from what we have already shown that $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$ is the unique LU factorization of $\mathbf{A}_{[k]}$. The matrix $\mathbf{U}_{[k]}$ is singular since $\mathbf{A}_{[k]}$ is singular and $\mathbf{L}_{[k]}$ is nonsingular. By block multiplication in (3.2) we have $\mathbf{C}_k = \mathbf{M}_k\mathbf{U}_{[k]}$ or $\mathbf{U}_{[k]}^T\mathbf{M}_k^T = \mathbf{C}_k^T$. This can be written as $n - k$ linear systems for the columns of \mathbf{M}_k^T . By assumption \mathbf{M}_k^T exists, but since $\mathbf{U}_{[k]}^T$ is singular \mathbf{M}_k is not unique, a contradiction. \square

A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ can have an LU factorization even if $\mathbf{A}_{[k]}$ is singular for some $k < n$. By Theorem 3.6 such an LU factorization cannot be unique.

Remark 3.7 (LU of upper triangular matrix)

An LU factorization of an upper triangular matrix \mathbf{A} is $\mathbf{A} = \mathbf{IA}$ so it always exists even if \mathbf{A} has zeros somewhere on the diagonal. By Lemma 1.9, if some a_{kk} is zero then $\mathbf{A}_{[k]}$ is singular and the LU factorization cannot be unique. In particular, for the zero matrix any unit lower triangular matrix can be used as \mathbf{L} in an LU factorization.

Remark 3.8 (PLU factorization)

We have shown that a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ has a unique LU factorization if and only if the leading principle submatrices $\mathbf{A}_{[k]}$ are nonsingular for $k = 1, \dots, n-1$. This condition seems fairly restrictive. However, for a nonsingular matrix \mathbf{A} there always is a permutation of the rows so that the permuted matrix has an LU factorization. We obtain a factorization of the form $\mathbf{P}^T \mathbf{A} = \mathbf{LU}$ or equivalently $\mathbf{A} = \mathbf{PLU}$, where \mathbf{P} is a permutation matrix, \mathbf{L} is unit lower triangular, and \mathbf{U} is upper triangular. We call this a **PLU factorization** of \mathbf{A} . (Cf. Theorem 1.29).

Exercise 3.9 (Row interchange)

Show that $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ has a unique LU factorization. Note that we have only interchanged rows in Example 3.1.

Exercise 3.10 (LU of singular matrix)

Find an LU factorization of the singular matrix $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$. Is it unique?

Exercise 3.11 (LU and determinant)

Suppose \mathbf{A} has an LU factorization $\mathbf{A} = \mathbf{LU}$. Show that $\det(\mathbf{A}_{[k]}) = u_{11}u_{22} \cdots u_{kk}$ for $k = 1, \dots, n$.

Exercise 3.12 (Diagonal elements in U)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n-1$. Use Exercise 3.11 to show that the diagonal elements u_{kk} in the LU factorization are

$$u_{11} = a_{11}, \quad u_{kk} = \frac{\det(\mathbf{A}_{[k]})}{\det(\mathbf{A}_{[k-1]})}, \text{ for } k = 2, \dots, n. \quad (3.4)$$

3.2 Block LU Factorization

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a block matrix of the form

$$\mathbf{A} := \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1m} \\ \vdots & & \vdots \\ \mathbf{A}_{m1} & \cdots & \mathbf{A}_{mm} \end{bmatrix}, \quad (3.5)$$

where each (diagonal) block \mathbf{A}_{ii} is square. We call the factorization

$$\mathbf{A} = \mathbf{L}\mathbf{U} = \begin{bmatrix} \mathbf{I} & & & \\ \mathbf{L}_{21} & \mathbf{I} & & \\ \vdots & & \ddots & \\ \mathbf{L}_{m1} & \cdots & \mathbf{L}_{m,m-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{11} & \cdots & \mathbf{U}_{1m} \\ & \mathbf{U}_{21} & \cdots & \mathbf{U}_{2m} \\ & & \ddots & \vdots \\ & & & \mathbf{U}_{mm} \end{bmatrix} \quad (3.6)$$

a **block LU factorization of \mathbf{A}** . Here the i th diagonal blocks \mathbf{I} and \mathbf{U}_{ii} in \mathbf{L} and \mathbf{U} have the same order as \mathbf{A}_{ii} , the i th diagonal block in \mathbf{A} .

The results for elementwise LU factorization carry over to block LU factorization as follows.

Theorem 3.13 (Block LU theorem)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a block matrix of the form (3.5). Then \mathbf{A} has a unique block LU factorization (3.6) if and only if the leading principal block submatrices

$$\mathbf{A}_{\{k\}} := \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1k} \\ \vdots & & \vdots \\ \mathbf{A}_{k1} & \cdots & \mathbf{A}_{kk} \end{bmatrix}$$

are nonsingular for $k = 1, \dots, m - 1$.

Proof. Suppose $\mathbf{A}_{\{k\}}$ is nonsingular for $k = 1, \dots, m - 1$. Following the proof in Theorem 3.6 suppose $\mathbf{A}_{\{m-1\}}$ has a unique block LU factorization $\mathbf{A}_{\{m-1\}} = \mathbf{L}_{\{m-1\}}\mathbf{U}_{\{m-1\}}$, and that $\mathbf{A}_{\{1\}}, \dots, \mathbf{A}_{\{m-1\}}$ are nonsingular. Then $\mathbf{L}_{\{m-1\}}$ and $\mathbf{U}_{\{m-1\}}$ are nonsingular and

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} \mathbf{A}_{\{m-1\}} & \mathbf{B} \\ \mathbf{C}^T & \mathbf{A}_{mm} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{L}_{\{m-1\}} & \mathbf{0} \\ \mathbf{C}^T \mathbf{U}_{\{m-1\}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{\{m-1\}} & \mathbf{L}_{\{m-1\}}^{-1} \mathbf{B} \\ 0 & \mathbf{A}_{mm} - \mathbf{C}^T \mathbf{U}_{\{m-1\}}^{-1} \mathbf{L}_{\{m-1\}}^{-1} \mathbf{B} \end{bmatrix}, \end{aligned} \quad (3.7)$$

is a block LU factorization of \mathbf{A} . It is unique by derivation. Conversely, suppose \mathbf{A} has a unique block LU factorization $\mathbf{A} = \mathbf{L}\mathbf{U}$. Then as in Lemma 3.5 it is easily seen that $\mathbf{A}_{\{k\}} = \mathbf{L}_{\{k\}}\mathbf{U}_{\{k\}}$ is the unique block LU factorization of $\mathbf{A}_{[k]}$ for $k = 1, \dots, m$. The rest of the proof is similar to the proof of Theorem 3.6. \square

Remark 3.14 (Comparing LU and block LU)

The number of arithmetic operations for the block LU factorization is the same as for the ordinary LU factorization. An advantage of the block method is that it combines many of the operations into matrix operations.

Remark 3.15 (A block LU is not an LU)

Note that (3.6) is not an LU factorization of \mathbf{A} since the \mathbf{U}_{ii} 's are not upper triangular in general. To relate the block LU factorization to the usual LU factorization we assume that each \mathbf{U}_{ii} has an LU factorization $\mathbf{U}_{ii} = \tilde{\mathbf{L}}_{ii}\tilde{\mathbf{U}}_{ii}$. Then $\mathbf{A} = \hat{\mathbf{L}}\hat{\mathbf{U}}$, where $\hat{\mathbf{L}} := \mathbf{L} \operatorname{diag}(\tilde{\mathbf{L}}_{ii})$ and $\hat{\mathbf{U}} := \operatorname{diag}(\tilde{\mathbf{L}}_{ii}^{-1})\mathbf{U}$, and this is an ordinary LU factorization of \mathbf{A} .

Exercise 3.16 (Making a block LU into an LU)

Show that $\hat{\mathbf{L}}$ is unit lower triangular and $\hat{\mathbf{U}}$ is upper triangular.

3.3 The Symmetric LU Factorization

We consider next LU factorization of a real symmetric matrix.

Definition 3.17 (Symmetric LU)

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$. A factorization $\mathbf{A} = \mathbf{LDL}^T$, where \mathbf{L} is unit lower triangular and \mathbf{D} is diagonal is called a **symmetric LU factorization** or an **LDLT factorization** of \mathbf{A} .

A matrix which has a symmetric LU factorization must be symmetric since $\mathbf{A}^T = (\mathbf{LDL}^T)^T = \mathbf{LDL}^T = \mathbf{A}$.

Example 3.18 (2×2 symmetric LU)

Let $a, b, c \in \mathbb{R}$. A symmetric LU factorization of $\mathbf{A} := \begin{bmatrix} a & b \\ b & c \end{bmatrix}$ must satisfy the equations

$$\begin{bmatrix} a & b \\ b & c \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_1 & 1 \end{bmatrix} \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix} \begin{bmatrix} 1 & l_1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} d_1 & d_1 l_1 \\ l_1 d_1 & l_1^2 d_1 + d_2 \end{bmatrix}$$

for the unknowns l_1 in \mathbf{L} and d_1, d_2 in \mathbf{D} . The equations are

$$d_1 = a, \quad l_1 a = b, \quad d_2 = c - al_1^2.$$

As in the nonsymmetric case the main problem is the nonlinear equation. Again there are three cases

1. $a \neq 0$: The matrix has a unique symmetric LU factorization with $l_1 = b/a$.
2. $a = 0, b \neq 0$: No symmetric LU factorization exists.
3. $a = b = 0$: The LU factorization exists, but it is not unique. Any value for l_1 can be used.

Consider the four matrices

$$\mathbf{A}_1 := \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_3 := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_4 := \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}.$$

Then the symmetric LU factorization is unique for \mathbf{A}_1 and \mathbf{A}_3 , is not unique for \mathbf{A}_4 and does not exist for \mathbf{A}_2 .

In view of this example it might come as no surprise that Theorem 3.6 carries over to the symmetric case. Again we start with an lemma.

Lemma 3.19 (Symmetric LU of leading principal sub matrices)

Suppose $\mathbf{A} = \mathbf{LDL}^T$ is a symmetric LU factorization of $\mathbf{A} \in \mathbb{R}^{n \times n}$. For $k = 1, \dots, n$ let $\mathbf{A}_{[k]}, \mathbf{L}_{[k]}, \mathbf{D}_{[k]}$ be the leading principal submatrices of $\mathbf{A}, \mathbf{L}, \mathbf{D}$, respectively. Then $\mathbf{A}_{[k]} = \mathbf{L}_{[k]} \mathbf{D}_{[k]} \mathbf{L}_{[k]}^T$ is a symmetric LU factorization of $\mathbf{A}_{[k]}$ for $k = 1, \dots, n$.

Proof. For $k = 1, \dots, n-1$ we partition $\mathbf{A} = \mathbf{LDL}^T$ as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[k]} & \mathbf{B}_k \\ \mathbf{C}_k & \mathbf{F}_k \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{[k]} & \mathbf{0} \\ \mathbf{M}_k & \mathbf{N}_k \end{bmatrix} \begin{bmatrix} \mathbf{D}_{[k]} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_k \end{bmatrix} \begin{bmatrix} \mathbf{L}_{[k]}^T & \mathbf{M}_k^T \\ \mathbf{0} & \mathbf{N}_k^T \end{bmatrix} = \mathbf{LDL}^T, \quad (3.8)$$

where $\mathbf{F}_k, \mathbf{N}_k, \mathbf{E}_k \in \mathbb{R}^{n-k, n-k}$. Block multiplication gives $\mathbf{A}_{[k]} = \mathbf{L}_{[k]} \mathbf{D}_{[k]} \mathbf{L}_{[k]}^T$. Since $\mathbf{L}_{[k]}$ is unit lower triangular and $\mathbf{D}_{[k]}$ is diagonal this is a symmetric LU factorization of $\mathbf{A}_{[k]}$. \square

Theorem 3.20 [Symmetric LU Theorem]

The matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has a unique symmetric LU factorization if and only if $\mathbf{A} = \mathbf{A}^T$ and $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n-1$.

Proof. If \mathbf{A} is nonsingular then \mathbf{D} is nonsingular and it can be shown that the theorem is a simple corollary of the LU theorem. To prove the general case we repeat the proof of Theorem 3.6 incorporating the necessary changes. Suppose $\mathbf{A}^T = \mathbf{A}$ and that $\mathbf{A}_{[k]}$ is nonsingular for $k = 1, \dots, n-1$. Note that $\mathbf{A}_{[k]}^T = \mathbf{A}_{[k]}$ for $k = 1, \dots, n$. We use induction on n to show that \mathbf{A} has a unique symmetric LU factorization. The result is clearly true for $n = 1$, since the unique symmetric LU factorization of a 1-by-1 matrix is $[a_{11}] = [1][a_{11}][1]$. Suppose that $\mathbf{A}_{[n-1]}$ has a unique symmetric LU factorization $\mathbf{A}_{[n-1]} = \mathbf{L}_{n-1} \mathbf{D}_{n-1} \mathbf{L}_{n-1}^T$, and that $\mathbf{A}_{[1]}, \dots, \mathbf{A}_{[n-1]}$ are nonsingular. By block multiplication

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[n-1]} & \mathbf{b} \\ \mathbf{b}^T & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{n-1} & \mathbf{0} \\ \mathbf{x}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{D}_{n-1} & \mathbf{0} \\ 0 & d_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{n-1}^T & \mathbf{x} \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (3.9)$$

if and only if $\mathbf{b} = \mathbf{L}_{n-1}\mathbf{D}_{n-1}\mathbf{x}$ and $a_{nn} = d_{nn} + \mathbf{x}^T\mathbf{D}_{n-1}\mathbf{x}$. Thus we obtain a symmetric LU factorization of \mathbf{A} that is unique since \mathbf{L}_{n-1} and \mathbf{D}_{n-1} are nonsingular.

For the converse we use Lemma 3.19 in the same way as Lemma 3.5 was used to prove Theorem 3.6. \square

3.4 Positive Definite- and Positive Semidefinite Matrices

Symmetric positive definite matrices occur often in scientific computing. In this section we study some properties of positive (semi)definite matrices. We study only real matrices, but consider both the symmetric and nonsymmetric case.

3.4.1 Definition and Examples

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a square matrix. The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

is called a **quadratic form**. We say that \mathbf{A} is

- (i) **positive definite** if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all nonzero $\mathbf{x} \in \mathbb{R}^n$.
- (ii) **positive semidefinite** if $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.
- (iii) **negative (semi)definite** if $-\mathbf{A}$ is positive (semi)definite.
- (iv) **symmetric positive (semi)definite** if \mathbf{A} is symmetric in addition to being positive (semi)definite.
- (v) **symmetric negative (semi)definite** if \mathbf{A} is symmetric in addition to being negative (semi)definite.

We observe the following.

- A matrix is positive definite if it is positive semidefinite and in addition

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}. \quad (3.10)$$

- A positive definite matrix must be nonsingular. Indeed, if $\mathbf{A} \mathbf{x} = \mathbf{0}$ for some $\mathbf{x} \in \mathbb{R}^n$ then $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$ which by (3.10) implies that $\mathbf{x} = \mathbf{0}$.

Since $\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{x}$, the matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive (semi)definite if and only if the symmetric part $(\mathbf{A} + \mathbf{A}^T)/2$ of \mathbf{A} is symmetric positive (semi)definite.

The zero-matrix is symmetric positive semidefinite, while the unit matrix is symmetric positive definite.

Example 3.21 (2×2 positive definite)

The family of matrices

$$\mathbf{A}[a] := \begin{bmatrix} 2 & 2-a \\ a & 1 \end{bmatrix}, \quad a \in \mathbb{R}$$

is positive definite for any $a \in \mathbb{R}$. Indeed for any nonzero $\mathbf{x} = [x_1, x_2]^T \in \mathbb{R}^2$

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 2x_1^2 + (2-a)x_1x_2 + ax_2x_1 + x_2^2 = x_1^2 + (x_1 + x_2)^2 > 0.$$

Lemma 3.22 (\mathbf{T} is symmetric positive definite)

The second derivative matrix $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{n \times n}$ is symmetric positive definite.

Proof. Clearly \mathbf{T} is symmetric. For any $\mathbf{x} \in \mathbb{R}^n$

$$\begin{aligned} \mathbf{x}^T \mathbf{T} \mathbf{x} &= 2 \sum_{i=1}^n x_i^2 - \sum_{i=1}^{n-1} x_i x_{i+1} - \sum_{i=2}^n x_{i-1} x_i \\ &= \sum_{i=1}^{n-1} x_i^2 - 2 \sum_{i=1}^{n-1} x_i x_{i+1} + \sum_{i=1}^{n-1} x_{i+1}^2 + x_1^2 + x_n^2 \\ &= x_1^2 + x_n^2 + \sum_{i=1}^{n-1} (x_{i+1} - x_i)^2. \end{aligned}$$

Thus $\mathbf{x}^T \mathbf{T} \mathbf{x} \geq 0$ and if $\mathbf{x}^T \mathbf{T} \mathbf{x} = 0$ then $x_1 = x_n = 0$ and $x_i = x_{i+1}$ for $i = 1, \dots, n-1$ which implies that $\mathbf{x} = 0$. Hence \mathbf{T} is positive definite. \square

Symmetric positive definite matrices are important in nonlinear optimization.

Example 3.23 (Gradient and Hessian)

Consider (cf. (C.1)) the gradient ∇f and hessian $Hf := \nabla \nabla^T f$ of a function $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n, \quad \nabla \nabla^T f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

We assume that f has continuous first and second partial derivatives on Ω .

Under suitable conditions on the domain Ω it is shown in advanced calculus texts that if $\nabla f(\mathbf{x}) = \mathbf{0}$ and $\nabla \nabla^T f(\mathbf{x})$ is symmetric positive definite then \mathbf{x} is a local minimum for f . This can be shown using the second-order Taylor expansion (C.2). Moreover, \mathbf{x} is a local maximum if $\nabla f(\mathbf{x}) = \mathbf{0}$ and $\nabla \nabla^T f(\mathbf{x})$ is negative definite.

3.4.2 Principal Submatrices

Theorem 3.24 (A general criterium)

Let m, n be positive integers. If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is (symmetric) positive semidefinite then $\mathbf{B} := \mathbf{X}^T \mathbf{A} \mathbf{X} \in \mathbb{R}^{m \times m}$ is (symmetric) positive semidefinite for any $\mathbf{X} \in \mathbb{R}^{n \times m}$. If in addition \mathbf{A} is (symmetric) positive definite and \mathbf{X} has linearly independent columns then \mathbf{B} is (symmetric) positive definite.

Proof. Let $\mathbf{y} \in \mathbb{R}^m$ and set $\mathbf{x} := \mathbf{X}\mathbf{y} \in \mathbb{R}^n$. Then $\mathbf{y}^T \mathbf{B} \mathbf{y} = \mathbf{y}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{y} = \mathbf{x}^T \mathbf{A} \mathbf{x}$. This is nonnegative if \mathbf{A} is positive semidefinite and positive if \mathbf{A} is positive definite and \mathbf{X} has linearly independent columns. For then \mathbf{x} is nonzero if \mathbf{y} is nonzero. If \mathbf{A} is symmetric then \mathbf{B} is symmetric and the statements about symmetry follows. \square

Theorem 3.25 (Principal submatrices)

Suppose $\mathbf{A} \in \mathbb{R}^{n,n}$ is a positive (semi) definite matrix and let \mathbf{B} be a principal submatrix. Then

1. \mathbf{B} is positive (semi) definite,
2. the real eigenvalues of \mathbf{B} are positive (nonnegative),
3. $\det(\mathbf{B}) > 0 (\geq 0)$.

Proof.

1. Suppose the submatrix $\mathbf{B} = \mathbf{A}(\mathbf{r}, \mathbf{r})$ is defined by the rows and columns $\mathbf{r} = [r_1, \dots, r_k]^T$ of \mathbf{A} . Let $\mathbf{X} = [\mathbf{e}_{r_1}, \dots, \mathbf{e}_{r_k}] \in \mathbb{R}^{n \times k}$. Then $\mathbf{B} := \mathbf{X}^T \mathbf{A} \mathbf{X}$, and \mathbf{B} is positive (semi) definite by Theorem 3.24.
2. Suppose (λ, \mathbf{x}) is an eigenpair of \mathbf{A} and that λ is real. Since \mathbf{A} is real we can choose \mathbf{x} to be real. Multiplying $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ by \mathbf{x}^T and solving for λ we find $\lambda = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} > 0 (\geq 0)$.
3. The determinant of \mathbf{B} equals the product of its eigenvalues. The eigenvalues are either real and positive or occur in complex conjugate pairs. The product of two nonzero complex conjugate numbers is positive.

□

A nonsymmetric positive definite matrix can have complex eigenvalues. For example, the eigenvalues of $\mathbf{A}[a]$ in Example 3.21 are positive for $a \in [1 - \frac{\sqrt{5}}{2}, 1 + \frac{\sqrt{5}}{2}]$ and complex for other values of \mathbf{A} .

From Theorem 3.25 and the LU Theorem 3.6 we obtain

Corollary 3.26 (positive (semi)definite criteria)

Suppose \mathbf{A} is a positive (semi)definite matrix. Then

1. Any principal submatrix is positive (semi)definite.
2. A positive definite matrix has a unique LU factorization.
3. Real eigenvalues of \mathbf{A} are positive (nonnegative).
4. $\det(\mathbf{A}) > 0$ ($\det(\mathbf{A}) \geq 0$).

3.4.3 The Symmetric Positive Definite Case

Taking $\mathbf{A} := \mathbf{I}$ and $\mathbf{X} := \mathbf{A}$ in Theorem 3.24 gives $\mathbf{B} = \mathbf{A}^T \mathbf{I} \mathbf{A} = \mathbf{A}^T \mathbf{A}$, and we obtain

Corollary 3.27 ($\mathbf{A}^T \mathbf{A}$ is symmetric positive semidefinite)

Let m, n be positive integers. If $\mathbf{A} \in \mathbb{R}^{m,n}$ then $\mathbf{A}^T \mathbf{A}$ is symmetric positive semidefinite. If in addition \mathbf{A} has linearly independent columns then $\mathbf{A}^T \mathbf{A}$ is symmetric positive definite.

Consider the eigenvalues of a real symmetric positive definite matrix. Note that such a matrix is Hermitian.

Lemma 3.28 (Eigenvalues of a Hermitian matrix)

All eigenvalues of a Hermitian matrix are real.

Proof. Suppose $\mathbf{A}^* = \mathbf{A}$ and $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ with $\mathbf{x} \neq 0$. We multiply both sides of $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ from the left by \mathbf{x}^* and divide by $\mathbf{x}^*\mathbf{x}$ to obtain $\lambda = \frac{\mathbf{x}^*\mathbf{A}\mathbf{x}}{\mathbf{x}^*\mathbf{x}}$. Taking complex conjugates we find $\bar{\lambda} = \lambda^* = \frac{(\mathbf{x}^*\mathbf{A}\mathbf{x})^*}{(\mathbf{x}^*\mathbf{x})^*} = \frac{\mathbf{x}^*\mathbf{A}^*\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \frac{\mathbf{x}^*\mathbf{A}\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \lambda$, and λ is real. □

Lemma 3.29 (Symmetry and positive eigenvalues)

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive definite if and only if it is symmetric and all eigenvalues are positive.

Proof. By Lemma 3.28 all eigenvalues of \mathbf{A} are real, and by Theorem 3.26 all eigenvalues are positive. Suppose conversely that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric with positive eigenvalues $\lambda_1, \dots, \lambda_n$. By the spectral theorem (cf. Theorem 6.39) we have $\mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{D}$, where $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}$ and $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$. Let $\mathbf{x} \in \mathbb{R}^n$ be nonzero and define $\mathbf{c} := \mathbf{U}^T \mathbf{x} = (c_1, \dots, c_n)^T$. Then $\mathbf{c}^T \mathbf{c} = \mathbf{x}^T \mathbf{U} \mathbf{U}^T \mathbf{x} = \mathbf{x}^T \mathbf{x}$ so \mathbf{c} is nonzero. Since $\mathbf{x} = \mathbf{U} \mathbf{c}$ we find

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{U} \mathbf{c})^T \mathbf{A} \mathbf{U} \mathbf{c} = \mathbf{c}^T \mathbf{U}^T \mathbf{A} \mathbf{U} \mathbf{c} = \mathbf{c}^T \mathbf{D} \mathbf{c} = \sum_{j=1}^n \lambda_j c_j^2 > 0$$

and it follows that \mathbf{A} is positive definite. \square

Lemma 3.30 (Symmetric positive definite and symmetric LU)

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive definite if and only if it has a symmetric LU factorization $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^T$ with positive diagonal elements in \mathbf{D} .

Proof. Suppose \mathbf{A} is symmetric positive definite. By Theorem 3.26 the leading principal submatrices $\mathbf{A}_{[k]} \in \mathbb{R}^{k \times k}$ are nonsingular for $k = 1, \dots, n - 1$, and \mathbf{A} has a unique symmetric LU factorization $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^T$ by Theorem 3.20. The i th diagonal element in \mathbf{D} is positive, $d_{ii} = \mathbf{e}_i^T \mathbf{D} \mathbf{e}_i = \mathbf{e}_i^T \mathbf{L}^{-1} \mathbf{A} \mathbf{L}^{-T} \mathbf{e}_i = \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i > 0$. Indeed, \mathbf{L}^{-T} is nonsingular so $\mathbf{x}_i := \mathbf{L}^{-T} \mathbf{e}_i$ is nonzero.

Conversely, suppose \mathbf{A} has a symmetric LU factorization $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^T$ with positive diagonal elements in \mathbf{D} . Then \mathbf{A} is symmetric and for any nonzero $\mathbf{y} \in \mathbb{R}^n$ we have $\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{L} \mathbf{D} \mathbf{L}^T \mathbf{x} = \mathbf{y}^T \mathbf{A} \mathbf{y} > 0$, since $\mathbf{y} := \mathbf{L}^T \mathbf{x} \neq \mathbf{0}$. \square

The following characterizations can be used to decide if a matrix is symmetric positive definite. Recall that a **leading principal minor** is the determinant of a leading principal submatrix.

Theorem 3.31 (Symmetric positive definite characterization)

The following statements are equivalent for a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$.

1. \mathbf{A} is symmetric positive definite.
2. \mathbf{A} has only positive eigenvalues.
3. All leading principal minors are positive.
4. $\mathbf{A} = \mathbf{B} \mathbf{B}^T$ for a nonsingular $\mathbf{B} \in \mathbb{R}^{n \times n}$.

Proof. 1 \Leftrightarrow 2 was shown in Lemma 3.29. We show that 1 \Rightarrow 3 \Rightarrow 4 \Rightarrow 1.

1 \Rightarrow 3: This follows from Theorem 3.25

3 \Rightarrow 4: Since all principal minors of \mathbf{A} are positive the principal submatrices

$A_{[k]}$ are nonsingular for all k and therefore \mathbf{A} has a symmetric LU factorization $\mathbf{A} = \mathbf{LDL}^T$ with positive diagonal elements in \mathbf{D} . But then $\mathbf{A} = \mathbf{BB}^T$, where $\mathbf{B} := \mathbf{LD}^{1/2}$, with $\mathbf{D}^{1/2} := \text{diag}(d_{1,1}^{1/2}, \dots, d_{n \times n}^{1/2})$.

4 \Rightarrow 1: This follows from Corollary 3.27. \square

3.5 The Cholesky Factorization

From Lemma 3.30 it follows that if \mathbf{A} is symmetric positive definite if and only if it has a symmetric LU factorization, and from the proof of 3. implies 4 in that theorem we can write this in the form $\mathbf{A} = \mathbf{BB}^T$ where \mathbf{B} is lower triangular matrix with positive diagonal elements. Such a factorization has a special name.

Definition 3.32 (Cholesky)

A factorization $\mathbf{A} = \mathbf{LL}^T$ where \mathbf{L} is lower triangular with positive diagonal elements is called a **Cholesky factorization** of \mathbf{A} . The matrix \mathbf{L} is called a **Cholesky factor**.

From the discussion before the definition we have

Theorem 3.33 (Cholesky)

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has a Cholesky factorization if and only if it is symmetric positive definite. Moreover, the Cholesky factorization is unique.

Proof. We still need to show uniqueness. Suppose $\mathbf{LL}^T = \mathbf{SS}^T$ are two Cholesky factorizations of the symmetric positive definite matrix \mathbf{A} . Since \mathbf{A} is nonsingular both \mathbf{L} and \mathbf{S} are nonsingular. Then $\mathbf{S}^{-1}\mathbf{L} = \mathbf{S}^T\mathbf{L}^{-T}$, where by Lemma 1.9 $\mathbf{S}^{-1}\mathbf{L}$ is lower triangular and $\mathbf{S}^T\mathbf{L}^{-T}$ is upper triangular, with diagonal elements ℓ_{ii}/s_{ii} and s_{ii}/ℓ_{ii} , respectively. But then both matrices must be equal to the same diagonal matrix and $\ell_{ii}^2 = s_{ii}^2$. By positivity $\ell_{ii} = s_{ii}$ and we conclude that $\mathbf{S}^{-1}\mathbf{L} = \mathbf{I} = \mathbf{S}^T\mathbf{L}^{-T}$ which means that $\mathbf{L} = \mathbf{S}$. \square

A Cholesky factorization can also be written in the equivalent form $\mathbf{A} = \mathbf{R}^T\mathbf{R}$, where $\mathbf{R} = \mathbf{L}^T$ is upper triangular with positive diagonal elements. The matrix \mathbf{A} must be symmetric since \mathbf{LL}^T is symmetric.

Example 3.34 (2×2)

The matrix $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$ has a symmetric LU- and a Cholesky-factorization given by

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -\frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 1 & -\frac{1}{2} \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sqrt{2} & 0 \\ -1/\sqrt{2} & \sqrt{3/2} \end{bmatrix} \begin{bmatrix} \sqrt{2} & -1/\sqrt{2} \\ 0 & \sqrt{3/2} \end{bmatrix}.$$

Consider computing the Cholesky factorization directly. The equation $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ implies that

$$a_{ik} = \sum_{j=1}^n \ell_{ij} \ell_{kj} = \sum_{j=1}^{\min(i,k)} \ell_{ij} \ell_{kj}, \quad i, k = 1, \dots, n. \quad (3.11)$$

Suppose we have computed the $k - 1$ first columns of \mathbf{L} . The k th column can then be computed from (3.11). Indeed, letting $i = k$ and solving for ℓ_{kk} we find

$$\ell_{kk} = (a_{kk} - \sum_{j=1}^{k-1} \ell_{kj}^2)^{1/2}, \quad (3.12)$$

and similarly for $i > k$

$$\ell_{ik} = (a_{ik} - \sum_{j=1}^{k-1} \ell_{ij} \ell_{kj}) / \ell_{kk} \quad i = k + 1, \dots, n. \quad (3.13)$$

Since \mathbf{A} is symmetric positive definite the Cholesky factor \mathbf{L} exists, is unique, and real, and therefore the term under the square root in (3.12) must be positive. We note however that we can encounter problems in floating point computation if the term is very small.

It is easily seen that the Cholesky-factorization of an n -by- n matrix based on (3.12) and (3.13) requires $\frac{1}{2}G_n = n^3/3$ arithmetic operations. The halving of the count compared to Gaussian elimination is due to the symmetry of \mathbf{A} .

If \mathbf{A} is d -banded then the same is true for the Cholesky factor.

Lemma 3.35 (Banded Cholesky factor)

Suppose \mathbf{A} is symmetric positive definite with Cholesky-factor \mathbf{L} . If $a_{ik} = 0$ for $i > k + d$, then $\ell_{ik} = 0$ for $i > k + d$.

Proof. We show that if \mathbf{L} has bandwidth d in its first $k - 1$ columns then column k also has bandwidth d . The proof then follows by induction on k . Now, if $i > k + d$, then $a_{ik} = 0$, and if \mathbf{L} has bandwidth d in its first $k - 1$ columns then $\ell_{ij} = 0$ for $j = 1, \dots, k - 1$. By (3.13) $\ell_{ik} = 0$. \square

We obtain formulas for the Cholesky factorization of a band matrix by simply replacing the lower bound $j = 1$ by $j = \max(1, k - d)$ in (3.12) and (3.13) and letting i run from $k + 1$ to $\min(n, k + d)$ in (3.13). This leads to the following algorithm.

Algorithm 3.36 (bandcholesky)

Suppose \mathbf{A} is symmetric positive definite. An lower triangular matrix \mathbf{L} is computed in sparse form so that $\mathbf{A} = \mathbf{LL}^T$. Only the lower triangular part of \mathbf{A} is used.

```

1 function L=bandcholesky(A,d)
2 %L=bandcholesky(A,d)
3 n=length(A);
4 L=sparse(zeros(n,n));
5 for k=1:n
6     km=max(1,k-d); kp=min(n,k+d); s=L(k,km:k-1);
7     L(k,k)=sqrt(A(k,k)-s*s');
8     L(k+1:kp,k)=(A(k+1:kp,k) - ...
9         L(k+1:kp,km:k-1)*s')/L(k,k);
10 end

```

For a different algorithm based on outer products, which also works for symmetric positive semidefinite matrices, see Algorithm 3.42.

The leading term in an operation count for a band matrix is $O(d^2n)$. When d is small this is a considerable saving compared to the count $\frac{1}{2}G_n = n^3/3$ for a full matrix.

There is also a banded version of the symmetric LU factorization which requires approximately the same number of arithmetic operations as the Cholesky factorization. The choice between using a symmetric LU factorization or an \mathbf{LL}^T factorization depends on several factors. Usually an LU or a symmetric LU factorization is preferred for matrices with small bandwidth (tridiagonal, pentadiagonal), while the \mathbf{LL}^T factorization is restricted to symmetric positive definite matrices and is often used when the bandwidth is larger.

3.6 The Symmetric Positive Semidefinite Case

We start with the following necessary conditions for a matrix to be symmetric positive semidefinite. It shows that if a diagonal element a_{ii} is zero then all elements in row i and column i are zero.

Lemma 3.37 (Criteria symmetric semidefinite)

If \mathbf{A} is symmetric positive semidefinite then for all i, j

1. $|a_{ij}| \leq (a_{ii} + a_{jj})/2,$
2. $|a_{ij}| \leq \sqrt{a_{ii}a_{jj}},$
3. $\max_{i,j} |a_{ij}| = \max_i a_{ii},$
4. $a_{ii} = 0 \implies a_{ij} = a_{ji} = 0, \text{ fixed } i, \text{ all } j.$

Proof. 3. follows from 1. and 4. from 2. Now

$$0 \leq (\alpha \mathbf{e}_i + \beta \mathbf{e}_j)^T \mathbf{A} (\alpha \mathbf{e}_i + \beta \mathbf{e}_j) = \alpha^2 a_{ii} + \beta^2 a_{jj} + 2\alpha\beta a_{ij}, \text{ all } i, j, \alpha, \beta \in \mathbb{R}, \quad (3.14)$$

since \mathbf{A} is symmetric positive semidefinite. Taking $\alpha = 1, \beta = \pm 1$ we obtain $a_{ii} + a_{jj} \pm 2a_{ij} \geq 0$ and this implies 1. 2. follows trivially from 1. if $a_{ii} = a_{jj} = 0$. Suppose one of them, say a_{ii} is nonzero. Note that $a_{ii} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i > 0$. Taking $\alpha = -a_{ij}, \beta = a_{ii}$ in (3.14) we find

$$0 \leq a_{ij}^2 a_{ii} + a_{ii}^2 a_{jj} - 2a_{ij}^2 a_{ii} = a_{ii}(a_{ii}a_{jj} - a_{ij}^2).$$

But then $a_{ii}a_{jj} - a_{ij}^2 \geq 0$ and 2. follows. \square

As an illustration consider the matrices

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 1 & 2 \\ 2 & 2 \end{bmatrix}, \quad \mathbf{A}_3 = \begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}.$$

None of them is positive semidefinite, since neither 1. nor 2. hold.

A symmetric positive semidefinite matrix has a factorization that is similar to the Cholesky factorization.

Definition 3.38 (Semi-Cholesky factorization)

A factorization $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ where \mathbf{L} is lower triangular with nonnegative diagonal elements is called a **semi-Cholesky factorization**.

Note that a semi-Cholesky factorization of a symmetric positive definite matrix is necessarily a Cholesky factorization. For if \mathbf{A} is positive definite then it is nonsingular and then \mathbf{L} must be nonsingular. Thus the diagonal elements of \mathbf{L} cannot be zero.

Theorem 3.39 (Characterization, semi-Cholesky factorization)

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has a semi-Cholesky factorization $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ if and only if it is symmetric positive semidefinite.

Proof. If $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ is a semi-Cholesky factorization then \mathbf{A} is symmetric positive semidefinite by Corollary 3.27. For the converse we use induction on n . A positive semidefinite matrix of order one has a semi-Cholesky factorization since the one and only element in \mathbf{A} is nonnegative. Suppose any symmetric positive semidefinite matrix of order $n - 1$ has a semi-Cholesky factorization and suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive semidefinite. We partition \mathbf{A} as follows

$$\mathbf{A} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix}, \quad \alpha \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^{n-1}, \mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}. \quad (3.15)$$

There are two cases. Suppose first $\alpha = e_1^T \mathbf{A} e_1 > 0$. We claim that $\mathbf{C} := \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$ is symmetric positive semidefinite. \mathbf{C} is symmetric. To show that \mathbf{C} is positive semidefinite we consider any $\mathbf{y} \in \mathbb{R}^{n-1}$ and define $\mathbf{x}^T := [-\mathbf{v}^T \mathbf{y}/\alpha, \mathbf{y}^T] \in \mathbb{R}^n$. Then

$$\begin{aligned} 0 \leq \mathbf{x}^T \mathbf{A} \mathbf{x} &= [-\mathbf{v}^T \mathbf{y}/\alpha, \mathbf{y}^T] \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} \begin{bmatrix} -\mathbf{v}^T \mathbf{y}/\alpha \\ \mathbf{y} \end{bmatrix} \\ &= [0, -(\mathbf{v}^T \mathbf{y})\mathbf{v}^T/\alpha + \mathbf{y}^T \mathbf{B}] \begin{bmatrix} -\mathbf{v}^T \mathbf{y}/\alpha \\ \mathbf{y} \end{bmatrix} \\ &= -(\mathbf{v}^T \mathbf{y})(\mathbf{v}^T \mathbf{y})/\alpha + \mathbf{y}^T \mathbf{B} \mathbf{y} = \mathbf{y}^T \mathbf{C} \mathbf{y}, \end{aligned} \quad (3.16)$$

since $(\mathbf{v}^T \mathbf{y})\mathbf{v}^T \mathbf{y} = (\mathbf{v}^T \mathbf{y})^T \mathbf{v}^T \mathbf{y} = \mathbf{y}^T \mathbf{v} \mathbf{v}^T \mathbf{y}$. So $\mathbf{C} \in \mathbb{R}^{(n-1) \times (n-1)}$ is symmetric positive semidefinite and by the induction hypothesis it has a semi-Cholesky factorization $\mathbf{C} = \mathbf{L}_1 \mathbf{L}_1^T$. The matrix

$$\mathbf{L}^T := \begin{bmatrix} \beta & \mathbf{v}^T/\beta \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix}, \quad \beta := \sqrt{\alpha}, \quad (3.17)$$

is upper triangular with nonnegative diagonal elements and

$$\mathbf{L} \mathbf{L}^T = \begin{bmatrix} \beta & \mathbf{0} \\ \mathbf{v}/\beta & \mathbf{L}_1 \end{bmatrix} \begin{bmatrix} \beta & \mathbf{v}^T/\beta \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} = \mathbf{A}$$

is a semi-Cholesky factorization of \mathbf{A} .

If $\alpha = 0$ then it follows from 4. in Lemma 3.37 that $\mathbf{v} = \mathbf{0}$. Moreover, $\mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}$ in (3.15) is positive semidefinite and therefore has a semi-Cholesky factorization $\mathbf{B} = \mathbf{L}_1 \mathbf{L}_1^T$. But then $\mathbf{L} \mathbf{L}^T$, where $\mathbf{L} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1 \end{bmatrix}$ is a semi-Cholesky factorization of \mathbf{A} . Indeed, \mathbf{L} is lower triangular and

$$\mathbf{L} \mathbf{L}^T = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1 \end{bmatrix} \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B} \end{bmatrix} = \mathbf{A}.$$

□

Theorem 3.40 (Positive symmetric semidefinite characterization)
The following is equivalent for a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$.

1. \mathbf{A} is positive semidefinite.
2. \mathbf{A} has only nonnegative eigenvalues.
3. $\mathbf{A} = \mathbf{B}\mathbf{B}^T$ for some $\mathbf{B} \in \mathbb{R}^{n \times n}$.
4. All principal minors are nonnegative.

Proof. The proof of $1. \Leftrightarrow 2$ follows as in the proof of Theorem 3.31. $1. \Leftrightarrow 3$. follows from Theorem 3.39 while $1. \Rightarrow 4$. is a consequence of Theorem 3.26. To prove $4. \Rightarrow 1$. one first shows that $\epsilon\mathbf{I} + \mathbf{A}$ is symmetric positive definite for all $\epsilon > 0$ (Cf. page 567 of [20]). But then $\mathbf{x}^T \mathbf{A} \mathbf{x} = \lim_{\epsilon \rightarrow 0} \mathbf{x}^T (\epsilon\mathbf{I} + \mathbf{A}) \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$. \square

In 4. of Theorem 3.40 we require nonnegativity of all principal minors, while only positivity of leading principal minors was required for positive definite matrices (cf. Theorem 3.31). To see that nonnegativity of the leading principal minors is not enough consider the matrix $\mathbf{A} := \begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix}$. The leading principal minors are nonnegative, but \mathbf{A} is not positive semidefinite.

3.6.1 An Algorithm for SemiCholesky Factorization of a Banded Matrix

Recall that a matrix \mathbf{A} is d -banded if $a_{ij} = 0$ for $|i - j| > d$. A (semi)Cholesky factorization preserves bandwidth.

Theorem 3.41 (Bandwidth semi Cholesy factor)

The Cholesky factor \mathbf{L} given by (3.17) has the same bandwidth as \mathbf{A} .

Proof. Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ is d -banded. Then $\mathbf{v}^T = [\mathbf{u}^T, \mathbf{0}^T]$ in (3.15), where $\mathbf{u} \in \mathbb{R}^d$, and therefore $\mathbf{C} := \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$ differs from \mathbf{B} only in the upper left $d \times d$ corner. It follows that \mathbf{C} has the same bandwidth as \mathbf{B} and \mathbf{A} . By induction on n , $\mathbf{C} = \mathbf{L}_1 \mathbf{L}_1^T$, where \mathbf{L}_1^T has the same bandwidth as \mathbf{C} . But then \mathbf{L} in (3.17) has the same bandwidth as \mathbf{A} . \square

Consider now implementing an algorithm based on the previous discussion. Since \mathbf{A} is symmetric we only need to use the lower part of \mathbf{A} . The first column of \mathbf{L} is $[\beta, \mathbf{v}^T/\beta]^T$ if $\alpha > 0$. If $\alpha = 0$ then by 4 in Lemma 3.37 the first column of \mathbf{A} is zero and this is also the first column of \mathbf{L} . We obtain

$$\begin{aligned}
 &\text{if } A(1,1) > 0 \\
 &A(1,1) = \sqrt{A(1,1)} \\
 &A(2:n,1) = A(2:n,1)/A(1,1) \\
 &\text{for } j = 2:n \\
 &A(j:n,j) = A(j:n,j) - A(j,1) * A(j:n,1)
 \end{aligned}
 \tag{3.18}$$

Here we store the first column of \mathbf{L} in the first column of \mathbf{A} and the lower part of $\mathbf{C} = \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$ in the lower part of $A(2:n, 2:n)$.

The code can be made more efficient when \mathbf{A} is a d -banded matrix. We simply replace all occurrences of n by $\min(i + d, n)$.

Algorithm 3.42 (bandsemicholesky)

Suppose \mathbf{A} is symmetric positive semidefinite. A lower triangular matrix \mathbf{L} is computed so that $\mathbf{A} = \mathbf{LL}^T$. This is the Cholesky factorization of \mathbf{A} if \mathbf{A} is symmetric positive definite and a semi-Cholesky factorization of \mathbf{A} otherwise. The algorithm uses the Matlab command `tril`.

```

1 function L=bandsemicholesky(A,d)
2 n=length(A);
3 for k=1:n
4     if A(k,k)>0
5         kp=min(n,k+d);
6         A(k,k)=sqrt(A(k,k));
7         A(k+1:kp,k)=A(k+1:kp,k)/A(k,k);
8         for j=k+1:kp
9             A(j:kp,j)=A(j:kp,j)-A(j,k)*A(j:kp,k);
10        end
11    else
12        A(k:kp,k)=zeros(kp-k+1,1);
13    end
14 end
15 L=tril(A);

```

Figure 3.3: An algorithm for semi Cholesky factorization, see text.

Continuing the reduction we arrive at Algorithm 3.42 in Figure 3.3.

In the algorithm we overwrite the lower triangle of \mathbf{A} with the elements of \mathbf{L} . Column k of \mathbf{L} is zero for those k where $\ell_{kk} = 0$. We reduce round-off noise by forcing those rows to be zero. In the semidefinite case no update is necessary and we “do nothing”.

There are many versions of Cholesky factorizations, see [4]. Algorithm 3.36 is based on outer products \mathbf{vv}^T . An advantage of this formulation is that it can be extended to symmetric positive semidefinite matrices. However deciding when a diagonal element is zero is a problem in floating point arithmetic.

3.7 Review Questions

3.8.1 What is the content of

- the LU theorem?
- the block LU theorem?
- the symmetric LU theorem?

3.8.2 Is $\mathbf{A}^T \mathbf{A}$ symmetric positive definite?

- 3.8.3** • What class of matrices has a Cholesky factorization?
• what is the bandwidth of the Cholesky factor of a band matrix?
- 3.8.4** For a symmetric matrix give 3 conditions that are equivalent to positive definiteness.
- 3.8.5** What class of matrices has a semi-Cholesky factorization?

Chapter 4

The Kronecker Product

Matrices arising from 2D and 3D problems sometimes have a Kronecker product structure. Identifying a Kronecker structure can be very rewarding since it simplifies the study of such matrices.

4.1 Test Matrices

In this section we introduce some matrices which we will use to compare various algorithms in later chapters.

4.1.1 The 2D Poisson Problem

Let

$$\Omega := (0, 1)^2 = \{(x, y) : 0 < x, y < 1\}$$

be the open unit square with boundary $\partial\Omega$. Consider the problem

$$-\nabla^2 u := -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f \text{ on } \Omega, \quad (4.1)$$

$$u := 0 \text{ on } \partial\Omega.$$

Here the function f is given and continuous on Ω , and we seek a function $u = u(x, y)$ such that (4.1) holds and which is zero on $\partial\Omega$.

Let m be a positive integer. We solve the problem numerically by finding approximations $v_{j,k} \approx u(jh, kh)$ on a grid of points given by

$$\bar{\Omega}_h := \{(jh, kh) : j, k = 0, 1, \dots, m+1\}, \quad \text{where } h = 1/(m+1).$$

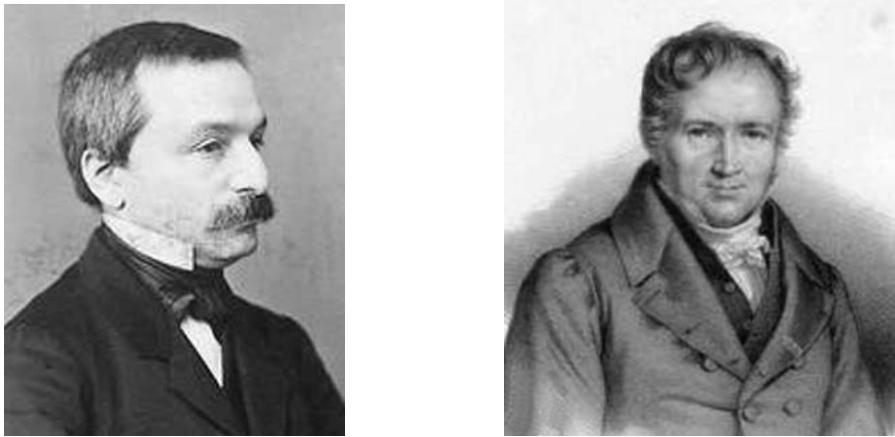


Figure 4.1: Leopold Kronecker, 1823-1891 (left), Siméon Denis Poisson, 1781-1840 (right).

The points $\Omega_h := \{(jh, kh) : j, k = 1, \dots, m\}$ are the interior points, while $\overline{\Omega}_h \setminus \Omega_h$ are the boundary points. The solution is zero at the boundary points. Using the difference approximation from Chapter 2 for the second derivative we obtain the following approximations for the partial derivatives

$$\frac{\partial^2 u(jh, kh)}{\partial x^2} \approx \frac{v_{j-1,k} - 2v_{j,k} + v_{j+1,k}}{h^2}, \quad \frac{\partial^2 u(jh, kh)}{\partial y^2} \approx \frac{v_{j,k-1} - 2v_{j,k} + v_{j,k+1}}{h^2}.$$

Inserting this in (4.1) and multiplying both sides by h^2 to obtain

$$(-v_{j-1,k} + 2v_{j,k} - v_{j+1,k}) + (-v_{j,k-1} + 2v_{j,k} - v_{j,k+1}) = h^2 f_{j,k} \quad (4.2)$$

or

$$4v_{j,k} - v_{j-1,k} - v_{j+1,k} - v_{j,k-1} - v_{j,k+1} = h^2 f_{j,k} := h^2 f(jh, kh). \quad (4.3)$$

From the boundary conditions we have in addition

$$v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0, \quad j, k = 0, 1, \dots, m+1. \quad (4.4)$$

The equations (4.3) and (4.4) define a set of linear equations for the unknowns $\mathbf{V} = [v_{jk}] \in \mathbb{R}^{m \times m}$.

Observe that (4.2) can be written as a matrix equation in the form

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2 \mathbf{F} \quad \text{with} \quad h = 1/(m+1), \quad (4.5)$$

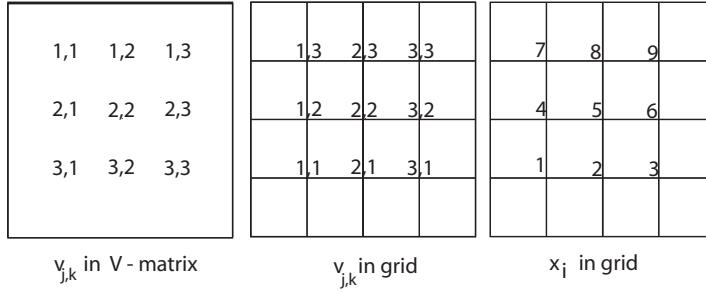


Figure 4.2: Numbering of grid points

where $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ is the second derivative matrix given by (2.2), $\mathbf{V} = (v_{jk}) \in \mathbb{R}^{m \times m}$, and $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m \times m}$. Indeed, the (j, k) element in $\mathbf{TV} + \mathbf{VT}$ is given by

$$\sum_{i=1}^m \mathbf{T}_{j,i} v_{i,k} + \sum_{i=1}^m v_{j,i} \mathbf{T}_{i,k},$$

and this is precisely the left hand side of (4.2).

To write (4.3) and (4.4) in standard form $\mathbf{Ax} = \mathbf{b}$ we need to order the unknowns $v_{j,k}$ in some way. The following operation of **vectorization** of a matrix gives one possible ordering.

Definition 4.1 (vec operation)

For any $\mathbf{B} \in \mathbb{R}^{m \times n}$ we define the vector

$$\text{vec}(\mathbf{B}) := [b_{11}, \dots, b_{m1}, b_{12}, \dots, b_{m2}, \dots, b_{1n}, \dots, b_{mn}]^T \in \mathbb{R}^{mn}$$

by stacking the columns of \mathbf{B} on top of each other.

Let $n = m^2$ and $\mathbf{x} := \text{vec}(\mathbf{V}) \in \mathbb{R}^n$. Note that forming \mathbf{x} by stacking the columns of \mathbf{V} on top of each other means an ordering of the grid points which for $m = 3$ is illustrated in Figure 4.2. We call this the **natural ordering**. The elements in (4.3) form a 5-point stencil, as shown in Figure 4.3.

To find the matrix \mathbf{A} we note that for values of j, k where the 5-point stencil does not touch the boundary, (4.3) takes the form

$$4x_i - x_{i-1} - x_{i+1} - x_{i-m} - x_{i+m} = b_i,$$

where $x_i = v_{jk}$ and $b_i = h^2 f_{jk}$. This must be modified close to the boundary. We obtain the linear system

$$\mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}, \quad \mathbf{b} \in \mathbb{R}^n, \quad n = m^2, \quad (4.6)$$

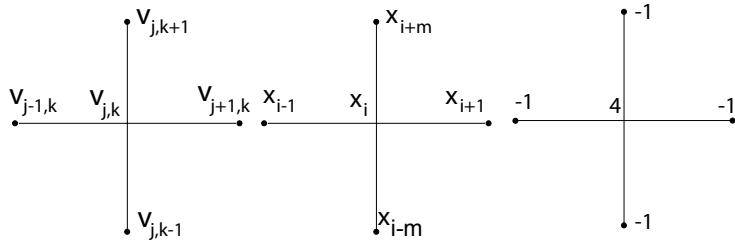


Figure 4.3: The 5-point stencil

where $\mathbf{x} = \text{vec}(\mathbf{V})$, $\mathbf{b} = h^2 \text{vec}(\mathbf{F})$ with $\mathbf{F} = (f_{jk}) \in \mathbb{R}^{m \times m}$, and \mathbf{A} is the **Poisson matrix** given by

$$\begin{aligned} a_{ii} &= 4, & i &= 1, \dots, n, \\ a_{i+1,i} = a_{i,i+1} &= -1, & i &= 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m, \\ a_{i+m,i} = a_{i,i+m} &= -1, & i &= 1, \dots, n-m, \\ a_{ij} &= 0, & & \text{otherwise.} \end{aligned} \quad (4.7)$$

For $m = 3$ we have the following matrix

$$\mathbf{A} = \begin{bmatrix} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{bmatrix}.$$

Exercise 4.2 (2 × 2 Poisson matrix)

Write down the Poisson matrix for $m = 2$ and show that it is strictly diagonally dominant.

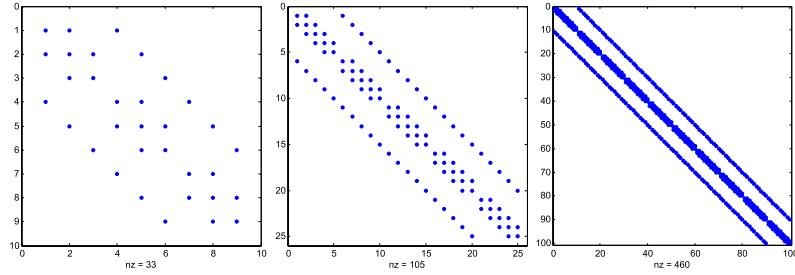


Figure 4.4: Band structure of the 2D test matrix, $n = 9$, $n = 25$, $n = 100$

4.1.2 The Test Matrices

The second derivative matrix $\mathbf{T} = \text{tridiag}(-1, 2, -1)$ is a special case of the tridiagonal matrix

$$\mathbf{T}_1 := \begin{bmatrix} d & a & 0 & & & \\ a & d & a & & & \\ 0 & \ddots & \ddots & \ddots & & \\ & & & & 0 & \\ & & & a & d & a \\ 0 & a & d & & & \end{bmatrix}, \quad (4.8)$$

where $a, d \in \mathbb{R}$. We call this the **1D test matrix**. It is symmetric and strictly diagonally dominant if $|d| > 2|a|$.

The (2-dimensional) Poisson matrix is a special case of the matrix $\mathbf{T}_2 = [a_{ij}] \in \mathbb{R}^{n \times n}$ with elements

$$\begin{aligned} a_{ii} &= 2d, \quad i = 1, \dots, n, \\ a_{i,i+1} = a_{i+1,i} &= a, \quad i = 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m, \\ a_{i,i+m} = a_{i+m,i} &= a, \quad i = 1, \dots, n-m, \\ a_{ij} &= 0, \quad \text{otherwise,} \end{aligned} \quad (4.9)$$

and where a, d are real numbers. We will refer to this matrix as simply the **2D test matrix**. The 2D test matrix is

- symmetric,
- a band matrix with bandwidth $m = \sqrt{n}$ (cf. Figure 4.4),
- strictly diagonally dominant if $|d| > 2|a|$,

- the Poisson matrix given by (4.7) when $a = -1$ and $d = 2$. This matrix is strictly diagonally dominant for $m = 2, n = 4$, but only diagonally dominant for $m > 2$.
- called the **averaging matrix** when $a = 1/9$ and $d = 5/18$. This matrix is strictly diagonally dominant for all n .

Properties of \mathbf{T}_2 can be derived from properties of \mathbf{T}_1 by using properties of the Kronecker product.

4.2 The Kronecker Product

Definition 4.3 (Kronecker Product)

For any positive integers p, q, r, s we define the **Kronecker product** of two matrices $\mathbf{A} \in \mathbb{R}^{p \times q}$ and $\mathbf{B} \in \mathbb{R}^{r \times s}$ as a matrix $\mathbf{C} \in \mathbb{R}^{pr \times qs}$ given in block form as

$$\mathbf{C} = \begin{bmatrix} \mathbf{Ab}_{1,1} & \mathbf{Ab}_{1,2} & \cdots & \mathbf{Ab}_{1,s} \\ \mathbf{Ab}_{2,1} & \mathbf{Ab}_{2,2} & \cdots & \mathbf{Ab}_{2,s} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Ab}_{r,1} & \mathbf{Ab}_{r,2} & \cdots & \mathbf{Ab}_{r,s} \end{bmatrix}.$$

We denote the Kronecker product of \mathbf{A} and \mathbf{B} by $\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$.

This definition of the Kronecker product is known more precisely as the **left Kronecker product**. In the literature one often finds the **right Kronecker product** which in our notation is given by $\mathbf{B} \otimes \mathbf{A}$.

As examples of Kronecker products which are relevant for our discussion, if

$$\mathbf{T}_1 = \begin{bmatrix} d & a \\ a & d \end{bmatrix} \quad \text{and} \quad \mathbf{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

then

$$\mathbf{T}_1 \otimes \mathbf{I} = \begin{bmatrix} d & a & 0 & 0 \\ a & d & 0 & 0 \\ 0 & 0 & d & a \\ 0 & 0 & a & d \end{bmatrix} \quad \text{and} \quad \mathbf{I} \otimes \mathbf{T}_1 = \begin{bmatrix} d & 0 & a & 0 \\ 0 & d & 0 & a \\ a & 0 & d & 0 \\ 0 & a & 0 & d \end{bmatrix}.$$

Also note that the Kronecker product $\mathbf{u} \otimes \mathbf{v} = [\mathbf{u}^T v_1, \dots, \mathbf{u}^T v_r]^T$ of two column vectors $\mathbf{u} \in \mathbb{R}^p$ and $\mathbf{v} \in \mathbb{R}^r$ is a column vector of length $p \times r$.

The 2D test matrix \mathbf{T}_2 can be written as a sum of two Kronecker products.

We see that

$$\mathbf{T}_2 = \begin{bmatrix} \mathbf{T}_1 & & & \\ & \mathbf{T}_1 & & \\ & & \ddots & \\ & & & \mathbf{T}_1 \\ & & & & \mathbf{T}_1 \end{bmatrix} + \begin{bmatrix} d\mathbf{I} & a\mathbf{I} & & \\ a\mathbf{I} & d\mathbf{I} & a\mathbf{I} & \\ & \ddots & \ddots & \ddots \\ & & a\mathbf{I} & d\mathbf{I} & a\mathbf{I} \\ & & & a\mathbf{I} & d\mathbf{I} \end{bmatrix} = \mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1.$$

Definition 4.4 (Kronecker sum)

For positive integers r, s, k , let $\mathbf{A} \in \mathbb{R}^{r \times r}$, $\mathbf{B} \in \mathbb{R}^{s \times s}$, and \mathbf{I}_k be the identity matrix of order k . The sum $\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}$ is known as the **Kronecker sum** of \mathbf{A} and \mathbf{B} .

In other words, the 2D test matrix is the Kronecker sum of two identical 1D test matrices.

The following simple arithmetic rules hold for Kronecker products. For scalars λ, μ and matrices $\mathbf{A}, \mathbf{A}_1, \mathbf{A}_2, \mathbf{B}, \mathbf{B}_1, \mathbf{B}_2, \mathbf{C}$ of dimensions such that the operations are defined, we have

$$\begin{aligned} (\lambda \mathbf{A}) \otimes (\mu \mathbf{B}) &= \lambda \mu (\mathbf{A} \otimes \mathbf{B}), \\ (\mathbf{A}_1 + \mathbf{A}_2) \otimes \mathbf{B} &= \mathbf{A}_1 \otimes \mathbf{B} + \mathbf{A}_2 \otimes \mathbf{B}, \\ \mathbf{A} \otimes (\mathbf{B}_1 + \mathbf{B}_2) &= \mathbf{A} \otimes \mathbf{B}_1 + \mathbf{A} \otimes \mathbf{B}_2, \\ (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} &= \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}), \\ (\mathbf{A} \otimes \mathbf{B})^T &= \mathbf{A}^T \otimes \mathbf{B}^T. \end{aligned} \tag{4.10}$$

Note however that in general we have $\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A}$, but it can be shown that there are permutation matrices \mathbf{P}, \mathbf{Q} such that $\mathbf{B} \otimes \mathbf{A} = \mathbf{P}(\mathbf{A} \otimes \mathbf{B})\mathbf{Q}$, see [13].

Exercise 4.5 (Properties of Kronecker products)

Prove (4.10).

The following **mixed product rule** is an essential tool for dealing with Kronecker products and sums.

Lemma 4.6 (Mixed Product Rule)

Suppose $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ are rectangular matrices with dimensions so that the products \mathbf{AC} and \mathbf{BD} are defined. Then the product $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D})$ is defined and

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}). \tag{4.11}$$

Proof. If $\mathbf{B} \in \mathbb{R}^{r,t}$ and $\mathbf{D} \in \mathbb{R}^{t,s}$ for some integers r, s, t , then

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \begin{bmatrix} \mathbf{Ab}_{1,1} & \cdots & \mathbf{Ab}_{1,t} \\ \vdots & & \vdots \\ \mathbf{Ab}_{r,1} & \cdots & \mathbf{Ab}_{r,t} \end{bmatrix} \begin{bmatrix} \mathbf{Cd}_{1,1} & \cdots & \mathbf{Cd}_{1,s} \\ \vdots & & \vdots \\ \mathbf{Cd}_{t,1} & \cdots & \mathbf{Cd}_{t,s} \end{bmatrix}.$$

Thus for all i, j

$$((\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}))_{i,j} = \mathbf{AC} \sum_{k=1}^t b_{i,k} d_{k,j} = (\mathbf{AC})(\mathbf{BD})_{i,j} = ((\mathbf{AC}) \otimes (\mathbf{BD}))_{i,j}.$$

□

The eigenvalues and eigenvectors of a Kronecker product can easily be determined if one knows the corresponding quantities for each of the factors in the product.

Lemma 4.7 (Eigenvalues of Kronecker products)

Suppose \mathbf{A} and \mathbf{B} are square matrices. Then the eigenvalues of $\mathbf{A} \otimes \mathbf{B}$ are products of eigenvalues of \mathbf{A} and \mathbf{B} , and the eigenvectors of $\mathbf{A} \otimes \mathbf{B}$ are Kronecker products of eigenvectors of \mathbf{A} and \mathbf{B} . More precisely, if $\mathbf{A} \in \mathbb{R}^{r \times r}$, $\mathbf{B} \in \mathbb{R}^{s \times s}$, and

$$\mathbf{Au}_i = \lambda_i \mathbf{u}_i, \quad i = 1, \dots, r, \quad \mathbf{Bv}_j = \mu_j \mathbf{v}_j, \quad j = 1, \dots, s,$$

then

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = \lambda_i \mu_j (\mathbf{u}_i \otimes \mathbf{v}_j), \quad i = 1, \dots, r, \quad j = 1, \dots, s. \quad (4.12)$$

Proof. Using (4.10) and (4.11) the proof is a oneliner. For all i, j

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\mathbf{Au}_i) \otimes (\mathbf{Bv}_j) = (\lambda_i \mathbf{u}_i) \otimes (\mu_j \mathbf{v}_j) = (\lambda_i \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j).$$

□

Consider next a Kronecker sum.

Lemma 4.8 (Eigenvalues of Kronecker sums)

For positive integers r, s let $\mathbf{A} \in \mathbb{R}^{r \times r}$ and $\mathbf{B} \in \mathbb{R}^{s \times s}$. Then the eigenvalues of the Kronecker sum $\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}$ are all sums of eigenvalues of \mathbf{A} and \mathbf{B} , and the eigenvectors of $\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}$ are all Kronecker products of eigenvectors of \mathbf{A} and \mathbf{B} . More precisely, if

$$\mathbf{Au}_i = \lambda_i \mathbf{u}_i, \quad i = 1, \dots, r, \quad \mathbf{Bv}_j = \mu_j \mathbf{v}_j, \quad j = 1, \dots, s,$$

then

$$(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\lambda_i + \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j), \quad i = 1, \dots, r, \quad j = 1, \dots, s. \quad (4.13)$$

Proof. Since $\mathbf{I}_s \mathbf{v}_j = \mathbf{v}_j$ for $j = 1, \dots, s$ and $\mathbf{I}_r \mathbf{u}_i = \mathbf{u}_i$ for $i = 1, \dots, r$, we obtain by Lemma 4.7 for all i, j

$$(\mathbf{A} \otimes \mathbf{I}_s)(\mathbf{u}_i \otimes \mathbf{v}_j) = \lambda_i(\mathbf{u}_i \otimes \mathbf{v}_j), \quad \text{and} \quad (\mathbf{I}_r \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = \mu_j(\mathbf{u}_i \otimes \mathbf{v}_j).$$

The result now follows by summing these relations. \square

In many cases the Kronecker product and sum inherit properties of their factors.

Lemma 4.9 (Kronecker product; inverse and positive definite)

1. *The matrices \mathbf{A} and \mathbf{B} are nonsingular if and only if $\mathbf{A} \otimes \mathbf{B}$ is nonsingular. In that case $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$.*
2. *If \mathbf{A} and \mathbf{B} are symmetric then $\mathbf{A} \otimes \mathbf{B}$ and $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$ are symmetric.*
3. *If one of \mathbf{A} , \mathbf{B} is symmetric positive definite and the other is symmetric positive semidefinite then $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$ is symmetric positive definite.*

Proof. Suppose that $\mathbf{A} \in \mathbb{R}^{r \times r}$ and $\mathbf{B} \in \mathbb{R}^{s \times s}$. 1. follows from the mixed product rule giving

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{A}^{-1} \otimes \mathbf{B}^{-1}) = (\mathbf{A}\mathbf{A}^{-1}) \otimes (\mathbf{B}\mathbf{B}^{-1}) = \mathbf{I}_r \otimes \mathbf{I}_s = \mathbf{I}_{rs}.$$

Thus $(\mathbf{A} \otimes \mathbf{B})$ is nonsingular with the indicated inverse. 2. and the symmetry part of 3. follow immediately from (4.10). Suppose \mathbf{A} is symmetric positive definite and \mathbf{B} is symmetric positive semidefinite. Then \mathbf{A} has positive eigenvalues and \mathbf{B} has nonnegative eigenvalues. By Lemma 4.8 the eigenvalues of $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$ are all positive and 3. follows. \square

In (4.5) we derived the matrix equation $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ for the unknowns \mathbf{V} in the discrete Poisson problem. With some effort we converted this matrix equation to a linear system in standard form $\mathbf{Ax} = \mathbf{b}$, where $\mathbf{A} = \mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}$, $\mathbf{x} = \text{vec}(\mathbf{V})$, and $\mathbf{b} = \text{vec}(\mathbf{F})$. This conversion could have been carried out with less effort using the following result.

Lemma 4.10 (Conversion Kronecker product to matrix equation)

Suppose $\mathbf{A} \in \mathbb{R}^{r \times r}$, $\mathbf{B} \in \mathbb{R}^{s \times s}$, and $\mathbf{F}, \mathbf{V} \in \mathbb{R}^{r \times s}$. Then we have

$$(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F}) \Leftrightarrow \mathbf{AVB}^T = \mathbf{F}, \quad (4.14)$$

$$(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F}) \Leftrightarrow \mathbf{AV} + \mathbf{VB}^T = \mathbf{F}. \quad (4.15)$$

Proof. We partition \mathbf{V} , \mathbf{F} , and \mathbf{B}^T by columns as $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_s]$, $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_s]$ and $\mathbf{B}^T = [\mathbf{b}_1, \dots, \mathbf{b}_s]$. Then we have

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\ \Leftrightarrow \quad &\left[\begin{array}{ccc} \mathbf{A}\mathbf{b}_{11} & \cdots & \mathbf{A}\mathbf{b}_{1s} \\ \vdots & & \vdots \\ \mathbf{A}\mathbf{b}_{s1} & \cdots & \mathbf{A}\mathbf{b}_{ss} \end{array} \right] \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_s \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1 \\ \vdots \\ \mathbf{f}_s \end{bmatrix} \\ \Leftrightarrow \quad &\mathbf{A} \left[\sum_j b_{1j} \mathbf{v}_j, \dots, \sum_j b_{sj} \mathbf{v}_j \right] = [\mathbf{f}_1, \dots, \mathbf{f}_s] \\ \Leftrightarrow \quad &\mathbf{A}[\mathbf{V}\mathbf{b}_1, \dots, \mathbf{V}\mathbf{b}_s] = \mathbf{F} \quad \Leftrightarrow \quad \mathbf{AVB}^T = \mathbf{F}. \end{aligned}$$

This proves (4.14). (4.15) follows immediately from (4.14) as follows

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\ \Leftrightarrow \quad &(\mathbf{AVI}_s^T + \mathbf{I}_r \mathbf{VB}^T) = \mathbf{F} \quad \Leftrightarrow \quad \mathbf{AV} + \mathbf{VB}^T = \mathbf{F}. \end{aligned}$$

□

For more on Kronecker products see [13].

4.3 Properties of the 1D and 2D Test Matrices

We can apply these results to the 2D test matrix \mathbf{T}_2 . We first consider the 1D test matrix. The eigenvectors of \mathbf{T}_1 are the columns of the sine matrix defined by

$$\mathbf{S} = \left[\sin \frac{jk\pi}{m+1} \right]_{j,k=1}^m \in \mathbb{R}^{m \times m}. \quad (4.16)$$

For $m = 3$,

$$\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3] = \begin{bmatrix} \sin \frac{\pi}{4} & \sin \frac{2\pi}{4} & \sin \frac{3\pi}{4} \\ \sin \frac{2\pi}{4} & \sin \frac{4\pi}{4} & \sin \frac{6\pi}{4} \\ \sin \frac{3\pi}{4} & \sin \frac{6\pi}{4} & \sin \frac{9\pi}{4} \end{bmatrix} = \begin{bmatrix} t & 1 & t \\ 1 & 0 & -1 \\ t & -1 & t \end{bmatrix}, \quad t := \frac{1}{\sqrt{2}}.$$

Lemma 4.11 (Eigenpairs of 2. derivative matrix)

Suppose $\mathbf{T}_1 = (t_{kj})_{k,j} = \text{tridiag}(a, d, a) \in \mathbb{R}^{m \times m}$ with $m \geq 2$, $a, d \in \mathbb{R}$, and let $h = 1/(m+1)$.

1. We have $\mathbf{T}_1 \mathbf{s}_j = \lambda_j \mathbf{s}_j$ for $j = 1, \dots, m$, where

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (4.17)$$

$$\lambda_j = d + 2a \cos(j\pi h). \quad (4.18)$$

2. The eigenvalues are distinct and the eigenvectors are orthogonal

$$\mathbf{s}_j^T \mathbf{s}_k = \frac{m+1}{2} \delta_{j,k}, \quad j, k = 1, \dots, m. \quad (4.19)$$

Proof. We find

$$\begin{aligned} (\mathbf{T}_1 \mathbf{s}_j)_k &= \sum_{l=1}^m t_{k,l} \sin(lj\pi h) = a [\sin((k-1)j\pi h) + \sin((k+1)j\pi h)] + d \sin(kj\pi h) \\ &= 2a \cos(j\pi h) \sin(kj\pi h) + d \sin(kj\pi h) = \lambda_j s_{k,j}, \end{aligned}$$

and 1. follows. Since $j\pi h = j\pi/(m+1) \in (0, \pi)$ for $j = 1, \dots, m$ and the cosine function is strictly monotone decreasing on $(0, \pi)$ the eigenvalues are distinct, and since \mathbf{T}_1 is symmetric it follows from Lemma 4.12 below that the eigenvectors \mathbf{s}_j are orthogonal. To finish the proof of (4.19) we compute the square of the Euclidian norm of each \mathbf{s}_j as follows:

$$\begin{aligned} \mathbf{s}_j^T \mathbf{s}_j &= \sum_{k=1}^m \sin^2(kj\pi h) = \sum_{k=0}^m \sin^2(kj\pi h) = \frac{1}{2} \sum_{k=0}^m (1 - \cos(2kj\pi h)) \\ &= \frac{m+1}{2} - \frac{1}{2} \sum_{k=0}^m \cos(2kj\pi h) = \frac{m+1}{2}, \end{aligned}$$

since the last cosine sum is zero. We show this by summing a geometric series of complex exponentials. With $i = \sqrt{-1}$ we find

$$\sum_{k=0}^m \cos(2kj\pi h) + i \sum_{k=0}^m \sin(2kj\pi h) = \sum_{k=0}^m e^{2ikj\pi h} = \frac{e^{2i(m+1)j\pi h} - 1}{e^{2ij\pi h} - 1} = 0,$$

and (4.19) follows. \square

Lemma 4.12 (Eigenpairs of a Hermitian matrix)

The eigenvalues of a Hermitian matrix are real. Moreover, eigenvectors corresponding to distinct eigenvalues are orthogonal.

Proof. The first part was shown in Lemma 3.28. Suppose that (λ, \mathbf{x}) and (μ, \mathbf{y}) are two eigenpairs for \mathbf{A} with $\mu \neq \lambda$. Multiplying $\mathbf{Ax} = \lambda \mathbf{x}$ by \mathbf{y}^* gives

$$\lambda \mathbf{y}^* \mathbf{x} = \mathbf{y}^* \mathbf{Ax} = (\mathbf{x}^* \mathbf{A}^* \mathbf{y})^* = (\mathbf{x}^* \mathbf{Ay})^* = (\mu \mathbf{x}^* \mathbf{y})^* = \mu \mathbf{y}^* \mathbf{x},$$

using that μ is real. Since $\lambda \neq \mu$ it follows that $\mathbf{y}^* \mathbf{x} = 0$, which means that \mathbf{x} and \mathbf{y} are orthogonal. \square

It is now easy to find the eigenpairs of the 2D test matrix and determine when it is positive definite.

Theorem 4.13 (Eigenpairs of 2D test matrix)

For fixed $m \geq 2$ let \mathbf{T}_2 be the matrix given by (4.9) and let $h = 1/(m+1)$.

1. We have $\mathbf{T}_2 \mathbf{x}_{j,k} = \lambda_{j,k} \mathbf{x}_{j,k}$ for $j, k = 1, \dots, m$, where

$$\mathbf{x}_{j,k} = \mathbf{s}_j \otimes \mathbf{s}_k, \quad (4.20)$$

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (4.21)$$

$$\lambda_{j,k} = 2d + 2a \cos(j\pi h) + 2a \cos(k\pi h). \quad (4.22)$$

2. The eigenvectors are orthogonal

$$\mathbf{x}_{j,k}^T \mathbf{x}_{p,q} = \frac{1}{4h^2} \delta_{j,p} \delta_{k,q}, \quad j, k, p, q = 1, \dots, m. \quad (4.23)$$

3. \mathbf{T}_2 is symmetric positive definite if $d > 0$ and $d \geq 2|a|$.

4. The Poisson and averaging matrices are symmetric positive definite.

Proof. 1. follows from Lemmas 4.8 and 4.11, since $\mathbf{T}_2 = \mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1$. Using the transpose rule, the mixed product rule and (4.19). we find for $j, k, p, q = 1, \dots, m$

$$(\mathbf{s}_j \otimes \mathbf{s}_k)^T (\mathbf{s}_p \otimes \mathbf{s}_q) = (\mathbf{s}_j^T \otimes \mathbf{s}_k^T)(\mathbf{s}_p \otimes \mathbf{s}_q) = (\mathbf{s}_j^T \mathbf{s}_p) \otimes (\mathbf{s}_k^T \mathbf{s}_q) = \frac{1}{4h^2} \delta_{j,p} \delta_{k,q}$$

and 2. follows. Since \mathbf{T}_2 is symmetric, 3. will follow if the eigenvalues are positive. But this is true if $d > 0$ and $d \geq 2|a|$, which hold for both choices $a = -1$, $d = 2$ and $a = 1/5$, $d = 5/18$. Thus the matrices in 4. are positive definite. \square

Exercise 4.14 (2. derivative matrix is positive definite)

Write down the eigenvalues of $\mathbf{T} = \text{tridiag}(-1, 2, -1)$ using Lemma 4.11 and conclude that \mathbf{T} is symmetric positive definite.

Exercise 4.15 (1D test matrix is positive definite?)

Use Lemma 4.11 to show that the matrix $\mathbf{T}_1 := \text{tridiag}(a, d, a) \in \mathbb{R}^{n \times n}$ is symmetric positive definite if $d > 0$ and $d \geq 2|a|$.

Exercise 4.16 (Eigenvalues 2 × 2 for 2D test matrix)

For $m = 2$ the matrix (4.9) is given by

$$\mathbf{A} = \begin{bmatrix} 2d & a & a & 0 \\ a & 2d & 0 & a \\ a & 0 & 2d & a \\ 0 & a & a & 2d \end{bmatrix}.$$

Show that $\lambda = 2a + 2d$ is an eigenvalue corresponding to the eigenvector $\mathbf{x} = [1, 1, 1, 1]^T$. Verify that apart from a scaling of the eigenvector this agrees with (4.22) and (4.21) for $j = k = 1$ and $m = 2$.

Exercise 4.17 (Nine point scheme for Poisson problem)

Consider the following 9 point difference approximation to the Poisson problem $-\nabla^2 u = f$, $u = 0$ on the boundary of the unit square (cf. (4.1))

$$\begin{aligned} \text{(a)} \quad -(\square_h v)_{j,k} &= (\mu f)_{j,k} & j, k = 1, \dots, m \\ \text{(b)} \quad 0 &= v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1}, & j, k = 0, 1, \dots, m+1, \\ \text{(c)} \quad -(\square_h v)_{j,k} &= [20v_{j,k} - 4v_{j-1,k} - 4v_{j,k-1} - 4v_{j+1,k} - 4v_{j,k+1} \\ &\quad - v_{j-1,k-1} - v_{j+1,k-1} - v_{j-1,k+1} - v_{j+1,k+1}]/(6h^2), \\ \text{(d)} \quad (\mu f)_{j,k} &= [8f_{j,k} + f_{j-1,k} + f_{j,k-1} + f_{j+1,k} + f_{j,k+1}]/12. \end{aligned} \quad (4.24)$$

- a) Write down the 4-by-4 system we obtain for $m = 2$.
- b) Find $v_{j,k}$ for $j, k = 1, 2$, if $f(x, y) = 2\pi^2 \sin(\pi x) \sin(\pi y)$ and $m = 2$. Answer: $v_{j,k} = 5\pi^2/66$.

It can be shown that (4.24) defines an $O(h^4)$ approximation to (4.1).

Exercise 4.18 (Matrix equation for nine point scheme)

Consider the nine point difference approximation to (4.1) given by (4.24) in Problem 4.17.

- a) Show that (4.24) is equivalent to the matrix equation

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} - \frac{1}{6}\mathbf{T}\mathbf{V}\mathbf{T} = h^2\mu\mathbf{F}. \quad (4.25)$$

Here $\mu\mathbf{F}$ has elements $(\mu f)_{j,k}$ given by (4.24d).

- b) Show that the standard form of the matrix equation (4.25) is $\mathbf{A}\mathbf{x} = \mathbf{b}$, where $\mathbf{A} = \mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T} - \frac{1}{6}\mathbf{T} \otimes \mathbf{T}$, $\mathbf{x} = \text{vec}(\mathbf{V})$, and $\mathbf{b} = h^2 \text{vec}(\mu\mathbf{F})$.

Exercise 4.19 (Biharmonic equation)

Consider the biharmonic equation

$$\begin{aligned} \nabla^4 u(s, t) &:= \nabla^2(\nabla^2 u(s, t)) = f(s, t) & (s, t) \in \Omega, \\ u(s, t) = 0, \quad \nabla^2 u(s, t) &= 0 & (s, t) \in \partial\Omega. \end{aligned} \quad (4.26)$$

Here Ω is the open unit square. The condition $\nabla^2 u = 0$ is called the Navier boundary condition. Moreover, $\nabla^4 u = u_{xxxx} + 2u_{xxyy} + u_{yyyy}$.

a) Let $v = -\nabla^2 u$. Show that (4.26) can be written as a system

$$\begin{aligned}-\nabla^2 v(s, t) &= f(s, t) & (s, t) \in \Omega \\ -\nabla^2 u(s, t) &= v(s, t) & (s, t) \in \Omega \\ u(s, t) &= v(s, t) = 0 & (s, t) \in \partial\Omega.\end{aligned}\quad (4.27)$$

b) Discretizing, using (4.2), with $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$, $h = 1/(m+1)$, and $\mathbf{F} = (f(jh, kh))_{j,k=1}^m$ we get two matrix equations

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2 \mathbf{F}, \quad \mathbf{T}\mathbf{U} + \mathbf{U}\mathbf{T} = h^2 \mathbf{V}.$$

Show that

$$(\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}) \text{vec}(\mathbf{V}) = h^2 \text{vec}(\mathbf{F}), \quad (\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}) \text{vec}(\mathbf{U}) = h^2 \text{vec}(\mathbf{V}).$$

and hence $\mathbf{A} = (\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T})^2$ is the matrix for the standard form of the discrete biharmonic equation.

c) Show that with $n = m^2$ the vector form and standard form of the systems in b) can be written

$$\mathbf{T}^2 \mathbf{U} + 2\mathbf{T}\mathbf{U}\mathbf{T} + \mathbf{U}\mathbf{T}^2 = h^4 \mathbf{F} \quad \text{and} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (4.28)$$

where $\mathbf{A} = \mathbf{T}^2 \otimes \mathbf{I} + 2\mathbf{T} \otimes \mathbf{T} + \mathbf{I} \otimes \mathbf{T}^2 \in \mathbb{R}^{n \times n}$, $\mathbf{x} = \text{vec}(\mathbf{U})$, and $\mathbf{b} = h^4 \text{vec}(\mathbf{F})$.

- d) Determine the eigenvalues and eigenvectors of the matrix \mathbf{A} in c) and show that it is symmetric positive definite. Also determine the bandwidth of \mathbf{A} .
- e) Suppose we want to solve the standard form equation $\mathbf{A}\mathbf{x} = \mathbf{b}$. We have two representations for the matrix \mathbf{A} , the product one in b) and the one in c). Which one would you prefer for the basis of an algorithm? Why?

4.4 Review Questions

4.4.1 Consider the Poisson matrix.

- Write this matrix as a Kronecker sum,
- how are its eigenvalues and eigenvectors related to the second derivative matrix?
- is it symmetric? positive definite?

4.4.2 What are the eigenpairs of $\mathbf{T}_1 := \text{tridiagonal}(a, d, a)$?

4.4.3 What are the inverse and transpose of a Kronecker product?

- 4.4.4**
- give a economical general way to solve the linear system $(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$?
 - Same for $(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$.

Chapter 5

Fast Direct Solution of a Large Linear System

5.1 Algorithms for a Banded Positive Definite System

In this chapter we present a fast method for solving $\mathbf{A}\mathbf{x} = \mathbf{b}$, where \mathbf{A} is the Poisson matrix (4.7). Thus, for $n = 3$

$$\begin{aligned}\mathbf{A} &= \left[\begin{array}{ccc|ccc|ccc} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ \hline -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ \hline 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{array} \right] \\ &= \left[\begin{array}{ccc} \mathbf{T} + 2\mathbf{I} & -\mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{T} + 2\mathbf{I} & -\mathbf{I} \\ \mathbf{0} & -\mathbf{I} & \mathbf{T} + 2\mathbf{I} \end{array} \right],\end{aligned}$$

where $\mathbf{T} = \text{tridiag}(-1, 2, -1)$. For the matrix \mathbf{A} we know by now that

1. It is symmetric positive definite.
2. It is banded.
3. It is block-tridiagonal.

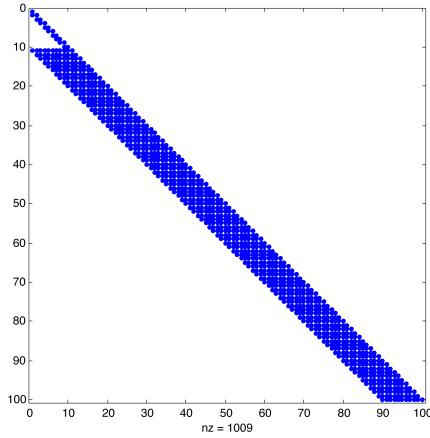


Figure 5.1: Fill-inn in the Cholesky factor of the Poisson matrix ($n = 100$).

4. We know the eigenvalues and eigenvectors of \mathbf{A} .
5. The eigenvectors are orthogonal.

5.1.1 Cholesky Factorization

Since \mathbf{A} is symmetric positive definite we can use a Cholesky factorization $\mathbf{A} = \mathbf{L}\mathbf{L}^T$, with \mathbf{L} lower triangular, to solve $\mathbf{A}\mathbf{x} = \mathbf{b}$. Since \mathbf{A} is banded with bandwidth $d = \sqrt{n}$ the matrix \mathbf{L} has bandwidth $d = \sqrt{n}$ (cf. Lemma 3.35) and the complexity of this factorization is $nd^2 = n^2$. We need to store \mathbf{A} possibly in sparse form.

The nonzero elements in \mathbf{L} are shown in Figure 5.1 for $n = 100$. Note that most of the zeros between the diagonals in \mathbf{A} have become nonzero in \mathbf{L} . This is known as **fill-inn**.

5.1.2 Block LU Factorization of a Block Tridiagonal Matrix

The Poisson matrix has a block tridiagonal structure. Consider finding the block LU factorization of a block tridiagonal matrix. We are looking for a factorization of the form

$$\begin{bmatrix} D_1 & C_1 \\ A_1 & D_2 & C_2 \\ & \ddots & \ddots & \ddots \\ & & A_{m-2} & D_{m-1} & C_{m-1} \\ & & A_{m-1} & D_m \end{bmatrix} = \begin{bmatrix} I & & & \\ L_1 & I & & \\ & \ddots & \ddots & \\ & & L_{m-1} & I \end{bmatrix} \begin{bmatrix} U_1 & C_1 \\ & \ddots & \ddots \\ & & U_{m-1} & C_{m-1} \\ & & & U_m \end{bmatrix}. \quad (5.1)$$

Here D_1, \dots, D_m and U_1, \dots, U_m are square matrices while A_1, \dots, A_{m-1} and C_1, \dots, C_{m-1} can be rectangular.

Using block multiplication the formulas (2.4) generalize to

$$\mathbf{U}_1 = \mathbf{D}_1, \quad \mathbf{L}_k = \mathbf{A}_k \mathbf{U}_k^{-1}, \quad \mathbf{U}_{k+1} = \mathbf{D}_{k+1} - \mathbf{L}_k \mathbf{C}_k, \quad k = 1, 2, \dots, m-1. \quad (5.2)$$

To solve the system $\mathbf{Ax} = \mathbf{b}$ we partition \mathbf{b} conformally with \mathbf{A} in the form $\mathbf{b}^T = [\mathbf{b}_1^T, \dots, \mathbf{b}_m^T]$. The formulas for solving $\mathbf{Ly} = \mathbf{b}$ and $\mathbf{Ux} = \mathbf{y}$ are as follows:

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{b}_1, & \mathbf{y}_k &= \mathbf{b}_k - \mathbf{L}_{k-1} \mathbf{y}_{k-1}, & k &= 2, 3, \dots, m, \\ \mathbf{x}_m &= \mathbf{U}_m^{-1} \mathbf{y}_m, & \mathbf{x}_k &= \mathbf{U}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_{k+1}), & k &= m-1, \dots, 2, 1. \end{aligned} \quad (5.3)$$

The solution is then $\mathbf{x}^T = [\mathbf{x}_1^T, \dots, \mathbf{x}_m^T]$. To find \mathbf{L}_k in (5.2) we solve the linear systems $\mathbf{L}_k \mathbf{U}_k = \mathbf{A}_k$. Similarly we need to solve a linear system to find \mathbf{x}_k in (5.3).

The number of arithmetic operations using block factorizations is $O(n^2)$, asymptotically the same as for Cholesky factorization. However we only need to store the $m \times m$ blocks and using matrix operations can be an advantage.

5.1.3 Other Methods

Other methods include

- Iterative methods. We study this in Chapters 9 and 10.
- Multigrid. See [7].
- Fast solvers based on diagonalization and the Fast Fourier Transform. See Sections 5.2, 5.3.

5.2 A Fast Poisson Solver based on Diagonalization

The algorithm we now derive will only require $O(n^{3/2})$ arithmetic operations and we only need to work with matrices of order m . Using the Fast Fourier Transform the number of arithmetic operations can be reduced further to $O(n \log n)$.

To start we recall that $\mathbf{Ax} = \mathbf{b}$ can be written as a matrix equation in the form (cf. (4.5))

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2 \mathbf{F} \quad \text{with} \quad h = 1/(m+1),$$

where $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ is the second derivative matrix, $\mathbf{V} = (v_{jk}) \in \mathbb{R}^{m \times m}$ are the unknowns, and $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m \times m}$ contains function values.

Recall that the eigenpairs of \mathbf{T} are given by

$$\begin{aligned}\mathbf{T}\mathbf{s}_j &= \lambda_j \mathbf{s}_j, \quad j = 1, \dots, m, \\ \mathbf{s}_j &= [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \\ \lambda_j &= 2 - 2 \cos(j\pi h) = 4 \sin^2(j\pi h/2), \quad h = 1/(m+1), \\ \mathbf{s}_j^T \mathbf{s}_k &= \delta_{jk}/(2h) \text{ for all } j, k.\end{aligned}$$

Let

$$\mathbf{S} := [\mathbf{s}_1, \dots, \mathbf{s}_m] = [\sin(jk\pi h)]_{j,k=1}^m \in \mathbb{R}^{m \times m}, \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_m). \quad (5.4)$$

Then $\mathbf{T}\mathbf{S} = \mathbf{SD}$ and $\mathbf{S}^T \mathbf{S} = \mathbf{S}^2 = \mathbf{I}/(2h)$. Define $\mathbf{X} \in \mathbb{R}^{m \times m}$ by $\mathbf{V} = \mathbf{SXS}$, where \mathbf{V} is the solution of $\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}$. Then

$$\begin{aligned}\mathbf{TV} + \mathbf{VT} &= h^2 \mathbf{F} \\ \xrightleftharpoons{\mathbf{V} = \mathbf{SXS}} \mathbf{TSXS} + \mathbf{SXST} &= h^2 \mathbf{F} \\ \xrightleftharpoons{\mathbf{S}(\)\mathbf{S}} \mathbf{STSXS}^2 + \mathbf{S}^2 \mathbf{XSTS} &= h^2 \mathbf{SFS} \\ \xrightleftharpoons{\mathbf{TS} = \mathbf{SD}} \mathbf{S}^2 \mathbf{DXS}^2 + \mathbf{S}^2 \mathbf{XS}^2 \mathbf{D} &= h^2 \mathbf{SFS} \\ \xrightleftharpoons{\mathbf{S}^2 = \mathbf{I}/(2h)} \mathbf{DX} + \mathbf{XD} &= 4h^4 \mathbf{SFS}.\end{aligned}$$

An equation of the form $\mathbf{DX} + \mathbf{XD} = \mathbf{B}$, where \mathbf{D} is diagonal is easy to solve. If $\mathbf{D} = \text{diag}(\lambda_j)$ we obtain for each element the equation $\lambda_j x_{jk} + x_{jk} \lambda_k = b_{jk}$ so $x_{jk} = b_{jk}/(\lambda_j + \lambda_k)$ for all j, k .

We now get the following algorithm to find the exact solution of $\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}$.

Algorithm 5.1 (Fast Poisson Solver)

We solve the Poisson problem $-\nabla^2 u = f$ on $\Omega = (0, 1)^2$ and $u = 0$ on $\partial\Omega$ using the 5-point scheme, i.e., let $m \in \mathbb{N}$, $h = 1/(m+1)$, and $\mathbf{F} = (f(jh, kh)) \in \mathbb{R}^{m \times m}$. We compute $\mathbf{V} \in \mathbb{R}^{m \times m}$, where $v_{jk} \approx u(jh, kh)$ by solving the equation $\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}$ using diagonalization of $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$.

```

1 function V=fastpoisson(F)
2 m=length(F); h=1/(m+1); hv=pi*h*(1:m)';
3 sigma=sin(hv/2).^2;
4 S=sin(hv*(1:m));
5 G=S*S;
6 X=h.^4*G./((sigma*ones(1,m)+ones(m,1)*sigma)');
7 V=zeros(m+2,m+2);
8 V(2:m+1,2:m+1)=S*X*S;
```

The formulas are fully vectorized and for convenience we have used $\sigma_j := \lambda_j/4$ instead of λ_j . Since the 6th line in Algorithm 5.1 only requires $O(m^2)$

arithmetic operations the complexity of this algorithm is for large m determined by the 4 m -by- m matrix multiplications and is given by $O(4 \times 2m^3) = O(8n^{3/2})$.⁹

5.3 A Fast Poisson Solver based on the Discrete Sine and Fourier Transforms

In Algorithm 5.1 we need to compute the product of the sine matrix $\mathbf{S} \in \mathbb{R}^{m \times m}$ given by (5.4) and a matrix $\mathbf{A} \in \mathbb{R}^{m \times m}$. Since the matrices are m -by- m this will normally require $O(m^3)$ operations. In this section we show that it is possible to calculate the products \mathbf{SA} and \mathbf{AS} in $O(m^2 \log_2 m)$ operations.

We need to discuss certain transforms known as the **Discrete Sine Transform**, the **Discrete Fourier Transform** and the **Fast Fourier Transform**. In addition we have the **Discrete Cosine Transform** which will not be discussed here. These transforms are of independent interest. They have applications to signal processing and image analysis, and are often used when one is dealing with discrete samples of data on a computer.

5.3.1 The Discrete Sine Transform (DST)

Given $\mathbf{v} = [v_1, \dots, v_m]^T \in \mathbb{R}^m$ we say that the vector $\mathbf{w} = [w_1, \dots, w_m]^T$ given by

$$w_j = \sum_{k=1}^m \sin\left(\frac{jk\pi}{m+1}\right) v_k, \quad j = 1, \dots, m$$

is the **Discrete Sine Transform** (DST) of \mathbf{v} . In matrix form we can write the DST as the matrix times vector $\mathbf{w} = \mathbf{Sv}$, where \mathbf{S} is the sine matrix given by (5.4). We can then identify the matrix $\mathbf{B} = \mathbf{SA}$ as the DST of $\mathbf{A} \in \mathbb{R}^{m,n}$, i.e. as the DST of the columns of \mathbf{A} . The product $\mathbf{B} = \mathbf{AS}$ can also be interpreted as a DST. Indeed, since \mathbf{S} is symmetric we have $\mathbf{B} = (\mathbf{SA}^T)^T$ which means that \mathbf{B} is the transpose of the DST of the rows of \mathbf{A} . It follows that we can compute the unknowns \mathbf{V} in Algorithm 5.1 by carrying out Discrete Sine Transforms on 4 m -by- m matrices in addition to the computation of \mathbf{X} .

5.3.2 The Discrete Fourier Transform (DFT)

The fast computation of the DST is based on its relation to the Discrete Fourier Transform (DFT) and the fact that the DFT can be computed by a technique known as the Fast Fourier Transform (FFT). To define the DFT let for $N \in \mathbb{N}$

$$\omega_N = \exp^{-2\pi i/N} = \cos(2\pi/N) - i \sin(2\pi/N), \quad (5.5)$$

⁹It is possible to compute \mathbf{V} using only two matrix multiplications and hence reduce the complexity to $O(4n^{3/2})$. This is detailed in Problem 5.8.



Figure 5.2: Jean Baptiste Joseph Fourier, 1768 - 1830.

where $i = \sqrt{-1}$ is the imaginary unit. Given $\mathbf{y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N$ we say that $\mathbf{z} = [z_1, \dots, z_N]^T$ given by

$$z_j = \sum_{k=1}^N \omega_N^{(j-1)(k-1)} y_k, \quad j = 1, \dots, N$$

is the **Discrete Fourier Transform** (DFT) of \mathbf{y} . We can write this as a matrix times vector product $\mathbf{z} = \mathbf{F}_N \mathbf{y}$, where the matrix \mathbf{F}_N is given by

$$\mathbf{F}_N = \left(\omega_N^{(j-1)(k-1)} \right)_{j,k=1}^N \in \mathbb{C}^{N \times N}. \quad (5.6)$$

This matrix is known as the **Fourier matrix**. If $\mathbf{A} \in \mathbb{R}^{N \times m}$ we say that $\mathbf{B} = \mathbf{F}_N \mathbf{A}$ is the DFT of \mathbf{A} .

As an example, since

$$\omega_4 = \exp^{-2\pi i / 4} = \cos(\pi/2) - i \sin(\pi/2) = -i$$

we find

$$\mathbf{F}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4 & \omega_4^2 & \omega_4^3 \\ 1 & \omega_4^2 & \omega_4^4 & \omega_4^6 \\ 1 & \omega_4^3 & \omega_4^6 & \omega_4^9 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -i & -1 & i \\ 1 & -1 & 1 & -1 \\ 1 & i & -1 & -i \end{bmatrix}. \quad (5.7)$$

The following lemma shows how the Discrete Sine Transform of order m can be computed from the Discrete Fourier Transform of order $2m + 2$.

Lemma 5.2 (Sine transform as Fourier transform)

Given a positive integer m and a vector $\mathbf{x} \in \mathbb{R}^m$. Component k of \mathbf{Sx} is equal to $i/2$ times component $k+1$ of $\mathbf{F}_{2m+2}\mathbf{z}$ where

$$\mathbf{z} = [0, x_1, \dots, x_m, 0, -x_m, -x_{m-1}, \dots, -x_1]^T \in \mathbb{R}^{2m+2}.$$

In symbols

$$(\mathbf{Sx})_k = \frac{i}{2} (\mathbf{F}_{2m+2}\mathbf{z})_{k+1}, \quad k = 1, \dots, m.$$

Proof. Let $\omega = \omega_{2m+2} = e^{-2\pi i/(2m+2)} = e^{-\pi i/(m+1)}$. Component $k+1$ of $\mathbf{F}_{2m+2}\mathbf{z}$ is given by

$$\begin{aligned} (\mathbf{F}_{2m+2}\mathbf{z})_{k+1} &= \sum_{j=1}^m x_j \omega^{jk} - \sum_{j=1}^m x_j \omega^{(2m+2-j)k} \\ &= \sum_{j=1}^m x_j (\omega^{jk} - \omega^{-jk}) \\ &= -2i \sum_{j=1}^m x_j \sin\left(\frac{jk\pi}{m+1}\right) = -2i(\mathbf{S}_m \mathbf{x})_k. \end{aligned}$$

Dividing both sides by $-2i$ proves the lemma. \square

It follows that we can compute the DST of length m by extracting m components from the DFT of length $N = 2m + 2$.

5.3.3 The Fast Fourier Transform (FFT)

From a linear algebra viewpoint the Fast Fourier Transform is a quick way to compute the matrix- vector product $\mathbf{F}_N \mathbf{y}$. Suppose N is even. The key to the FFT is a connection between \mathbf{F}_N and $\mathbf{F}_{N/2}$ which makes it possible to compute the FFT of order N as two FFT's of order $N/2$. By repeating this process we can reduce the number of arithmetic operations to compute a DFT from $O(N^2)$ to $O(N \log_2 N)$.

Suppose N is even. The connection between \mathbf{F}_N and $\mathbf{F}_{N/2}$ involves a permutation matrix $\mathbf{P}_N \in \mathbb{R}^{N \times N}$ given by

$$\mathbf{P}_N = [\mathbf{e}_1, \mathbf{e}_3, \dots, \mathbf{e}_{N-1}, \mathbf{e}_2, \mathbf{e}_4, \dots, \mathbf{e}_N],$$

where the $\mathbf{e}_k = (\delta_{j,k})$ are unit vectors. If \mathbf{A} is a matrix with N columns $[\mathbf{a}_1, \dots, \mathbf{a}_N]$ then

$$\mathbf{A}\mathbf{P}_N = [\mathbf{a}_1, \mathbf{a}_3, \dots, \mathbf{a}_{N-1}, \mathbf{a}_2, \mathbf{a}_4, \dots, \mathbf{a}_N],$$

i.e. post multiplying \mathbf{A} by \mathbf{P}_N permutes the columns of \mathbf{A} so that all the odd-indexed columns are followed by all the even-indexed columns. For example we have from (5.7)

$$\mathbf{P}_4 = [\mathbf{e}_1 \ \mathbf{e}_3 \ \mathbf{e}_2 \ \mathbf{e}_4] = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{F}_4 \mathbf{P}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & -i & i \\ 1 & 1 & -1 & -1 \\ 1 & -1 & i & -i \end{bmatrix},$$

where we have indicated a certain block structure of $\mathbf{F}_4 \mathbf{P}_4$. These blocks can be related to the 2-by-2 matrix \mathbf{F}_2 . We define the diagonal scaling matrix \mathbf{D}_2 by

$$\mathbf{D}_2 = \text{diag}(1, \omega_4) = \begin{bmatrix} 1 & 0 \\ 1 & -i \end{bmatrix}.$$

Since $\omega_2 = \exp^{-2\pi i/2} = -1$ we find

$$\mathbf{F}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \quad \mathbf{D}_2 \mathbf{F}_2 = \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix},$$

and we see that

$$\mathbf{F}_4 \mathbf{P}_4 = \left[\begin{array}{c|c} \mathbf{F}_2 & \mathbf{D}_2 \mathbf{F}_2 \\ \hline \mathbf{F}_2 & -\mathbf{D}_2 \mathbf{F}_2 \end{array} \right].$$

This result holds in general.

Theorem 5.3 (Fast Fourier Transform)

If $N = 2m$ is even then

$$\mathbf{F}_{2m} \mathbf{P}_{2m} = \left[\begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m \mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m \mathbf{F}_m \end{array} \right], \quad (5.8)$$

where

$$\mathbf{D}_m = \text{diag}(1, \omega_N, \omega_N^2, \dots, \omega_N^{m-1}). \quad (5.9)$$

Proof. Fix integers j, k with $0 \leq j, k \leq m-1$ and set $p = j+1$ and $q = k+1$. Since $\omega_m^m = 1$, $\omega_N^2 = \omega_m$, and $\omega_N^m = -1$ we find by considering elements in the four sub-blocks in turn

$$\begin{aligned} (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p,q} &= \omega_N^{j(2k)} &= \omega_m^{jk} &= (\mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p+m,q} &= \omega_N^{(j+m)(2k)} &= \omega_m^{(j+m)k} &= (\mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p,q+m} &= \omega_N^{j(2k+1)} &= \omega_N^j \omega_m^{jk} &= (\mathbf{D}_m \mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p+m,q+m} &= \omega_N^{(j+m)(2k+1)} &= -\omega_N^{j(2k+1)} &= (-\mathbf{D}_m \mathbf{F}_m)_{p,q}. \end{aligned}$$

It follows that the four m -by- m blocks of $\mathbf{F}_{2m} \mathbf{P}_{2m}$ have the required structure.

□

Using Theorem 5.3 we can carry out the DFT as a block multiplication. Let $\mathbf{y} \in \mathbb{R}^{2m}$ and set $\mathbf{w} = \mathbf{P}_{2m}^T \mathbf{y} = [\mathbf{w}_1, \mathbf{w}_2]^T$, where $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^m$. Then

$$\begin{aligned} \mathbf{F}_{2m} \mathbf{y} &= \mathbf{F}_{2m} \mathbf{P}_{2m} \mathbf{P}_{2m}^T \mathbf{y} = \mathbf{F}_{2m} \mathbf{P}_{2m} \mathbf{w} \\ &= \left[\begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m \mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m \mathbf{F}_m \end{array} \right] \left[\begin{array}{c} \mathbf{w}_1 \\ \mathbf{w}_2 \end{array} \right] = \left[\begin{array}{c} \mathbf{q}_1 + \mathbf{q}_2 \\ \mathbf{q}_1 - \mathbf{q}_2 \end{array} \right], \end{aligned}$$

where

$$\mathbf{q}_1 = \mathbf{F}_m \mathbf{w}_1, \quad \text{and} \quad \mathbf{q}_2 = \mathbf{D}_m(\mathbf{F}_m \mathbf{w}_2).$$

In order to compute $\mathbf{F}_{2m}\mathbf{y}$ we need to compute $\mathbf{F}_m\mathbf{w}_1$ and $\mathbf{F}_m\mathbf{w}_2$. Note that $\mathbf{w}_1^T = [y_1, y_3, \dots, y_{N-1}]$, while $\mathbf{w}_2^T = [y_2, y_4, \dots, y_N]$. This follows since $\mathbf{w}^T = [\mathbf{w}_1^T, \mathbf{w}_2^T] = \mathbf{y}^T \mathbf{P}_{2m}$ and post multiplying a vector by \mathbf{P}_{2m} moves odd indexed components to the left of all the even indexed components.

We have seen that by combining two FFT's of order m we obtain an FFT of order $2m$. If $N = 2^k$ then this process can be applied recursively as in the following Matlab function:

Algorithm 5.4 (Recursive FFT)

For $\mathbf{y} \in \mathbb{C}^n$ we compute the Fourier transform $\mathbf{z} = \mathbf{F}_n\mathbf{y}$.

```

1 function z=fftrec(y)
2 n=length(y);
3 if n==1
4   z=y;
5 else
6   q1=fftrec(y(1:2:n-1));
7   q2=exp(-2*pi*i/n).^(0:n/2-1).*fftrec(y(2:2:n));
8   z=[q1+q2 q1-q2];
9 end
```

Such a recursive version of FFT is useful for testing purposes, but is much too slow for large problems. A challenge for FFT code writers is to develop nonrecursive versions and also to handle efficiently the case where N is not a power of two. We refer to [28] for further details.

The complexity of the FFT is given by $\gamma N \log_2 N$ for some constant γ independent of N . To show this for the special case when N is a power of two let x_k be the complexity (the number of arithmetic operations) when $N = 2^k$. Since we need two FFT's of order $N/2 = 2^{k-1}$ and a multiplication with the diagonal matrix $\mathbf{D}_{N/2}$, it is reasonable to assume that $x_k = 2x_{k-1} + \gamma 2^k$ for some constant γ independent of k . Since $x_0 = 0$ we obtain by induction on k that $x_k = \gamma k 2^k$. Indeed, this holds for $k = 0$ and if $x_{k-1} = \gamma(k-1)2^{k-1}$ then $x_k = 2x_{k-1} + \gamma 2^k = 2\gamma(k-1)2^{k-1} + \gamma 2^k = \gamma k 2^k$. Reasonable implementations of FFT typically have $\gamma \approx 5$, see [28].

The efficiency improvement using the FFT to compute the DFT is spectacular for large N . The direct multiplication $\mathbf{F}_N\mathbf{y}$ requires $O(8n^2)$ arithmetic operations since complex arithmetic is involved. Assuming that the FFT uses $5N \log_2 N$ arithmetic operations we find for $N = 2^{20} \approx 10^6$ the ratio

$$\frac{8N^2}{5N \log_2 N} \approx 84000.$$

Thus if the FFT takes one second of computing time and the computing time is

proportional to the number of arithmetic operations then the direct multiplication would take something like 84000 seconds or 23 hours.

5.3.4 A Poisson Solver based on the FFT

We now have all the ingredients to compute the matrix products $\mathbf{S}\mathbf{A}$ and $\mathbf{A}\mathbf{S}$ using FFT's of order $2m + 2$ where m is the order of \mathbf{S} and \mathbf{A} . This can then be used for quick computation of the exact solution \mathbf{V} of the discrete Poisson problem in Algorithm 5.1. We first compute $\mathbf{H} = \mathbf{SF}$ using Lemma 5.2 and m FFT's, one for each of the m columns of \mathbf{F} . We then compute $\mathbf{G} = \mathbf{HS}$ by m FFT's, one for each of the rows of \mathbf{H} . After \mathbf{X} is determined we compute $\mathbf{Z} = \mathbf{SX}$ and $\mathbf{V} = \mathbf{ZS}$ by another $2m$ FFT's. In total the work amounts to $4m$ FFT's of order $2m + 2$. Since one FFT requires $O(\gamma(2m + 2) \log_2(2m + 2))$ arithmetic operations the $4m$ FFT's amount to

$$8\gamma m(m + 1) \log_2(2m + 2) \approx 8\gamma m^2 \log_2 m = 4\gamma n \log_2 n,$$

where $n = m^2$ is the size of the linear system $\mathbf{Ax} = \mathbf{b}$ we would be solving if Cholesky factorization was used. This should be compared to the $O(8n^{3/2})$ arithmetic operations used in Algorithm 5.1 requiring 4 straightforward matrix multiplications with \mathbf{S} . What is faster will depend heavily on the programming of the FFT and the size of the problem. We refer to [28] for other efficient ways to implement the DST.

Exercise 5.5 (Fourier matrix)

Show that the Fourier matrix \mathbf{F}_4 is symmetric, but not Hermitian.

Exercise 5.6 (Sine transform as Fourier transform)

Verify Lemma 5.2 directly when $m = 1$.

Exercise 5.7 (Explicit solution of the discrete Poisson equation)

Show that the exact solution of the discrete Poisson equation (4.3) and (4.4) can be written $\mathbf{V} = (v_{i,j})_{i,j=1}^m$, where

$$v_{ij} = \frac{1}{(m+1)^4} \sum_{p=1}^m \sum_{r=1}^m \sum_{k=1}^m \sum_{l=1}^m \frac{\sin\left(\frac{ip\pi}{m+1}\right) \sin\left(\frac{jr\pi}{m+1}\right) \sin\left(\frac{kp\pi}{m+1}\right) \sin\left(\frac{lr\pi}{m+1}\right)}{\left[\sin\left(\frac{p\pi}{2(m+1)}\right)\right]^2 + \left[\sin\left(\frac{r\pi}{2(m+1)}\right)\right]^2} f_{p,r}.$$

Exercise 5.8 (Improved version of Algorithm 5.1)

Algorithm 5.1 involves multiplying a matrix by \mathbf{S} four times. In this problem we show that it is enough to multiply by \mathbf{S} two times. We achieve this by diagonalizing only the second \mathbf{T} in $\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}$.

(a) Show that

$$\mathbf{T}\mathbf{X} + \mathbf{X}\mathbf{D} = \mathbf{C}, \text{ where } \mathbf{X} = \mathbf{V}\mathbf{S}, \text{ and } \mathbf{C} = h^2\mathbf{F}\mathbf{S}.$$

(b) Show that

$$(\mathbf{T} + \lambda_j \mathbf{I})\mathbf{x}_j = \mathbf{c}_j \quad j = 1, \dots, m, \quad (5.10)$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ and $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_m]$ and $\lambda_j = 4\sin^2(j\pi h/2)$. Thus we can find \mathbf{X} by solving m linear systems, one for each of the columns of \mathbf{X} . Recall that a tridiagonal $m \times m$ system can be solved by (2.4) and (2.5) in $8m - 7$ arithmetic operations. Give an algorithm to find \mathbf{X} which only requires $O(\delta m^2)$ arithmetic operations for some constant δ independent of m .

(c) Describe a method to compute \mathbf{V} which only requires $O(4m^3) = O(4n^{3/2})$ arithmetic operations.

(d) Describe a method based on the Fast Fourier Transform which requires $O(2\gamma n \log_2 n)$ where γ is the same constant as mentioned at the end of the last section.

Exercise 5.9 (Fast solution of 9 point scheme)

Consider the equation

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} - \frac{1}{6}\mathbf{T}\mathbf{V}\mathbf{T} = h^2\mu\mathbf{F},$$

that was derived in Exercise 4.18 for the 9-point scheme. Define the matrix \mathbf{X} by $\mathbf{V} = \mathbf{S}\mathbf{X}\mathbf{S} = (x_{j,k})$ where \mathbf{V} is the solution of (4.25). Show that

$$\mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} - \frac{1}{6}\mathbf{D}\mathbf{X}\mathbf{D} = 4h^4\mathbf{G}, \text{ where } \mathbf{G} = \mathbf{S}\mu\mathbf{F}\mathbf{S},$$

and that

$$x_{j,k} = \frac{h^4 g_{j,k}}{\sigma_j + \sigma_k - \frac{2}{3}\sigma_j\sigma_k}, \text{ where } \sigma_j = \sin^2((j\pi h)/2) \text{ for } j, k = 1, 2, \dots, m.$$

Show that $\sigma_j + \sigma_k - \frac{2}{3}\sigma_j\sigma_k > 0$ for $j, k = 1, 2, \dots, m$. Conclude that the matrix \mathbf{A} in Exercise 4.18 b) is symmetric positive definite and that (4.24) always has a solution \mathbf{V} .

Exercise 5.10 (Algorithm for fast solution of 9 point scheme)

Derive an algorithm for solving (4.24) which for large m requires essentially the same number of operations as in Algorithm 5.1. (We assume that $\mu\mathbf{F}$ already has been formed).

Exercise 5.11 (Fast solution of biharmonic equation)

For the biharmonic problem we derived in Exercise 4.19 the equation

$$\mathbf{T}^2 \mathbf{U} + 2\mathbf{T}\mathbf{U}\mathbf{T} + \mathbf{U}\mathbf{T}^2 = h^4 \mathbf{F}.$$

Define the matrix $\mathbf{X} = (x_{j,k})$ by $\mathbf{U} = \mathbf{S}\mathbf{X}\mathbf{S}$ where \mathbf{U} is the solution of (4.28). Show that

$$\mathbf{D}^2 \mathbf{X} + 2\mathbf{D}\mathbf{X}\mathbf{D} + \mathbf{X}\mathbf{D}^2 = 4h^6 \mathbf{G}, \text{ where } \mathbf{G} = \mathbf{S}\mathbf{F}\mathbf{S},$$

and that

$$x_{j,k} = \frac{h^6 g_{j,k}}{4(\sigma_j + \sigma_k)^2}, \text{ where } \sigma_j = \sin^2((j\pi h)/2) \text{ for } j, k = 1, 2, \dots, m.$$

Exercise 5.12 (Algorithm for fast solution of biharmonic equation)

Use Exercise 5.11 to derive an algorithm

```
function U=simplefastbiharmonic(F)
```

which requires only $O(\delta n^{3/2})$ operations to find \mathbf{U} in Problem 4.19. Here δ is some constant independent of n .

Exercise 5.13 (Check algorithm for fast solution of biharmonic equation)

In Exercise 5.12 compute the solution \mathbf{U} corresponding to $\mathbf{F} = \text{ones}(m, m)$. For some small m 's check that you get the same solution obtained by solving the standard form $\mathbf{Ax} = \mathbf{b}$ in (4.28). You can use $\mathbf{x} = \mathbf{A} \setminus \mathbf{b}$ for solving $\mathbf{Ax} = \mathbf{b}$. Use $\mathbf{F}(:)$ to vectorize a matrix and $\text{reshape}(\mathbf{x}, m, m)$ to turn a vector $\mathbf{x} \in \mathbb{R}^{m^2}$ into an $m \times m$ matrix. Use the Matlab command `surf(U)` for plotting U for, say, $m = 50$. Compare the result with Exercise 5.12 by plotting the difference between both matrices.

Exercise 5.14 (Fast solution of biharmonic equation using 9 point rule)

Repeat Exercises 4.19, 5.12 and 5.13 using the nine point rule (4.24) to solve the system (4.27).

5.4 Review Questions

5.4.1 Consider the Poisson matrix.

- What is the bandwidth of its Cholesky factor?
- approximately how many arithmetic operations does it take to find the Cholesky factor?

- same question for block LU,
- same question for the fast Poisson solver with and without FFT.

5.4.2 What is the discrete sine transform and discrete Fourier transform of a vector?

Part II

Some Matrix Theory

Chapter 6

Matrix Reduction by Similarity Transformations

A basic problem in numerical linear algebra is to compute eigenvalues and eigenvectors of a matrix \mathbf{A} . Before attempting to find eigenvalues and eigenvectors of \mathbf{A} (exceptions are made for certain sparse matrices), it should be reduced by similarity transformations to a simpler form. For example, if the simpler form is a triangular matrix \mathbf{R} , then the diagonal elements of \mathbf{R} are the eigenvalues of \mathbf{A} . The contents of this chapter is mainly theoretical, but the results are useful in numerical analysis.

6.1 Orthonormal, Unitary, and Similar Matrices

Orthogonal and unitary similarity transformations are particularly important since they are insensitive to noise in the elements of the matrix. We start by reviewing some basic facts about matrices with orthonormal columns.

6.1.1 Orthonormal and Unitary Matrices

Definition 6.1 (Orthonormal matrix)

A matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is said to be orthonormal if $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$.

Warning: An orthonormal matrix is often called an “orthogonal matrix” in the literature.

Theorem 6.2 (Orthonormal matrix)

Suppose $\mathbf{Q} \in \mathbb{R}^{n \times n}$ and let $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{y}$ be the standard inner product on \mathbb{R}^n . The following are equivalent:

1. \mathbf{Q} is orthonormal,

2. $\mathbf{Q}^{-1} = \mathbf{Q}^T$,
3. $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}$,
4. $\|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \mathbb{R}^n$.
5. $\langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.
6. the columns of \mathbf{Q} form an orthonormal basis for \mathbb{R}^n ,
7. the rows of \mathbf{Q} form an orthonormal basis for \mathbb{R}^n .

Moreover, the product of two orthonormal matrices is orthonormal.

Proof. Let $\mathbf{q}_1, \dots, \mathbf{q}_n$ be the columns of \mathbf{Q} .

1 \Rightarrow 2 Since \mathbf{Q} is square its left inverse \mathbf{Q}^T is the inverse.

2 \Rightarrow 3 \mathbf{Q}^T is a right inverse since it is an inverse.

3 \Rightarrow 1 \mathbf{Q}^T is a leftinverse since it is a right inverse.

1 \Rightarrow 4 $\|\mathbf{Q}\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q}\mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2$ for all $\mathbf{x} \in \mathbb{R}^n$.

4 \Rightarrow 5 We have

$$\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x} + \mathbf{y}\|_2^2 - \|\mathbf{x}\|_2^2 - \|\mathbf{y}\|_2^2 = \|\mathbf{Q}\mathbf{x} + \mathbf{Q}\mathbf{y}\|_2^2 - \|\mathbf{Q}\mathbf{x}\|_2^2 - \|\mathbf{Q}\mathbf{y}\|_2^2 = 2\langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{y} \rangle.$$

5 \Rightarrow 1 Taking $\mathbf{x} = \mathbf{e}_i$ and $\mathbf{y} = \mathbf{e}_j$ we find $(\mathbf{Q}^T \mathbf{Q})_{ij} = \langle \mathbf{q}_i, \mathbf{q}_j \rangle = \langle \mathbf{Q}\mathbf{e}_i, \mathbf{Q}\mathbf{e}_j \rangle = \langle \mathbf{e}_i, \mathbf{e}_j \rangle = \delta_{ij}$ for all $i, j = 1, \dots, n$.

1 \Leftrightarrow 6 Since $(\mathbf{Q}^T \mathbf{Q})_{ij} = \langle \mathbf{q}_i, \mathbf{q}_j \rangle$ we see that $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ if and only if $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ is an orthonormal set and hence linearly independent. The spanning property holds since the number of columns of \mathbf{Q} equals the dimension of \mathbb{R}^n .

3 \Leftrightarrow 7 The transpose of the rows of \mathbf{Q} are the columns of \mathbf{Q}^T and $(\mathbf{Q}^T)^T = \mathbf{Q}$. Therefore 3 \Leftrightarrow 7 follows from 1 \Leftrightarrow 6 applied to \mathbf{Q}^T .

We have shown that 1, 2, ..., 7 are equivalent. Finally, suppose \mathbf{Q}_1 and \mathbf{Q}_2 are orthonormal. Then $(\mathbf{Q}_1 \mathbf{Q}_2)^T \mathbf{Q}_1 \mathbf{Q}_2 = \mathbf{Q}_2^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{Q}_2 = \mathbf{I}$ so the product $\mathbf{Q}_1 \mathbf{Q}_2$ is orthonormal. . \square

Consider now the complex case.

Definition 6.3 (Unitary matrix)

A matrix $\mathbf{U} \in \mathbb{C}^{n \times n}$ is said to be **unitary** if $\mathbf{U}^* \mathbf{U} = \mathbf{I}$.

Note that a real unitary matrix is orthonormal.

Theorem 6.4 (Unitary matrix)

Suppose $\mathbf{U} \in \mathbb{C}^{n \times n}$. The following are equivalent:

1. \mathbf{U} is unitary,
2. $\mathbf{U}^{-1} = \mathbf{U}^*$,
3. $\mathbf{U}\mathbf{U}^* = \mathbf{I}$,
4. $\|\mathbf{U}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \mathbb{C}^n$.
5. $\langle \mathbf{U}\mathbf{x}, \mathbf{U}\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$,
6. the columns of \mathbf{U} form an orthonormal basis for \mathbb{C}^n ,
7. the rows of \mathbf{U} form an orthonormal basis for \mathbb{C}^n .

The product $\mathbf{U}_1\mathbf{U}_2$ of two unitary matrices \mathbf{U}_1 and \mathbf{U}_2 is unitary.

Proof. The proof is almost identical to the proof of the real case. For the proof of 4 \Rightarrow 5 we use (25). \square

Exercise 6.5 (Unitary matrix)

Provide the details of the proof of Theorem 6.4.

6.1.2 Similarity Transformations

Row operations are used in Gaussian elimination to reduce a matrix to triangular form, but row operations change the eigenvalues of a matrix. We need a transformation which can be used to simplify a matrix without changing the eigenvalues.

Definition 6.6 (Similar matrices)

Two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ are said to be **similar** if there is a nonsingular matrix $\mathbf{S} \in \mathbb{C}^{n \times n}$ such that $\mathbf{B} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}$. The transformation $\mathbf{A} \rightarrow \mathbf{B}$ is called a **similarity transformation**. It is called a **unitary similarity transformation** if \mathbf{S} is unitary and an **orthonormal similarity transformation** if $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ and $\mathbf{S} \in \mathbb{R}^{n \times n}$ is orthonormal.

Note that $\mathbf{S}^{-1} = \mathbf{S}^*$ for a unitary similarity transformation.

Theorem 6.7 (Eigenvalues of similar matrices)

Similar matrices have the same characteristic polynomial and therefore the same eigenvalues.

Proof. Let $\mathbf{B} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}$. By properties of determinants

$$\begin{aligned}\pi_{\mathbf{B}}(\lambda) &= \det(\mathbf{S}^{-1}\mathbf{A}\mathbf{S} - \lambda\mathbf{I}) = \det(\mathbf{S}^{-1}(\mathbf{A} - \lambda\mathbf{I})\mathbf{S}) \\ &= \det(\mathbf{S}^{-1}) \det(\mathbf{A} - \lambda\mathbf{I}) \det(\mathbf{S}) = \det(\mathbf{S}^{-1}\mathbf{S}) \det(\mathbf{A} - \lambda\mathbf{I}) = \pi_{\mathbf{A}}(\lambda).\end{aligned}$$

But then \mathbf{A} and \mathbf{B} have the same characteristic polynomial and hence the same eigenvalues. \square

As a corollary we have the following useful result.

Corollary 6.8 (Spectra of AB and BA)

For any $A \in \mathbb{C}^{m \times n}$ and $B \in \mathbb{C}^{n \times m}$ the matrices AB and BA have the same spectrum apart from a possible zero eigenvalue. More precisely,

$$\lambda^n \pi_{AB}(\lambda) = \lambda^m \pi_{BA}(\lambda), \quad \lambda \in \mathbb{C}.$$

Proof. Define block matrices of order $n + m$ by

$$E = \begin{bmatrix} AB & \mathbf{0} \\ B & \mathbf{0} \end{bmatrix}, \quad F = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ B & BA \end{bmatrix}, \quad S = \begin{bmatrix} I & A \\ \mathbf{0} & I \end{bmatrix}.$$

By Property 6. of Theorem 0.66 we have $\pi_E(\lambda) = \lambda^n \pi_{AB}(\lambda)$ and $\pi_F(\lambda) = \lambda^m \pi_{BA}(\lambda)$. But $ES = SF$ so E and F are similar and have the same characteristic polynomial by Theorem 6.7. \square

The case where $S^{-1}AS$ is a diagonal matrix is of special interest.

Theorem 6.9 (Eigenvectors of nondefective matrices)

If $S^{-1}AS = \text{diag}(\lambda_1, \dots, \lambda_n)$ then $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A , and the columns of S are linearly independent eigenvectors of A .

Proof. By Theorem 6.7 the diagonal elements of the diagonal matrix are the eigenvalues of A . We find $AS = S \text{diag}(\lambda_1, \dots, \lambda_n)$ or $As_i = \lambda_i s_i$, $i = 1, \dots, n$, where $S = [s_1, \dots, s_n]$. Since S has an inverse it is nonsingular and the columns of S are linearly independent eigenvectors of A . \square

A matrix with a full set of linearly independent eigenvectors has a special name.

Definition 6.10 (Defective and nondefective matrix)

A matrix is **nondefective** if the eigenvectors are linearly independent. and **defective** otherwise.

If $A \in \mathbb{C}^{n \times n}$ is nondefective with linearly independent eigenvectors v_1, \dots, v_n then these eigenvectors form a basis for \mathbb{C}^n and an $x \in \mathbb{C}^n$ can be written $x = \sum_{j=1}^n c_j v_j$ for some scalars c_1, \dots, c_n . We call this an **eigenvector expansion** of x .

Theorem 6.9 implies

Corollary 6.11 (Diagonalizable matrix)

A matrix is diagonalizable if and only if it is nondefective.

Proof. If \mathbf{A} is diagonalizable then it is nondefective by Theorem 6.9. Conversely, let the columns of $\mathbf{S} := [\mathbf{s}_1, \dots, \mathbf{s}_n]$ be the eigenvectors of \mathbf{A} . Since \mathbf{A} is nondefective \mathbf{S} is nonsingular, and by the proof of Theorem 6.9 $\mathbf{S}^{-1}\mathbf{AS}$ is diagonal. \square

We also state

Corollary 6.12 (Eigenvectors of \mathbf{A}^*)

If \mathbf{A} is nondefective then \mathbf{A}^* is nondefective, and if $\mathbf{S}^{-1}\mathbf{AS} = \text{diag}(\lambda_1, \dots, \lambda_n)$ then the columns of \mathbf{S}^{-*} are the eigenvectors of \mathbf{A}^* . Moreover, the eigenvalues of \mathbf{A}^* are the complex conjugates of the eigenvalues of \mathbf{A} .

Proof. If $\mathbf{S}^{-1}\mathbf{AS} = \text{diag}(\lambda_1, \dots, \lambda_n)$ then $\mathbf{A}^*\mathbf{S}^{-*} = \mathbf{S}^{-*}\text{diag}(\overline{\lambda_1}, \dots, \overline{\lambda_n})$. Since \mathbf{S}^{-*} is nonsingular the results follow. \square

6.2 Linear Independence of Eigenvectors

For distinct eigenvalues we have the following result.

Theorem 6.13 (Distinct eigenvalues)

Eigenvectors corresponding to distinct eigenvalues are linearly independent.

Proof. Let m be the smallest positive integer so that $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly dependent eigenvectors of \mathbf{A} . Clearly $m \geq 2$ since $\mathbf{x}_1 \neq \mathbf{0}$. For some nonzero $[c_1, \dots, c_m]$ we have

$$\sum_{j=1}^m c_j \mathbf{x}_j = \mathbf{0}. \quad (6.1)$$

Applying \mathbf{A} to (6.1) we obtain by linearity $\sum_{j=1}^m c_j \lambda_j \mathbf{x}_j = \mathbf{0}$. From this relation we subtract λ_m times (6.1) and find $\sum_{j=1}^{m-1} c_j (\lambda_j - \lambda_m) \mathbf{x}_j = \mathbf{0}$. But since $\lambda_j - \lambda_m \neq 0$ for $j = 1, \dots, m-1$ and at least one $c_j \neq 0$ for $j < m$ we see that $\{\mathbf{x}_1, \dots, \mathbf{x}_{m-1}\}$ is linearly dependent, contradicting the minimality of m . \square

Corollary 6.14 (Nondefective matrix)

A matrix with distinct eigenvalues is nondefective.

Proof. By the previous theorem the eigenvectors are linearly independent. \square

6.2.1 Algebraic and Geometric Multiplicity of Eigenvalues

A defective matrix must necessarily have one or more multiple eigenvalues, but as the following example shows this is not sufficient.

Example 6.15 (Two upper triangular matrices)

Consider the 2 matrices of order 3

$$\mathbf{A}_1 := \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

Both matrices are upper triangular and have an eigenvalue $\lambda = 1$ of (algebraic) multiplicity 3.

1. The eigenvectors of \mathbf{A}_1 are the unit vectors $\mathbf{x}_i = \mathbf{e}_i$, $i = 1, 2, 3$. Thus \mathbf{A}_1 is nondefective.
2. An eigenvector $\mathbf{x} = [x_1, x_2, x_3]^T$ of \mathbf{A}_2 must be a solution of the homogenous triangular linear system

$$(\mathbf{A} - \mathbf{I})\mathbf{x} = \mathbf{0} \text{ or } \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

But then $x_2 = x_3 = 0$ and any eigenvector must be a multiple of \mathbf{e}_1 . We conclude that \mathbf{A}_2 is defective.

For multiple eigenvalues we need to distinguish between two kinds of multiplicities. We will show that a matrix is nondefective if and only if the two multiplicities are the same for all eigenvalues.

The eigenvalues $\lambda_1, \dots, \lambda_n$ of $\mathbf{A} \in \mathbb{C}^{n \times n}$ are the roots of the characteristic polynomial

$$\pi_{\mathbf{A}}(\lambda) := \det(\mathbf{A} - \lambda\mathbf{I}) = (\lambda_1 - \lambda) \cdots (\lambda_n - \lambda),$$

or in terms of multiplicities a_i of distinct eigenvalues μ_i

$$\pi_{\mathbf{A}}(\lambda) = (\mu_1 - \lambda)^{a_1} \cdots (\mu_r - \lambda)^{a_r}, \quad \mu_i \neq \mu_j, \quad i \neq j, \quad \sum_{i=1}^r a_i = n, \quad (6.2)$$

Thus, $\sigma(\mathbf{A}) = \{\lambda_1, \dots, \lambda_n\} = \{\mu_1, \dots, \mu_r\}$ and the positive integer $a_i = a(\mu_i) = a_{\mathbf{A}}(\mu_i)$ is called the **algebraic multiplicity** of μ_i . An eigenvalue λ is simple (double, triple) if a_i is equal to one (two, three).

To define a second kind of multiplicity we consider for each $\lambda \in \sigma(\mathbf{A})$ the nullspace

$$\ker(\mathbf{A} - \lambda\mathbf{I}) := \{\mathbf{x} \in \mathbb{C}^n : (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}\} \quad (6.3)$$

of $\mathbf{A} - \lambda\mathbf{I}$. The nullspace is a subspace of \mathbb{C}^n consisting of all eigenvectors of \mathbf{A} corresponding to the eigenvalue λ . The dimension of the subspace must be at least one since $\mathbf{A} - \lambda\mathbf{I}$ is singular.

Definition 6.16 (Geometric multiplicity)

The geometric multiplicity $g = g(\lambda) = g_{\mathbf{A}}(\lambda)$ of an eigenvalue λ of \mathbf{A} is the dimension of the nullspace $\ker(\mathbf{A} - \lambda\mathbf{I})$.

Example 6.17 (Geometric multiplicity)

The $n \times n$ identity matrix \mathbf{I} has the eigenvalue $\lambda = 1$ with $\pi_{\mathbf{I}}(\lambda) = (1 - \lambda)^n$. Since $\mathbf{I} - \lambda\mathbf{I}$ is the zero matrix when $\lambda = 1$, the nullspace of $\mathbf{I} - \lambda\mathbf{I}$ is all of n -space and it follows that $a = g = n$. On the other hand we saw in Example 6.15 that the 3×3 matrix $\mathbf{A} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$ has the eigenvalue $\lambda = 1$ with $a = 3$ and $g = 1$.

The geometric multiplicity of an eigenvalue is always bounded above by the algebraic multiplicity of the eigenvalue.

Theorem 6.18 ($g \leq a$)

For any square matrix \mathbf{A} and any $\lambda \in \sigma(\mathbf{A})$ we have $g_{\mathbf{A}}(\lambda) \leq a_{\mathbf{A}}(\lambda)$.

Proof. Let $\{\mathbf{v}_1, \dots, \mathbf{v}_g\}$ with $g := g_{\mathbf{A}}(\lambda)$ be an orthonormal basis for $\ker(\mathbf{A} - \lambda\mathbf{I})$, and extend this set to an orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ for \mathbb{C}^n . By Theorem 6.4 the matrix $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n \times n}$ is unitary and $\mathbf{V}^{-1} = \mathbf{V}^*$. Partition \mathbf{V} as $\mathbf{V} = [\mathbf{V}_1, \mathbf{V}_2]$, where $\mathbf{V}_1 := [\mathbf{v}_1, \dots, \mathbf{v}_g]$ and $\mathbf{V}_2 := [\mathbf{v}_{g+1}, \dots, \mathbf{v}_n]$. Then $\mathbf{A}\mathbf{V}_1 = \lambda\mathbf{V}_1$, $\mathbf{V}_1^*\mathbf{V}_1 = \mathbf{I}$, $\mathbf{V}_2^*\mathbf{V}_1 = \mathbf{0}$, and

$$\mathbf{B} := \mathbf{V}^*\mathbf{A}\mathbf{V} = \begin{bmatrix} \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} \mathbf{A} \begin{bmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{V}_1^*\mathbf{A}\mathbf{V}_1 & \mathbf{V}_1^*\mathbf{A}\mathbf{V}_2 \\ \mathbf{V}_2^*\mathbf{A}\mathbf{V}_1 & \mathbf{V}_2^*\mathbf{A}\mathbf{V}_2 \end{bmatrix} = \begin{bmatrix} \lambda\mathbf{I} & \mathbf{V}_1^*\mathbf{A}\mathbf{V}_2 \\ \mathbf{0} & \mathbf{C} \end{bmatrix},$$

where $\mathbf{C} := \mathbf{V}_2^*\mathbf{A}\mathbf{V}_2$. Since \mathbf{B} is block triangular, Property 6. of Theorem 6.66 implies that $\pi_{\mathbf{B}}(z) = (z - \lambda)^g \pi_{\mathbf{C}}(z)$. But then $a_{\mathbf{B}}(\lambda) \geq g$. Since \mathbf{A} and \mathbf{B} are similar they have the same characteristic polynomial, and it follows that $a_{\mathbf{A}}(\lambda) = a_{\mathbf{B}}(\lambda) \geq g_{\mathbf{A}}(\lambda)$. \square

Definition 6.19 (Defective eigenvalue)

An eigenvalue λ satisfying $g_{\mathbf{A}}(\lambda) < a_{\mathbf{A}}(\lambda)$ is said to be **defective**.

By Example 6.17 the eigenvalue $\lambda = 1$ of $\begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$ is defective.

Theorem 6.20 (The number of linearly independent eigenvectors)

The number of linearly independent eigenvectors of a matrix equals the sum of the geometric multiplicities of the eigenvalues.

Proof. Let \mathbf{A} have distinct eigenvalues μ_1, \dots, μ_r with algebraic multiplicities a_1, \dots, a_r and geometric multiplicities g_1, \dots, g_r . It clearly holds for $r = 1$. Suppose $\{\mathbf{v}_{j,1}, \dots, \mathbf{v}_{j,g_j}\}$ is a basis for $\ker(\mathbf{A} - \mu_j \mathbf{I})$ for $j = 1, \dots, r$. We claim that the combined set $\{\mathbf{v}_{j,k}\}_{k=1,j=1}^{g_j,r}$ is linearly independent. We show this using induction on r . Suppose $\{\mathbf{v}_{jk}\}_{k=1,j=1}^{g_j,r-1}$ is linearly independent and assume

$$\sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} \mathbf{v}_{jk} = \mathbf{0} \text{ for some scalars } a_{jk}. \quad (6.4)$$

Now multiply this equation by $(\mathbf{A} - \mu_r \mathbf{I})$. Since $\mathbf{A}\mathbf{v}_{j,k} = \mu_j \mathbf{v}_{j,k}$ we obtain

$$\mathbf{0} = \sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} (\mathbf{A} - \mu_r \mathbf{I}) \mathbf{v}_{jk} = \sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} (\mu_j - \mu_r) \mathbf{v}_{jk} = \sum_{j=1}^{r-1} \sum_{k=1}^{g_j} a_{jk} (\mu_j - \mu_r) \mathbf{v}_{jk}.$$

By the induction hypothesis $a_{jk}(\mu_j - \mu_r) = 0$ and hence $a_{jk} = 0$ for $j < r$. In (6.4) we are left with $\sum_{k=1}^{g_r} a_{rk} \mathbf{v}_{rk} = \mathbf{0}$. But since $\{\mathbf{v}_{j,r}\}$ are linearly independent we also have $a_{rk} = 0$ for $k = 1, \dots, g_r$. Thus $\{\mathbf{v}_{jk}\}_{k=1,j=1}^{g_j,r}$ is linearly independent and it follows that the number of linearly independent eigenvectors is equal to $\sum_j g_j$. \square

Corollary 6.21 (Linearly independent eigenvectors characterization)

A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ has n linearly independent eigenvectors if and only if the algebraic and geometric multiplicity of all eigenvalues are the same.

Proof. Since $g_j \leq a_j$ for all j and $\sum_j a_j = n$ we have $\sum_j g_j = n$ if and only if $a_j = g_j$ for $j = 1, \dots, r$. \square

Theorem 6.22 (Geometric multiplicity of similar matrices)

Similar matrices have the same eigenvalues with the same algebraic and geometric multiplicities.

Proof. Similar matrices have the same characteristic polynomials and only the invariance of geometric multiplicity needs to be shown. Suppose $\lambda \in \sigma(\mathbf{A})$, $\dim \ker(\mathbf{S}^{-1} \mathbf{AS} - \lambda \mathbf{I}) = k$, and $\dim \ker(\mathbf{A} - \lambda \mathbf{I}) = \ell$. We need to show that $k = \ell$. Suppose $\mathbf{v}_1, \dots, \mathbf{v}_k$ is a basis for $\ker(\mathbf{S}^{-1} \mathbf{AS} - \lambda \mathbf{I})$. Then $\mathbf{S}^{-1} \mathbf{AS} \mathbf{v}_i = \lambda \mathbf{v}_i$ or $\mathbf{AS} \mathbf{v}_i = \lambda \mathbf{S} \mathbf{v}_i$, $i = 1, \dots, k$. But then $\{\mathbf{S} \mathbf{v}_1, \dots, \mathbf{S} \mathbf{v}_k\} \subset \ker(\mathbf{A} - \lambda \mathbf{I})$, which implies that $k \leq \ell$. \square

Exercise 6.23 (Find eigenpair example)

Find eigenvalues and eigenvectors of $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 2 \end{bmatrix}$.

Exercise 6.24 (Idempotent matrix)

Let $\lambda \in \sigma(\mathbf{A})$ where $\mathbf{A}^2 = \mathbf{A} \in \mathbb{C}^{n \times n}$. Show that $\lambda = 0$ or $\lambda = 1$. (A matrix is called **idempotent** if $\mathbf{A}^2 = \mathbf{A}$).

Exercise 6.25 (Idempotent matrix)

Let $\lambda \in \sigma(\mathbf{A})$ where $\mathbf{A}^k = 0$ for some $k \in \mathbb{N}$. Show that $\lambda = 0$. (A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ such that $\mathbf{A}^k = 0$ for some $k \in \mathbb{N}$ is called **nilpotent**).

Exercise 6.26 (Eigenvalues of a unitary matrix)

Let $\lambda \in \sigma(\mathbf{A})$, where $\mathbf{A}^* \mathbf{A} = \mathbf{I}$. Show that $|\lambda| = 1$.

Exercise 6.27 (Nonsingular approximation of a singular matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is singular. Then we can find $\epsilon_0 > 0$ such that $\mathbf{A} + \epsilon \mathbf{I}$ is nonsingular for all $\epsilon \in (0, \epsilon_0)$. Hint: $\det(\mathbf{A}) = \lambda_1 \lambda_2 \cdots \lambda_n$, where λ_i are the eigenvalues of \mathbf{A} .

Exercise 6.28 (Companion matrix)

For $q_0, \dots, q_{n-1} \in \mathbb{C}$ let $f(\lambda) = \lambda^n + q_{n-1}\lambda^{n-1} + \cdots + q_0$ be a polynomial of degree n in λ . We derive two matrices which have $(-1)^n f$ as its characteristic polynomial.

a) Show that $f = (-1)^n \pi_{\mathbf{A}}$ where

$$\mathbf{A} = \begin{bmatrix} -q_{n-1} & -q_{n-2} & \cdots & -q_1 & -q_0 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}.$$

\mathbf{A} is called the **companion matrix** of f .

b) Show that $f = (-1)^n \pi_{\mathbf{A}'}$ where

$$\mathbf{A}' = \begin{bmatrix} 0 & 0 & \cdots & 0 & -q_0 \\ 1 & 0 & \cdots & 0 & -q_1 \\ 0 & 1 & \cdots & 0 & -q_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -q_{n-1} \end{bmatrix}.$$

Thus \mathbf{A}' can also be regarded as a companion matrix for f .

6.3 Normal Matrices

In this section we consider the reduction of a matrix to triangular, or almost triangular form, and characterize matrices with orthonormal eigenvectors.

6.3.1 The Schur Decomposition

Although not every matrix can be diagonalized it can be brought into triangular form by a *unitary* similarity transformation.

Theorem 6.29 (Schur decomposition)

For each $\mathbf{A} \in \mathbb{C}^{n \times n}$ there exists a unitary matrix $\mathbf{U} \in \mathbb{C}^{n \times n}$ such that $\mathbf{R} := \mathbf{U}^* \mathbf{A} \mathbf{U}$ is upper triangular.

Proof. We use induction on n . For $n = 1$ the matrix \mathbf{U} is the 1×1 identity matrix. Assume that the theorem is true for matrices of order k and suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$, where $n := k + 1$. Let $(\lambda_1, \mathbf{v}_1)$ be an eigenpair for \mathbf{A} with $\|\mathbf{v}_1\|_2 = 1$. By Theorem 0.39 we can extend \mathbf{v}_1 to an orthonormal basis $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ for \mathbb{C}^n . The matrix $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n \times n}$ is unitary, and the first column of the product $\mathbf{V}^* \mathbf{A} \mathbf{V}$ is

$$\mathbf{V}^* \mathbf{A} \mathbf{V} e_1 = \mathbf{V}^* \mathbf{A} \mathbf{v}_1 = \lambda_1 \mathbf{V}^* \mathbf{v}_1 = \lambda_1 e_1.$$

It follows that

$$\mathbf{V}^* \mathbf{A} \mathbf{V} = \left[\begin{array}{c|c} \lambda_1 & \mathbf{x}^* \\ \mathbf{0} & \mathbf{M} \end{array} \right], \text{ for some } \mathbf{M} \in \mathbb{C}^{k \times k} \text{ and } \mathbf{x} \in \mathbb{C}^k. \quad (6.5)$$

By the induction hypothesis there is a unitary matrix $\mathbf{W}_1 \in \mathbb{C}^{k \times k}$ such that $\mathbf{W}_1^* \mathbf{M} \mathbf{W}_1$ is upper triangular. Define

$$\mathbf{W} = \left[\begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1 \end{array} \right] \text{ and } \mathbf{U} = \mathbf{V} \mathbf{W}.$$

Then \mathbf{W} and \mathbf{U} (cf. Theorem 6.2) are unitary and

$$\begin{aligned} \mathbf{U}^* \mathbf{A} \mathbf{U} &= \mathbf{W}^* (\mathbf{V}^* \mathbf{A} \mathbf{V}) \mathbf{W} = \left[\begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1^* \end{array} \right] \left[\begin{array}{c|c} \lambda_1 & \mathbf{x}^* \\ \mathbf{0} & \mathbf{M} \end{array} \right] \left[\begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1 \end{array} \right] \\ &= \left[\begin{array}{c|c} \lambda_1 & \mathbf{x}^* \mathbf{W}_1 \\ \mathbf{0} & \mathbf{W}_1^* \mathbf{M} \mathbf{W}_1 \end{array} \right], \end{aligned}$$

is upper triangular. \square

By using the unitary transformation \mathbf{V} on the $n \times n$ matrix \mathbf{A} , we obtain a matrix \mathbf{M} of order $n - 1$. \mathbf{M} has the same eigenvalues as \mathbf{A} except λ . Thus we can find another eigenvalue of \mathbf{A} by working with a smaller matrix \mathbf{M} . This is an example of a **deflation** technique which is very useful in numerical work.

Example 6.30 (Deflation example)

The matrix $\mathbf{T} := \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$ has an eigenpair $(2, \mathbf{x}_1)$, where $\mathbf{x}_1 = [-1, 0, 1]^T$.

We can extend \mathbf{x}_1 to a basis $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ for \mathbb{R}^3 by defining $\mathbf{x}_2 = [0, 1, 0]^T$, $\mathbf{x}_3 = [1, 0, 1]^T$. This is already an orthogonal basis and normalizing we obtain the orthonormal matrix

$$\mathbf{V} = \begin{bmatrix} -\frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

We obtain (6.5) with $\lambda = 2$ and

$$\mathbf{C} = \begin{bmatrix} 2 & -\sqrt{2} \\ -\sqrt{2} & 2 \end{bmatrix}.$$

We can now find the remaining eigenvalues of \mathbf{A} from the 2×2 matrix \mathbf{C} .

If \mathbf{A} has complex eigenvalues then \mathbf{U} will be complex even if \mathbf{A} is real. The following is a real version of Theorem 6.29.

Theorem 6.31 (Schur form, real eigenvalues)

For each $\mathbf{A} \in \mathbb{R}^{n \times n}$ with real eigenvalues there exists a matrix $\mathbf{U} \in \mathbb{R}^{n \times n}$ with $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, such that $\mathbf{U}^T \mathbf{A} \mathbf{U}$ is upper triangular.

Proof. Consider the proof of Theorem 6.29. Since \mathbf{A} and λ_1 are real the eigenvector \mathbf{v}_1 is real and the matrix \mathbf{W} is real and $\mathbf{W}^T \mathbf{W} = \mathbf{I}$. By the induction hypothesis \mathbf{V} is real and $\mathbf{V}^T \mathbf{V} = \mathbf{I}$. But then also $\mathbf{U} = \mathbf{V} \mathbf{W}$ is real and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$. \square

Exercise 6.32 (Schur decomposition example)

Show that a Schur decomposition of $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 2 \end{bmatrix}$ is $\mathbf{U}^T \mathbf{A} \mathbf{U} = \begin{bmatrix} -1 & -1 \\ 0 & 4 \end{bmatrix}$, where $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$.

The matrices \mathbf{U} and \mathbf{R} in the Schur decomposition are called **Schur factors**.

A real matrix with complex eigenvalues cannot be reduced to triangular form by a real unitary similarity transformation. For example, the matrix $\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ has eigenvalues $\pm i$ and is not similar to a real triangular matrix. How far can we reduce a real matrix \mathbf{A} by a real unitary similarity transformation? To study this we note that the complex eigenvalues of \mathbf{A} occur in conjugate pairs, $\lambda = \mu + i\nu$, $\bar{\lambda} = \mu - i\nu$, where μ, ν are real. The real 2×2 matrix

$$\mathbf{M} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix} \tag{6.6}$$

has eigenvalues $\lambda = \mu + i\nu$ and $\bar{\lambda} = \mu - i\nu$.

Definition 6.33 (Quasi-triangular matrix)

We say that a matrix is **quasi-triangular** if it is block triangular with only 1×1 and 2×2 blocks on the diagonal. Moreover, no 2×2 block should have real eigenvalues.

As an example consider

$$\mathbf{R} = \left[\begin{array}{cc|cc|cc} 2 & 1 & 3 & 4 & 5 \\ -1 & 2 & 4 & 3 & 2 \\ \hline 0 & 0 & 1 & 2 & 3 \\ \hline 0 & 0 & 0 & 3 & 2 \\ 0 & 0 & 0 & -1 & 1 \end{array} \right].$$

\mathbf{R} has three diagonal blocks:

$$\mathbf{D}_1 = \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{D}_2 = [1], \quad \mathbf{D}_3 = \begin{bmatrix} 3 & 2 \\ -1 & 1 \end{bmatrix}.$$

By Theorem 0.66 the eigenvalues of \mathbf{R} are the union of the eigenvalues of \mathbf{D}_1 , \mathbf{D}_2 and \mathbf{D}_3 . The eigenvalues of \mathbf{D}_1 are $2+i$ and $2-i$, while \mathbf{D}_2 has eigenvalue 1, and \mathbf{D}_3 has the same eigenvalues as \mathbf{D}_1 . Thus \mathbf{R} has one real eigenvalue 1 corresponding to the 1×1 block and complex eigenvalues $2+i, 2-i$ with multiplicity 2 corresponding to the two 2×2 blocks.

Any $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be reduced to quasi-triangular form by a real orthonormal similarity transformation. A proof is given in Section 6.7.

6.3.2 Matrices with Orthonormal Eigenvectors

It is possible to characterize matrices that have a diagonal Schur factorization.

Definition 6.34 (Normal Matrix)

A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is said to be **normal** if $\mathbf{A}\mathbf{A}^* = \mathbf{A}^*\mathbf{A}$.

Examples of normal matrices are

1. $\mathbf{A}^* = \mathbf{A}$, (Hermitian)
2. $\mathbf{A}^* = -\mathbf{A}$, (Skew-Hermitian)
3. $\mathbf{A}^* = \mathbf{A}^{-1}$, (Unitary)
4. $\mathbf{A} = \mathbf{D}$. (Diagonal)

The 2. derivative matrix \mathbf{T} in (2.2) and the discrete Poisson matrix (cf. Lemma 4.13) are examples of normal matrices.

Exercise 6.35 (Skew-Hermitian matrix)

Suppose $\mathbf{C} = \mathbf{A} + i\mathbf{B}$, where $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. Show that \mathbf{C} is skew-Hermitian if and only if $\mathbf{A}^T = -\mathbf{A}$ and $\mathbf{B}^T = \mathbf{B}$.

Exercise 6.36 (Eigenvalues of a skew-Hermitian matrix)

Show that any eigenvalue of a skew-Hermitian matrix is purely imaginary.

The following theorem says that a matrix has orthonormal eigenpairs if and only if it is normal.

Theorem 6.37 (Orthonormal eigenpairs characterization)

A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is unitary similar with a diagonal matrix if and only if it is normal.

Proof. If $\mathbf{B} = \mathbf{U}^* \mathbf{A} \mathbf{U}$, with \mathbf{B} diagonal, and $\mathbf{U}^* \mathbf{U} = \mathbf{I}$, then

$$\begin{aligned}\mathbf{A} \mathbf{A}^* &= (\mathbf{U} \mathbf{B} \mathbf{U}^*)(\mathbf{U} \mathbf{B}^* \mathbf{U}^*) = \mathbf{U} \mathbf{B} \mathbf{B}^* \mathbf{U}^* \text{ and} \\ \mathbf{A}^* \mathbf{A} &= (\mathbf{U} \mathbf{B}^* \mathbf{U}^*)(\mathbf{U} \mathbf{B} \mathbf{U}^*) = \mathbf{U} \mathbf{B}^* \mathbf{B} \mathbf{U}^*.\end{aligned}$$

Now $\mathbf{B} \mathbf{B}^* = \mathbf{B}^* \mathbf{B}$ since \mathbf{B} is diagonal, and \mathbf{A} is normal.

Suppose $\mathbf{A}^* \mathbf{A} = \mathbf{A} \mathbf{A}^*$. By Theorem 6.29 we can find \mathbf{U} with $\mathbf{U}^* \mathbf{U} = \mathbf{I}$ such that $\mathbf{B} = \mathbf{U}^* \mathbf{A} \mathbf{U}$ is upper triangular. Since \mathbf{A} is normal \mathbf{B} is normal. Indeed,

$$\mathbf{B} \mathbf{B}^* = \mathbf{U}^* \mathbf{A} \mathbf{U} \mathbf{U}^* \mathbf{A}^* \mathbf{U} = \mathbf{U}^* \mathbf{A} \mathbf{A}^* \mathbf{U} = \mathbf{U}^* \mathbf{A}^* \mathbf{A} \mathbf{U} = \mathbf{B}^* \mathbf{B}.$$

The proof is complete if we can show that an upper triangular normal matrix \mathbf{B} must be diagonal. The diagonal elements in $\mathbf{E} := \mathbf{B}^* \mathbf{B}$ and $\mathbf{F} := \mathbf{B} \mathbf{B}^*$ are given by

$$e_{ii} = \sum_{k=1}^n \bar{b}_{ki} b_{ki} = \sum_{k=1}^i |b_{ki}|^2 \text{ and } f_{ii} = \sum_{k=1}^n b_{ik} \bar{b}_{ik} = \sum_{k=i}^n |b_{ik}|^2.$$

The result now follows by equating e_{ii} and f_{ii} for $i = 1, 2, \dots, n$. In particular for $i = 1$ we have $|b_{11}|^2 = |b_{11}|^2 + |b_{12}|^2 + \dots + |b_{1n}|^2$, so $b_{1k} = 0$ for $k = 2, 3, \dots, n$. Suppose $b_{jk} = 0$ for $j = 1, \dots, i-1$, $k = j+1, \dots, n$. Then

$$e_{ii} = \sum_{k=1}^i |b_{ki}|^2 = |b_{ii}|^2 = \sum_{k=i}^n |b_{ik}|^2 = f_{ii}$$

so $b_{ik} = 0$, $k = i+1, \dots, n$. By induction on the rows we see that \mathbf{B} is diagonal.

□

6.4 Hermitian Matrices

The special cases where \mathbf{A} is Hermitian, or real and symmetric, deserve special attention.

Theorem 6.38 (Spectral theorem, complex form)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is Hermitian. Then \mathbf{A} has real eigenvalues $\lambda_1, \dots, \lambda_n$. Moreover, there is a unitary matrix $\mathbf{U} \in \mathbb{C}^{n \times n}$ such that $\mathbf{U}^* \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n)$. For any such \mathbf{U} the columns $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ of \mathbf{U} are orthonormal eigenvectors of \mathbf{A} and $\mathbf{A}\mathbf{u}_j = \lambda_j \mathbf{u}_j$ for $j = 1, \dots, n$.

Proof. That the eigenvalues are real was shown in Lemma 4.12. The rest follows from Theorem 6.37. \square

There is also a real version.

Theorem 6.39 (Spectral Theorem (real form))

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{A}^T = \mathbf{A}$. Then \mathbf{A} has real eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Moreover, there is an orthonormal matrix $\mathbf{U} \in \mathbb{R}^{n \times n}$ such that $\mathbf{U}^T \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. For any such \mathbf{U} the columns $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ of \mathbf{U} are orthonormal eigenvectors of \mathbf{A} and $\mathbf{A}\mathbf{u}_j = \lambda_j \mathbf{u}_j$ for $j = 1, \dots, n$.

Proof. Since a real symmetric matrix has real eigenvalues and eigenvectors this follows from Theorem 6.38. \square

Example 6.40 The orthonormal diagonalization of $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$ is $\mathbf{U}^T \mathbf{A} \mathbf{U} = \text{diag}(1, 3)$, where $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$.

6.4.1 The Rayleigh Quotient

The Rayleigh quotient is an important tool when studying eigenvalues.

Definition 6.41 (Rayleigh quotient)

For $\mathbf{A} \in \mathbb{C}^{n \times n}$ and a nonzero \mathbf{x} the number

$$R(\mathbf{x}) = R_{\mathbf{A}}(\mathbf{x}) := \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}}$$

is called a **Rayleigh quotient**.

If (λ, \mathbf{x}) is an eigenpair for \mathbf{A} then $R(\mathbf{x}) = \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} = \lambda$.

Equation (6.7) in the following lemma shows that the Rayleigh quotient of a normal matrix is a **convex combination** of its eigenvalues.

Lemma 6.42 (Convex combination of the eigenvalues)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is normal with orthonormal eigenpairs $(\lambda_j, \mathbf{u}_j)$, $j = 1, 2, \dots, n$. Then the Rayleigh quotient is a convex combination of the eigenvalues of \mathbf{A}

$$R_{\mathbf{A}}(\mathbf{x}) = \frac{\sum_{i=1}^n \lambda_i |c_i|^2}{\sum_{j=1}^n |c_j|^2}, \quad \mathbf{x} \neq \mathbf{0}, \quad \mathbf{x} = \sum_{j=1}^n c_j \mathbf{u}_j. \quad (6.7)$$

Proof. By orthonormality of the eigenvectors $\mathbf{x}^* \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n \bar{c}_i \bar{u}_i c_j u_j = \sum_{j=1}^n |c_j|^2$. Similarly, $\mathbf{x}^* \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n \bar{c}_i \bar{u}_i c_j \lambda_j u_j = \sum_{i=1}^n \lambda_i |c_i|^2$. and (6.7) follows. This is clearly a combination of nonnegative quantities and a convex combination since $\sum_{i=1}^n |c_i|^2 / \sum_{j=1}^n |c_j|^2 = 1$. \square

6.4.2 Minmax Theorems

There are some useful characterizations of the eigenvalues of a Hermitian matrix. First we show

Theorem 6.43 (Minmax)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is Hermitian with eigenvalues $\lambda_1, \dots, \lambda_n$, ordered so that $\lambda_1 \geq \dots \geq \lambda_n$. Let $1 \leq k \leq n$. For any subspace \mathcal{S} of \mathbb{C}^n of dimension $n - k + 1$

$$\lambda_k \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad (6.8)$$

with equality for $\mathcal{S} = \tilde{\mathcal{S}} := \text{span}(\mathbf{u}_k, \dots, \mathbf{u}_n)$ and $\mathbf{x} = \mathbf{u}_k$. Here $(\lambda_j, \mathbf{u}_j)$, $1 \leq j \leq n$ are orthonormal eigenpairs for \mathbf{A} .

Proof. Let \mathcal{S} be any subspace of \mathbb{C}^n of dimension $n - k + 1$ and define $\mathcal{S}' := \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$. We need to find $\mathbf{y} \in \mathcal{S}$ so that $R(\mathbf{y}) \geq \lambda_k$. Now $\mathcal{S} + \mathcal{S}' := \{s + s' : s \in \mathcal{S}, s' \in \mathcal{S}'\}$ is a subspace of \mathbb{C}^n and by (7)

$$\dim(\mathcal{S} \cap \mathcal{S}') = \dim(\mathcal{S}) + \dim(\mathcal{S}') - \dim(\mathcal{S} + \mathcal{S}') \geq (n - k + 1) + k - n = 1.$$

It follows that $\mathcal{S} \cap \mathcal{S}'$ is nonempty. Let $\mathbf{y} \in \mathcal{S} \cap \mathcal{S}' = \sum_{j=1}^k c_j \mathbf{u}_j$ with $\sum_{j=1}^k |c_j|^2 = 1$. Defining $c_j = 0$ for $k+1 \leq j \leq n$, we obtain by Lemma 6.42

$$\max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \geq R(\mathbf{y}) = \sum_{j=1}^n \lambda_j |c_j|^2 = \sum_{j=1}^k \lambda_j |c_j|^2 \geq \sum_{j=1}^k \lambda_k |c_j|^2 = \lambda_k,$$

and (6.8) follows. If $\mathbf{y} \in \tilde{\mathcal{S}}$, say $\mathbf{y} = \sum_{j=k}^n d_j \mathbf{u}_j$ with $\sum_{j=k}^n |d_j|^2 = 1$ then again by Lemma 6.42 $R(\mathbf{y}) = \sum_{j=k}^n \lambda_j |d_j|^2 \leq \lambda_k$, and since $\mathbf{y} \in \tilde{\mathcal{S}}$ is arbitrary we have

$\max_{\substack{\mathbf{x} \in \tilde{\mathcal{S}} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \leq \lambda_k$ and equality in (6.8) follows for $\mathcal{S} = \tilde{\mathcal{S}}$. Moreover, $R(\mathbf{u}_k) = \lambda_k$. \square

There is also a maxmin version of this result.

Theorem 6.44 (Maxmin)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is Hermitian with eigenvalues $\lambda_1, \dots, \lambda_n$, ordered so that $\lambda_1 \geq \dots \geq \lambda_n$. Let $1 \leq k \leq n$. For any subspace \mathcal{S} of \mathbb{C}^n of dimension k

$$\lambda_k \geq \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad (6.9)$$

with equality for $\mathcal{S} = \tilde{\mathcal{S}} := \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ and $\mathbf{x} = \mathbf{u}_k$. Here $(\lambda_j, \mathbf{u}_j)$, $1 \leq j \leq n$ are orthonormal eigenpairs for \mathbf{A} .

Proof. The proof is very similar to the proof of Theorem 6.43. We define $\mathcal{S}' := \text{span}(\mathbf{u}_k, \dots, \mathbf{u}_n)$ and show that $R(\mathbf{y}) \leq \lambda_k$ for some $\mathbf{y} \in \mathcal{S} \cap \mathcal{S}'$. It is easy to see that $R(\mathbf{y}) \geq \lambda_k$ for any $\mathbf{y} \in \tilde{\mathcal{S}}$. \square

These theorems immediately lead to classical minmax and maxmin characterizations.

Corollary 6.45 (The Courant-Fischer Theorem)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is Hermitian with eigenvalues $\lambda_1, \dots, \lambda_n$, ordered so that $\lambda_1 \geq \dots \geq \lambda_n$. Then

$$\lambda_k = \min_{\dim(\mathcal{S})=n-k+1} \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) = \max_{\dim(\mathcal{S})=k} \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad k = 1, \dots, n. \quad (6.10)$$

Using Theorem 6.43 we can prove inequalities of eigenvalues without knowing the eigenvectors and we can get both upper and lower bounds.

Theorem 6.46 (Eigenvalue perturbation for Hermitian matrices)

Let $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ be Hermitian with eigenvalues $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ and $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$. Then

$$\alpha_k + \varepsilon_n \leq \beta_k \leq \alpha_k + \varepsilon_1, \quad \text{for } k = 1, \dots, n, \quad (6.11)$$

where $\varepsilon_1 \geq \varepsilon_2 \geq \dots \geq \varepsilon_n$ are the eigenvalues of $\mathbf{E} := \mathbf{B} - \mathbf{A}$.

Proof. Since \mathbf{E} is a sum of Hermitian matrices it is Hermitian and the eigenvalues are real. Let (α_j, \mathbf{u}_j) , $j = 1, \dots, n$ be orthonormal eigenpairs for \mathbf{A} and let $\mathcal{S} := \text{span}\{\mathbf{u}_k, \dots, \mathbf{u}_n\}$. By Theorem 6.43 we obtain

$$\beta_k \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{B}}(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{A}}(\mathbf{x}) + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{E}}(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{A}}(\mathbf{x}) + \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{E}}(\mathbf{x}) = \alpha_k + \varepsilon_1,$$

and this proves the upper inequality. For the lower one we define $\mathbf{D} := -\mathbf{E}$ and observe that $-\varepsilon_n$ is the largest eigenvalue of \mathbf{D} . Since $\mathbf{A} = \mathbf{B} + \mathbf{D}$ it follows from the result just proved that $\alpha_k \leq \beta_k - \varepsilon_n$, which is the same as the lower inequality. \square

In many applications of this result the eigenvalues of the matrix \mathbf{E} will be small and then the theorem states that the eigenvalues of \mathbf{B} are close to those of \mathbf{A} . Moreover, it associates a unique eigenvalue of \mathbf{A} with each eigenvalue of \mathbf{B} .

Exercise 6.47 (Eigenvalue perturbation for Hermitian matrices)

Show that in Theorem 6.46, if \mathbf{E} is symmetric positive semidefinite then $\beta_i \geq \alpha_i$.

6.4.3 The Hoffman-Wielandt Theorem

We can also give a bound involving all eigenvalues. The following theorem shows that the eigenvalue problem for a normal matrix is well conditioned.

Theorem 6.48 (Hoffman-Wielandt Theorem)

Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ are both normal matrices with eigenvalues $\lambda_1, \dots, \lambda_n$ and μ_1, \dots, μ_n , respectively. Then there is a permutation i_1, \dots, i_n of $1, 2, \dots, n$ such that

$$\sum_{j=1}^n |\mu_{i_j} - \lambda_j|^2 \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2. \quad (6.12)$$

For a proof of this theorem see [[24], p. 190]. For a Hermitian matrix we can use the identity permutation if we order both set of eigenvalues in nonincreasing or nondecreasing order.

Exercise 6.49 (Hoffman-Wielandt)

Show that (6.12) does not hold for the matrices $\mathbf{A} := \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix}$ and $\mathbf{B} := \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}$. Why does this not contradict the Hoffman-Wielandt theorem?

6.5 Left Eigenvectors

Definition 6.50 (Left eigenpair)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a square matrix, $\lambda \in \mathbb{C}$ and $\mathbf{y} \in \mathbb{C}^n$. We say that (λ, \mathbf{y}) is a left eigenpair for \mathbf{A} if $\mathbf{y}^* \mathbf{A} = \lambda \mathbf{y}^*$ and \mathbf{y} is nonzero.

Since $\mathbf{A}^* \mathbf{y} = \bar{\lambda} \mathbf{y}$ Theorem 0.66 implies that λ is an eigenvalue of \mathbf{A} , while a left eigenvector is an eigenvector of \mathbf{A}^* . Thus left and right eigenvalues are identical, but left and right eigenvectors are in general different. For an Hermitian matrix the right and left eigenpairs are the same

Left- and right eigenvectors corresponding to distinct eigenvalues are orthogonal.

Theorem 6.51 (Biorthogonality)

Suppose (μ, \mathbf{y}) and (λ, \mathbf{x}) are left and right eigenpairs of $\mathbf{A} \in \mathbb{C}^{n \times n}$. If $\lambda \neq \mu$ then $\mathbf{y}^* \mathbf{x} = 0$.

Proof. Using the eigenpair relation in two ways we obtain $\mathbf{y}^* \mathbf{A} \mathbf{x} = \lambda \mathbf{y}^* \mathbf{x} = \mu \mathbf{y}^* \mathbf{x}$ and we conclude that $\mathbf{y}^* \mathbf{x} = 0$. \square

Right and left eigenvectors corresponding to the same eigenvalue are sometimes orthogonal, sometimes not.

Theorem 6.52 (Simple eigenvalue)

Suppose (λ, \mathbf{x}) and (λ, \mathbf{y}) are right and left eigenpairs of $\mathbf{A} \in \mathbb{C}^{n \times n}$. If λ has algebraic multiplicity one then $\mathbf{y}^* \mathbf{x} \neq 0$.

Proof. Assume that $\|\mathbf{x}\|_2 = 1$. We have (cf. (6.5))

$$\mathbf{V}^* \mathbf{A} \mathbf{V} = \left[\begin{array}{c|c} \lambda & \mathbf{z}^* \\ \mathbf{0} & \mathbf{M} \end{array} \right],$$

where \mathbf{V} is unitary and $\mathbf{V} \mathbf{e}_1 = \mathbf{x}$. We show that if $\mathbf{y}^* \mathbf{x} = 0$ then λ is also an eigenvalue of \mathbf{M} contradicting the multiplicity assumption of λ . Let $\mathbf{u} := \mathbf{V}^* \mathbf{y}$. Then

$$(\mathbf{V}^* \mathbf{A}^* \mathbf{V}) \mathbf{u} = \mathbf{V}^* \mathbf{A}^* \mathbf{y} = \bar{\lambda} \mathbf{V}^* \mathbf{y} = \bar{\lambda} \mathbf{u},$$

so $(\bar{\lambda}, \mathbf{u})$ is an eigenpair of $\mathbf{V}^* \mathbf{A}^* \mathbf{V}$. But then $\mathbf{y}^* \mathbf{x} = \mathbf{u}^* \mathbf{V}^* \mathbf{V} \mathbf{e}_1$. Suppose that $\mathbf{u}^* \mathbf{e}_1 = 0$, i.e., $\mathbf{u} = [\mathbf{v}]$ for some nonzero $\mathbf{v} \in \mathbb{C}^{n-1}$. Then

$$\mathbf{V}^* \mathbf{A}^* \mathbf{V} \mathbf{u} = \left[\begin{array}{c|c} \bar{\lambda} & \mathbf{0}^* \\ \mathbf{z} & \mathbf{M}^* \end{array} \right] \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{M}^* \mathbf{v} \end{bmatrix} = \bar{\lambda} \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix}$$

and by Theorem 0.66 it follows that λ is an eigenvalue of \mathbf{M} . \square

The case with multiple eigenvalues is more complicated. For example, the matrix $\mathbf{A} := \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ has one eigenvalue $\lambda = 1$ of algebraic multiplicity two, one right eigenvector $\mathbf{x} = \mathbf{e}_1$ and one left eigenvector $\mathbf{y} = \mathbf{e}_2$. Thus \mathbf{x} and \mathbf{y} are orthogonal.

Theorem 6.53 (Biorthogonal eigenvector expansion)

If $\mathbf{A} \in \mathbb{C}^{n \times n}$ has linearly independent right eigenvectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ then there exists a set of left eigenvectors $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ with $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$. Conversely, if $\mathbf{A} \in \mathbb{C}^{n \times n}$ has linearly independent left eigenvectors $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ then there exists a set of right eigenvectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ with $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$. For any scaling of these sets we have the eigenvector expansions

$$\mathbf{v} = \sum_{j=1}^n \frac{\mathbf{y}_j^* \mathbf{v}}{\mathbf{y}_j^* \mathbf{x}_j} \mathbf{x}_j = \sum_{k=1}^n \frac{\mathbf{x}_k^* \mathbf{v}}{\mathbf{y}_k^* \mathbf{x}_k} \mathbf{y}_k. \quad (6.13)$$

Proof. For any right eigenpairs $(\lambda_1, \mathbf{x}_1), \dots, (\lambda_n, \mathbf{x}_n)$ and left eigenpairs $(\lambda_1, \mathbf{y}_1), \dots, (\lambda_n, \mathbf{y}_n)$ of \mathbf{A} we have $\mathbf{AX} = \mathbf{XD}$, $\mathbf{Y}^* \mathbf{A} = \mathbf{DY}^*$, where

$$\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n], \quad \mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n], \quad \mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n).$$

If \mathbf{X} is nonsingular then $\mathbf{X}^{-1} \mathbf{A} = \mathbf{DX}^{-1}$ and it follows that $\mathbf{Y}^* := \mathbf{X}^{-1}$ contains a collection of left eigenvectors such that $\mathbf{Y}^* \mathbf{X} = \mathbf{I}$. Thus the columns of \mathbf{Y} are linearly independent and $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$. Similarly, if \mathbf{Y} is nonsingular then $\mathbf{AY}^{-*} = \mathbf{Y}^{-*} \mathbf{D}$ and it follows that $\mathbf{X} := \mathbf{Y}^{-*}$ contains a collection of linearly independent right eigenvectors such that $\mathbf{Y}^* \mathbf{X} = \mathbf{I}$. If $\mathbf{v} = \sum_{j=1}^n c_j \mathbf{x}_j$ then $\mathbf{y}_i^* \mathbf{v} = \sum_{j=1}^n c_j \mathbf{y}_i^* \mathbf{x}_j = c_i \mathbf{y}_i^* \mathbf{x}_i$, so $c_i = \mathbf{y}_i^* \mathbf{v} / \mathbf{y}_i^* \mathbf{x}_i$ for $i = 1, \dots, n$ and the first expansion in (6.13) follows. The second expansion follows similarly. \square

For an Hermitian matrix the right eigenvectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ are also left eigenvectors and (6.13) takes the form

$$\mathbf{v} = \sum_{j=1}^n \frac{\mathbf{x}_j^* \mathbf{v}}{\mathbf{x}_j^* \mathbf{x}_j} \mathbf{x}_j. \quad (6.14)$$

Exercise 6.54 (Biorthogonal expansion)

Determine right and left eigenpairs for the matrix $\mathbf{A} := \begin{bmatrix} 3 & 1 \\ 2 & 2 \end{bmatrix}$ and the two expansions in (6.13) for any $\mathbf{v} \in \mathbb{R}^2$.

Exercise 6.55 (Generalized Rayleigh quotient)

For $\mathbf{A} \in \mathbb{C}^{n \times n}$ and any $\mathbf{y}, \mathbf{x} \in \mathbb{C}^n$ with $\mathbf{y}^* \mathbf{x} \neq 0$ the quantity $R(\mathbf{y}, \mathbf{x}) = R_{\mathbf{A}}(\mathbf{y}, \mathbf{x}) := \frac{\mathbf{y}^* \mathbf{Ax}}{\mathbf{y}^* \mathbf{x}}$ is called a **generalized Rayleigh quotient** for \mathbf{A} . Show that if (λ, \mathbf{x}) is a right eigenpair for \mathbf{A} then $R(\mathbf{y}, \mathbf{x}) = \lambda$ for any \mathbf{y} with $\mathbf{y}^* \mathbf{x} \neq 0$. Also show that if (λ, \mathbf{y}) is a left eigenpair for \mathbf{A} then $R(\mathbf{y}, \mathbf{x}) = \lambda$ for any \mathbf{x} with $\mathbf{y}^* \mathbf{x} \neq 0$.

6.6 The Jordan Form and the Minimal Polynomial

We have seen that any square matrix can be triangularized by a unitary similarity transformation. Moreover, any nondefective matrix can be diagonalized. The following question arises. How close to a diagonal matrix can we reduce a defective matrix by a similarity transformation? The main result is Theorem 6.57. For a proof, see for example [12].

Definition 6.56 (Jordan block)

A **Jordan block**, denoted $\mathbf{J}_m(\lambda)$ is an $m \times m$ matrix of the form

$$\mathbf{J}_m(\lambda) := \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda & 1 & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \ddots & \\ 0 & 0 & 0 & \cdots & \lambda & 1 \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix}$$

A 3×3 Jordan block has the form $\mathbf{J}_3(\lambda) = \begin{bmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{bmatrix}$. We see that λ is an eigenvalue of $\mathbf{J}_m(\lambda)$ and any eigenvector must be a multiple of e_1 . Thus, the eigenvectors of $\mathbf{J}_m(\lambda)$ have algebraic multiplicity m and geometric multiplicity one.

The Jordan canonical form is a decomposition of a matrix into Jordan blocks. See [12] for a proof.

Theorem 6.57 (The Jordan form of a matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ has k distinct eigenvalues $\lambda_1, \dots, \lambda_k$ of algebraic multiplicities a_1, \dots, a_k and geometric multiplicities g_1, \dots, g_k . There is a nonsingular matrix $\mathbf{S} \in \mathbb{C}^{n \times n}$ such that

$$\mathbf{J} := \mathbf{S}^{-1} \mathbf{A} \mathbf{S} = \text{diag}(\mathbf{U}_1, \dots, \mathbf{U}_k), \text{ with } \mathbf{U}_i \in \mathbb{C}^{a_i, a_i}, \quad (6.15)$$

where each \mathbf{U}_i is block diagonal having g_i Jordan blocks along the diagonal

$$\mathbf{U}_i = \text{diag}(\mathbf{J}_{m_{i,1}}(\lambda_i), \dots, \mathbf{J}_{m_{i,g_i}}(\lambda_i)). \quad (6.16)$$

Here $m_{i,1}, \dots, m_{i,g_i}$ are unique integers so that $m_{i,1} \geq m_{i,2} \geq \dots \geq m_{i,g_i}$ and $a_i = \sum_{j=1}^{g_i} m_{i,j}$ for all i .

The matrix \mathbf{J} in (6.15) is called the **Jordan form** of \mathbf{A} . As an example consider the Jordan form

$$\mathbf{J} := \text{diag}(\mathbf{U}_1, \mathbf{U}_2) = \begin{bmatrix} \begin{smallmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{smallmatrix} & & & \\ & \begin{smallmatrix} 2 & 1 \\ 0 & 2 \\ 0 & 2 \end{smallmatrix} & & \\ & & \begin{smallmatrix} 2 \\ 0 \\ 2 \end{smallmatrix} & \\ & & & \begin{smallmatrix} 3 & 1 \\ 0 & 3 \end{smallmatrix} \end{bmatrix} \in \mathbb{R}^{8,8}. \quad (6.17)$$

The eigenvalues together with their algebraic and geometric multiplicities can be read off directly from the Jordan form.

- $\mathbf{U}_1 = \text{diag}(\mathbf{J}_3(2), \mathbf{J}_2(2), \mathbf{J}_1(2))$ and $\mathbf{U}_2 = \mathbf{J}_2(3)$.
- 2 is an eigenvalue of algebraic multiplicity 6 and geometric multiplicity 3.
- 3 is an eigenvalue of algebraic multiplicity 2 and geometric multiplicity 1.

Each \mathbf{U}_i is upper triangular with the eigenvalue λ_i on the diagonal and consists of g_i Jordan blocks. These Jordan blocks can be taken in any order and it is customary to refer to any such block diagonal matrix as the Jordan form of \mathbf{A} . Thus in the example the matrix

$$\mathbf{J} := \begin{bmatrix} \begin{smallmatrix} 3 & 1 \\ 0 & 3 \end{smallmatrix} & & & \\ & \begin{smallmatrix} 2 & 1 \\ 0 & 2 \\ 0 & 2 \end{smallmatrix} & & \\ & & \begin{smallmatrix} 2 \\ 0 \\ 2 \end{smallmatrix} & \\ & & & \begin{smallmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{smallmatrix} \end{bmatrix}$$

is also a Jordan form of \mathbf{A} . In any Jordan form of this \mathbf{A} the sizes of the 4 Jordan blocks $\mathbf{J}_3(2), \mathbf{J}_2(2), \mathbf{J}_1(2), \mathbf{J}_2(3)$ are uniquely given.

The columns of \mathbf{S} are called **principal vectors**. They satisfy the matrix equation $\mathbf{AS} = \mathbf{SJ}$. As an example, in (6.17) we have $\mathbf{S} = [s_1, \dots, s_8]$ and we find

$$\begin{aligned}\mathbf{As}_1 &= 2s_1, & \mathbf{As}_2 &= 2s_2 + s_1, & \mathbf{As}_3 &= 2s_3 + s_2, \\ \mathbf{As}_4 &= 2s_4, & \mathbf{As}_5 &= 2s_5 + s_4, \\ \mathbf{As}_6 &= 2s_6, \\ \mathbf{As}_7 &= 3s_7, & \mathbf{As}_8 &= 3s_8 + s_7,\end{aligned}$$

We see that the first principal vector in each Jordan block is an eigenvector of \mathbf{A} . The remaining principal vectors are not eigenvectors.

Exercise 6.58 (Jordan example)

For the Jordan form of the matrix $\mathbf{A} = \begin{bmatrix} 3 & 0 & -1 \\ -4 & 1 & 0 \\ -4 & 0 & -1 \end{bmatrix}$ we have $\mathbf{J} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$. Find \mathbf{S} .

Exercise 6.59 (Big Jordan example)

Find the Jordan form of the matrix

$$\mathbf{A} = \frac{1}{9} \begin{bmatrix} 10 & 16 & -8 & -5 & 6 & 1 & -3 & 4 \\ -7 & 32 & -7 & -10 & 12 & 2 & -6 & 8 \\ -6 & 12 & 12 & -15 & 18 & 3 & -9 & 12 \\ -5 & 10 & -5 & -2 & 24 & 4 & -12 & 16 \\ -4 & 8 & -4 & -16 & 30 & 14 & -15 & 20 \\ -3 & 6 & -3 & -12 & 9 & 24 & -9 & 24 \\ -2 & 4 & -2 & -8 & 6 & -2 & 15 & 28 \\ -1 & 2 & -1 & -4 & 3 & -1 & -6 & 41 \end{bmatrix}. \quad (6.18)$$

The following lemma is useful when studying powers of matrices.

Lemma 6.60 (Properties of the Jordan form)

Let \mathbf{J} be the Jordan form of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ as given in Theorem 6.57. Then for $r = 0, 1, 2, \dots, m = 2, 3, \dots$, and any $\lambda \in \mathbb{C}$

1. $\mathbf{A}^r = \mathbf{SJ}^r \mathbf{S}^{-1}$,
2. $\mathbf{J}^r = \text{diag}(\mathbf{U}_1^r, \dots, \mathbf{U}_k^r)$,
3. $\mathbf{U}_i^r = \text{diag}(\mathbf{J}_{m_{i,1}}(\lambda_i)^r, \dots, \mathbf{J}_{m_{i,g_i}}(\lambda_i)^r)$,
4. $\mathbf{E}_m^r = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{m-r} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ for $1 \leq r \leq m-1$, where $\mathbf{E}_m := \mathbf{J}_m(\lambda) - \lambda \mathbf{I}_m$,
5. $\mathbf{E}_m^m = \mathbf{0}$.
6. $\mathbf{J}_m(\lambda)^r = (\mathbf{E}_m + \lambda \mathbf{I}_m)^r = \sum_{k=0}^{\min\{r,m-1\}} \binom{r}{k} \lambda^{r-k} \mathbf{E}_m^k$

Proof.

1. We have $\mathbf{A}^2 = \mathbf{SJS}^{-1}\mathbf{SJS}^{-1} = \mathbf{SJ}^2\mathbf{S}^{-1}$ and 1. follows by induction on r .
2. This follows since \mathbf{J} is block diagonal.
3. Each $\mathbf{J}_{m_i,j}$ is block diagonal.
4. We have

$$\mathbf{E}_m = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \ddots & \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{m-1} \\ 0 & \mathbf{0}^T \end{bmatrix}. \quad (6.19)$$

The result follow for $r = 1$ and for general $r \leq m - 1$ by induction.

5. $\mathbf{E}_m^m = \mathbf{E}_m^{m-1}\mathbf{E}_m = \mathbf{0}$.
6. This follows from the binomial theorem since \mathbf{I}_m and \mathbf{E}_m commute and $\mathbf{E}^m = \mathbf{0}$.

□

Exercise 6.61 (Jordan block example)

Determine \mathbf{J}_3^r for $r \geq 1$.

Exercise 6.62 (Powers of a Jordan block)

Find \mathbf{J}^{100} and \mathbf{A}^{100} for the matrix in Exercise 6.58.

6.6.1 The Minimal Polynomial

Let \mathbf{J} be the Jordan form of \mathbf{A} given in Theorem 6.57. Since \mathbf{A} and \mathbf{J} are similar they have the same characteristic polynomial, and since the Jordan form of \mathbf{A} is upper triangular with the eigenvalues of \mathbf{A} on the diagonal we have

$$\pi_{\mathbf{A}}(\lambda) = \pi_{\mathbf{J}}(\lambda) = \prod_{i=1}^k \prod_{j=1}^{g_i} (\lambda_i - \lambda)^{m_{ij}}.$$

The polynomials $p_{ij}(\lambda) := (\lambda_i - \lambda)^{m_{ij}}$ are called the **elementary divisors** of \mathbf{A} . They divide the characteristic polynomial.

Definition 6.63 (Minimal polynomial of a matrix)

Suppose $\mathbf{A} = \mathbf{SJS}^{-1}$ is the Jordan canonical form of \mathbf{A} . The polynomial

$$\mu(z) := \prod_{i=1}^k (\lambda_i - z)^{m_i} \text{ where } m_i := \max_{1 \leq j \leq g_i} m_{ij},$$

is called the **minimal polynomial** of \mathbf{A} .

Since each factor in $\mu(z)$ is also a factor in $\pi_{\mathbf{A}}(z)$, we have the factorization $\pi_{\mathbf{A}}(z) = \mu(z)\nu(z)$ for some polynomial $\nu(z)$.

Exercise 6.64 (Minimal polynomial example)

What is the characteristic polynomial and the minimal polynomial of the matrix \mathbf{J} in (6.17)?

To see in what way the minimal polynomial is minimal, we consider two matrices defined from the characteristic polynomial $\pi_{\mathbf{A}}$ and the minimal polynomial. Substituting a matrix for the independent variable in these polynomial we obtain

$$\pi_{\mathbf{A}}(\mathbf{A}) := \prod_{i=1}^k \prod_{j=1}^{g_i} (\lambda_i \mathbf{I} - \mathbf{A})^{m_{ij}}, \quad \mu(\mathbf{A}) := \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{A})^{m_i}. \quad (6.20)$$

By induction it is easy to see that $\mu(\mathbf{A})$ and $\pi_{\mathbf{A}}(\mathbf{A})$ are polynomials in the matrix \mathbf{A} . Moreover, $\mu(\mathbf{A}) = \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{S} \mathbf{J} \mathbf{S}^{-1})^{m_i} = \mathbf{S} \mu(\mathbf{J}) \mathbf{S}^{-1}$, so that $\mu(\mathbf{A}) = \mathbf{0}$ if and only if $\mu(\mathbf{J}) = \mathbf{0}$. Now,

$$\begin{aligned} \mu(\mathbf{J}) &= \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{J})^{m_i} = \prod_{i=1}^k \text{diag}((\lambda_i \mathbf{I} - \mathbf{U}_1)^{m_i}, \dots, (\lambda_i \mathbf{I} - \mathbf{U}_k)^{m_i}) \\ &= \text{diag}\left(\prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{U}_1)^{m_i}, \dots, \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{U}_k)^{m_i}\right) = \mathbf{0}, \end{aligned}$$

since $(\lambda_r \mathbf{I} - \mathbf{U}_r)^{m_r} = \mathbf{0}$ for $r = 1, \dots, k$. To show the latter we observe that

$$\begin{aligned} (\lambda_r \mathbf{I} - \mathbf{U}_r)^{m_r} &= \text{diag}((\lambda_r \mathbf{I} - \mathbf{J}_{m_{r1}})^{m_r}, \dots, (\lambda_r \mathbf{I} - \mathbf{J}_{m_{rg_r}})^{m_r}) \\ &= \text{diag}(\mathbf{E}_{m_{r1}}^{m_r}, \dots, \mathbf{E}_{m_{rg_r}}^{m_r}) = \mathbf{0}, \end{aligned}$$

by Lemma 6.60 and the maximality of m_r .

We have shown that a matrix satisfies its minimal polynomial equation $\mu(\mathbf{A}) = \mathbf{0}$. Moreover, the degree of any polynomial p such that $p(\mathbf{A}) = \mathbf{0}$ is at least as large as the degree $d = \sum_{i=1}^k m_i$ of the minimal polynomial μ . This follows from the proof since any such polynomial must contain the elementary divisors $(\lambda_i - \lambda)^{m_i}$ for $i = 1, \dots, k$. Since the minimal polynomial divides the characteristic polynomial we obtain as a corollary the **Cayley-Hamilton Theorem** which says that a matrix satisfies its characteristic equation $\pi_{\mathbf{A}}(\mathbf{A}) = \mathbf{0}$.

Exercise 6.65 (Similar matrix polynomials)

Show that $p(\mathbf{B}) = \mathbf{S}^{-1} p(\mathbf{A}) \mathbf{S}$ for any polynomial p and any similar matrices $\mathbf{B} = \mathbf{S}^{-1} \mathbf{A} \mathbf{S}$.

Exercise 6.66 (Minimal polynomial of a diagonalizable matrix)

What is the minimal polynomial of the unit matrix and more generally of a diagonalizable matrix?

6.7 Proof of the Real Schur Form

In this section we prove the following theorem.

Theorem 6.67 (Proof of real Schur form)

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then we can find $\mathbf{U} \in \mathbb{R}^{n \times n}$ with $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ such that $\mathbf{U}^T \mathbf{A} \mathbf{U}$ is quasi-triangular.

Proof. If \mathbf{A} has only real eigenvalues, Theorem 6.31 gives the result. Suppose $\lambda = \mu + i\nu$, $\mu, \nu \in \mathbb{R}$, is an eigenvalue of \mathbf{A} with $\nu \neq 0$. Let $\mathbf{z} = \mathbf{x} + i\mathbf{y}$, $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, be an eigenvector of \mathbf{A} corresponding to λ . Since

$$\mathbf{A}\mathbf{z} = \mathbf{A}(\mathbf{x} + i\mathbf{y}) = (\mu + i\nu)(\mathbf{x} + i\mathbf{y}) = \mu\mathbf{x} - \nu\mathbf{y} + i(\nu\mathbf{x} + \mu\mathbf{y}),$$

we find by comparing real and imaginary parts

$$\mathbf{A}\mathbf{x} = \mu\mathbf{x} - \nu\mathbf{y}, \quad \mathbf{A}\mathbf{y} = \nu\mathbf{x} + \mu\mathbf{y}. \quad (6.21)$$

We claim that \mathbf{x} and \mathbf{y} are linearly independent. First we note that $\nu \neq 0$ implies $\mathbf{x} \neq \mathbf{0}$, $\mathbf{y} \neq \mathbf{0}$. For if $\mathbf{x} = \mathbf{0}$ then (6.21) implies that $\mathbf{0} = -\nu\mathbf{y}$, and hence $\mathbf{y} = \mathbf{0}$ as well, contradicting the nonzeroness of the eigenvector. Similarly, if $\mathbf{y} = \mathbf{0}$ then $\mathbf{0} = \nu\mathbf{x}$, again resulting in a zero eigenvector. Suppose $\mathbf{y} = \alpha\mathbf{x}$ for some α . Replacing \mathbf{y} by $\alpha\mathbf{x}$ in (6.21), we find $\mathbf{A}\mathbf{x} = (\mu - \alpha\nu)\mathbf{x}$ and $\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{y}/\alpha = (\mu + \nu/\alpha)\mathbf{x}$. But then $\mu - \alpha\nu = \mu + \nu/\alpha$ or $\alpha^2 = -1$. Since \mathbf{x} and \mathbf{y} are real, we cannot have both $\mathbf{y} = \alpha\mathbf{x}$ and $\alpha^2 = -1$. We conclude that \mathbf{x} and \mathbf{y} are linearly independent.

(6.21) can be written in matrix form as

$$\mathbf{A}\mathbf{X}_1 = \mathbf{X}_1 \mathbf{M}, \quad \mathbf{X}_1 = [\mathbf{x}, \mathbf{y}] \in \mathbb{R}^{n,2}, \quad \mathbf{M} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix}. \quad (6.22)$$

By Theorem 11.12 we can find an orthonormal matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{Q}\mathbf{X}_1 = \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix}$$

where $\mathbf{R} \in \mathbb{R}^{2,2}$ is upper triangular. Since \mathbf{X}_1 has linearly independent columns, \mathbf{R} is nonsingular. Let $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$ and define

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2] = [\mathbf{x}, \mathbf{y}, \mathbf{q}_3, \dots, \mathbf{q}_n].$$

We find

$$\mathbf{Q}\mathbf{X} = [\mathbf{Q}\mathbf{X}_1, \mathbf{Q}\mathbf{q}_3, \dots, \mathbf{Q}\mathbf{q}_n] = \begin{bmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix}.$$

Since \mathbf{R} is nonsingular, $\mathbf{Q}\mathbf{X}$ and \mathbf{X} are nonsingular. Moreover, using (6.22)

$$\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = [\mathbf{X}^{-1} \mathbf{A} \mathbf{X}_1, \mathbf{X}^{-1} \mathbf{A} \mathbf{X}_2] = [\mathbf{X}^{-1} \mathbf{X}_1 \mathbf{M}, \mathbf{X}^{-1} \mathbf{A} \mathbf{X}_2] = \begin{bmatrix} \mathbf{M} & \mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix}$$

for some matrices $\mathbf{B} \in \mathbb{R}^{2,n-2}$, $\mathbf{C} \in \mathbb{R}^{n-2,n-2}$. Now

$$\mathbf{Q}\mathbf{A}\mathbf{Q}^T = (\mathbf{Q}\mathbf{X})\mathbf{X}^{-1}\mathbf{A}\mathbf{X}(\mathbf{Q}\mathbf{X})^{-1} = \begin{bmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{R}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix},$$

or

$$\mathbf{Q}\mathbf{A}\mathbf{Q}^T = \begin{bmatrix} \mathbf{R}\mathbf{M}\mathbf{R}^{-1} & \mathbf{R}\mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix}. \quad (6.23)$$

By Theorem 6.7 the 2×2 matrix $\mathbf{R}\mathbf{M}\mathbf{R}^{-1}$ has the same eigenvalues λ and $\bar{\lambda}$ as \mathbf{M} . The remaining $n-2$ eigenvalues of \mathbf{A} are the eigenvalues of \mathbf{C} .

To complete the proof we use induction on n . The theorem is trivially true for $n = 1$ and $n = 2$. Suppose $n \geq 3$ and it holds for matrices of order $\leq n-1$. Let

$$\mathbf{V} = \begin{bmatrix} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{V}} \end{bmatrix}$$

where $\hat{\mathbf{V}} \in \mathbb{R}^{n-2,n-2}$, $\hat{\mathbf{V}}^T \hat{\mathbf{V}} = \mathbf{I}_{n-2}$ and $\hat{\mathbf{V}}^T \mathbf{C} \hat{\mathbf{V}}$ is quasi-triangular. Let $\mathbf{U} = \mathbf{Q}\mathbf{V}$. Then $\mathbf{U} \in \mathbb{R}^{n \times n}$, $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ and $\mathbf{U}^T \mathbf{A} \mathbf{U}$ is quasi-triangular. \square

6.8 Conclusions

Consider the eigenpair problem for some classes of matrices $\mathbf{A} \in \mathbb{C}^{n \times n}$.

Diagonal Matrices. The eigenpairs are easily determined. Since $\mathbf{A}\mathbf{e}_i = a_{ii}\mathbf{e}_i$ the eigenpairs are $(\lambda_i, \mathbf{e}_i)$, where $\lambda_i = a_{ii}$ for $i = 1, \dots, n$. Moreover, $a(\lambda_i) = g(\lambda_i)$ for all i , since the eigenvectors of \mathbf{A} are linearly independent.

Triangular Matrices Suppose \mathbf{A} is upper or lower triangular. Consider finding the eigenvalues Since $\det(\mathbf{A} - \lambda\mathbf{I}) = \prod_{i=1}^n (a_{ii} - \lambda)$ the eigenvalues are $\lambda_i = a_{ii}$ for $i = 1, \dots, n$, the diagonal elements of \mathbf{A} . To determine the eigenvectors can be challenging as Example 6.17 indicates.

Block Diagonal Matrices Suppose

$$\mathbf{A} = \text{diag}(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_r), \quad \mathbf{A}_i \in \mathbb{C}^{m_i \times m_i}.$$

Here the eigenpair problem reduces to r smaller problems. Let $\mathbf{A}_i \mathbf{X}_i = \mathbf{X}_i \mathbf{D}_i$ define the eigenpairs of \mathbf{A}_i for $i = 1, \dots, r$ and let $\mathbf{X} := \text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_r)$, $\mathbf{D} := \text{diag}(\mathbf{D}_1, \dots, \mathbf{D}_r)$. Then the eigenpairs for \mathbf{A} are given by

$$\begin{aligned} \mathbf{AD} &= \text{diag}(\mathbf{A}_1, \dots, \mathbf{A}_r) \text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_r) = \text{diag}(\mathbf{A}_1 \mathbf{X}_1, \dots, \mathbf{A}_r \mathbf{X}_r) \\ &= \text{diag}(\mathbf{X}_1 \mathbf{D}_1, \dots, \mathbf{X}_r \mathbf{D}_r) = \mathbf{XD}. \end{aligned}$$

Block Triangular matrices Matrices Let $A_{11}, A_{22}, \dots, A_{rr}$ be the diagonal blocks of A . By Property 8. of determinants

$$\det(A - \lambda I) = \prod_{i=1}^r \det(A_{ii} - \lambda I)$$

and the eigenvalues are found from the eigenvalues of the diagonal blocks.

6.9 Review Questions

- 6.9.1 Does A and A^T , A and A^* have the same eigenvalues? What about A^*A and AA^* ?
- 6.9.2 Can the geometric multiplicity of an eigenvalue be bigger than the algebraic multiplicity?
- 6.9.3 What are the eigenvalues of a diagonal matrix?
- 6.9.4 What are the Schur factors of a matrix?
- 6.9.5 What is a quasi-triangular matrix?
- 6.9.6 Give some classes of normal matrices. Why are normal matrices important?.
- 6.9.7 State the Courant-Fischer theorem.
- 6.9.8 State the Hoffman-Wieland theorem for Hermitian matrices.
- 6.9.9 What is a left eigenvector of a matrix.

Chapter 7

The Singular Value Decomposition

The singular value decomposition is useful both for theory and practice. Some of its applications include solving over-determined equations, principal component analysis in statistics, numerical determination of the rank of a matrix, algorithms used in search engines, and the theory of matrices.

7.1 Singular Values and Singular Vectors

We know from Theorem 6.37 that a square matrix \mathbf{A} can be diagonalized by a unitary similarity transformation if and only if it is normal, that is $\mathbf{A}^*\mathbf{A} = \mathbf{A}\mathbf{A}^*$. In particular, if $\mathbf{A} \in \mathbb{C}^{n \times n}$ is normal then it has a set of orthonormal eigenpairs $(\lambda_1, \mathbf{u}_1), \dots, (\lambda_n, \mathbf{u}_n)$. Letting $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{C}^{n \times n}$ and $\mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n)$ we have the spectral decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^*, \text{ where } \mathbf{U}^*\mathbf{U} = \mathbf{I}. \quad (7.1)$$

In this chapter we show that any matrix, even a rectangular one, can be diagonalized provided we allow two unitary matrices. A factorization of $\mathbf{A} \in \mathbb{C}^{m \times n}$ of the form

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*, \quad \mathbf{U} \in \mathbb{C}^{m \times m}, \mathbf{V} \in \mathbb{C}^{n \times n}, \Sigma \in \mathbb{R}^{m \times n}, \quad (7.2)$$

where \mathbf{U}, \mathbf{V} are unitary and Σ is a nonnegative diagonal matrix, i.e., $\Sigma_{i,j} = 0$ for all $i \neq j$, $\Sigma_{i,i} \geq 0$ for $i = \min(m, n)$, is called a **singular value decomposition**. The diagonal elements of Σ are called **singular values** and the columns $\mathbf{u}_1, \dots, \mathbf{u}_m$ of \mathbf{U} and $\mathbf{v}_1, \dots, \mathbf{v}_n$ of \mathbf{V} are called (left and right) **singular vectors**.

For a Hermitian matrix with nonnegative eigenvalues the spectral decomposition (7.1) is also a singular value decomposition. For example, the matrix

$\mathbf{A} := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ has eigenpairs $(2, [\begin{smallmatrix} 1 \\ 1 \end{smallmatrix}])$ and $(0, [\begin{smallmatrix} 1 \\ -1 \end{smallmatrix}])$. The factorization

$$\mathbf{A} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \mathbf{U}\Sigma\mathbf{V}^T, \quad \mathbf{V} = \mathbf{U},$$

is both a spectral decomposition and a singular value decomposition.

In general the singular values are not eigenvalues of \mathbf{A} . The singular vectors are eigenvectors, but of different matrices.

7.1.1 SVD and SVF

In this section we show that every matrix has a singular value decomposition (SVD) and a reduced form called the singular value factorization (SVF). Eigenpairs of $\mathbf{A}^* \mathbf{A}$ play a central role.

Lemma 7.1 (Eigenpairs of $\mathbf{A}^* \mathbf{A}$)

Suppose $m, n \in \mathbb{N}$ and $\mathbf{A} \in \mathbb{C}^{m \times n}$. The matrix $\mathbf{A}^* \mathbf{A}$ has eigenpairs $(\lambda_j, \mathbf{v}_j)$ for $j = 1, \dots, n$, where $\mathbf{v}_j^* \mathbf{v}_k = \delta_{jk}$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$. Moreover,

$$\sigma_j := \sqrt{\lambda_j} = \|\mathbf{A}\mathbf{v}_j\|_2, \text{ for } j = 1, \dots, n. \quad (7.3)$$

Proof. The matrix $\mathbf{A}^* \mathbf{A} \in \mathbb{C}^{n \times n}$ is Hermitian, and by Theorem 6.38 it has real eigenvalues λ_j and orthonormal eigenvectors \mathbf{v}_j for $j = 1, \dots, n$. For each j $\|\mathbf{A}\mathbf{v}_j\|_2^2 = (\mathbf{A}\mathbf{v}_j)^* \mathbf{A}\mathbf{v}_j = \mathbf{v}_j^* \mathbf{A}^* \mathbf{A}\mathbf{v}_j = \lambda_j \mathbf{v}_j^* \mathbf{v}_j = \lambda_j$, since $\mathbf{v}_j^* \mathbf{v}_j = 1$. Thus $\lambda_j \geq 0$, and furthermore (7.3) follows. \square

Theorem 7.2 (Orthogonal bases for column- and null space of \mathbf{A})

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ and let $(\lambda_j, \mathbf{v}_j)$ for $j = 1, \dots, n$ be orthonormal eigenpairs for $\mathbf{A}^* \mathbf{A}$. If $\lambda_j > 0$, $j = 1, \dots, r$ and $\lambda_j = 0$, $j = r+1, \dots, n$ then $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ is an orthogonal basis for the column space $\text{span}(\mathbf{A}) := \{\mathbf{A}\mathbf{y} \in \mathbb{C}^m : \mathbf{y} \in \mathbb{C}^n\}$ and $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ is an orthonormal basis for the nullspace $\ker(\mathbf{A}) := \{\mathbf{y} \in \mathbb{C}^n : \mathbf{A}\mathbf{y} = \mathbf{0}\}$.

Proof. By orthonormality of $\mathbf{v}_1, \dots, \mathbf{v}_n$ $(\mathbf{A}\mathbf{v}_j)^* \mathbf{A}\mathbf{v}_k = \mathbf{v}_j^* \mathbf{A}^* \mathbf{A}\mathbf{v}_k = \lambda_k \mathbf{v}_j^* \mathbf{v}_k = 0, j \neq k$, showing that $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_n$ are orthogonal vectors. Moreover, (7.3) implies that $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$ are nonzero and $\mathbf{A}\mathbf{v}_j = \mathbf{0}$ for $j = r+1, \dots, n$. In particular, the elements of $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ and $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ are linearly independent vectors in $\text{span}(\mathbf{A})$ and $\ker(\mathbf{A})$, respectively. The proof will be complete once it is shown that

$$\text{span}(\mathbf{A}) \subset \text{span}(\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r), \quad \ker(\mathbf{A}) \subset \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n). \quad (7.4)$$

Suppose $\mathbf{x} \in \text{span}(\mathbf{A})$. Then $\mathbf{x} = \mathbf{A}\mathbf{y}$ for some $\mathbf{y} \in \mathbb{C}^n$. Let $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{v}_j$ be an eigenvector expansion of \mathbf{y} . Since $\mathbf{A}\mathbf{v}_j = \mathbf{0}$ for $j = r+1, \dots, n$ we obtain

$$\mathbf{x} = \mathbf{A}\mathbf{y} = \sum_{j=1}^n c_j \mathbf{A}\mathbf{v}_j = \sum_{j=1}^r c_j \mathbf{A}\mathbf{v}_j \in \text{span}(\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r).$$

Finally, if $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{v}_j \in \ker(\mathbf{A})$ then $\mathbf{A}\mathbf{y} = \sum_{j=1}^r c_j \mathbf{A}\mathbf{v}_j = \mathbf{0}$, and $c_1 = \dots = c_r = 0$ since $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$ are linearly independent. But then $\mathbf{y} = \sum_{j=r+1}^n c_j \mathbf{v}_j \in \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n)$. \square

Definition 7.3 (Singular values)

We define the singular values $\sigma_1, \sigma_2, \dots, \sigma_n$ of $\mathbf{A} \in \mathbb{C}^{m \times n}$ to be the nonnegative square roots of the eigenvalues of $\mathbf{A}^* \mathbf{A}$. We will assume they are ordered so that for some integer r with $0 \leq r \leq n$,

$$\sigma_1 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots = \sigma_n. \quad (7.5)$$

Theorem 7.4 (rank=#positive singular values) The rank of a matrix is equal to the number of positive singular values.

Proof. Suppose \mathbf{A} has r positive singular values. By Theorem 7.2 $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ is an orthogonal basis for $\text{span}(\mathbf{A})$ and this means that \mathbf{A} has rank r . Conversely suppose $\text{rank}(\mathbf{A}) = r$ and \mathbf{A} has $s \neq r$ positive singular values. Again by Theorem 7.2 the set $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_s\}$ is a basis for $\text{span}(\mathbf{A})$ and $\text{rank}(\mathbf{A}) = s$ a contradiction. Thus $s = r$. \square

Every matrix has a singular value decomposition.

Theorem 7.5 (Existence of singular value decomposition)

Let $m, n \in \mathbb{N}$ and suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$. Then \mathbf{A} has a singular value decomposition of the form $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$, where $\mathbf{U} \in \mathbb{C}^{m \times m}$, $\mathbf{V} \in \mathbb{C}^{n \times n}$ are unitary, and $\Sigma \in \mathbb{R}^{m \times n}$ satisfies $\Sigma_{i,j} = 0$ for $i \neq j$ and $\Sigma_{i,i} = \sigma_i$, $i = 1, \dots, \min(m, n)$, with σ_i ordered singular values of \mathbf{A} . If \mathbf{A} is real then $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$, and \mathbf{U} , \mathbf{V} are real and orthonormal.

Proof. The existence proof is constructive. Let $(\lambda_j, \mathbf{v}_j)$, $j = 1, \dots, n$ be orthonormal eigenpairs for $\mathbf{A}^* \mathbf{A}$ ordered so that $\lambda_1 \geq \dots \geq \lambda_r > 0 = \lambda_{r+1} = \dots = \lambda_n$ for some r . Define $\sigma_i := \sqrt{\lambda_i}$, $i = 1, \dots, n$, Σ as in the theorem, and $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n]$. To define \mathbf{U} we start with

$$\mathbf{u}_j := \frac{1}{\sigma_j} \mathbf{A}\mathbf{v}_j, \text{ for } j = 1, \dots, r. \quad (7.6)$$

They are orthonormal since $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$ are orthogonal and $\sigma_j = \|\mathbf{A}\mathbf{v}_j\|_2 > 0$, $j = 1, \dots, r$ by Lemma 7.1. Extend $\mathbf{u}_1, \dots, \mathbf{u}_r$ in any way to an orthonormal basis $\mathbf{u}_1, \dots, \mathbf{u}_m$ for \mathbb{C}^m and define $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_m]$. Now

$$\mathbf{U}\Sigma = \mathbf{U}[\sigma_1\mathbf{e}_1, \dots, \sigma_r\mathbf{e}_r, \overbrace{\mathbf{0}, \dots, \mathbf{0}}^{n-r}] = [\sigma_1\mathbf{u}_1, \dots, \sigma_r\mathbf{u}_r, \mathbf{0}, \dots, \mathbf{0}] = [\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_n].$$

Thus $\mathbf{U}\Sigma = \mathbf{AV}$ and since \mathbf{V} is unitary we find $\mathbf{U}\Sigma\mathbf{V}^* = \mathbf{AVV}^* = \mathbf{A}$.

For a matrix with real elements the eigenvectors of $\mathbf{A}^T\mathbf{A}$ are real and the singular value decomposition takes the form $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$. \square

Suppose \mathbf{A} has rank r . The matrices $\Sigma, \mathbf{U}, \mathbf{V}$ in the SVD of \mathbf{A} can be written in block form

$$\begin{aligned} \Sigma &:= \begin{bmatrix} \Sigma_1 & \mathbf{0}_{r,n-r} \\ \mathbf{0}_{m-r,r} & \mathbf{0}_{m-r,n-r} \end{bmatrix} \in \mathbb{R}^{m \times n}, \text{ where } \Sigma_1 := \text{diag}(\sigma_1, \dots, \sigma_r), \\ \mathbf{U} &= [\mathbf{u}_1, \dots, \mathbf{u}_m] = [\mathbf{U}_1, \mathbf{U}_2], \quad \mathbf{U}_1 \in \mathbb{C}^{m,r}, \quad \mathbf{U}_2 \in \mathbb{C}^{m,m-r}, \\ \mathbf{V} &= [\mathbf{v}_1, \dots, \mathbf{v}_n] = [\mathbf{V}_1, \mathbf{V}_2], \quad \mathbf{V}_1 \in \mathbb{C}^{n,r}, \quad \mathbf{V}_2 \in \mathbb{C}^{n,n-r}. \end{aligned} \tag{7.7}$$

Here, for $k, l \geq 0$ the symbol $\mathbf{0}_{k,l} = []$ denotes the empty matrix if $k = 0$ or $l = 0$, and the zero matrix with k rows and l columns otherwise.

By block multiplication

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^* = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} = \mathbf{U}_1\Sigma_1\mathbf{V}_1^* = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*.$$

and we obtain a reduced factorization called the **singular value factorization (SVF)** and an outer product form of this factorization.

Corollary 7.6 (Singular value factorization)

Suppose $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$ is a singular value decomposition of a rank r matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$. Then \mathbf{A} has the singular value factorization

$$\mathbf{A} = \mathbf{U}_1\Sigma_1\mathbf{V}_1^* = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*, \quad \mathbf{U}_1 \in \mathbb{C}^{m,r}, \quad \Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_r), \quad \mathbf{V}_1 \in \mathbb{C}^{n,r},$$

where \mathbf{U}_1 and \mathbf{V}_1 have orthonormal columns.

A nonsingular square matrix has full rank and only positive singular values. Thus the SVD and SVF are the same for a nonsingular matrix.

The singular value decomposition is most often not unique. The matrix Σ is unique, but this is not true in general for \mathbf{U} and \mathbf{V} . Indeed, not even \mathbf{U}_1 and \mathbf{V}_1 are necessarily unique.

7.1.2 Examples

Example 7.7 (Nonsingular matrix)

Derive the following SVD.

$$\mathbf{A} := \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}. \quad (7.8)$$

Discussion: Eigenpairs of $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 97 & 96 \\ 96 & 153 \end{bmatrix} / 25$ are given by

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 9 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Taking square roots and normalizing we find $\sigma_1 = 3$, $\sigma_2 = 1$, $\mathbf{v}_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix} / 5$, and $\mathbf{v}_2 = \begin{bmatrix} 4 \\ -3 \end{bmatrix} / 5$. Thus $\mathbf{u}_1 := \mathbf{A} \mathbf{v}_1 / \sigma_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix} / 5$, $\mathbf{u}_2 := \mathbf{A} \mathbf{v}_2 / \sigma_2 = \begin{bmatrix} -4 \\ 3 \end{bmatrix} / 5$, and this shows (7.8). Since $m = n = r$ we have $\mathbf{U}_1 = \mathbf{U}$, $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}$ and $\mathbf{V}_1 = \mathbf{V}$ and the SVD and SVF are the same. The matrix \mathbf{A} is normal so that the singular values of \mathbf{A} are equal to the absolute value of the eigenvalues of \mathbf{A} (cf. Exercise 7.13). The eigenvalues of \mathbf{A} are $\lambda_1 = 3$ and $\lambda_2 = -1$. Thus $\lambda_2 \neq \sigma_2$.

Example 7.8 (Full row rank)

Find the singular value decomposition of

$$\mathbf{A} := \frac{1}{15} \begin{bmatrix} 14 & 4 & 16 \\ 2 & 22 & 13 \end{bmatrix} \in \mathbb{R}^{2 \times 3}.$$

Discussion: Eigenpairs of $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 8 & 4 & 10 \\ 4 & 20 & 14 \\ 10 & 14 & 17 \end{bmatrix} / 9$ are given by

$$\mathbf{B} \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} = 0 \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix}.$$

Thus $r = 2$ and

$$\boldsymbol{\Sigma} := \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{V} := \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix}.$$

From (7.6) we find $\mathbf{u}_1 = \mathbf{A} \mathbf{v}_1 / \sigma_1 = [3, 4]^T / 5$, $\mathbf{u}_2 = \mathbf{A} \mathbf{v}_2 / \sigma_2 = [4, -3]^T / 5$ and therefore

$$\mathbf{U} = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}.$$

Since $r = 2$, $\text{rank}(\mathbf{A}) = 2$, $\{\mathbf{u}_1, \mathbf{u}_2\}$ is an orthonormal basis for $\text{span}(\mathbf{A})$ and $\{\mathbf{v}_3\}$ is an orthonormal basis for $\ker(\mathbf{A})$. The SVF and outer product form of \mathbf{A} are

$$\mathbf{A} = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ 4 & -3 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 1 & -2 & 2 \\ 2 & 1 & -2 \end{bmatrix} = 2 \frac{1}{15} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \end{bmatrix} + 1 \frac{1}{15} \begin{bmatrix} -4 \\ 3 \end{bmatrix} \begin{bmatrix} 2 & -2 & 1 \end{bmatrix}.$$

Example 7.9 (Full column rank)

Find the SVD of

$$\mathbf{A}_1 = \frac{1}{15} \begin{bmatrix} 14 & 2 \\ 4 & 22 \\ 16 & 13 \end{bmatrix} \in \mathbb{R}^{3,2}.$$

Since $\mathbf{A}_1 = \mathbf{A}^T$, where \mathbf{A} is the matrix in Example 7.8 we can find an SVD of \mathbf{A}_1 by simply transposing the SVD of \mathbf{A} . Thus

$$\mathbf{A}_1 = (\mathbf{U}\Sigma\mathbf{V}^T)^T = \mathbf{V}\Sigma^T\mathbf{U}^T = \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}. \quad (7.9)$$

Alternatively we can follow the recipe from the proof of Theorem 7.5. Eigenpairs of

$$\mathbf{B}_1 = \mathbf{A}_1^T \mathbf{A}_1 = \frac{1}{25} \begin{bmatrix} 52 & 36 \\ 36 & 73 \end{bmatrix}$$

are found from

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 4 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = 1 \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Thus $\sigma_1 = 2$, $\sigma_2 = 1$, and $\mathbf{U} = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}$ as before. Now

$$\mathbf{v}_1 = \mathbf{A}_1 \mathbf{u}_1 / \sigma_1 = [1, 2, 2]^T / 3, \quad \mathbf{v}_2 = \mathbf{A}_1 \mathbf{u}_2 / \sigma_2 = [2, -2, 1]^T / 3.$$

Since $m = 3$ we also need \mathbf{v}_3 which should be orthogonal to \mathbf{v}_1 and \mathbf{v}_2 . $\mathbf{v}_3 = [2, 1, -2]^T$ is such a vector and normalizing \mathbf{v}_3 we obtain (7.9).

Example 7.10 ($r < n < m$)

Find the SVD of

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}.$$

Discussion: Eigenpairs of

$$\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

are derived from

$$\mathbf{B} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0 \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

and we find $\sigma_1 = 2$, $\sigma_2 = 0$. Thus $r = 1$, $m = 3$, $n = 2$ and

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \Sigma_1 = [2], \quad \mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

The equation (7.6) implies $\mathbf{u}_1 = \mathbf{Av}_1/\sigma_1 = \mathbf{s}_1/\sqrt{2}$, where $\mathbf{s}_1 = [1, 1, 0]^T$. To find the other columns of \mathbf{U} we can extend \mathbf{s}_1 to a basis $\{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3\}$ for \mathbb{R}^3 , apply the Gram-Schmidt orthogonalization process to $\{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3\}$, and then normalize. Choosing the basis

$$\mathbf{s}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{s}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{s}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

we find from (29)

$$\mathbf{w}_1 = \mathbf{s}_1, \quad \mathbf{w}_2 = \mathbf{s}_2 - \frac{\mathbf{s}_2^T \mathbf{w}_1}{\mathbf{w}_1^T \mathbf{w}_1} \mathbf{w}_1 = \begin{bmatrix} -1/2 \\ 1/2 \\ 0 \end{bmatrix}, \quad \mathbf{w}_3 = \mathbf{s}_3 - \frac{\mathbf{s}_3^T \mathbf{w}_1}{\mathbf{w}_1^T \mathbf{w}_1} \mathbf{w}_1 - \frac{\mathbf{s}_3^T \mathbf{w}_2}{\mathbf{w}_2^T \mathbf{w}_2} \mathbf{w}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Normalizing the \mathbf{w}_i 's we obtain $\mathbf{u}_1 = \mathbf{w}_1/\|\mathbf{w}_1\|_2 = [1/\sqrt{2}, 1/\sqrt{2}, 0]^T$, $\mathbf{u}_2 = \mathbf{w}_2/\|\mathbf{w}_2\|_2 = [-1/\sqrt{2}, 1/\sqrt{2}, 0]^T$, and $\mathbf{u}_3 = \mathbf{w}_3/\|\mathbf{w}_3\|_2 = [0, 0, 1]^T$. Therefore, $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$, where

$$\mathbf{U} := \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{3,3}, \quad \Sigma := \begin{bmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{3,2}, \quad \mathbf{V} := \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} \in \mathbb{R}^{2,2}.$$

Exercise 7.11 (SVD examples)

Find the singular value decomposition of the following matrices

$$(a) \quad A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}.$$

$$(b) \quad A = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 2 & 2 \end{bmatrix}.$$

Exercise 7.12 (More SVD examples)

Find the singular value decomposition of the following matrices

$$(a) \quad A = \mathbf{e}_1 \text{ the first unit vector in } \mathbb{R}^m.$$

$$(b) \quad A = \mathbf{e}_n^T \text{ the last unit vector in } \mathbb{R}^n.$$

$$(c) \quad A = \begin{bmatrix} -1 & 0 \\ 0 & 3 \end{bmatrix}.$$

Exercise 7.13 (Singular values of a normal matrix)

Show that the singular values of a normal matrix are the absolute values of its eigenvalues.¹⁰

The method we used to find the singular value decomposition in the previous examples and exercises can be suitable for hand calculation with small matrices, but it is not appropriate as a basis for a general purpose numerical method. In particular, the Gram-Schmidt orthogonalization process is not numerically stable, and forming $\mathbf{A}^*\mathbf{A}$ can lead to extra errors in the computation. Standard computer implementations of the singular value decomposition ([25]) first reduces \mathbf{A} to bidiagonal form and then use an adapted version of the QR algorithm where the matrix $\mathbf{A}^*\mathbf{A}$ is not formed. The QR algorithm is discussed in Chapter 14.

7.2 SVD and the Four Fundamental Subspaces

The columns of the singular vectors form orthonormal bases for the four fundamental subspaces $\text{span}(\mathbf{A})$, $\ker(\mathbf{A})$, $\text{span}(\mathbf{A}^*)$, and $\ker(\mathbf{A}^*)$. To show this we will use the following lemma.

Lemma 7.14 (SVD of $\mathbf{A}^*\mathbf{A}$ and $\mathbf{A}\mathbf{A}^*$)

In any singular value decomposition $\mathbf{U}\Sigma\mathbf{V}^*$ of \mathbf{A} the columns of \mathbf{U} and \mathbf{V} are eigenvectors of the matrices $\mathbf{A}\mathbf{A}^*$ and $\mathbf{A}^*\mathbf{A}$, respectively. Moreover, $\mathbf{A}\mathbf{A}^*$ and $\mathbf{A}^*\mathbf{A}$ have the same nonzero eigenvalues with the same algebraic multiplicity.

Proof. We find

$$\begin{aligned}\mathbf{A}\mathbf{A}^* &= (\mathbf{U}\Sigma\mathbf{V}^*)(\mathbf{V}\Sigma^T\mathbf{U}^*) = \mathbf{U}\mathbf{D}_1\mathbf{U}^*, \quad \mathbf{D}_1 = \text{diag}(\sigma_1^2, \dots, \sigma_r^2, 0, \dots, 0) \in \mathbb{R}^{m \times m}, \\ \mathbf{A}^*\mathbf{A} &= (\mathbf{V}\Sigma^T\mathbf{U}^*)(\mathbf{U}\Sigma\mathbf{V}^*) = \mathbf{V}\mathbf{D}_2\mathbf{V}^*, \quad \mathbf{D}_2 = \text{diag}(\sigma_1^2, \dots, \sigma_r^2, 0, \dots, 0) \in \mathbb{R}^{n \times n}.\end{aligned}$$

But these are spectral decompositions of $\mathbf{A}\mathbf{A}^*$ and $\mathbf{A}^*\mathbf{A}$ and the result follows¹¹ \square

Theorem 7.15 (Singular vectors and orthonormal bases)

For positive integers m, n let $\mathbf{A} \in \mathbb{C}^{m \times n}$ have rank r and a singular value decomposition $\mathbf{A} = [\mathbf{u}_1, \dots, \mathbf{u}_m]\Sigma[\mathbf{v}_1, \dots, \mathbf{v}_n]^* = \mathbf{U}\Sigma\mathbf{V}^*$. Then the singular vectors satisfy

$$\begin{aligned}\mathbf{A}\mathbf{v}_i &= \sigma_i\mathbf{u}_i, \quad i = 1, \dots, r, \quad \mathbf{A}\mathbf{v}_i = 0, \quad i = r+1, \dots, n, \\ \mathbf{A}^*\mathbf{u}_i &= \sigma_i\mathbf{v}_i, \quad i = 1, \dots, r, \quad \mathbf{A}^*\mathbf{u}_i = 0, \quad i = r+1, \dots, m.\end{aligned}\tag{7.10}$$

¹⁰hint: Compute $\mathbf{A}^*\mathbf{A}$ using the representation for \mathbf{A} in Theorem 6.37.

¹¹That the two matrices have the same eigenvalues also follows from a more general result. Indeed, the characteristic polynomials are related by $\lambda^n \rho_{\mathbf{A}\mathbf{A}^*}(\lambda) = \lambda^m \rho_{\mathbf{A}^*\mathbf{A}}(\lambda)$ for $\lambda \in \mathbb{C}$. (cf. Theorem 6.8).

Moreover,

1. $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ is an orthonormal basis for $\text{span}(\mathbf{A})$,
 2. $\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$ is an orthonormal basis for $\ker(\mathbf{A}^*)$,
 3. $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ is an orthonormal basis for $\text{span}(\mathbf{A}^*)$,
 4. $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ is an orthonormal basis for $\ker(\mathbf{A})$.
- (7.11)

Proof. If $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$ then $\mathbf{AV} = \mathbf{U}\Sigma$, or in terms of the block partition (7.7) $\mathbf{A}[\mathbf{V}_1, \mathbf{V}_2] = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$. But then $\mathbf{AV}_1 = \mathbf{U}_1\Sigma$, $\mathbf{AV}_2 = \mathbf{0}$, and this implies the first part of (7.10). Taking transpose of $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$ gives $\mathbf{A}^* = \mathbf{V}\Sigma^T\mathbf{U}^*$ or $\mathbf{A}^*\mathbf{U} = \mathbf{V}\Sigma^T$. Using the block partition as before we obtain the last part of (7.10).

Since $\mathbf{v}_1, \dots, \mathbf{v}_n$ are orthonormal eigenvectors for $\mathbf{A}^*\mathbf{A}$ (cf. Lemma 7.14) it follows from (7.10) and Theorem 7.2 that $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ is an orthonormal basis for $\ker(\mathbf{A})$ and $\{\mathbf{Av}_1, \dots, \mathbf{Av}_r\}$ is an orthogonal basis for $\text{span}(\mathbf{A})$. By (7.10) $\mathbf{u}_j := \frac{1}{\sigma_j} \mathbf{Av}_j$, $j = 1, \dots, r$ and this implies that $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ is an orthonormal basis for $\text{span}(\mathbf{A})$. Also by Lemma 7.14 $\mathbf{u}_1, \dots, \mathbf{u}_m$ are orthonormal eigenvectors for \mathbf{AA}^* with the same nonzero eigenvalues as $\mathbf{A}^*\mathbf{A}$. Applying (7.10) and Theorem 7.2 to \mathbf{A}^* it follows that $\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$ is an orthonormal basis for $\ker(\mathbf{A}^*)$ and $\{\mathbf{A}^*\mathbf{u}_1, \dots, \mathbf{A}^*\mathbf{u}_r\}$ form an orthogonal basis for $\text{span}(\mathbf{A}^*)$. By (7.10) $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ is an orthonormal basis for $\text{span}(\mathbf{A}^*)$. \square

By counting the number of elements in the four bases in (7.11) we obtain from Theorem 7.15 a new proof of a fundamental result in matrix analysis. Recall that $\text{null}(\mathbf{A})$ is defined as the dimension of $\ker(\mathbf{A})$.

Corollary 7.16 (Counting dimensions of fundamental subspaces)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$. Then

1. $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}) = n$,
2. $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}^*) = m$,
3. $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^*)$.

The next result can also be useful.

Theorem 7.17 (Rank and nullity relations)

For any $\mathbf{A} \in \mathbb{C}^{m \times n}$ we have

1. $\text{rank } \mathbf{A} = \text{rank}(\mathbf{A}^*\mathbf{A}) = \text{rank}(\mathbf{AA}^*)$,
2. $\text{null}(\mathbf{A}^*\mathbf{A}) = \text{null } \mathbf{A}$, and $\text{null}(\mathbf{AA}^*) = \text{null}(\mathbf{A}^*)$,

Proof. The three matrices \mathbf{A} , $\mathbf{A}^*\mathbf{A}$, and $\mathbf{A}\mathbf{A}^*$ have the same rank since they have the same number of positive singular values (cf. Lemma 7.14). The result about the dimension of the null spaces follows from Corollary 7.16. Indeed, $\text{null}(\mathbf{A}^*\mathbf{A}) = n - \text{rank}(\mathbf{A}^*\mathbf{A}) = n - \text{rank}(\mathbf{A}) = \text{null}(\mathbf{A})$ and $\text{null}(\mathbf{A}\mathbf{A}^*) = m - \text{rank}(\mathbf{A}\mathbf{A}^*) = m - \text{rank}(\mathbf{A}) = \text{null}(\mathbf{A}^*)$. \square

Exercise 7.18 (Orthonormal bases example)

Let \mathbf{A} and \mathbf{B} be as in Example 7.9. Give orthonormal bases for $\text{span}(\mathbf{B})$ and $\ker(\mathbf{B})$ and explain why $\text{span}(\mathbf{B}) \oplus \ker(\mathbf{A})$ is an orthogonal decomposition of \mathbb{R}^3 .

Exercise 7.19 (Some spanning sets)

Show for any $\mathbf{A} \in \mathbb{C}^{m \times n}$ that $\text{span}(\mathbf{A}^*\mathbf{A}) = \text{span}(\mathbf{V}_1) = \text{span}(\mathbf{A}^*)$

Exercise 7.20 (Singular values and eigenpair of composite matrix)

Let $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m \geq n$ have singular values $\sigma_1, \dots, \sigma_n$, left singular vectors $\mathbf{u}_1, \dots, \mathbf{u}_m \in \mathbb{C}^m$, and right singular vectors $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{C}^n$. Show that the matrix

$$\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix}$$

has the $n+m$ eigenpairs

$$\{(\sigma_1, \mathbf{p}_1), \dots, (\sigma_n, \mathbf{p}_n), (-\sigma_1, \mathbf{q}_1), \dots, (-\sigma_n, \mathbf{q}_n), (0, \mathbf{r}_{n+1}), \dots, (0, \mathbf{r}_m)\},$$

where

$$\mathbf{p}_i = \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix}, \quad \mathbf{q}_i = \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix}, \quad \mathbf{r}_j = \begin{bmatrix} \mathbf{u}_j \\ \mathbf{0} \end{bmatrix}, \text{ for } i = 1, \dots, n \text{ and } j = n+1, \dots, m.$$

7.3 A Geometric Interpretation

The singular value decomposition gives insight into the geometry of a linear transformation. Consider the linear transformation $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ given by $\mathbf{z} \rightarrow \mathbf{Az}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$. The function \mathbf{T} maps the unit sphere $\mathcal{S} := \{\mathbf{z} \in \mathbb{R}^n : \|\mathbf{z}\|_2 = 1\}$ onto an ellipsoid in \mathbb{R}^m . The singular values are the length of the semiaxes. We describe this in the square nonsingular case. Suppose $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ is a singular value decomposition of $\mathbf{A} \in \mathbb{R}^{n \times n}$. Since \mathbf{A} has rank n we have $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$, with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$ and $\mathbf{A}^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$. Let $\mathcal{E} := \mathbf{A}\mathcal{S} \subset \mathbb{R}^n$ be the image of \mathcal{S} under the transformation \mathbf{T} . If $\mathbf{x} \in \mathcal{E}$ then $\mathbf{x} = \mathbf{Az}$ for some $\mathbf{z} \in \mathcal{S}$ and we find

$$\begin{aligned} 1 &= \|\mathbf{z}\|_2^2 = \|\mathbf{A}^{-1}\mathbf{Az}\|_2^2 = \|\mathbf{A}^{-1}\mathbf{x}\|_2^2 = \|\mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{x}\|_2^2 \\ &= \|\Sigma^{-1}\mathbf{U}^T\mathbf{x}\|_2^2 = \|\Sigma^{-1}\mathbf{y}\|_2^2 = \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2}, \end{aligned}$$

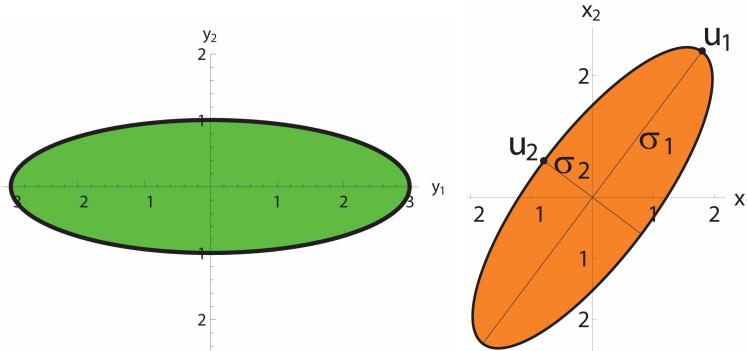


Figure 7.1: The ellipse $y_1^2/9 + y_2^2 = 1$ (left) and the rotated ellipse \mathbf{AS} (right).

where $\mathbf{y} := \mathbf{U}^T \mathbf{x}$ and we used $\|\mathbf{V}\mathbf{v}\|_2 = \|\mathbf{v}\|_2$ for a vector \mathbf{v} . The equation $1 = \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2}$ describes an ellipsoid in \mathbb{R}^n with semiaxes of length σ_j along the unit vectors \mathbf{e}_j for $j = 1, \dots, n$. Since the orthonormal transformation $\mathbf{U}\mathbf{y} \rightarrow \mathbf{x}$ preserves length, the image $\mathcal{E} = \mathbf{AS}$ is an ellipsoid with semiaxes along the left singular vectors $\mathbf{u}_j = \mathbf{U}\mathbf{e}_j$ of length σ_j . Since $\mathbf{A}\mathbf{v}_j = \sigma_j \mathbf{u}_j$, the right singular vectors defines points in \mathcal{S} that are mapped onto the semiaxes of \mathcal{E} .

Example 7.21 (Ellipse)

Consider the transformation $\mathbf{A} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by the matrix

$$\mathbf{A} := \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix}$$

in Example 7.7. Recall that $\sigma_1 = 3$, $\sigma_2 = 1$, $\mathbf{u}_1 = [3, 4]^T/5$ and $\mathbf{u}_2 = [-4, 3]^T/5$. The ellipses $y_1^2/\sigma_1^2 + y_2^2/\sigma_2^2 = 1$ and $\mathcal{E} = \mathbf{AS}$ are shown in Figure 7.1. Since $\mathbf{y} = \mathbf{U}^T \mathbf{x} = [3/5x_1 + 4/5x_2, -4/5x_1 + 3/5x_2]^T$, the equation for the ellipse on the right is

$$\frac{(\frac{3}{5}x_1 + \frac{4}{5}x_2)^2}{9} + \frac{(-\frac{4}{5}x_1 + \frac{3}{5}x_2)^2}{1} = 1,$$

7.4 Determining the Rank of a Matrix Numerically

In many elementary linear algebra courses a version of Gaussian elimination, called Gauss-Jordan elimination, is used to determine the rank of a matrix. To carry this out by hand for a large matrix can be a Herculean task and using a computer and floating point arithmetic the result will not be reliable. Entries, which in the final result should have been zero, will have nonzero values because of round-off errors.

As an alternative we can use the singular value decomposition to determine rank. Although success is not at all guaranteed, the result will be more reliable than if Gauss-Jordan elimination is used.

By Theorem 7.5 the rank of a matrix is equal to the number of nonzero singular values and if we have computed the singular values, then all we have to do is to count the nonzero ones. The problem however is the same as for Gaussian elimination. Due to round-off errors none of the computed singular values are likely to be zero.

7.4.1 The Frobenius Norm

This norm will be used in a discussion of how many of the computed singular values can possibly be considered to be zero. The **Frobenius norm**, of a matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ is defined by

$$\|\mathbf{A}\|_F := \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}. \quad (7.12)$$

There is a relation between the Frobenius norm of a matrix and its singular values. First we derive some elementary properties of this norm. A systematic study of matrix norms is given in the next chapter.

Lemma 7.22 (Frobenius norm properties)

For any $m, n \in \mathbb{N}$ and any matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$

1. $\|\mathbf{A}^*\|_F = \|\mathbf{A}\|_F$,
2. $\|\mathbf{A}\|_F^2 = \sum_{j=1}^n \|\mathbf{a}_{:j}\|_2^2$,
3. $\|\mathbf{U}\mathbf{A}\|_F = \|\mathbf{A}\mathbf{V}\|_F = \|\mathbf{A}\|_F$ for any unitary matrices $\mathbf{U} \in \mathbb{C}^{m \times m}$ and $\mathbf{V} \in \mathbb{C}^{n \times n}$,
4. $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$ for any $\mathbf{B} \in \mathbb{C}^{n,k}$, $k \in \mathbb{N}$,
5. $\|\mathbf{Ax}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$, for all $\mathbf{x} \in \mathbb{C}^n$.

Proof.

1. $\|\mathbf{A}^*\|_F^2 = \sum_{j=1}^n \sum_{i=1}^m |\bar{a}_{ij}|^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 = \|\mathbf{A}\|_F^2$.
2. This follows since the Frobenius norm is the Euclidian norm of a vector, $\|\mathbf{A}\|_F := \|\text{vec}(\mathbf{A})\|_2$, where $\text{vec}(\mathbf{A}) \in \mathbb{C}^{mn}$ is the vector obtained by stacking the columns of \mathbf{A} on top of each other.

3. Recall that $\|\mathbf{U}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \mathbb{C}^n$ if $\mathbf{U}^*\mathbf{U} = I$. Applying this to each column $\mathbf{a}_{:j}$ of \mathbf{A} we find $\|\mathbf{U}\mathbf{A}\|_F^2 \stackrel{2}{=} \sum_{j=1}^n \|\mathbf{U}\mathbf{a}_{:j}\|_2^2 = \sum_{j=1}^n \|\mathbf{a}_{:j}\|_2^2 \stackrel{2}{=} \|\mathbf{A}\|_F^2$. Similarly, since $\mathbf{V}\mathbf{V}^* = I$ we find $\|\mathbf{A}\mathbf{V}\|_F \stackrel{1}{=} \|\mathbf{V}^*\mathbf{A}^*\|_F = \|\mathbf{A}^*\|_F \stackrel{1}{=} \|\mathbf{A}\|_F$.
4. Using the Cauchy-Schwarz inequality and 2. we obtain

$$\|\mathbf{AB}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^k |\mathbf{a}_{:i}^T \mathbf{b}_{:j}|^2 \leq \sum_{i=1}^m \sum_{j=1}^k \|\mathbf{a}_{:i}\|_2^2 \|\mathbf{b}_{:j}\|_2^2 = \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.$$

5. Since $\|\mathbf{v}\|_F = \|\mathbf{v}\|_2$ for a vector this follows by taking $k = 1$ and $\mathbf{B} = \mathbf{x}$ in 4.

□

Theorem 7.23 (Frobenius norm and singular values)

We have $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_n^2}$, where $\sigma_1, \dots, \sigma_n$ are the singular values of \mathbf{A} .

Proof. Using Lemma 7.22 we find $\|\mathbf{A}\|_F \stackrel{3}{=} \|\mathbf{U}^*\mathbf{AV}\|_F = \|\Sigma\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_n^2}$. □

7.4.2 Low Rank Approximation

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m \geq n$ has singular value decomposition $\mathbf{A} = \mathbf{U} \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*$, where $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$. We choose $\epsilon > 0$ and let $1 \leq r \leq n$ be the smallest integer such that $\sigma_{r+1}^2 + \dots + \sigma_n^2 < \epsilon^2$. Define $\mathbf{A}' := \mathbf{U} \begin{bmatrix} \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*$, where $\mathbf{D}' := \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{n \times n}$. By Lemma 7.22

$$\|\mathbf{A} - \mathbf{A}'\|_F = \|\mathbf{U} \begin{bmatrix} \mathbf{D} - \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*\|_F = \|\begin{bmatrix} \mathbf{D} - \mathbf{D}' \\ \mathbf{0} \end{bmatrix}\|_F = \sqrt{\sigma_{r+1}^2 + \dots + \sigma_n^2} < \epsilon.$$

Thus, if ϵ is small then \mathbf{A} is near a matrix \mathbf{A}' of rank r . This can be used to determine rank numerically. We choose an r such that $\sqrt{\sigma_{r+1}^2 + \dots + \sigma_n^2}$ is “small”. Then we postulate that $\text{rank}(\mathbf{A}) = r$ since \mathbf{A} is close to a matrix of rank r .

The following theorem shows that of all $m \times n$ matrices of rank r , \mathbf{A}' is closest to \mathbf{A} measured in the Frobenius norm.

Theorem 7.24 (Best low rank approximation)

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ has singular values $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. For any $r \leq \text{rank}(\mathbf{A})$ we have

$$\|\mathbf{A} - \mathbf{A}'\|_F = \min_{\substack{\mathbf{B} \in \mathbb{R}^{m \times n} \\ \text{rank}(\mathbf{B})=r}} \|\mathbf{A} - \mathbf{B}\|_F = \sqrt{\sigma_{r+1}^2 + \dots + \sigma_n^2}.$$

For the proof of this theorem we refer to p. 322 of [25].

Exercise 7.25 (Rank example)

Consider the singular value decomposition

$$\mathbf{A} := \begin{bmatrix} 0 & 3 & 3 \\ 4 & 1 & -1 \\ 4 & 1 & -1 \\ 0 & 3 & 3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 6 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ \frac{2}{3} & -\frac{1}{3} & -\frac{2}{3} \\ \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \end{bmatrix}$$

- (a) Give orthonormal bases for $\text{span}(\mathbf{A})$, $\text{span}(\mathbf{A}^T)$, $\ker(\mathbf{A})$, $\ker(\mathbf{A}^T)$ and $\text{span}(\mathbf{A})^\perp$.
- (b) Explain why for all matrices $\mathbf{B} \in \mathbb{R}^{4,3}$ of rank one we have $\|\mathbf{A} - \mathbf{B}\|_F \geq 6$.
- (c) Give a matrix \mathbf{A}_1 of rank one such that $\|\mathbf{A} - \mathbf{A}_1\|_F = 6$.

Exercise 7.26 (Another rank example)

Let \mathbf{A} be the $n \times n$ matrix that for $n = 4$ takes the form

$$\mathbf{A} = \begin{bmatrix} 1 & -1 & -1 & -1 \\ 0 & 1 & -1 & -1 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Thus \mathbf{A} is upper triangular with diagonal elements one and all elements above the diagonal equal to -1 . Let \mathbf{B} be the matrix obtained from \mathbf{A} by changing the $(n, 1)$ element from zero to -2^{2-n} .

- (a) Show that $\mathbf{Bx} = \mathbf{0}$, where $\mathbf{x} := [2^{n-2}, 2^{n-3}, \dots, 2^0, 1]^T$. Conclude that \mathbf{B} is singular, $\det(\mathbf{A}) = 1$, and $\|\mathbf{A} - \mathbf{B}\|_F = 2^{2-n}$. Thus even if $\det(\mathbf{A})$ is not small Frobenius norm of $\mathbf{A} - \mathbf{B}$ is small for large n , and the matrix \mathbf{A} is very close to being singular for large n .
- (b) Use Theorem 7.24 to show that the smallest singular value σ_n of \mathbf{A} is bounded above by 2^{2-n} .

7.5 The Minmax Theorem for Singular Values and the Hoffman-Wielandt Theorem

We have a minmax and maxmin characterization for singular values.

Theorem 7.27 (The Courant-Fischer Theorem for Singular Values)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ has singular values $\sigma_1, \sigma_2, \dots, \sigma_n$ ordered so that $\sigma_1 \geq \dots \geq \sigma_n$. Then for $k = 1, \dots, n$

$$\sigma_k = \min_{\dim(S)=n-k+1} \max_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} = \max_{\dim(S)=k} \min_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}. \quad (7.13)$$

Proof. Since

$$\frac{\|\mathbf{A}\mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} = \frac{(\mathbf{A}\mathbf{x})^*(\mathbf{A}\mathbf{x})}{\mathbf{x}^*\mathbf{x}} = \frac{\mathbf{x}^*(\mathbf{A}^*\mathbf{A})\mathbf{x}}{\mathbf{x}^*\mathbf{x}}$$

is the Rayleigh quotient $R_{\mathbf{A}^*\mathbf{A}}(\mathbf{x})$ of $\mathbf{A}^*\mathbf{A}$, and since the singular values of \mathbf{A} are the nonnegative square roots of the eigenvalues of $\mathbf{A}^*\mathbf{A}$, the results follow from the Courant-Fischer Theorem for eigenvalues, see Theorem 6.45. \square

By taking $k = 1$ and $k = n$ in (7.13) we obtain for any $\mathbf{A} \in \mathbb{C}^{m \times n}$

$$\sigma_1 = \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}, \quad \sigma_n = \min_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}. \quad (7.14)$$

This follows since the only subspace of \mathbb{C}^n of dimension n is \mathbb{C}^n itself.

The Hoffman-Wielandt Theorem, see Theorem 6.48, for eigenvalues of Hermitian matrices can be written

$$\sum_{j=1}^n |\mu_j - \lambda_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2 := \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2, \quad (7.15)$$

where $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ are both Hermitian matrices with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ and $\mu_1 \geq \dots \geq \mu_n$, respectively.

For singular values we have a similar result, see also Section 12.6.

Theorem 7.28 (Hoffman-Wielandt Theorem for singular values)

For any $m, n \in \mathbb{N}$ and $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m \times n}$ we have

$$\sum_{j=1}^n |\beta_j - \alpha_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2. \quad (7.16)$$

where $\alpha_1 \geq \dots \geq \alpha_n$ and $\beta_1 \geq \dots \geq \beta_n$ are the singular values of \mathbf{A} and \mathbf{B} , respectively.

7.5.1 Proof of the Hoffman-Wielandt theorem for singular values

We apply the Hoffman-Wielandt Theorem for eigenvalues to the Hermitian matrices

$$\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \text{ and } \mathbf{D} := \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ \mathbf{B}^* & \mathbf{0} \end{bmatrix} \in \mathbb{C}^{m+n \times m+n}.$$

If \mathbf{C} and \mathbf{D} has eigenvalues $\lambda_1 \geq \dots \geq \lambda_{m+n}$ and $\mu_1 \geq \dots \geq \mu_{m+n}$, respectively then

$$\sum_{j=1}^{m+n} |\lambda_j - \mu_j|^2 \leq \|\mathbf{C} - \mathbf{D}\|_F^2. \quad (7.17)$$

Suppose \mathbf{A} has rank r and singular value decomposition $\mathbf{U}\Sigma\mathbf{V}^*$. We use (7.10) and determine the eigenpairs of \mathbf{C} as follows.

$$\begin{aligned} \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} \mathbf{A}\mathbf{v}_i \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} \alpha_i \mathbf{u}_i \\ \alpha_i \mathbf{v}_i \end{bmatrix} = \alpha_i \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix}, \quad i = 1, \dots, r, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} -\mathbf{A}\mathbf{v}_i \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} -\alpha_i \mathbf{u}_i \\ \alpha_i \mathbf{v}_i \end{bmatrix} = -\alpha_i \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix}, \quad i = 1, \dots, r, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{0} \end{bmatrix} &= \begin{bmatrix} \mathbf{0} \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} = 0 \begin{bmatrix} \mathbf{u}_i \\ \mathbf{0} \end{bmatrix}, \quad i = r+1, \dots, m, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} \mathbf{A}\mathbf{v}_i \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} = 0 \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_i \end{bmatrix}, \quad i = r+1, \dots, n. \end{aligned}$$

Thus \mathbf{C} has the $2r$ eigenvalues $\alpha_1, -\alpha_1, \dots, \alpha_r, -\alpha_r$ and $m+n-2r$ additional zero eigenvalues. Similarly, if \mathbf{B} has rank s then \mathbf{D} has the $2s$ eigenvalues $\beta_1, -\beta_1, \dots, \beta_s, -\beta_s$ and $m+n-2s$ additional zero eigenvalues. Let

$$t := \max(r, s).$$

Then

$$\begin{aligned} \lambda_1 \geq \dots \geq \lambda_{m+n} &= \alpha_1 \geq \dots \geq \alpha_t \geq 0 = \dots = 0 \geq -\alpha_t \geq \dots \geq -\alpha_1, \\ \mu_1 \geq \dots \geq \mu_{m+n} &= \beta_1 \geq \dots \geq \beta_t \geq 0 = \dots = 0 \geq -\beta_t \geq \dots \geq -\beta_1. \end{aligned}$$

We find $\sum_{j=1}^{m+n} |\lambda_j - \mu_j|^2 = 2 \sum_{i=1}^t |\alpha_i - \beta_i|^2$ and

$$\|\mathbf{C} - \mathbf{D}\|_F^2 = \left\| \begin{bmatrix} \mathbf{0} & \mathbf{A} - \mathbf{B} \\ \mathbf{A}^* - \mathbf{B}^* & \mathbf{0} \end{bmatrix} \right\|_F^2 = \|\mathbf{B} - \mathbf{A}\|_F^2 + \|(\mathbf{B} - \mathbf{A})^*\|_F^2 = 2\|\mathbf{B} - \mathbf{A}\|_F^2.$$

But then (7.17) implies $\sum_{i=1}^t |\alpha_i - \beta_i|^2 \leq \|\mathbf{B} - \mathbf{A}\|_F^2$. Since $t \leq n$ and $\alpha_i = \beta_i = 0$ for $i = t+1, \dots, n$ we obtain (7.16).

7.6 Review Questions

7.6.1 Consider an SVD and an SVF of a matrix \mathbf{A} .

- What are the singular values of \mathbf{A} ?
- how is the SVD defined?
- how can we find an SVF if we know an SVD?
- how can we find an SVD if we know an SVF?
- what are the relations between the singular vectors?
- which singular vectors form bases for $\text{span}(\mathbf{A})$ and $\ker(\mathbf{A}^*)$?

7.6.2 How are the Frobenius norm and singular values related?

7.6.3 State the Courant-Fischer theorem for singular values.

7.6.4 State the Hoffman-Wieland theorem for singular values.

Chapter 8

Matrix Norms

To measure the size of a matrix we can use a matrix norm. In this chapter we initiate a systematic study of matrix norms.

8.1 Matrix Norms

For simplicity we consider only matrix norms on the vector space $(\mathbb{C}^{m \times n}, \mathbb{C})$. All results also holds for $(\mathbb{R}^{m \times n}, \mathbb{R})$.

Definition 8.1 (Matrix Norms)

Suppose m, n are positive integers. A function $\|\cdot\|: \mathbb{C}^{m \times n} \rightarrow \mathbb{R}$ is called a **matrix norm** on $\mathbb{C}^{m \times n}$ if for all $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m \times n}$ and all $c \in \mathbb{C}$

1. $\|\mathbf{A}\| \geq 0$ with equality if and only if $\mathbf{A} = 0$. (positivity)
2. $\|c\mathbf{A}\| = |c| \|\mathbf{A}\|$. (homogeneity)
3. $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$. (subadditivity)

A matrix norm is simply a vector norm on the finite dimensional vector space $(\mathbb{C}^{m \times n}, \mathbb{C})$ of $m \times n$ matrices. Adapting Theorem 0.20 to this situation gives

Theorem 8.2 (Matrix norm equivalence)

All matrix norms are equivalent. Thus, if $\|\cdot\|$ and $\|\cdot\|'$ are two matrix norms on $\mathbb{C}^{m \times n}$ then there are positive constants μ and M such that

$$\mu \|\mathbf{A}\| \leq \|\mathbf{A}\|' \leq M \|\mathbf{A}\|$$

holds for all $\mathbf{A} \in \mathbb{C}^{m \times n}$. Moreover, a matrix norm is a continuous function.

From any vector norm $\|\cdot\|_V$ on \mathbb{C}^{mn} we can define a matrix norm on $\mathbb{C}^{m \times n}$ by $\|\mathbf{A}\| := \|\text{vec}(\mathbf{A})\|_V$, where $\text{vec}(\mathbf{A}) \in \mathbb{C}^{mn}$ is the vector obtained by stacking the columns of \mathbf{A} on top of each other. In particular, to the p vector norms for $p = 1, 2, \infty$, we have the corresponding **sum norm**, **Frobenius norm**, and **max norm** defined by

$$\|\mathbf{A}\|_S := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|, \quad \|\mathbf{A}\|_F := \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}, \quad \|\mathbf{A}\|_M := \max_{i,j} |a_{ij}|. \quad (8.1)$$

Of these norms the Frobenius norm is the most useful. Some of its properties were derived in Lemma 7.22 and Theorem 7.23.

8.1.1 Consistent and Subordinate Matrix Norms

Since matrices can be multiplied it is useful to have an analogue of subadditivity for matrix multiplication. For square matrices the product \mathbf{AB} is defined in a fixed space $\mathbb{C}^{n \times n}$, while in the rectangular case matrix multiplication combines matrices in different spaces. The following definition captures this distinction.

Definition 8.3 (Consistent Matrix Norms)

A matrix norm is called **consistent** on $\mathbb{C}^{n \times n}$ if

$$4. \quad \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (\text{submultiplicativity})$$

holds for all $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$. A matrix norm is **consistent** if it is defined on $\mathbb{C}^{m \times n}$ for all $m, n \in \mathbb{N}$, and 4. holds for all matrices \mathbf{A}, \mathbf{B} for which the product \mathbf{AB} is defined.

Clearly the three norms in (8.1) are defined for all $m, n \in \mathbb{N}$. From Lemma 7.22 it follows that the Frobenius norm is consistent.

Exercise 8.4 (Consistency of sum norm?)

Show that the sum norm is consistent.

Exercise 8.5 (Consistency of max norm?)

Show that the max norm is not consistent by considering $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$.

Exercise 8.6 (Consistency of modified max norm?)

(a) Show that the norm

$$\|\mathbf{A}\| := \sqrt{mn} \|\mathbf{A}\|_M, \quad \mathbf{A} \in \mathbb{C}^{m \times n}$$

is a consistent matrix norm.

(b) Show that the constant \sqrt{mn} can be replaced by m and by n .

For a consistent matrix norm on $\mathbb{C}^{n \times n}$ we have the inequality

$$\|\mathbf{A}^k\| \leq \|\mathbf{A}\|^k \text{ for } k \in \mathbb{N}. \quad (8.2)$$

When working with norms one often has to bound the vector norm of a matrix times a vector by the norm of the matrix times the norm of the vector. This leads to the following definition.

Definition 8.7 (Subordinate Matrix Norms)

Suppose $m, n \in \mathbb{N}$ are given, let $\|\cdot\|_\alpha$ on \mathbb{C}^m and $\|\cdot\|_\beta$ on \mathbb{C}^n be vector norms, and let $\|\cdot\|$ be a matrix norm on $\mathbb{C}^{m \times n}$. We say that the matrix norm $\|\cdot\|$ is **subordinate** to the vector norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ if $\|\mathbf{Ax}\|_\alpha \leq \|\mathbf{A}\| \|\mathbf{x}\|_\beta$ for all $\mathbf{A} \in \mathbb{C}^{m \times n}$ and all $\mathbf{x} \in \mathbb{C}^n$. If $\|\cdot\|_\alpha = \|\cdot\|_\beta$ then we say that $\|\cdot\|$ is subordinate to $\|\cdot\|_\alpha$.

By Lemma 7.22 we have $\|\mathbf{Ax}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$, for all $\mathbf{x} \in \mathbb{C}^n$. Thus the Frobenius norm is subordinate to the Euclidian vector norm.

Exercise 8.8 (The sum norm is subordinate to?)

Show that the sum norm is subordinate to the l_1 -norm.

Exercise 8.9 (The max norm is subordinate to?)

- (a) Show that the max norm is subordinate to the ∞ and 1 norm, i. e., $\|\mathbf{Ax}\|_\infty \leq \|\mathbf{A}\|_M \|\mathbf{x}\|_1$ holds for all $\mathbf{A} \in \mathbb{C}^{m \times n}$ and all $\mathbf{x} \in \mathbb{C}^n$.
- (b) Show that $\|\mathbf{Ae}_l\|_\infty = \|\mathbf{A}\|_M \|\mathbf{e}_l\|_1$, where $\|\mathbf{A}\|_M = |a_{kl}|$ for some k .
- (c) Show that $\|\mathbf{A}\|_M = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|_\infty}{\|\mathbf{x}\|_1}$.

8.1.2 Operator Norms

Corresponding to vector norms on \mathbb{C}^n and \mathbb{C}^m there is an induced matrix norm on $\mathbb{C}^{m \times n}$ which we call the **operator norm**.

Definition 8.10 (Operator Norm)

Suppose $m, n \in \mathbb{N}$ are given and let $\|\cdot\|_\alpha$ be a vector norm on \mathbb{C}^m and $\|\cdot\|_\beta$ a vector norm on \mathbb{C}^n . For $\mathbf{A} \in \mathbb{C}^{m \times n}$ we define

$$\|\mathbf{A}\| := \|\mathbf{A}\|_{\alpha, \beta} := \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta}. \quad (8.3)$$

We call this the (α, β) **operator norm**, the (α, β) -norm, or simply the α -norm if $\alpha = \beta$.

Before we show that the (α, β) -norm is a matrix norm we make some observations.

1. It is enough to take the max over subsets of \mathbb{C}^n . For example

$$\|\mathbf{A}\|_{\alpha, \beta} = \max_{\mathbf{x} \notin \ker(\mathbf{A})} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta} = \max_{\|\mathbf{x}\|_\beta=1} \|\mathbf{Ax}\|_\alpha. \quad (8.4)$$

It is obvious that only \mathbf{x} 's outside the null space $\ker(\mathbf{A})$ of \mathbf{A} need to be considered. It is enough to take the max over the β -norm unit sphere in \mathbb{C}^n since

$$\max_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta} = \max_{\mathbf{x} \neq 0} \left\| \mathbf{A} \left(\frac{\mathbf{x}}{\|\mathbf{x}\|_\beta} \right) \right\|_\alpha = \max_{\|\mathbf{x}\|_\beta=1} \|\mathbf{Ax}\|_\alpha.$$

2. The operator norm $\|\mathbf{A}\|$ is subordinate to the vector norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$. Thus

$$\|\mathbf{Ax}\|_\alpha \leq \|\mathbf{A}\| \|\mathbf{x}\|_\beta \text{ for all } \mathbf{A} \in \mathbb{C}^{m \times n} \text{ and } \mathbf{x} \in \mathbb{C}^n. \quad (8.5)$$

3. We can use max instead of sup in (8.3). This follows by the following compactness argument. The unit sphere $S_\beta = \{\mathbf{x} \in \mathbb{C}^n : \|\mathbf{x}\|_\beta = 1\}$ is bounded. It is also finite dimensional and closed, and hence compact. Moreover, since the vector norm $\|\cdot\|_\alpha$ is a continuous function, it follows that the function $f : S_\beta \rightarrow \mathbb{R}$ given by $f(\mathbf{x}) = \|\mathbf{Ax}\|_\alpha$ is continuous. But then f attains its max and min and we have

$$\|\mathbf{A}\|_{\alpha, \beta} = \|\mathbf{Ax}^*\|_\alpha \text{ for some } \mathbf{x}^* \in \mathbb{C}^n \text{ with } \|\mathbf{x}^*\|_\beta = 1. \quad (8.6)$$

Lemma 8.11 (The operator norm is a matrix norm)

The operator norm given by (8.3) is a matrix norm on $\mathbb{C}^{m \times n}$. The operator norm is consistent if the vector norm $\|\cdot\|_\alpha$ is defined for all $m \in \mathbb{N}$ and $\|\cdot\|_\beta = \|\cdot\|_\alpha$.

Proof. We use (8.4). In 2. and 3. below we take the max over the unit sphere S_β .

1. Nonnegativity is obvious. If $\|\mathbf{A}\| = 0$ then $\|\mathbf{Ay}\|_\beta = 0$ for each $\mathbf{y} \in \mathbb{C}^n$. In particular, each column \mathbf{Ae}_j in \mathbf{A} is zero. Hence $\mathbf{A} = 0$.

2. $\|c\mathbf{A}\| = \max_{\mathbf{x}} \|c\mathbf{A}\mathbf{x}\|_\alpha = \max_{\mathbf{x}} |c| \|\mathbf{A}\mathbf{x}\|_\alpha = |c| \|\mathbf{A}\|.$
3. $\|\mathbf{A} + \mathbf{B}\| = \max_{\mathbf{x}} \|(\mathbf{A} + \mathbf{B})\mathbf{x}\|_\alpha \leq \max_{\mathbf{x}} \|\mathbf{A}\mathbf{x}\|_\alpha + \max_{\mathbf{x}} \|\mathbf{B}\mathbf{x}\|_\alpha = \|\mathbf{A}\| + \|\mathbf{B}\|.$
4.
$$\begin{aligned} \|\mathbf{AB}\| &= \max_{\mathbf{Bx} \neq \mathbf{0}} \frac{\|\mathbf{ABx}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \max_{\mathbf{Bx} \neq \mathbf{0}} \frac{\|\mathbf{ABx}\|_\alpha}{\|\mathbf{Bx}\|_\alpha} \frac{\|\mathbf{Bx}\|_\alpha}{\|\mathbf{x}\|_\alpha} \\ &\leq \max_{\mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{Ay}\|_\alpha}{\|\mathbf{y}\|_\alpha} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Bx}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \|\mathbf{A}\| \|\mathbf{B}\|. \end{aligned}$$

□

For any α -norm of the $n \times n$ identity matrix we find

$$\|\mathbf{I}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ix}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \max_{\mathbf{x} \neq \mathbf{0}} 1 = 1.$$

For the Frobenius norm we find $\|\mathbf{I}\|_F = \sqrt{n}$, and this shows that the Frobenius norm is not an operator norm for $n > 1$.

8.1.3 The Operator p -Norms

Recall that the p or ℓ_p vector norms (10) are given by

$$\|\mathbf{x}\|_p := \left(\sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad p \geq 1, \quad \|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|.$$

The operator norms $\|\cdot\|_p$ defined from these p -vector norms are used quite frequently for $p = 1, 2, \infty$. We define for any $1 \leq p \leq \infty$

$$\|\mathbf{A}\|_p := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p} = \max_{\|\mathbf{y}\|_p=1} \|\mathbf{Ay}\|_p. \quad (8.7)$$

For $p = 1, 2, \infty$ we have explicit expressions for these norms.

Theorem 8.12 (onetwoinfnorms)

For $\mathbf{A} \in \mathbb{C}^{m \times n}$ we have

$$\begin{aligned} \|\mathbf{A}\|_1 &:= \max_{1 \leq j \leq n} \|\mathbf{Ae}_j\|_1 = \max_{1 \leq j \leq n} \sum_{k=1}^m |a_{k,j}|, && (\text{max column sum}) \\ \|\mathbf{A}\|_2 &:= \sigma_1, && (\text{largest singular value of } \mathbf{A}) \\ \|\mathbf{A}\|_\infty &= \max_{1 \leq k \leq m} \|\mathbf{e}_k^T \mathbf{A}\|_1 = \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{k,j}|, && (\text{max row sum}). \end{aligned} \quad (8.8)$$

The **two-norm** $\|\mathbf{A}\|_2$ is also called the **spectral norm** of \mathbf{A} .

Proof. The result for $p = 2$ follows from the minmax theorem for singular values. Indeed, by (7.14) we have $\sigma_1 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}$. For $p = 1, \infty$ we do the following:

(a) We derive a constant K_p such that $\|\mathbf{A}\mathbf{x}\|_p \leq K_p$ for any $\mathbf{x} \in \mathbb{C}^n$ with $\|\mathbf{x}\|_p = 1$.

(b) We give an extremal vector $\mathbf{y}^* \in \mathbb{C}^n$ with $\|\mathbf{y}^*\|_p = 1$ so that $\|\mathbf{A}\mathbf{y}^*\|_p = K_p$. It then follows from (8.7) that $\|\mathbf{A}\|_p = \|\mathbf{A}\mathbf{y}^*\|_p = K_p$.

1-norm: Define K_1 , c and \mathbf{y}^* by $K_1 := \max_{1 \leq j \leq n} \sum_{k=1}^m |a_{kj}| =: \sum_{k=1}^m |a_{kc}|$ and $\mathbf{y}^* := \mathbf{e}_c$, a unit vector. Then $\|\mathbf{y}^*\|_1 = 1$ and we obtain

(a)

$$\|\mathbf{A}\mathbf{x}\|_1 = \sum_{k=1}^m \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \sum_{k=1}^m \sum_{j=1}^n |a_{kj}| |x_j| = \sum_{j=1}^n \left(\sum_{k=1}^m |a_{kj}| \right) |x_j| \leq K_1.$$

(b) $\|\mathbf{A}\mathbf{y}^*\|_1 = K_1$.

∞ -norm: Define K_∞ , r and \mathbf{y}^* by $K_\infty := \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{kj}| =: \sum_{j=1}^n |a_{rj}|$ and $\mathbf{y}^* := [e^{-i\theta_1}, \dots, e^{-i\theta_n}]^T$, where $a_{rj} = |a_{rj}| e^{i\theta_j}$ for $j = 1, \dots, n$.

(a) $\|\mathbf{A}\mathbf{x}\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{kj}| |x_j| \leq K_\infty$.

(b) $\|\mathbf{A}\mathbf{y}^*\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} e^{-i\theta_j} \right| = K_\infty$.

The last equality is correct because $\left| \sum_{j=1}^n a_{kj} e^{-i\theta_j} \right| \leq \sum_{j=1}^n |a_{kj}| \leq K_\infty$ with equality for $k = r$.

□

Example 8.13 (Compare onetwoinfnorms)

In Example 7.8 we found that the largest singular value of the matrix $\mathbf{A} := \frac{1}{15} \begin{bmatrix} 14 & 4 & 16 \\ 2 & 22 & 13 \end{bmatrix}$, is $\sigma_1 = 2$. We find

$$\|\mathbf{A}\|_1 = \frac{29}{15}, \quad \|\mathbf{A}\|_2 = 2, \quad \|\mathbf{A}\|_\infty = \frac{37}{15}, \quad \|\mathbf{A}\|_F = \sqrt{5}.$$

We observe that the values of these norms do not differ by much.

In some cases the spectral norm is equal to an eigenvalue of the matrix.

Theorem 8.14 (Spectral norm)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ has singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ and eigenvalues $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$. Then

$$\|\mathbf{A}\|_2 = \sigma_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\sigma_n}, \tag{8.9}$$

$$\|\mathbf{A}\|_2 = \lambda_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\lambda_n}, \quad \text{if } \mathbf{A} \text{ is symmetric positive definite,} \tag{8.10}$$

$$\|\mathbf{A}\|_2 = |\lambda_1| \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{|\lambda_n|}, \quad \text{if } \mathbf{A} \text{ is normal.} \tag{8.11}$$

For the norms of \mathbf{A}^{-1} we assume of course that \mathbf{A} is nonsingular.

Proof. Since $1/\sigma_n$ is the largest singular value of \mathbf{A}^{-1} , (8.9) follows. As shown in Exercise 7.13 the singular values of a symmetric positive definite matrix (normal matrix) are equal to the eigenvalues (absolute value of the eigenvalues). This implies (8.10) and (8.11). \square

Exercise 8.15 (Spectral norm of the inverse)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular. Use (8.9) and (7.14) to show that

$$\|\mathbf{A}^{-1}\|_2 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{x}\|_2}{\|\mathbf{Ax}\|_2}.$$

Exercise 8.16 (p -norm example)

Let

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}.$$

Compute $\|\mathbf{A}\|_p$ and $\|\mathbf{A}^{-1}\|_p$ for $p = 1, 2, \infty$.

The following result is sometimes useful.

Theorem 8.17 (Spectral norm bound)

For any $\mathbf{A} \in \mathbb{C}^{m \times n}$ we have $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty$.

Proof. Let $(\sigma_1^2, \mathbf{v}_1)$ be an eigenpair for $\mathbf{A}^* \mathbf{A}$ corresponding to the largest singular value of \mathbf{A} . Then

$$\|\mathbf{A}\|_2^2 \|\mathbf{v}_1\|_1 = \sigma_1^2 \|\mathbf{v}_1\|_1 = \|\sigma_1^2 \mathbf{v}_1\|_1 = \|\mathbf{A}^* \mathbf{A} \mathbf{v}_1\|_1 \leq \|\mathbf{A}^*\|_1 \|\mathbf{A}\|_1 \|\mathbf{v}_1\|_1.$$

Observing that $\|\mathbf{A}^*\|_1 = \|\mathbf{A}\|_\infty$ by Theorem 8.12 and canceling $\|\mathbf{v}_1\|_1$ proves the result. \square

8.1.4 Unitary Invariant Matrix Norms

Definition 8.18 (Unitary invariant norm)

A matrix norm $\|\cdot\|$ on $\mathbb{C}^{m \times n}$ is called **unitary invariant** if $\|\mathbf{U} \mathbf{A} \mathbf{V}\| = \|\mathbf{A}\|$ for any $\mathbf{A} \in \mathbb{C}^{m \times n}$ and any unitary matrices $\mathbf{U} \in \mathbb{C}^{m \times m}$ and $\mathbf{V} \in \mathbb{C}^{n \times n}$.

When an unitary invariant matrix norm is used, the size of a perturbation is not increased by a unitary transformation. Thus if \mathbf{U} and \mathbf{V} are unitary then $\mathbf{U}(\mathbf{A} + \mathbf{E})\mathbf{V} = \mathbf{U}\mathbf{A}\mathbf{V} + \mathbf{F}$, where $\|\mathbf{F}\| = \|\mathbf{E}\|$.

It follows from Lemma 7.22 that the Frobenius norm is unitary invariant. We show here that this also holds for the spectral norm. It can be shown that the spectral norm is the only unitary invariant operator norm, see [12] p. 308.

Theorem 8.19 (Unitary invariant norms)

The Frobenius norm and the spectral norm are unitary invariant. Moreover $\|\mathbf{A}^\|_F = \|\mathbf{A}\|_F$ and $\|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2$.*

Proof. The results for the Frobenius norm follow from Lemma 7.22. Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ and let $\mathbf{U} \in \mathbb{C}^{m \times m}$ and $\mathbf{V} \in \mathbb{C}^{n \times n}$ be unitary. Since the 2-vector norm is unitary invariant we obtain

$$\|\mathbf{U}\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{U}\mathbf{A}\mathbf{x}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{A}\|_2.$$

Now \mathbf{A} and \mathbf{A}^* have the same nonzero singular values, and it follows from Theorem 8.12 that $\|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2$. Moreover \mathbf{V}^* is unitary. Using these facts we find

$$\|\mathbf{A}\mathbf{V}\|_2 = \|(\mathbf{A}\mathbf{V})^*\|_2 = \|\mathbf{V}^*\mathbf{A}^*\|_2 = \|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2.$$

□

Exercise 8.20 (Univariance of spectral norm)

Show that $\|\mathbf{V}\mathbf{A}\|_2 = \|\mathbf{A}\|_2$ holds even for a rectangular \mathbf{V} as long as $\mathbf{V}^*\mathbf{V} = \mathbf{I}$.

Exercise 8.21 ($\|\mathbf{AU}\|_2$ rectangular \mathbf{A})

Find $\mathbf{A} \in \mathbb{R}^{2 \times 2}$ and $\mathbf{U} \in \mathbb{R}^{2 \times 1}$ with $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ such that $\|\mathbf{AU}\|_2 < \|\mathbf{A}\|_2$. Thus, in general, $\|\mathbf{AU}\|_2 = \|\mathbf{A}\|_2$ does not hold for a rectangular \mathbf{U} even if $\mathbf{U}^*\mathbf{U} = \mathbf{I}$.

Exercise 8.22 (p -norm of diagonal matrix)

Show that $\|\mathbf{A}\|_p = \rho(\mathbf{A}) := \max |\lambda_i|$ (the largest eigenvalue of \mathbf{A}), $1 \leq p \leq \infty$, when \mathbf{A} is a diagonal matrix.

Exercise 8.23 (spectral norm of a column vector)

A vector $\mathbf{a} \in \mathbb{C}^m$ can also be considered as a column vector $\mathbf{A} \in \mathbb{C}^{m,1}$.

(a) Show that the spectral matrix norm (2-norm) of \mathbf{A} equals the Euclidean vector norm of \mathbf{a} .

(b) Show that $\|\mathbf{A}\|_p = \|\mathbf{a}\|_p$ for $1 \leq p \leq \infty$.

Exercise 8.24 (Norm of absolute value matrix)

If $\mathbf{A} \in \mathbb{C}^{m \times n}$ has elements a_{ij} , let $|\mathbf{A}| \in \mathbb{R}^{m \times n}$ be the matrix with elements $|a_{ij}|$.

- (a) Compute $|\mathbf{A}|$ if $\mathbf{A} = \begin{bmatrix} 1+i & -2 \\ 1 & 1-i \end{bmatrix}$, $i = \sqrt{-1}$.
- (b) Show that for any $\mathbf{A} \in \mathbb{C}^{m \times n}$ $\|\mathbf{A}\|_F = \||\mathbf{A}|\|_F$, $\|\mathbf{A}\|_p = \||\mathbf{A}|\|_p$ for $p = 1, \infty$.
- (c) Show that for any $\mathbf{A} \in \mathbb{C}^{m \times n}$ $\|\mathbf{A}\|_2 \leq \||\mathbf{A}|\|_2$.
- (d) Find a real symmetric 2×2 matrix \mathbf{A} such that $\|\mathbf{A}\|_2 < \||\mathbf{A}|\|_2$.

Exercise 8.25 (Spectral norm)

Let $m, n \in \mathbb{N}$ and $\mathbf{A} \in \mathbb{C}^{m \times n}$. Show that

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1} |\mathbf{y}^* \mathbf{A} \mathbf{x}|.$$

8.1.5 Absolute and Monotone Norms

A vector norm on \mathbb{C}^n is called an **absolute norm** if $\|\mathbf{x}\| = \|\mathbf{|x|}\|$ for all $\mathbf{x} \in \mathbb{C}^n$. Here $\mathbf{|x|} := [|x_1|, \dots, |x_n|]^T$, the absolute values of the components of \mathbf{x} . Clearly the vector p norms are absolute norms. We state without proof (see Theorem 5.5.10 of [12]) that a vector norm on \mathbb{C}^n is an absolute norm if and only if it is a **monotone norm**, i.e.,

$$|x_i| \leq |y_i|, \quad i = 1, \dots, n \implies \|\mathbf{x}\| \leq \|\mathbf{y}\|, \text{ for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

Absolute and monotone matrix norms are defined as for vector norms.

Exercise 8.26 (Absolute norms)

Show that the Frobenius norm and the $1, \infty$ operator norms are absolute norms.

Exercise 8.27 (Is the spectral norm an absolute norm?)

Show that the spectral norm is not an absolute norm.

The study of matrix norms will be continued in Chapter 9.

8.2 The Condition Number with Respect to Inversion

Consider the system of two linear equations

$$\begin{array}{rcl} x_1 & + & x_2 = 20 \\ x_1 & + & (1 - 10^{-16})x_2 = 20 - 10^{-15} \end{array}$$

whose exact solution is $x_1 = x_2 = 10$. If we replace the second equation by

$$x_1 + (1 + 10^{-16})x_2 = 20 - 10^{-15},$$

the exact solution changes to $x_1 = 30$, $x_2 = -10$. Here a small change in one of the coefficients, from $1 - 10^{-16}$ to $1 + 10^{-16}$, changed the exact solution by a large amount.

A mathematical problem in which the solution is very sensitive to changes in the data is called **ill-conditioned**. Such problems are difficult to solve on a computer.

In this section we consider what effect a small change (perturbation) in the data \mathbf{A}, \mathbf{b} has on the solution \mathbf{x} of a linear system $\mathbf{Ax} = \mathbf{b}$. Suppose \mathbf{y} solves $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b} + \mathbf{e}$ where \mathbf{E} is a (small) $n \times n$ matrix and \mathbf{e} a (small) vector. How large can $\mathbf{y} - \mathbf{x}$ be? To measure this we use vector and matrix norms. In this section $\|\cdot\|$ will denote a vector norm on \mathbb{C}^n and also a submultiplicative matrix norm on $\mathbb{C}^{n \times n}$ which in addition is subordinate to the vector norm. Thus for any $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ and any $\mathbf{x} \in \mathbb{C}^n$ we have

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \text{ and } \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|.$$

This is satisfied if the matrix norm is the operator norm corresponding to the given vector norm, but is also satisfied for the Frobenius matrix norm and the Euclidian vector norm. This follows from Lemma 7.22.

Suppose \mathbf{x} and \mathbf{y} are vectors in \mathbb{C}^n that we want to compare. The difference $\|\mathbf{y} - \mathbf{x}\|$ measures the **absolute error** in \mathbf{y} as an approximation to \mathbf{x} , while $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$ and $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$ are measures for the **relative error**.

We consider first a perturbation in the right-hand side \mathbf{b} .

Theorem 8.28 (Perturbation in the right-hand side)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular, $\mathbf{b}, \mathbf{e} \in \mathbb{C}^n$, $\mathbf{b} \neq \mathbf{0}$ and $\mathbf{Ax} = \mathbf{b}$, $\mathbf{Ay} = \mathbf{b} + \mathbf{e}$. Then

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (8.12)$$

Proof. Subtracting $\mathbf{Ax} = \mathbf{b}$ from $\mathbf{Ay} = \mathbf{b} + \mathbf{e}$ we have $\mathbf{A}(\mathbf{y} - \mathbf{x}) = \mathbf{e}$ or $\mathbf{y} - \mathbf{x} = \mathbf{A}^{-1}\mathbf{e}$. Combining $\|\mathbf{y} - \mathbf{x}\| = \|\mathbf{A}^{-1}\mathbf{e}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{e}\|$ and $\|\mathbf{b}\| = \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$ we obtain the upper bound in (8.12). Combining $\|\mathbf{e}\| \leq \|\mathbf{A}\| \|\mathbf{y} - \mathbf{x}\|$ and $\|\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{b}\|$ we obtain the lower bound. \square

Consider (8.12). $\|\mathbf{e}\|/\|\mathbf{b}\|$ is a measure of the size of the perturbation \mathbf{e} relative to the size of \mathbf{b} . The upper bound says that $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$ in the worst case can be

$$K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$$

times as large as $\|e\|/\|b\|$. $K(\mathbf{A})$ is called the **condition number with respect to inversion of a matrix**, or just the condition number, if it is clear from the context that we are talking about solving linear systems or inverting a matrix. The condition number depends on the matrix \mathbf{A} and on the norm used. If $K(\mathbf{A})$ is large, \mathbf{A} is called **ill-conditioned** (with respect to inversion). If $K(\mathbf{A})$ is small, \mathbf{A} is called **well-conditioned** (with respect to inversion). We always have $K(\mathbf{A}) \geq 1$. For since $\|\mathbf{x}\| = \|\mathbf{I}\mathbf{x}\| \leq \|\mathbf{I}\|\|\mathbf{x}\|$ for any \mathbf{x} , by subordinance we have $\|\mathbf{I}\| \geq 1$ and therefore by submultiplicativity $\|\mathbf{A}\|\|\mathbf{A}^{-1}\| \geq \|\mathbf{A}\mathbf{A}^{-1}\| = \|\mathbf{I}\| \geq 1$.

Since all matrix norms are equivalent, the dependence of $K(\mathbf{A})$ on the norm chosen is less important than the dependence on \mathbf{A} . Sometimes one chooses the spectral norm when discussing properties of the condition number, and the ℓ_1 , ℓ_∞ , or Frobenius norm when one wishes to compute it or estimate it.

Explicit expressions for the 2-norm condition number follow from Theorem 8.14.

Theorem 8.29 (Spectral condition number)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$ and eigenvalues $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| > 0$. Then $K_2(\mathbf{A}) := \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \sigma_1/\sigma_n$. Moreover,

$$K_2(\mathbf{A}) = \begin{cases} \lambda_1/\lambda_n, & \text{if } \mathbf{A} \text{ is symmetric positive definite,} \\ |\lambda_1|/|\lambda_n|, & \text{if } \mathbf{A} \text{ is normal.} \end{cases} \quad (8.13)$$

It follows that \mathbf{A} is ill-conditioned with respect to inversion if and only if σ_1/σ_n is large, or λ_1/λ_n is large when \mathbf{A} is symmetric positive definite.

Consider next the effect of a perturbation in the coefficient matrix. Suppose $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$ with \mathbf{A} nonsingular. We like to compare the solution \mathbf{x} and \mathbf{y} of the systems $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$. We expect $\mathbf{A} + \mathbf{E}$ to be nonsingular if the elements of \mathbf{E} are sufficiently small and we need to address this question. Consider first the case where $\mathbf{A} = \mathbf{I}$.

Theorem 8.30 (Nonsingularity of perturbation of identity)

Suppose $\mathbf{B} \in \mathbb{C}^{n \times n}$ and $\|\mathbf{B}\| < 1$ for some operator norm $\|\cdot\|$. Then $\mathbf{I} - \mathbf{B}$ is nonsingular and

$$\frac{1}{1 + \|\mathbf{B}\|} \leq \|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}. \quad (8.14)$$

Proof. Suppose $(\mathbf{I} - \mathbf{B})\mathbf{x} = \mathbf{0}$ for some nonzero $\mathbf{x} \in \mathbb{C}^n$. Then $\mathbf{x} = \mathbf{B}\mathbf{x}$ so that $\|\mathbf{x}\| = \|\mathbf{B}\mathbf{x}\| \leq \|\mathbf{B}\|\|\mathbf{x}\|$. But then $\|\mathbf{B}\| \geq 1$, a contradiction. It follows that $\mathbf{I} - \mathbf{B}$ is nonsingular. Next, since

$$\begin{aligned} \|\mathbf{I}\| &= \|(\mathbf{I} - \mathbf{B})(\mathbf{I} - \mathbf{B})^{-1}\| \leq \|\mathbf{I} - \mathbf{B}\| \|(\mathbf{I} - \mathbf{B})^{-1}\| \\ &\leq (\|\mathbf{I}\| + \|\mathbf{B}\|) \|(\mathbf{I} - \mathbf{B})^{-1}\| \end{aligned}$$

we obtain

$$\frac{\|\mathbf{I}\|}{\|\mathbf{I}\| + \|\mathbf{B}\|} \leq \|(\mathbf{I} - \mathbf{B})^{-1}\|. \quad (8.15)$$

Since $\|\mathbf{I}\| = 1$ for an operator norm the lower bound in (8.15) is equal to the lower bound in (8.14).

Taking norms and using the inverse triangle inequality in

$$\mathbf{I} = (\mathbf{I} - \mathbf{B})(\mathbf{I} - \mathbf{B})^{-1} = (\mathbf{I} - \mathbf{B})^{-1} - \mathbf{B}(\mathbf{I} - \mathbf{B})^{-1}$$

implies

$$1 = \|\mathbf{I}\| \geq \|(\mathbf{I} - \mathbf{B})^{-1}\| - \|\mathbf{B}(\mathbf{I} - \mathbf{B})^{-1}\| \geq (1 - \|\mathbf{B}\|) \|(\mathbf{I} - \mathbf{B})^{-1}\|$$

and the upper bound follows. \square

We show in Section 9.2 that Theorem 8.30 holds for the Frobenius norm, and more generally for any consistent matrix norm on $\mathbb{C}^{n \times n}$. In particular, since $\|\mathbf{I}\| \geq 1$ for any consistent matrix norm on $\mathbb{C}^{n \times n}$ the lower bound in (8.14) follows from the lower bound in (8.15).

Theorem 8.31 (Nonsingularity of perturbation)

Suppose $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$, $\mathbf{b} \in \mathbb{C}^n$ with \mathbf{A} invertible and $\mathbf{b} \neq \mathbf{0}$. If $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1/2$ for some matrix norm consistent on $\mathbb{C}^{n \times n}$ then $\mathbf{A} + \mathbf{E}$ is invertible. If $\mathbf{Ax} = \mathbf{b}$ and $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$ then

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{y}\|} \leq \|\mathbf{A}^{-1}\mathbf{E}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}, \quad (8.16)$$

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (8.17)$$

Proof. Since $r < 1$ Theorem 8.30 implies that the matrix $\mathbf{I} - \mathbf{B} := \mathbf{I} + \mathbf{A}^{-1}\mathbf{E}$ is nonsingular and then $\mathbf{A} + \mathbf{E} = \mathbf{A}(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})$ is nonsingular. Subtracting $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$ from $\mathbf{Ax} = \mathbf{b}$ gives $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{E}\mathbf{y}$ or $\mathbf{x} - \mathbf{y} = \mathbf{A}^{-1}\mathbf{E}\mathbf{y}$. Taking norms and dividing by $\|\mathbf{y}\|$ proves (8.16). Solving $\mathbf{x} - \mathbf{y} = \mathbf{A}^{-1}\mathbf{E}\mathbf{y}$ for \mathbf{y} we obtain $\mathbf{y} = (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\mathbf{x}$. By (8.14)

$$\|\mathbf{y}\| \leq \|(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\| \|\mathbf{x}\| \leq \frac{\|\mathbf{x}\|}{1 - \|\mathbf{A}^{-1}\mathbf{E}\|} \leq 2\|\mathbf{x}\|.$$

But then (8.17) follows from (8.16). \square

In Theorem 8.31 we gave a bound for the relative error in \mathbf{x} as an approximation to \mathbf{y} , (8.16), and the relative error in \mathbf{y} as an approximation to \mathbf{x} , (8.17).

$\|\mathbf{E}\|/\|\mathbf{A}\|$ is a measure for the size of the perturbation \mathbf{E} in \mathbf{A} relative to the size of \mathbf{A} . The condition number again plays a crucial role. $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$ can be as large as $K(\mathbf{A})$ times $\|\mathbf{E}\|/\|\mathbf{A}\|$. It can be shown that the upper bound can be attained for any \mathbf{A} and any \mathbf{b} . In deriving the upper bound we used the inequality $\|\mathbf{A}^{-1}\mathbf{E}\mathbf{y}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{E}\| \|\mathbf{y}\|$. For a more or less random perturbation \mathbf{E} this is not a severe overestimate for $\|\mathbf{A}^{-1}\mathbf{E}\mathbf{y}\|$. In the situation where \mathbf{E} is due to round-off errors (8.16) can give a fairly realistic estimate for $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$.

Suppose we have computed an approximate solution \mathbf{y} to $\mathbf{Ax} = \mathbf{b}$. The vector $\mathbf{r}(\mathbf{y}) := \mathbf{Ay} - \mathbf{b}$ is called the **residual vector**, or just the residual. We can bound $\mathbf{x} - \mathbf{y}$ in term of \mathbf{r} .

Theorem 8.32 (Perturbation and residual)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mathbf{b} \in \mathbb{C}^n$, \mathbf{A} is nonsingular and $\mathbf{b} \neq \mathbf{0}$. Let $\mathbf{r}(\mathbf{y}) = \mathbf{Ay} - \mathbf{b}$ for each $\mathbf{y} \in \mathbb{C}^n$. If $\mathbf{Ax} = \mathbf{b}$ then

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|}. \quad (8.18)$$

Proof. We simply take $\mathbf{e} = \mathbf{r}(\mathbf{y})$ in Theorem 8.28. \square

If \mathbf{A} is well-conditioned, (8.18) says that $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\| \approx \|\mathbf{r}(\mathbf{y})\|/\|\mathbf{b}\|$. In other words, the accuracy in \mathbf{y} is about the same order of magnitude as the residual as long as $\|\mathbf{b}\| \approx 1$. If \mathbf{A} is ill-conditioned, anything can happen. We can for example have an accurate solution even if the residual is large.

We end this section with a perturbation result for the inverse matrix. Again the condition number plays a crucial role.

Theorem 8.33 (Perturbation of inverse matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular and let $\|\cdot\|$ be a consistent matrix norm on $\mathbb{C}^{n \times n}$. If $\mathbf{E} \in \mathbb{C}^{n \times n}$ is so small that $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1$ then $\mathbf{A} + \mathbf{E}$ is nonsingular and

$$\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - r}. \quad (8.19)$$

If $r < 1/2$ then

$$\frac{\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\|}{\|\mathbf{A}^{-1}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (8.20)$$

Proof. We showed in Theorem 8.31 that $\mathbf{A} + \mathbf{E}$ is nonsingular and since $(\mathbf{A} + \mathbf{E})^{-1} = (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\mathbf{A}^{-1}$ we obtain $\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \|(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\| \|\mathbf{A}^{-1}\|$ and (8.19) follows from (8.14) and the remark after the Theorem. From the identity

$$(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1} = -\mathbf{A}^{-1}\mathbf{E}(\mathbf{A} + \mathbf{E})^{-1}$$

we obtain by (8.19)

$$\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{E}\| \|(\mathbf{A} + \mathbf{E})^{-1}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\| \|\mathbf{A}^{-1}\|}{\|\mathbf{A}\|} \frac{\|\mathbf{A}^{-1}\|}{1 - r}.$$

Dividing by $\|\mathbf{A}^{-1}\|$ and setting $r = 1/2$ proves (8.20). \square

Exercise 8.34 (Sharpness of perturbation bounds)

The upper and lower bounds for $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$ given by (8.12) can be attained for any matrix \mathbf{A} , but only for special choices of \mathbf{b} . Suppose $\mathbf{y}_\mathbf{A}$ and $\mathbf{y}_{\mathbf{A}^{-1}}$ are vectors with $\|\mathbf{y}_\mathbf{A}\| = \|\mathbf{y}_{\mathbf{A}^{-1}}\| = 1$ and $\|\mathbf{A}\| = \|\mathbf{A}\mathbf{y}_\mathbf{A}\|$ and $\|\mathbf{A}^{-1}\| = \|\mathbf{A}^{-1}\mathbf{y}_{\mathbf{A}^{-1}}\|$.

(a) Show that the upper bound in (8.12) is attained if $\mathbf{b} = \mathbf{A}\mathbf{y}_\mathbf{A}$ and $\mathbf{e} = \mathbf{y}_{\mathbf{A}^{-1}}$.

(b) Show that the lower bound is attained if $\mathbf{b} = \mathbf{y}_{\mathbf{A}^{-1}}$ and $\mathbf{e} = \mathbf{A}\mathbf{y}_\mathbf{A}$.

Exercise 8.35 (Condition number of 2. derivative matrix)

In this exercise we will show that for $m \geq 1$

$$\frac{4}{\pi^2}(m+1)^2 - 2/3 < \text{cond}_p(\mathbf{T}) \leq \frac{1}{2}(m+1)^2, \quad p = 1, 2, \infty, \quad (8.21)$$

where $\mathbf{T} := \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ and $\text{cond}_p(\mathbf{T}) := \|\mathbf{T}\|_p \|\mathbf{T}^{-1}\|_p$ is the p -norm condition number of \mathbf{T} . The p matrix norm is given by (8.7). You will need the explicit inverse of \mathbf{T} given by (2.10) and the eigenvalues given in Lemma 4.11. As usual we define $h := 1/(m+1)$.

a) Show that for $m \geq 3$

$$\text{cond}_1(\mathbf{T}) = \text{cond}_\infty(\mathbf{T}) = \frac{1}{2} \begin{cases} h^{-2}, & m \text{ odd}, \\ h^{-2} - 1, & m \text{ even}. \end{cases} \quad (8.22)$$

and that $\text{cond}_1(\mathbf{T}) = \text{cond}_\infty(\mathbf{T}) = 3$ for $m = 2$.

b) Show that for $p = 2$ and $m \geq 1$ we have

$$\text{cond}_2(\mathbf{T}) = \cot^2\left(\frac{\pi h}{2}\right) = 1/\tan^2\left(\frac{\pi h}{2}\right).$$

c) Show the bounds

$$\frac{4}{\pi^2}h^{-2} - \frac{2}{3} < \text{cond}_2(\mathbf{T}) < \frac{4}{\pi^2}h^{-2}. \quad (8.23)$$

Hint: For the upper bound use the inequality $\tan x > x$ valid for $0 < x < \pi/2$. For the lower bound we use the inequality $\cot^2 x > \frac{1}{x^2} - \frac{2}{3}$ for $x > 0$. This can be derived for $0 < x < \pi$ by first showing that the second derivative of $\cot^2 x$ is positive and then use Taylor's theorem.

d) Show (8.21).

8.3 Proof that the p -Norms are Norms

We want to show

Theorem 8.36 (The p norms are norms)

Let for $1 \leq p \leq \infty$ and $\mathbf{x} \in \mathbb{C}^n$

$$\|\mathbf{x}\|_p := \left(\sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad \|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|.$$

Then for all $1 \leq p \leq \infty$, $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ and all $a \in \mathbb{C}$

1. $\|\mathbf{x}\|_p \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$. (positivity)
2. $\|a\mathbf{x}\|_p = |a| \|\mathbf{x}\|_p$. (homogeneity)
3. $\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$. (subadditivity)

Positivity and homogeneity follows immediately. To show the subadditivity we need some elementary properties of convex functions.

Definition 8.37 (Convex function)

Let $I \subset \mathbb{R}$ be an interval. A function $f : I \rightarrow \mathbb{R}$ is called convex if

$$f((1 - \lambda)x_1 + \lambda x_2) \leq (1 - \lambda)f(x_1) + \lambda f(x_2) \quad (8.24)$$

for all $x_1 < x_2 \in I$ and all $\lambda \in [0, 1]$. The sum $\sum_{j=1}^n \lambda_j x_j$ is called a **convex combination** of x_1, \dots, x_n if $\lambda_j \geq 0$ for $j = 1, \dots, n$ and $\sum_{j=1}^n \lambda_j = 1$.

The convexity condition is illustrated in Figure 8.1.

Lemma 8.38 (A sufficient condition for convexity)

If $f \in C^2[a, b]$ and $f''(x) \geq 0$ for $x \in [a, b]$ then f is convex.

Proof. We recall the formula for linear interpolation with remainder, (cf a book on numerical methods) For any $a \leq x_1 \leq x \leq x_2 \leq b$ there is a $c \in [x_1, x_2]$ such that

$$f(x) = \frac{x_2 - x}{x_2 - x_1} f(x_1) + \frac{x - x_1}{x_2 - x_1} f(x_2) + (x - x_1)(x - x_2) f''(c)/2.$$

Thus $f(x) = (1 - \lambda)f(x_1) + \lambda f(x_2) + (x_2 - x_1)^2 \lambda(\lambda - 1) f''(c)/2$, where $\lambda := \frac{x - x_1}{x_2 - x_1}$. Since $\lambda \in [0, 1]$ the remainder term is negative. Moreover,

$$x = \frac{x_2 - x}{x_2 - x_1} x_1 + \frac{x - x_1}{x_2 - x_1} x_2 = (1 - \lambda)x_1 + \lambda x_2$$

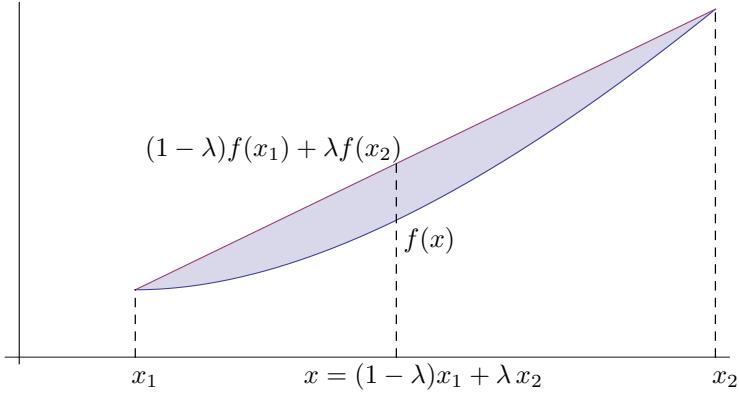


Figure 8.1: A convex function.

so that (8.24) holds, and f is convex. \square

The following inequality is elementary, but can be used to prove many non-trivial inequalities.

Theorem 8.39 (Jensen's Inequality)

Suppose $I \in \mathbb{R}$ is an interval and $f : I \rightarrow \mathbb{R}$ is convex. Then for all $n \in \mathbb{N}$, all $\lambda_1, \dots, \lambda_n$ with $\lambda_j \geq 0$ for $j = 1, \dots, n$ and $\sum_{j=1}^n \lambda_j = 1$, and all $z_1, \dots, z_n \in I$

$$f\left(\sum_{j=1}^n \lambda_j\right) \leq \sum_{j=1}^n \lambda_j f(z_j).$$

Proof. We use induction on n . The result is trivial for $n = 1$. Let $n \geq 2$, assume the inequality holds for $k = n - 1$, and let λ_j, z_j for $j = 1, \dots, n$ be given as in the theorem. Since $n \geq 2$ we have $\lambda_i < 1$ for at least one i so assume without loss of generality that $\lambda_1 < 1$. Define u by $u := \sum_{j=2}^n \frac{\lambda_j}{1-\lambda_1} z_j$. Since $\sum_{j=2}^n \lambda_j = 1 - \lambda_1$ this is a convex combination of k terms and the induction hypothesis implies that $f(u) \leq \sum_{j=2}^n \frac{\lambda_j}{1-\lambda_1} f(z_j)$. But then by the convexity of f

$$f\left(\sum_{j=1}^n \lambda_j\right) = f(\lambda_1 z_1 + (1 - \lambda_1) u) \leq \lambda_1 f(z_1) + (1 - \lambda_1) f(u) \leq \sum_{j=1}^n \lambda_j f(z_j)$$

and the inequality holds for $k + 1 = n$. \square

Corollary 8.40 (Weighted geometric/arithmetic mean inequality)

Suppose $\sum_{j=1}^n \lambda_j a_j$ is a convex combination of nonnegative numbers a_1, \dots, a_n . Then

$$a_1^{\lambda_1} a_2^{\lambda_2} \cdots a_n^{\lambda_n} \leq \sum_{j=1}^n \lambda_j a_j, \quad (8.25)$$

where $0^0 := 0$.

Proof. The result is trivial if one or more of the a_j 's are zero so assume $a_j > 0$ for all j . Consider the function $f : (0, \infty)$ given by $f(x) = -\log x$. Since $f'(x) = -1/x < 0$ and $f''(x) = 1/x^2 > 0$ for $x \in (0, \infty)$, this function is monotone and convex. By Jensen's inequality

$$-\log \left(\sum_{j=1}^n \lambda_j a_j \right) \leq -\sum_{j=1}^n \lambda_j \log(a_j) = -\log \left(a_1^{\lambda_1} \cdots a_n^{\lambda_n} \right)$$

and the inequality follows by monotonicity of f . \square

Taking $\lambda_j = \frac{1}{n}$ for all j in (8.25) we obtain the classical **geometric/arithmetic mean inequality**

$$(a_1 a_2 \cdots a_n)^{\frac{1}{n}} \leq \frac{1}{n} \sum_{j=1}^n a_j. \quad (8.26)$$

Corollary 8.41 (Hölder's inequality)

For $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ and $1 \leq p \leq \infty$

$$\sum_{j=1}^n |x_j y_j| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \text{ where } \frac{1}{p} + \frac{1}{q} = 1.$$

Proof. We leave the proof for $p = 1$ and $p = \infty$ as an exercise so assume $1 < p < \infty$. For any $a, b \geq 0$ the weighted arithmetic/geometric mean inequality implies that

$$a^{\frac{1}{p}} b^{\frac{1}{q}} \leq \frac{1}{p} a + \frac{1}{q} b, \text{ where } \frac{1}{p} + \frac{1}{q} = 1. \quad (8.27)$$

If $\mathbf{x} = \mathbf{0}$ or $\mathbf{y} = \mathbf{0}$ there is nothing to prove so assume that both \mathbf{x} and \mathbf{y} are nonzero. Using 8.27 on each term we obtain

$$\frac{1}{\|\mathbf{x}\|_p \|\mathbf{y}\|_q} \sum_{j=1}^n |x_j y_j| = \sum_{j=1}^n \left(\frac{|x_j|^p}{\|\mathbf{x}\|_p^p} \right)^{\frac{1}{p}} \left(\frac{|y_j|^q}{\|\mathbf{y}\|_q^q} \right)^{\frac{1}{q}} \leq \sum_{j=1}^n \left(\frac{1}{p} \frac{|x_j|^p}{\|\mathbf{x}\|_p^p} + \frac{1}{q} \frac{|y_j|^q}{\|\mathbf{y}\|_q^q} \right) = 1$$

and the proof of the inequality is complete. \square

Corollary 8.42 (Minkowski's inequality)

For $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ and $1 \leq p \leq \infty$

$$\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p.$$

Proof. We leave the proof for $p = 1$ and $p = \infty$ as an exercise so assume $1 < p < \infty$. We write

$$\|\mathbf{x} + \mathbf{y}\|_p^p = \sum_{j=1}^n |x_j + y_j|^p \leq \sum_{j=1}^n |x_j| |x_j + y_j|^{p-1} + \sum_{j=1}^n |y_j| |x_j + y_j|^{p-1}.$$

We apply Hölder's inequality with exponent p and q to each sum. In view of the relation $(p-1)q = p$ the result is

$$\|\mathbf{x} + \mathbf{y}\|_p^p \leq \|\mathbf{x}\|_p \|\mathbf{x} + \mathbf{y}\|_p^{p/q} + \|\mathbf{y}\|_p \|\mathbf{x} + \mathbf{y}\|_p^{p/q}.$$

Since $p - \frac{p}{q} = 1$ the inequality follows. \square

Exercise 8.43 (p norm for $p = 1$ and $p = \infty$)

Show that $\|\cdot\|_p$ is a vector norm in \mathbb{R}^n for $p = 1, p = \infty$.

Exercise 8.44 (The p - norm unit sphere)

The set

$$S_p = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_p = 1\}$$

is called the unit sphere in \mathbb{R}^n with respect to p . Draw S_p for $p = 1, 2, \infty$ for $n = 2$.

Exercise 8.45 (Sharpness of p -norm inequaltiy)

Let $1 \leq p$. Produce a vector \mathbf{x}_l such that $\|\mathbf{x}_l\|_\infty = \|\mathbf{x}_l\|_p$ and another vector \mathbf{x}_u such that $\|\mathbf{x}_u\|_p = n^{1/p} \|\mathbf{x}_u\|_p \infty$. Thus the inequalities in (13) are sharp.

Exercise 8.46 (p -norm inequaltiies for arbitrary p)

If $1 \leq q \leq p \leq \infty$ then

$$\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_q \leq n^{1/q-1/p} \|\mathbf{x}\|_p, \quad \mathbf{x} \in \mathbb{C}^n.$$

Hint: For the rightmost inequality use Jensen's inequality Cf. Theorem 8.39 with $f(z) = z^{p/q}$ and $z_i = |x_i|^q$. For the left inequality consider first $y_i = x_i/\|\mathbf{x}\|_\infty$, $i = 1, 2, \dots, n$.

8.4 Review Questions

- 8.4.1** • What is a consistent matrix norm?
• what is a subordinate matrix norm?
• is an operator norm consistent?
• why is the Frobenius norm not an operator norm?
• what is the spectral norm of a matrix?
• how do we compute $\|A\|_\infty$?
• what is the spectral condition number of a symmetric positive definite matrix?
- 8.4.2** Why is $\|A\|_2 \leq \|A\|_F$ for any matrix A ?
- 8.4.3** What is the spectral norm of the inverse of a normal matrix?

Part III

Iterative Methods for Large Linear Systems

Chapter 9

The Classical Iterative Methods

Gaussian elimination and Cholesky factorization are **direct methods**. In absence of rounding errors they find the exact solution using a finite number of arithmetic operations. In an **iterative method** we start with an approximation \mathbf{x}_0 to the exact solution \mathbf{x} and then compute a sequence $\{\mathbf{x}_k\}$ such that hopefully $\mathbf{x}_k \rightarrow \mathbf{x}$. Iterative methods are mainly used for large sparse systems, i.e., where many of the elements in the coefficient matrix are zero. The main advantages of iterative methods are reduced storage requirements and ease of implementation. In an iterative method the main work in each iteration is a matrix times vector multiplication, an operation which often does not need storing the matrix, not even in sparse form.

We consider the classical iterative methods of Jacobi, Gauss-Seidel, and an accelerated version of Gauss-Seidel's method called Successive OverRelaxation (SOR). David Young developed in his thesis a beautiful theory describing the convergence rate of SOR, see [31].

We give the main points of this theory specialized to the average- and discrete Poisson matrix. With a careful choice of an acceleration parameter the amount of work using SOR on the discrete Poisson problem is the same as for the fast Poisson solver without FFT. Moreover, SOR is not restricted to constant coefficient methods on a rectangle. However, to obtain fast convergence using SOR it is necessary to have a good estimate for the acceleration parameter.

To study these methods we prove three theorems Theorems 9.5, 9.8 and 9.10 which are basic matrix analysis results.



Figure 9.1: David M. Young Jr., 1923-2008

9.1 Classical Iterative Methods; Component Form

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular with nonzero diagonal elements and let $\mathbf{b} \in \mathbb{C}^n$. Solving the i th equation of $\mathbf{Ax} = \mathbf{b}$ for $\mathbf{x}(i)$, we obtain a fixed-point form of $\mathbf{Ax} = \mathbf{b}$

$$\mathbf{x}(i) = \left(- \sum_{j=1}^{i-1} a_{ij} \mathbf{x}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}(j) + b_i \right) / a_{ii}, \quad i = 1, 2, \dots, n. \quad (9.1)$$

Suppose we know an approximation $\mathbf{x}_k = [\mathbf{x}_k(1), \dots, \mathbf{x}_k(n)]^T$ to the exact solution \mathbf{x} of $\mathbf{Ax} = \mathbf{b}$.

1. In **Jacobi's method (J method)** we substitute \mathbf{x}_k into the right hand side of (9.1) and compute a new approximation by

$$\mathbf{x}_{k+1}(i) = \left(- \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_k(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii}, \quad \text{for } i = 1, 2, \dots, n. \quad (9.2)$$

2. **Gauss-Seidel's method (GS method)** is a modification of Jacobi's method, where we use the new $\mathbf{x}_{k+1}(i)$ immediately after it has been computed.

$$\mathbf{x}_{k+1}(i) = \left(- \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii}, \quad \text{for } i = 1, 2, \dots, n. \quad (9.3)$$

3. The **Successive Over Relaxation method (SOR method)** is obtained by introducing an acceleration parameter $0 < \omega < 2$ in the GS method. We write $\mathbf{x}(i) = \omega \mathbf{x}(i) + (1 - \omega) \mathbf{x}(i)$ and this leads to the method

$$\mathbf{x}_{k+1}(i) = \omega \left(- \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii} + (1 - \omega) \mathbf{x}_k(i). \quad (9.4)$$

The SOR method reduces to the Gauss-Seidel method for $\omega = 1$. Denoting the right hand side of (9.3) by \mathbf{x}_{k+1}^{gs} we can write (9.4) as $\mathbf{x}_{k+1} = \omega \mathbf{x}_{k+1}^{gs} +$

$(1 - \omega)\mathbf{x}_k$, and we see that \mathbf{x}_{k+1} is located on the straight line passing through the two points \mathbf{x}_{k+1}^{gs} and \mathbf{x}_k . The restriction $0 < \omega < 2$ is necessary for convergence (cf. Theorem 9.30). Normally we choose the relaxation parameter ω in the range $1 \leq \omega < 2$ and then \mathbf{x}_{k+1} is computed by linear extrapolation, i.e., it is not located between \mathbf{x}_{k+1}^{gs} and \mathbf{x}_k .

4. We mention also briefly the Symmetric Successive Over Relaxation method **SSOR**. One iteration in SSOR consists of two SOR sweeps. A forward SOR sweep (9.4), computing an approximation denoted $\mathbf{x}_{k+1/2}$ instead of \mathbf{x}_{k+1} , is followed by a back SOR sweep computing

$$\mathbf{x}_{k+1}(i) = \omega \left(-\sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1/2}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_{k+1}(j) + b_i \right) / a_{ii} + (1 - \omega) \mathbf{x}_{k+1/2}(i) \quad (9.5)$$

in the order $i = n, n-1, \dots, 1$. The method is slower and more complicated than the SOR method. Its main use is as a symmetric preconditioner. For if \mathbf{A} is symmetric then SSOR combines the two SOR steps in such a way that the resulting iteration matrix is similar to a symmetric matrix. We will not discuss this method any further here and refer to Section 10.7 for an alternative example of a preconditioner.

We will refer to the J,GS, and SOR methods as the **classical (iteration) methods**.

9.1.1 The Discrete Poisson System

Consider the classical methods applied to the discrete Poisson matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ given by (4.7). Let $n = m^2$ and set $h = 1/(m+1)$. In component form the linear system $\mathbf{Ax} = \mathbf{b}$ can be written (cf. (4.3))

$$4\mathbf{u}(i,j) - \mathbf{u}(i-1,j) - \mathbf{u}(i+1,j) - \mathbf{u}(i,j-1) - \mathbf{u}(i,j+1) = h^2 \mathbf{f}(i,j), \quad i, j = 1, \dots, m,$$

with homogenous boundary conditions (4.4). Solving for $\mathbf{u}(i,j)$ we obtain

$$\mathbf{u}(i,j) = (\mathbf{u}(i-1,j) + \mathbf{u}(i+1,j) + \mathbf{u}(i,j-1) + \mathbf{u}(i,j+1) + h^2 \mathbf{f}(i,j)) / 4. \quad (9.6)$$

	k_{100}	k_{2500}	$k_{10\ 000}$	$k_{40\ 000}$	$k_{160\ 000}$
J	385	8386			
GS	194	4194			
SOR	35	164	324	645	1286

Table 9.1: The number of iterations k_n to solve the $n \times n$ discrete Poisson problem using the methods of Jacobi, Gauss-Seidel, and SOR (see text) with a tolerance 10^{-8} .

The J, GS , and SOR methods can now be written

$$\begin{aligned}
 J : \quad & \mathbf{v}_{k+1}(i, j) = \left(\mathbf{v}_k(i-1, j) + \mathbf{v}_k(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \right. \\
 & \quad \left. + h^2 \mathbf{f}(i, j) \right) / 4 \\
 GS : \quad & \mathbf{v}_{k+1}(i, j) = \left(\mathbf{v}_{k+1}(i-1, j) + \mathbf{v}_{k+1}(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \right. \\
 & \quad \left. + h^2 \mathbf{f}(i, j) \right) / 4 \\
 SOR : \quad & \mathbf{v}_{k+1}(i, j) = \omega \left(\mathbf{v}_{k+1}(i-1, j) + \mathbf{v}_{k+1}(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \right. \\
 & \quad \left. + h^2 \mathbf{f}(i, j) \right) / 4 + (1 - \omega) \mathbf{v}_k(i, j).
 \end{aligned} \tag{9.7}$$

For GS and SOR we use the **natural ordering** i.e., with i, j in increasing order $i, j = 1, \dots, m$, while for J any ordering can be used.

In Algorithm 9.1 we give a Matlab program to test the convergence of Jacobi's method on the discrete Poisson problem. We carry out Jacobi iterations on the linear system (9.6) with $\mathbf{F} = (f_{ij}) \in \mathbb{R}^{m \times m}$, starting with $\mathbf{V}^{(0)} = \mathbf{0} \in \mathbb{R}^{(m+2) \times (m+2)}$. The output is the number of iterations k , to obtain $\|\mathbf{V}^{(k)} - \mathbf{U}\|_M := \max_{i,j} |v^{(k)}_{ij} - u_{ij}| < tol$. Here $(u_{ij}) \in \mathbb{R}^{(m+2) \times (m+2)}$ is the "exact" solution of (9.6) computed using the fast Poisson solver in Algorithm 5.1. We set $k = K + 1$ if convergence is not obtained in K iterations. In Table 9.1 we show the output $k = k_n$ from this algorithm using $\mathbf{F} = \text{ones}(m, m)$ for $m = 10, 50$, $K = 10^4$, and $tol = 10^{-8}$. We also show the number of iterations for Gauss-Seidel and SOR with a value of ω known as the optimal acceleration parameter $\omega = 2/(1 + \sin(\frac{\pi}{m+1}))$. We will derive this value later.

Algorithm 9.1 (Jacobi)

```

1 function k=jdp(F,K,tol)
2 m=length(F); U=fastpoisson(F);
3 V=zeros(m+2,m+2); W=V; E=F/(m+1)^2;
4 for k=1:K
5   for i=2:m+1
6     for j=2:m+1
7       W(i,j)=(V(i-1,j)+V(i+1,j)+V(i,j-1)...
8         +V(i,j+1)+E(i-1,j-1))/4;
9     end
10    end
11    if max(max(abs(W-U)))<tol, return
12  end
13  V=W;
14 end
15 k=K+1;

```

For the GS and SOR methods we have used Algorithm 9.2. This is the analog of Algorithm 9.1 using GS and SOR instead of J to solve the discrete Poisson problem. w is an acceleration parameter with $0 < w < 2$. For $w = 1$ we obtain Gauss-Seidel's method.

Algorithm 9.2 (SOR)

```

1 function k=sordp(F,K,w,tol)
2 m=length(F); U=fastpoisson(F);
3 V=zeros(m+2,m+2); E=F/(m+1)^2;
4 for k=1:K
5   for i=2:m+1
6     for j=2:m+1
7       V(i,j)=w*(V(i-1,j)+V(i+1,j)+V(i,j-1)...
8         +V(i,j+1)+E(i-1,j-1))/4+(1-w)*V(i,j);
9     end
10    end
11    if max(max(abs(V-U)))<tol, return
12  end
13 end
14 k=K+1;

```

We make several remarks about these programs and the results in Table 9.1.

1. The rate (speed) of convergence is quite different for the three methods. The J and GS method converge, but rather slowly. The J method needs about twice as many iterations as the GS method. The improvement using the SOR method with optimal ω is rather spectacular.
2. We show in Section 9.4.1 that the number of iterations k_n for a size n problem

is $k_n = O(n)$ for the J and GS method and $k_n = O(\sqrt{n})$ for SOR with optimal ω . The choice of tol will only influence the constants multiplying n or \sqrt{n} .

3. From (9.7) it follows that each iteration requires $O(n)$ arithmetic operations. Thus the number of arithmetic operations to achieve a given tolerance is $O(k_n \times n)$. Therefore the number of arithmetic operations for the J and GS method is $O(n^2)$, while it is only $O(n^{3/2})$ for the SOR method with optimal ω . Asymptotically, for J and GS this is the same as using banded Cholesky, while SOR competes with the fast method (without FFT).
4. We do not need to store the coefficient matrix so the storage requirements for these methods on the discrete Poisson problem is $O(n)$, asymptotically the same as for the fast methods. For the GS and SOR method we can store the new $\mathbf{v}_{k+1}(i, j)$ in the same location as $\mathbf{v}_k(i, j)$. For Jacobi's method we need an extra array. (w in Algorithm 9.1).
5. Jacobi's method has the advantage that it can be easily parallelized.

9.1.2 Matrix Formulations of the Classical Methods

To study convergence it is convenient to use matrix formulations of the classical methods. In general we can construct an iterative method by choosing a nonsingular matrix \mathbf{M} and write $\mathbf{A}\mathbf{x} = \mathbf{b}$ in the equivalent form $\mathbf{M}^{-1}\mathbf{A}\mathbf{x} = \mathbf{M}^{-1}\mathbf{b}$. This system can be written $\mathbf{x} = \mathbf{x} - \mathbf{M}^{-1}\mathbf{A}\mathbf{x} + \mathbf{M}^{-1}\mathbf{b} = (\mathbf{I} - \mathbf{M}^{-1}\mathbf{A})\mathbf{x} + \mathbf{M}^{-1}\mathbf{b}$, and we obtain $\mathbf{A}\mathbf{x} = \mathbf{b}$ in a **fixed point form**

$$\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}, \quad \mathbf{G} = \mathbf{I} - \mathbf{M}^{-1}\mathbf{A}, \quad \mathbf{c} = \mathbf{M}^{-1}\mathbf{b}. \quad (9.8)$$

The corresponding iterative method is given by

$$\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}. \quad (9.9)$$

Different choices of \mathbf{M} leads to different iterative methods. The matrix \mathbf{M} can be interpreted in two ways. It is a **preconditioning matrix** since a good choice of \mathbf{M} can lead to a system $\mathbf{M}^{-1}\mathbf{A}\mathbf{x} = \mathbf{M}^{-1}\mathbf{b}$ with smaller condition number. It is also known as a **splitting matrix**, since if we split \mathbf{A} in the form $\mathbf{A} = \mathbf{M} + (\mathbf{A} - \mathbf{M})$ then $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be written $\mathbf{M}\mathbf{x} = (\mathbf{M} - \mathbf{A})\mathbf{x} + \mathbf{b}$ and this leads to the iterative method

$$\mathbf{M}\mathbf{x}_{k+1} = (\mathbf{M} - \mathbf{A})\mathbf{x}_k + \mathbf{b} \quad (9.10)$$

which is equivalent to (9.9).

To study convergence we subtract (9.8) from (9.9). The vector \mathbf{c} cancels and we obtain

$$\boldsymbol{\epsilon}_{k+1} = \mathbf{G}\boldsymbol{\epsilon}_k, \quad \boldsymbol{\epsilon}_j := \mathbf{x}_j - \mathbf{x}, \quad j = 0, 1, 2, \dots \quad (9.11)$$

By induction on k

$$\epsilon_k = \mathbf{G}^k \epsilon_0, \quad k = 0, 1, 2, \dots \quad (9.12)$$

Thus the convergence depends on the behavior of the powers \mathbf{G}^k as k increases. The matrix \mathbf{M} should be chosen so that all elements in \mathbf{G}^k converges quickly to zero and such that the linear system (9.10) is easy to solve for \mathbf{x}_{k+1} . We will see that these are conflicting demands. \mathbf{M} should be an approximation to \mathbf{A} to obtain a \mathbf{G} with small elements, but then (9.10) might not be easy to solve for \mathbf{x}_{k+1} .

9.1.3 The Splitting Matrices for the Classical Methods

Before continuing the study of convergence we derive \mathbf{M} for the classical methods. It is convenient to write \mathbf{A} as a sum of three matrices, $\mathbf{A} = \mathbf{D} - \mathbf{A}_L - \mathbf{A}_R$, where $-\mathbf{A}_L$, \mathbf{D} , and $-\mathbf{A}_R$ are the lower, diagonal, and upper part of \mathbf{A} , respectively. Thus $\mathbf{D} := \text{diag}(a_{11}, \dots, a_{nn})$,

$$\mathbf{A}_L := \begin{bmatrix} 0 & & & & \\ -a_{21} & 0 & & & \\ \vdots & \ddots & \ddots & & \\ -a_{n,1} & \cdots & -a_{n,n-1} & 0 & \end{bmatrix}, \quad \mathbf{A}_R := \begin{bmatrix} 0 & -a_{12} & \cdots & -a_{1n} \\ & \ddots & \ddots & \vdots \\ & 0 & -a_{n-1,n} & 0 \end{bmatrix}. \quad (9.13)$$

Proposition 9.3 (Splitting matrices for J, GS, and SOR)

The splitting matrices $\mathbf{M}_J, \mathbf{M}_1, \mathbf{M}_\omega$ for the J, GS, and SOR method are given by

$$\mathbf{M}_J = \mathbf{D}, \quad \mathbf{M}_1 = \mathbf{D} - \mathbf{A}_L, \quad \mathbf{M}_\omega = \omega^{-1} \mathbf{D} - \mathbf{A}_L. \quad (9.14)$$

Proof. The equation $\mathbf{Ax} = \mathbf{b}$ can be written $\mathbf{Dx} - \mathbf{A}_L \mathbf{x} - \mathbf{A}_R \mathbf{x} = \mathbf{b}$ or $\mathbf{Dx} = \mathbf{A}_L \mathbf{x} + \mathbf{A}_R \mathbf{x} + \mathbf{b}$. This leads to

$$\begin{aligned} J: \quad & \mathbf{Dx}_{k+1} = \mathbf{A}_L \mathbf{x}_k + \mathbf{A}_R \mathbf{x}_k + \mathbf{b}, \text{ or} \\ & \mathbf{M}_J \mathbf{x}_{k+1} = (\mathbf{A}_L + \mathbf{A}_R) \mathbf{x}_k + \mathbf{b}, \\ GS: \quad & \mathbf{Dx}_{k+1} = \mathbf{A}_L \mathbf{x}_{k+1} + \mathbf{A}_R \mathbf{x}_k + \mathbf{b}, \text{ or} \\ & \mathbf{M}_1 \mathbf{x}_{k+1} = \mathbf{A}_R \mathbf{x}_k + \mathbf{b}, \\ SOR: \quad & \mathbf{Dx}_{k+1} = \omega(\mathbf{A}_L \mathbf{x}_{k+1} + \mathbf{A}_R \mathbf{x}_k + \mathbf{b}) + (1 - \omega) \mathbf{Dx}_k, \text{ or} \\ & \mathbf{M}_\omega \mathbf{x}_{k+1} = (\mathbf{A}_R + (\omega^{-1} - 1) \mathbf{D}) \mathbf{x}_k + \mathbf{b}. \end{aligned} \quad (9.15)$$

These expressions are of the form (9.10). \square

Example 9.4 (Splitting matrices)

For the system

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

we find

$$\mathbf{A}_L = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \mathbf{A}_R = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix},$$

and

$$\mathbf{M}_J = \mathbf{D} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \mathbf{M}_\omega = \omega^{-1} \mathbf{D} - \mathbf{A}_L = \begin{bmatrix} 2\omega^{-1} & 0 \\ -1 & 2\omega^{-1} \end{bmatrix}.$$

The iteration matrix $\mathbf{G}_\omega = \mathbf{I} - \mathbf{M}_\omega^{-1} \mathbf{A}$ is given by

$$\mathbf{G}_\omega = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} \omega/2 & 0 \\ \omega^2/4 & \omega/2 \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 1-\omega & \omega/2 \\ \omega(1-\omega)/2 & 1-\omega+\omega^2/4 \end{bmatrix}. \quad (9.16)$$

For the J and GS method we have

$$\mathbf{G}_J = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A} = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix}, \quad \mathbf{G}_1 = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix}. \quad (9.17)$$

We could have derived these matrices directly from the component form of the iteration. For example, for the GS method we have the component form

$$\mathbf{x}_{k+1}(1) = \frac{1}{2} \mathbf{x}_k(2) + \frac{1}{2}, \quad \mathbf{x}_{k+1}(2) = \frac{1}{2} \mathbf{x}_{k+1}(1) + \frac{1}{2}.$$

Substituting the value of $\mathbf{x}_{k+1}(1)$ from the first equation into the second equation we find

$$\mathbf{x}_{k+1}(2) = \frac{1}{2} \left(\frac{1}{2} \mathbf{x}_k(2) + \frac{1}{2} \right) + \frac{1}{2} = \frac{1}{4} \mathbf{x}_k(2) + \frac{3}{4}.$$

Thus

$$\mathbf{x}_{k+1} = \begin{bmatrix} \mathbf{x}_{k+1}(1) \\ \mathbf{x}_{k+1}(2) \end{bmatrix} = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix} \begin{bmatrix} \mathbf{x}_k(1) \\ \mathbf{x}_k(2) \end{bmatrix} + \begin{bmatrix} 1/2 \\ 3/4 \end{bmatrix} = \mathbf{G}_1 \mathbf{x}_k + \mathbf{c}.$$

9.2 Convergence and Spectral Radius

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a square matrix. In this section we consider the special matrix sequence $\{\mathbf{A}^k\}$ of powers of \mathbf{A} . Such a sequence occurs in iterative methods (cf (9.12)), in Markov processes in statistics, and in many other applications. We want to know when this sequence converges to the zero matrix.

To start we define the **spectral radius** of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ as the maximum absolute value of its eigenvalues.

$$\rho(\mathbf{A}) := \max_{\lambda \in \sigma(\mathbf{A})} |\lambda|. \quad (9.18)$$

In this section we show the following theorem.

Theorem 9.5 (When is $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$?)

For any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have

$$\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0} \iff \rho(\mathbf{A}) < 1.$$

Clearly $\rho(\mathbf{A}) < 1$ is a necessary condition for $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$. For if (λ, \mathbf{x}) is an eigenpair of \mathbf{A} with $|\lambda| \geq 1$ then $\mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}$, and it follows that $\mathbf{A}^k \mathbf{x}$ does not tend to zero. But then we cannot have $\mathbf{A}^k \rightarrow \mathbf{0}$.

The sufficiency condition is much harder to show. It is enough to find a matrix norm such that $\|\mathbf{A}\| < 1$. Moreover the norm should be consistent on $\mathbb{C}^{n \times n}$, i.e., $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$ for all $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$. For then $\|\mathbf{A}^k\| \leq \|\mathbf{A}\|^k \rightarrow 0$.

We first show.

Theorem 9.6 (Any consistent norm majorizes the spectral radius)

For any matrix norm $\|\cdot\|$ which is consistent on $\mathbb{C}^{n \times n}$ and any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$.

Proof. Let (λ, \mathbf{x}) be an eigenpair for \mathbf{A} and define $\mathbf{X} := [\mathbf{x}, \dots, \mathbf{x}] \in \mathbb{C}^{n \times n}$. Then $\lambda \mathbf{X} = \mathbf{AX}$, which implies $|\lambda| \|\mathbf{X}\| = \|\lambda \mathbf{X}\| = \|\mathbf{AX}\| \leq \|\mathbf{A}\| \|\mathbf{X}\|$. Since $\|\mathbf{X}\| \neq 0$ we obtain $|\lambda| \leq \|\mathbf{A}\|$. \square

The next theorem shows that if $\rho(\mathbf{A}) < 1$ then $\|\mathbf{A}\| < 1$ for some consistent matrix norm on $\mathbb{C}^{n \times n}$, thus completing the proof of Theorem 9.5.

Theorem 9.7 (The spectral radius can be approximated by a norm)

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\epsilon > 0$ be given. There is a consistent matrix norm $\|\cdot\|$ on $\mathbb{C}^{n \times n}$ such that $\rho(\mathbf{A}) \leq \|\mathbf{A}\| \leq \rho(\mathbf{A}) + \epsilon$.

Proof. Let \mathbf{A} have eigenvalues $\lambda_1, \dots, \lambda_n$. By the Schur Triangulation Theorem 6.29 there is a unitary matrix \mathbf{U} and an upper triangular matrix $\mathbf{R} = [r_{ij}]$ such that $\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{R}$. For $t > 0$ we define $\mathbf{D}_t := \text{diag}(t, t^2, \dots, t^n) \in \mathbb{R}^{n \times n}$, and note that the (i, j) element in $\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}$ is given by $t^{i-j} r_{ij}$ for all i, j . For $n = 3$

$$\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1} = \begin{bmatrix} \lambda_1 & t^{-1} r_{12} & t^{-2} r_{13} \\ 0 & \lambda_2 & t^{-1} r_{23} \\ 0 & 0 & \lambda_3 \end{bmatrix}.$$

For each $\mathbf{B} \in \mathbb{C}^{n \times n}$ and $t > 0$ we define $\|\mathbf{B}\|_t := \|\mathbf{D}_t \mathbf{U}^* \mathbf{B} \mathbf{U} \mathbf{D}_t^{-1}\|_1$. We leave it as an exercise to show that $\|\cdot\|_t$ is a consistent matrix norm on $\mathbb{C}^{n \times n}$. We define $\|\mathbf{B}\| := \|\mathbf{B}\|_t$, where t is chosen so large that the sum of the absolute values of all

off-diagonal elements in $\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}$ is less than ϵ . Then

$$\begin{aligned}\|\mathbf{A}\| &= \|\mathbf{D}_t \mathbf{U}^* \mathbf{A} \mathbf{U} \mathbf{D}_t^{-1}\|_1 = \|\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |(\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1})_{ij}| \\ &\leq \max_{1 \leq j \leq n} (|\lambda_j| + \epsilon) = \rho(\mathbf{A}) + \epsilon.\end{aligned}$$

□

A consistent matrix norm of a matrix can be much larger than the spectral radius. However the following result holds.

Theorem 9.8 (Spectral radius convergence)

For any consistent matrix norm $\|\cdot\|$ on $\mathbb{C}^{n \times n}$ and any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have

$$\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A}). \quad (9.19)$$

Proof. By Theorems 0.66 and 9.6 we obtain $\rho(\mathbf{A})^k = \rho(\mathbf{A}^k) \leq \|\mathbf{A}^k\|$ for any $k \in \mathbb{N}$ so that $\rho(\mathbf{A}) \leq \|\mathbf{A}^k\|^{1/k}$. Let $\epsilon > 0$ and consider the matrix $\mathbf{B} := (\rho(\mathbf{A}) + \epsilon)^{-1} \mathbf{A}$. Then $\rho(\mathbf{B}) = \rho(\mathbf{A}) / (\rho(\mathbf{A}) + \epsilon) < 1$ and $\|\mathbf{B}^k\| \rightarrow 0$ by Theorem 9.5 as $k \rightarrow \infty$. Choose $N \in \mathbb{N}$ such that $\|\mathbf{B}^k\| < 1$ for all $k \geq N$. Then for $k \geq N$

$$\|\mathbf{A}^k\| = \|(\rho(\mathbf{A}) + \epsilon)\mathbf{B}\|^k = (\rho(\mathbf{A}) + \epsilon)^k \|\mathbf{B}^k\| < (\rho(\mathbf{A}) + \epsilon)^k.$$

We have shown that $\rho(\mathbf{A}) \leq \|\mathbf{A}^k\|^{1/k} \leq \rho(\mathbf{A}) + \epsilon$ for $k \geq N$. Since ϵ is arbitrary the result follows. □

Exercise 9.9 (Slow spectral radius convergence)

The convergence $\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A})$ can be quite slow. Consider

$$\mathbf{A} := \begin{bmatrix} \lambda & a & 0 & \cdots & 0 & 0 \\ 0 & \lambda & a & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \ddots & \\ 0 & 0 & 0 & \cdots & \lambda & a \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

If $|\lambda| = \rho(\mathbf{A}) < 1$ then $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$ for any $a \in \mathbb{R}$. We show below that the $(1, n)$ element of \mathbf{A}^k is given by $f(k) := \binom{k}{n-1} a^{n-1} \lambda^{k-n+1}$ for $k \geq n-1$.

- (a) Pick an n , e.g. $n = 5$, and make a plot of $f(k)$ for $\lambda = 0.9$, $a = 10$, and $n-1 \leq k \leq 200$. Your program should also compute $\max_k f(k)$. Use your program to determine how large k must be before $f(k) < 10^{-8}$.

- (b) We can determine the elements of \mathbf{A}^k explicitly for any k . Let $\mathbf{E} := (\mathbf{A} - \lambda\mathbf{I})/a$. Show by induction that $\mathbf{E}^k = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{n-k} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ for $1 \leq k \leq n-1$ and that $\mathbf{E}^n = \mathbf{0}$.
- (c) We have $\mathbf{A}^k = (a\mathbf{E} + \lambda\mathbf{I})^k = \sum_{j=0}^{\min\{k,n-1\}} \binom{k}{j} a^j \lambda^{k-j} \mathbf{E}^j$ and conclude that the $(1, n)$ element is given by $f(k)$ for $k \geq n-1$.

9.2.1 Neumann Series

Let \mathbf{B} be a square matrix. In this section we consider the **Neumann Series** $\sum_{k=0}^{\infty} \mathbf{B}^k$ which is a matrix analogue of a geometric series of numbers.

Consider an infinite series $\sum_{m=0}^{\infty} \mathbf{y}_m$ of vectors in \mathbb{C}^n . We say that the series converges if the sequence of partial sums $\{\mathbf{x}_k\}$ given by $\mathbf{x}_k = \sum_{m=0}^k \mathbf{y}_m$ converges. A sufficient condition for convergence is that $\sum_{m=0}^{\infty} \|\mathbf{y}_m\|$ converges for some vector norm. We say that the series converges **absolutely** if this is the case. Convergence and absolute convergence of a series of matrices is defined analogously.

Theorem 9.10 (Neumann Series)

Suppose $\mathbf{B} \in \mathbb{C}^{n \times n}$. Then

1. The series $\sum_{k=0}^{\infty} \mathbf{B}^k$ converges if and only if $\rho(\mathbf{B}) < 1$.
2. If $\rho(\mathbf{B}) < 1$ then $(\mathbf{I} - \mathbf{B})$ is nonsingular and $(\mathbf{I} - \mathbf{B})^{-1} = \sum_{k=0}^{\infty} \mathbf{B}^k$.
3. If $\|\mathbf{B}\| < 1$ for some consistent matrix norm $\|\cdot\|$ on $\mathbb{C}^{n \times n}$ then

$$\|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}. \quad (9.20)$$

Proof.

1. Suppose $\rho(\mathbf{B}) < 1$. We show that the sequence $\{\mathbf{A}_m\}$ of partial sums $\mathbf{A}_m := \sum_{k=0}^m \mathbf{B}^k$ is a Cauchy sequence and hence convergent. Let $\epsilon > 0$. By Theorem 9.7 there is a consistent matrix norm $\|\cdot\|$ on $\mathbb{C}^{n \times n}$ such that $\|\mathbf{B}\| < 1$. Then for $l > m$

$$\|\mathbf{A}_l - \mathbf{A}_m\| = \left\| \sum_{k=m+1}^l \mathbf{B}^k \right\| \leq \sum_{k=m+1}^l \|\mathbf{B}\|^k \leq \frac{\|\mathbf{B}\|^{m+1}}{1 - \|\mathbf{B}\|} \leq \epsilon$$

provided $m \geq N$ and N is such that $\frac{\|\mathbf{B}\|^{N+1}}{1 - \|\mathbf{B}\|} \leq \epsilon$. Thus $\{\mathbf{A}_m\}$ is a Cauchy sequence and hence convergent.

Conversely, suppose (λ, \mathbf{x}) is an eigenpair for \mathbf{B} with $\lambda \geq 1$. Now for $l > m$

$$\|(\mathbf{A}_l - \mathbf{A}_m)\mathbf{x}\| = \left\| \sum_{k=m+1}^l \mathbf{B}^k \mathbf{x} \right\| = \left\| \sum_{k=m+1}^l \lambda^k \mathbf{x} \right\| = \|\mathbf{x}\| \sum_{k=m+1}^l |\lambda|^k \geq |\lambda|^{m+1} \|\mathbf{x}\|.$$

But then $\{\mathbf{A}_m\}$ cannot be a Cauchy sequence and hence not convergent.

2. By induction on m it follows that

$$\left(\sum_{k=0}^m \mathbf{B}^k \right) (\mathbf{I} - \mathbf{B}) = \mathbf{I} - \mathbf{B}^{m+1}. \quad (9.21)$$

For if $\left(\sum_{k=0}^{m-1} \mathbf{B}^k \right) (\mathbf{I} - \mathbf{B}) = \mathbf{I} - \mathbf{B}^m$ then

$$\left(\sum_{k=0}^m \mathbf{B}^k \right) (\mathbf{I} - \mathbf{B}) = \left(\sum_{k=0}^{m-1} \mathbf{B}^k + \mathbf{B}^m \right) (\mathbf{I} - \mathbf{B}) = \mathbf{I} - \mathbf{B}^m + \mathbf{B}^m - \mathbf{B}^{m+1} = \mathbf{I} - \mathbf{B}^{m+1}.$$

Since $\rho(\mathbf{B}) < 1$ we conclude that $\mathbf{B}^{m+1} \rightarrow 0$ and hence taking limits in (9.21) we obtain $\left(\sum_{k=0}^{\infty} \mathbf{B}^k \right) (\mathbf{I} - \mathbf{B}) = \mathbf{I}$ which completes the proof of 2.

3. By 2: $\|(\mathbf{I} - \mathbf{B})^{-1}\| = \left\| \sum_{k=0}^{\infty} \mathbf{B}^k \right\| \leq \sum_{k=0}^{\infty} \|\mathbf{B}\|^k = \frac{1}{1 - \|\mathbf{B}\|}$.

□

Exercise 9.11 (A special norm)

Show that $\|\mathbf{B}\|_t := \|\mathbf{D}_t \mathbf{U}^* \mathbf{B} \mathbf{U} \mathbf{D}_t^{-1}\|_1$ defined in the proof of Theorem 9.7 is a consistent matrix norm on $\mathbb{C}^{n \times n}$.

Exercise 9.12 (When is $\mathbf{A} + \mathbf{E}$ nonsingular?)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular and $\mathbf{E} \in \mathbb{C}^{n \times n}$. Show that $\mathbf{A} + \mathbf{E}$ is nonsingular if and only if $\rho(\mathbf{A}^{-1} \mathbf{E}) < 1$.

9.3 Convergence of Fixed-point Iteration

We have seen that the classical methods can be written in the form (9.9) for a suitable \mathbf{M} . Starting with \mathbf{x}_0 this defines a sequence $\{\mathbf{x}_k\}$ of vectors in \mathbb{C}^n . If $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}$ for some $\mathbf{x} \in \mathbb{C}^n$ then \mathbf{x} is a solution of $\mathbf{x} = \mathbf{Gx} + \mathbf{c}$ since

$$\mathbf{x} = \lim_{k \rightarrow \infty} \mathbf{x}_{k+1} = \lim_{k \rightarrow \infty} (\mathbf{Gx}_k + \mathbf{c}) = \mathbf{G} \lim_{k \rightarrow \infty} \mathbf{x}_k + \mathbf{c} = \mathbf{Gx} + \mathbf{c}.$$

For a general $\mathbf{G} \in \mathbb{C}^{n \times n}$ and $\mathbf{c} \in \mathbb{C}^n$ a solution of $\mathbf{x} = \mathbf{Gx} + \mathbf{c}$ is called a **fixed-point** and the iteration $\mathbf{x}_{k+1} = \mathbf{Gx}_k + \mathbf{c}$ a **fixed-point iteration**. The fixed-point is unique if $\mathbf{I} - \mathbf{G}$ is nonsingular.

Consider next convergence of fixed-point iteration.

Definition 9.13 (Convergence of fixed-point iteration)

We say that the iterative method $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$ converges if the sequence $\{\mathbf{x}_k\}$ converges for any starting vector \mathbf{x}_0 .

Lemma 9.14 (Convergence of an iterative method)

The iterative method $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$ converges if and only if $\lim_{k \rightarrow \infty} \mathbf{G}^k = \mathbf{0}$.

Proof. By (9.12) we have $\boldsymbol{\epsilon}_k = \mathbf{G}^k \boldsymbol{\epsilon}_0$ for $k \geq 0$, where $\boldsymbol{\epsilon}_j := \mathbf{x}_j - \mathbf{x}$ for $j \in \mathbb{N}$. Now $\boldsymbol{\epsilon}_k \rightarrow \mathbf{0}$ if $\mathbf{G}^k \rightarrow \mathbf{0}$. The converse follows by choosing \mathbf{x}_0 so that $\boldsymbol{\epsilon}_0 = \mathbf{e}_j$, the j th unit vector for $j = 1, \dots, n$. \square

Using Theorem 9.6 we obtain the following theorem:

Theorem 9.15 (When does an iterative method converge?)

Suppose $\mathbf{G} \in \mathbb{C}^{n \times n}$ and $\mathbf{c} \in \mathbb{C}^n$. The iteration $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$ converges if and only if $\rho(\mathbf{G}) < 1$.

Since $\rho(\mathbf{G}) < \|\mathbf{G}\|$ for any consistent matrix norm on $\mathbb{C}^{n \times n}$ (cf. Theorem 9.6) we obtain

Corollary 9.16 (Sufficient condition for convergence)

If $\|\mathbf{G}\| < 1$ for some consistent matrix norm, then the iteration $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$ converges.

Exercise 9.17 (Divergence example for J and GS)

Show that both Jacobi's method and Gauss-Seidel's method diverge for $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$.

Exercise 9.18 (J and GS on spline matrix)

Explain why J and GS converge for the cubic spline matrix \mathbf{N} in Chapter 2. (This is mainly of academic interest since each iteration requires $O(n)$ arithmetic operations the same as the complete Gaussian elimination process for a tridiagonal system.)

Exercise 9.19 (Strictly diagonally dominance; The J method)

Show that the J method converges if $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$ for $i = 1, \dots, n$.

Exercise 9.20 (Strictly diagonally dominance; The GS method)

Consider the GS method. Suppose $r := \max_i r_i < 1$, where $r_i = \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|}$. Show using induction on i that $|\boldsymbol{\epsilon}_{k+1}(j)| \leq r \|\boldsymbol{\epsilon}_k\|_\infty$ for $j = 1, \dots, i$. Conclude that Gauss-Seidel's method is convergent when \mathbf{A} is strictly diagonally dominant.

Consider next the **rate of convergence**. Suppose $\|\cdot\|$ is a matrix norm that is subordinate to a vector norm also denoted by $\|\cdot\|$. Taking norms in $\boldsymbol{\epsilon}_k = \mathbf{G}^k \boldsymbol{\epsilon}_0$ we obtain

$$\|\boldsymbol{\epsilon}_k\| = \|\mathbf{G}^k \boldsymbol{\epsilon}_0\| \leq \|\mathbf{G}^k\| \|\boldsymbol{\epsilon}_0\| \approx \rho(\mathbf{G})^k \|\boldsymbol{\epsilon}_0\|.$$

For the last formula we apply Theorem 9.8 which says that $\lim_{k \rightarrow \infty} \|\mathbf{G}^k\|^{1/k} = \rho(\mathbf{G})$. Thus for fast convergence we should use a \mathbf{G} with small spectral radius.

Lemma 9.21 (Number of iterations)

Suppose $\rho(\mathbf{G}) = 1 - \eta$ for some $0 < \eta < 1$, $\|\cdot\|$ a consistent matrix norm, and let $s \in \mathbb{N}$. Then

$$\tilde{k} := \frac{\log(10)s}{\eta} \quad (9.22)$$

is an estimate for the smallest number of iterations k so that $\rho(\mathbf{G})^k \leq 10^{-s}$.

Proof. \tilde{k} is an approximate solution of the equation $\rho(\mathbf{G})^k = 10^{-s}$. Indeed, taking logarithms we find $k \log \rho(\mathbf{G}) = -s \log 10$. Thus

$$k = -\frac{s \log(10)}{\log(1-\eta)} = \frac{s \log(10)}{\eta + O(\eta^2)} \approx \frac{\log(10)s}{\eta} = \tilde{k}.$$

□

Exercise 9.22 (Estimate in Lemma 9.21 can be exact)

Consider the iteration in Example 9.4. Show that $\rho(\mathbf{G}_J) = 1/2$. Then show that $\mathbf{x}_k(1) = \mathbf{x}_k(2) = 1 - 2^{-k}$ for $k \geq 0$. Thus the estimate in Lemma 9.21 is exact in this case.

The convergence $\lim_{k \rightarrow \infty} \|\mathbf{G}^k\|^{1/k} = \rho(\mathbf{G})$ can be quite slow, (cf. Exercise 9.9).

9.3.1 Stopping the Iteration

In Algorithms 9.1 and 9.2 we had access to the exact solution and could stop the iteration when the error was sufficiently small in the infinity norm. The decision when to stop is obviously more complicated when the exact solution is not known. One possibility is to choose a vector norm, keep track of $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|$, and stop when this number is sufficiently small. This must be applied with some care if $\|\mathbf{G}\|$ is close to one, as the following result indicates.

Lemma 9.23 (Be careful when stopping)

Suppose $\|\mathbf{G}\| < 1$ for some consistent matrix norm which is subordinate to a vector norm also denoted by $\|\cdot\|$. If $\mathbf{x}_k = \mathbf{G}\mathbf{x}_{k-1} + \mathbf{c}$ and $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$. Then

$$\|\mathbf{x}_k - \mathbf{x}\| \leq \frac{\|\mathbf{G}\|}{1 - \|\mathbf{G}\|} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|, \quad k \geq 1. \quad (9.23)$$

Proof. We find

$$\begin{aligned}\|\mathbf{x}_k - \mathbf{x}\| &= \|\mathbf{G}(\mathbf{x}_{k-1} - \mathbf{x})\| \leq \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}\| \\ &= \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}\| \leq \|\mathbf{G}\| (\|\mathbf{x}_{k-1} - \mathbf{x}_k\| + \|\mathbf{x}_k - \mathbf{x}\|).\end{aligned}$$

Thus $(1 - \|\mathbf{G}\|) \|\mathbf{x}_k - \mathbf{x}\| \leq \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}_k\|$ which implies (9.23). \square

Another possibility is to stop when the residual vector $\mathbf{r}_k := \mathbf{b} - \mathbf{A}\mathbf{x}_k$ is sufficiently small in some norm. To use the residual vector for stopping it is convenient to write the iterative method (9.9) in an alternative form. If \mathbf{M} is the splitting matrix of the method then by (9.10) we have $\mathbf{M}\mathbf{x}_{k+1} = \mathbf{M}\mathbf{x}_k - \mathbf{A}\mathbf{x}_k + \mathbf{b}$. This leads to

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{M}^{-1}\mathbf{r}_k, \quad \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k. \quad (9.24)$$

Testing on \mathbf{r}_k works fine if \mathbf{A} is well conditioned, but Theorem 8.32 shows that the relative error in the solution can be much larger than the relative error in \mathbf{r}_k if \mathbf{A} is ill-conditioned.

9.3.2 Richardson's Method (R method)

This method is based on the simple splitting $\mathbf{M}_R := \alpha\mathbf{I}$, where α is a nonzero scalar. By (9.24) we obtain Richardson's method in the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha^{-1}\mathbf{r}_k, \quad \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k. \quad (9.25)$$

If all eigenvalues of \mathbf{A} have positive real parts then the R method converges provided α is sufficiently large.

Proposition 9.24 (Convergence of Richardson's method)

Suppose all eigenvalues of \mathbf{A} have positive real parts and that α is real. Then there is an α_0 such that the R method converges for $\alpha > \alpha_0$. If \mathbf{A} has positive eigenvalues $0 < \lambda_n \leq \dots \leq \lambda_1$ then the spectral radius of

$$\mathbf{G}(\alpha) := \mathbf{I} - \alpha^{-1}\mathbf{A}$$

is uniquely minimized if $\alpha = \alpha^*$, where

$$\alpha^* := \frac{\lambda_1 + \lambda_n}{2}, \quad \text{and } \rho(\mathbf{G}(\alpha^*)) = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}. \quad (9.26)$$

Proof. The eigenvalues of $\mathbf{G}(\alpha)$ are

$$\mu_j(\alpha) = 1 - \lambda_j/\alpha, \quad j = 1, \dots, n,$$

and if $u_j := Re\lambda_j > 0$ then

$$|\mu_j(\alpha)|^2 = \left(1 - \frac{\lambda_j}{\alpha}\right)\left(1 - \frac{\bar{\lambda}_j}{\alpha}\right) = 1 - 2\frac{u_j}{\alpha} + \frac{|\lambda_j|^2}{\alpha^2} = 1 - \frac{|\lambda_j|^2}{\alpha^2} \left(\frac{2\alpha u_j}{|\lambda_j|^2} - 1\right) < 1$$

if $2\alpha > \max_j(|\lambda_j|^2/u_j)$ and the R method converges. We next show that $\rho(\mathbf{G}(\alpha)) > \rho(\mathbf{G}(\alpha^*))$ if $\alpha \neq \alpha^*$. Indeed, if $\alpha > \alpha^*$ then

$$\rho(\mathbf{G}(\alpha)) \geq \mu_n(\alpha) = 1 - \lambda_n/\alpha > 1 - \lambda_n/\alpha^* = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} = \rho(\mathbf{G}(\alpha^*)).$$

Next, if $\alpha < \alpha^*$ then

$$-\rho(\mathbf{G}(\alpha)) \leq \mu_1(\alpha) = 1 - \lambda_1/\alpha < 1 - \lambda_1/\alpha^* = -\frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} = -\rho(\mathbf{G}(\alpha^*)),$$

and again $\rho(\mathbf{G}(\alpha)) > \rho(\mathbf{G}(\alpha^*))$. \square

9.4 Convergence of the Classical Methods for the Discrete Poisson Matrix

The matrix \mathbf{A} in (4.7) is symmetric positive definite (cf. Theorem 4.13). We show in Theorem 9.31 that the SOR method converges for all $0 < \omega < 2$ if \mathbf{A} is symmetric positive definite. So the GS method converges, but the J method does not converge for all symmetric positive definite matrices.

Exercise 9.25 (The GS method converges, but not the J method)

Show (by finding its eigenvalues) that the matrix

$$\begin{bmatrix} 1 & a & a \\ a & 1 & a \\ a & a & 1 \end{bmatrix}$$

is symmetric positive definite for $-1/2 < a < 1$, but that the J method does not converge for $1/2 < a < 1$.

For the discrete Poisson problem we can determine explicitly the eigenvalues of the iteration matrices and thus not only show convergence, but also estimate the number of iterations necessary to achieve a given accuracy.

Recall that by (4.22) the eigenvalues $\lambda_{j,k}$ of \mathbf{A} given by (4.7) are

$$\lambda_{j,k} = 4 - 2\cos(j\pi h) - 2\cos(k\pi h), \quad j, k = 1, \dots, m, \quad h = 1/(m+1).$$

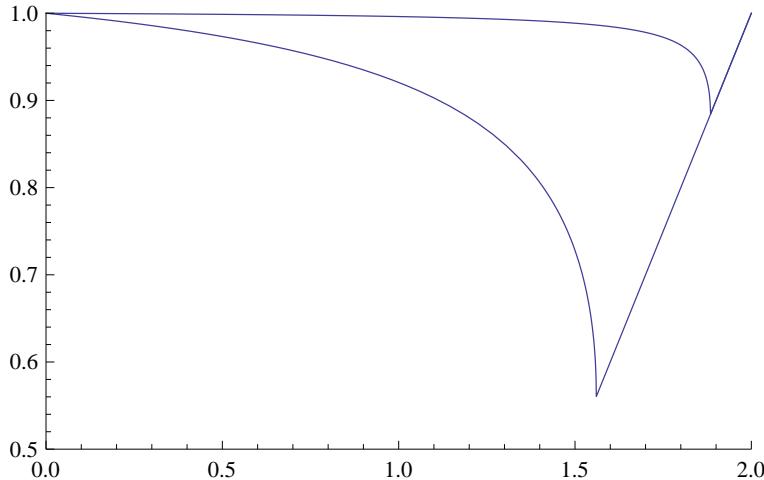


Figure 9.2: $\rho(\mathbf{G}_\omega)$ with $\omega \in [0, 2]$ for $n = 100$, (lower curve) and $n = 2500$ (upper curve).

Consider first Jacobi's method. The matrix $\mathbf{G}_J = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{I} - \mathbf{A}/4$ has eigenvalues

$$\mu_{j,k} = 1 - \frac{1}{4}\lambda_{j,k} = \frac{1}{2}\cos(j\pi h) + \frac{1}{2}\cos(k\pi h), \quad j, k = 1, \dots, m. \quad (9.27)$$

It follows that $\rho(\mathbf{G}_J) = \cos(\pi h) < 1$ and the J method converges for all starting values and all right hand sides.

For the SOR method it is possible to explicitly determine $\rho(\mathbf{G}_\omega)$ for any $\omega \in (0, 2)$. The following result will be shown in Section 9.5.

Theorem 9.26 (The spectral radius of SOR matrix)

Consider the SOR iteration (9.7), whith the natural ordering. The spectral radius of \mathbf{G}_ω is

$$\rho(\mathbf{G}_\omega) = \begin{cases} \frac{1}{4} \left(\omega\beta + \sqrt{(\omega\beta)^2 - 4(\omega - 1)} \right)^2, & \text{for } 0 < \omega \leq \omega^*, \\ \omega - 1, & \text{for } \omega^* < \omega < 2, \end{cases} \quad (9.28)$$

where $\beta := \rho(\mathbf{G}_J)$ and

$$\omega^* := \frac{2}{1 + \sqrt{1 - \beta^2}} > 1. \quad (9.29)$$

Moreover,

$$\rho(\mathbf{G}_\omega) > \rho(\mathbf{G}_{\omega^*}) \text{ for } \omega \in (0, 2) \setminus \{\omega^*\}. \quad (9.30)$$

	n=100	n=2500	k_{100}	k_{2500}
J	0.959493	0.998103	446	9703
GS	0.920627	0.99621	223	4852
SOR	0.56039	0.88402	32	150

Table 9.2: Spectral radia for \mathbf{G}_J , \mathbf{G}_1 , \mathbf{G}_{ω^*} and the smallest integer k_n such that $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$.

A plot of $\rho(\mathbf{G}_\omega)$ as a function of $\omega \in (0, 2)$ is shown in Figure 9.2 for $n = 100$ (lower curve) and $n = 2500$ (upper curve). As ω increases the spectral radius of \mathbf{G}_ω decreases monotonically to the minimum ω^* . Then it increases linearly to the value one for $\omega = 2$. We call ω^* the **optimal relaxation parameter**.

For the discrete Poisson problem we have $\beta = \cos(\pi h)$ and it follows from (9.28), (9.29) that

$$\omega^* = \frac{2}{1 + \sin(\pi h)}, \quad \rho(\mathbf{G}_{\omega^*}) = \omega^* - 1 = \frac{1 - \sin(\pi h)}{1 + \sin(\pi h)}, \quad h = \frac{1}{m+1}. \quad (9.31)$$

Letting $\omega = 1$ in (9.28) we find $\rho(\mathbf{G}_1) = \beta^2 = \rho(\mathbf{G}_J)^2 = \cos^2(\pi h)$. Thus, for the discrete Poisson problem the J method needs twice as many iterations as the GS method for a given accuracy.

The values of $\rho(\mathbf{G}_J)$, $\rho(\mathbf{G}_1)$, and $\rho(\mathbf{G}_{\omega^*}) = \omega^* - 1$ are shown in Table 9.2 for $n = 100$ and $n = 2500$. We also show the smallest integer k_n such that $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$. This is an estimate for the number of iteration needed to obtain an accuracy of 10^{-8} . These values are comparable to the exact values given in Table 9.1.

9.4.1 Number of Iterations

Let s be a positive integer. We can now estimate the number of iterations k_n to obtain $\rho(\mathbf{G})^{k_n} < 10^{-s}$ for the J, GS and SOR method with optimal ω . We use Lemma 9.21 that provided the estimate

$$\tilde{k}_n = \frac{\log(10)s}{\eta}, \quad \rho(\mathbf{G}) = 1 - \eta.$$

Note that $h = 1/(m+1) \approx n^{-1/2}$. The estimates we derive agree with those we found numerically in Section 9.1.1.

- J: $\rho(\mathbf{G}_J) = \cos(\pi h) = 1 - \eta$, $\eta = 1 - \cos(\pi h) = \frac{1}{2}\pi^2 h^2 + O(h^4) = \frac{\pi^2}{2}/n + O(n^{-2})$. Thus

$$\tilde{k}_n = \frac{2 \log(10)s}{\pi^2} n + O(n^{-1}) = O(n).$$

- GS: $\rho(\mathbf{G}_1) = \cos^2(\pi h) = 1 - \eta$, $\eta = 1 - \cos^2(\pi h) = \sin^2 \pi h = \pi^2 h^2 + O(h^4) = \pi^2/n + O(n^{-2})$. Thus

$$\tilde{k}_n = \frac{\log(10)s}{\pi^2} n + O(n^{-1}) = O(n).$$

- SOR: $\rho(\mathbf{G}_{\omega^*}) = \frac{1-\sin(\pi h)}{1+\sin(\pi h)} = 1 - 2\pi h + O(h^2)$. Thus

$$\tilde{k}_n = \frac{\log(10)s}{\pi^2} \sqrt{n} + O(n^{-1/2}) = O(\sqrt{n}).$$

Exercise 9.28 (Convergence example for fix point iteration)

Consider for $a \in \mathbb{C}$

$$\mathbf{x} := \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 & a \\ a & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1-a \\ 1-a \end{bmatrix} =: \mathbf{G}\mathbf{x} + \mathbf{c}.$$

Starting with $\mathbf{x}_0 = \mathbf{0}$ show by induction

$$\mathbf{x}_k(1) = \mathbf{x}_k(2) = 1 - a^k, \quad k \geq 0,$$

and conclude that the iteration converges to the fixed-point $\mathbf{x} = [1, 1]^T$ for $|a| < 1$ and diverges for $|a| > 1$. Show that $\rho(\mathbf{G}) = 1 - \eta$ with $\eta = 1 - |a|$. Compute the estimate (9.22) for the rate of convergence for $a = 0.9$ and $s = 16$ and compare with the true number of iterations determined from $|a|^k \leq 10^{-16}$.

9.5 Convergence Analysis for SOR

The iteration matrix \mathbf{G}_ω for the SOR method can be written in two alternative forms that are both useful for the analysis.

Lemma 9.29 (SOR iteration matrix)

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{D} = \text{diag}(a_{11}, \dots, a_{nn})$ are both nonsingular. Then

$$\mathbf{G}_\omega = \mathbf{I} - (\omega^{-1} \mathbf{D} - \mathbf{A}_L)^{-1} \mathbf{A} = (\mathbf{I} - \omega \mathbf{L})^{-1} (\omega \mathbf{R} + (1 - \omega) \mathbf{I}), \quad (9.32)$$

where \mathbf{A}_L and \mathbf{A}_R are given by (9.13) and

$$\mathbf{L} := \mathbf{D}^{-1} \mathbf{A}_L, \quad \mathbf{R} := \mathbf{D}^{-1} \mathbf{A}_R, \quad \text{so that } \mathbf{D}^{-1} \mathbf{A} = \mathbf{I} - \mathbf{L} - \mathbf{R}. \quad (9.33)$$

Proof. For the first form see (9.9) and Proposition 9.3. Solving the SOR part of (9.15) for \mathbf{x}_{k+1} gives

$$\mathbf{x}_{k+1} = \omega(\mathbf{L}\mathbf{x}_{k+1} + \mathbf{R}\mathbf{x}_k + \mathbf{D}^{-1}\mathbf{b}) + (1 - \omega)\mathbf{x}_k,$$

or

$$(\mathbf{I} - \omega \mathbf{L})\mathbf{x}_{k+1} = (\omega \mathbf{R} + (1 - \omega)\mathbf{I})\mathbf{x}_k + \omega \mathbf{D}^{-1}\mathbf{b}.$$

Solving for \mathbf{x}_{k+1} we obtain $\mathbf{x}_{k+1} = \mathbf{G}_\omega \mathbf{x}_k + \mathbf{c}$, where \mathbf{G}_ω is given by the second form in (9.32). \square

We start with the following convergence result.

Theorem 9.30 (Necessary condition for convergence of SOR)

The SOR method diverges if ω is not in the interval $(0, 2)$.

Proof. Recall that the determinant of a product equals the product of determinants and that the determinant of a triangular matrix equals the product of the diagonal elements. From (9.32) we obtain

$$\det(\mathbf{G}_\omega) = \det((\mathbf{I} - \omega \mathbf{L})^{-1}) \det(\omega \mathbf{R} + (1 - \omega)\mathbf{I}).$$

Since $\mathbf{I} - \omega \mathbf{L}$ is lower triangular with ones on the diagonal it follows from Lemma 1.9 that the first determinant equals one. The matrix $\omega \mathbf{R} + (1 - \omega)\mathbf{I}$ is upper triangular with $1 - \omega$ on the diagonal and therefore its determinant equals $(1 - \omega)^n$. It follows that $\det(\mathbf{G}_\omega) = (1 - \omega)^n$.

Since the determinant of a matrix equals the product of its eigenvalues we must have $|\lambda| \geq |1 - \omega|$ for at least one eigenvalue λ of \mathbf{G}_ω . We conclude that $\rho(\mathbf{G}_\omega) \geq |1 - \omega|$. But then $\rho(\mathbf{G}_\omega) \geq 1$ if ω is not in the interval $(0, 2)$ and by Theorem 9.15 SOR diverges. \square

We next show that SOR converges for all $\omega \in (0, 2)$ if \mathbf{A} is symmetric positive definite.

Theorem 9.31 (SOR on positive definite matrix)

SOR converges for a symmetric positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ if and only if $0 < \omega < 2$. In particular, Gauss-Seidel's method converges for a symmetric positive definite matrix.

Proof. By Theorem 9.30 convergence implies $0 < \omega < 2$. Suppose $0 < \omega < 2$. The eigenpair equation $\mathbf{G}_\omega \mathbf{x} = \lambda \mathbf{x}$ can be written $\mathbf{x} - (\omega^{-1} \mathbf{D} - \mathbf{A}_L)^{-1} \mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ or

$$\mathbf{A} \mathbf{x} = (\omega^{-1} \mathbf{D} - \mathbf{A}_L) \mathbf{y}, \quad \mathbf{y} := (1 - \lambda) \mathbf{x}. \quad (9.34)$$

Since $\mathbf{A} = -\mathbf{A}_L + \mathbf{D} - \mathbf{A}_R$ we find

$$(\omega^{-1} \mathbf{D} - \mathbf{D} + \mathbf{A}_R) \mathbf{y} = (\omega^{-1} \mathbf{D} - \mathbf{A}_L - \mathbf{A}) \mathbf{y} \stackrel{(9.34)}{=} \mathbf{A} \mathbf{x} - \mathbf{A} \mathbf{y} = \lambda \mathbf{A} \mathbf{x},$$

so that by taking inner products and replacing \mathbf{A}_R^* by \mathbf{A}_L

$$\begin{aligned} \langle \mathbf{y}, \lambda \mathbf{A} \mathbf{x} \rangle &= \langle \mathbf{y}, (\omega^{-1} \mathbf{D} - \mathbf{D} + \mathbf{A}_R) \mathbf{y} \rangle = \langle (\omega^{-1} \mathbf{D} - \mathbf{D} + \mathbf{A}_R^*) \mathbf{y}, \mathbf{y} \rangle \\ &= \langle (\omega^{-1} \mathbf{D} - \mathbf{D} + \mathbf{A}_L) \mathbf{y}, \mathbf{y} \rangle. \end{aligned} \quad (9.35)$$

Taking inner product with \mathbf{y} in (9.34) and adding to (9.35) we obtain

$$\begin{aligned}\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \lambda \mathbf{A}\mathbf{x} \rangle &= \langle (\omega^{-1}\mathbf{D} - \mathbf{A}_L)\mathbf{y}, \mathbf{y} \rangle + \langle (\omega^{-1}\mathbf{D} - \mathbf{D} + \mathbf{A}_L)\mathbf{y}, \mathbf{y} \rangle \\ &= (2\omega^{-1} - 1)\langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle = (2\omega^{-1} - 1)(1 - \lambda)(1 - \bar{\lambda})\langle \mathbf{D}\mathbf{x}, \mathbf{x} \rangle \\ &= (2\omega^{-1} - 1)|1 - \lambda|^2\langle \mathbf{D}\mathbf{x}, \mathbf{x} \rangle.\end{aligned}$$

On the other hand, since \mathbf{A} is symmetric

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \lambda \mathbf{A}\mathbf{x} \rangle = (1 - \bar{\lambda})\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle + (1 - \lambda)\bar{\lambda}\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle = (1 - |\lambda|^2)\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle.$$

Thus,

$$(2\omega^{-1} - 1)|1 - \lambda|^2\langle \mathbf{D}\mathbf{x}, \mathbf{x} \rangle = (1 - |\lambda|^2)\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle. \quad (9.36)$$

Since \mathbf{A} is symmetric positive definite we observe that also \mathbf{D} is symmetric positive definite. Furthermore we cannot have $\lambda = 1$ for then $\mathbf{y} = \mathbf{0}$ which by (9.34) implies that \mathbf{A} is singular. Since $0 < \omega < 2$ implies $\omega^{-1} > 1/2$ the left side of (9.36) is positive and hence the right hand side is positive as well. We conclude that $|\lambda| < 1$. But then $\rho(\mathbf{G}_\omega) < 1$ and SOR converges. \square

9.6 The Optimal SOR Parameter ω

The following analysis holds both for the discrete Poisson matrix and the averaging matrix given by (4.9). A more general theory is presented in [31]. Consider first how the eigenvalues of \mathbf{G}_J and \mathbf{G}_ω are related.

Theorem 9.32 (The optimal ω)

Consider for $a, d \in \mathbb{R}$ the SOR method applied to the matrix (4.9), where we use the natural ordering. Moreover, assume $\omega \in (0, 2)$.

1. If $\lambda \neq 0$ is an eigenvalue of \mathbf{G}_ω then

$$\mu := \frac{\lambda + \omega - 1}{\omega\lambda^{1/2}} \quad (9.37)$$

is an eigenvalue of \mathbf{G}_J .

2. If μ is an eigenvalue of \mathbf{G}_J and λ satisfies the equation

$$\mu\omega\lambda^{1/2} = \lambda + \omega - 1 \quad (9.38)$$

then λ is an eigenvalue of \mathbf{G}_ω .

Proof. For simplicity of notation we assume that $a = -1$ and $d = 2$. The component equations in this proof hold for $i, j = 1, \dots, m$. Suppose (λ, \mathbf{w}) is an eigenpair for \mathbf{G}_ω . By (9.32) $(\mathbf{I} - \omega \mathbf{L})^{-1}(\omega \mathbf{R} + (1 - \omega)\mathbf{I})\mathbf{w} = \lambda \mathbf{w}$ or

$$(\omega \mathbf{R} + \lambda \omega \mathbf{L})\mathbf{w} = (\lambda + \omega - 1)\mathbf{w}. \quad (9.39)$$

Let $\mathbf{w} = \text{vec}(\mathbf{W})$, where $\mathbf{W} \in \mathbb{C}^{m \times m}$. Then (9.39) can be written

$$\frac{\omega}{4}(\lambda w_{i-1,j} + \lambda w_{i,j-1} + w_{i+1,j} + w_{i,j+1}) = (\lambda + \omega - 1)w_{i,j}, \quad (9.40)$$

where $w_{i,j} = 0$ if $i \in \{0, m+1\}$ or $j \in \{0, m+1\}$. We claim that (μ, \mathbf{v}) is an eigenpair for \mathbf{G}_J , where μ is given by (9.37) and $\mathbf{v} = \text{vec}(\mathbf{V})$ with

$$v_{i,j} := \lambda^{-(i+j)/2} w_{i,j}. \quad (9.41)$$

Indeed, replacing $w_{i,j}$ by $\lambda^{-(i+j)/2} v_{i,j}$ in (9.40) and cancelling the common factor $\lambda^{-(i+j)/2}$ we obtain

$$\frac{\omega}{4}(v_{i-1,j} + v_{i,j-1} + v_{i+1,j} + v_{i,j+1}) = \lambda^{-1/2}(\lambda + \omega - 1)v_{i,j}.$$

But then

$$\mathbf{G}_J \mathbf{v} = (\mathbf{L} + \mathbf{R})\mathbf{v} = \frac{\lambda + \omega - 1}{\omega \lambda^{1/2}} = \mu \mathbf{v}.$$

For the converse let (μ, \mathbf{v}) be an eigenpair for \mathbf{G}_J and let as before $\mathbf{v} = \text{vec}(\mathbf{V})$, $\mathbf{W} = \text{vec}(\mathbf{W})$ with $v_{i,j} = \lambda^{-(i+j)/2} w_{i,j}$. The equation $\mathbf{G}_J \mathbf{v} = \mu \mathbf{v}$ can be written

$$\frac{1}{4}(v_{i-1,j} + v_{i,j-1} + v_{i+1,j} + v_{i,j+1}) = \mu v_{i,j}.$$

Let λ be a solution of (9.38). Replacing $v_{i,j}$ by $\lambda^{-(i+j)/2} w_{i,j}$ and canceling $\lambda^{-(i+j)/2}$ we obtain

$$\frac{1}{4}(\lambda^{1/2} w_{i-1,j} + \lambda^{1/2} w_{i,j-1} + \lambda^{-1/2} w_{i+1,j} + \lambda^{-1/2} w_{i,j+1}) = \mu w_{i,j},$$

or, multiplying by $\omega \lambda^{1/2}$

$$\frac{\omega}{4}(\lambda w_{i-1,j} + \lambda w_{i,j-1} + w_{i+1,j} + w_{i,j+1}) = \omega \mu \lambda^{1/2} w_{i,j},$$

Thus, if $\omega \mu^{1/2} = \lambda + \omega - 1$ then by (9.40) (λ, \mathbf{w}) is an eigenpair for \mathbf{G}_ω . \square

Proof of Theorem 9.26

By (4.22) the eigenvalues of $\mathbf{G}_J = \mathbf{I} - \mathbf{A}/(2d)$ are given by

$$\mu_{j,k} = -a(\cos(j\pi h) + \cos(k\pi h))/(2d), \quad j, k = 1, \dots, m.$$

Thus the eigenvalues are real and if μ is an eigenvalue then $-\mu$ is also an eigenvalue. Thus it is enough to consider positive eigenvalues μ . For simplicity of notation let again $a = -1$ and $d = 2$. Solving (9.38) for λ gives

$$\lambda(\mu) := \frac{1}{4} \left(\omega\mu \pm \sqrt{(\omega\mu)^2 - 4(\omega - 1)} \right)^2. \quad (9.42)$$

Both roots $\lambda(\mu)$ are eigenvalues of \mathbf{G}_ω . The discriminant

$$d(\omega) := (\omega\mu)^2 - 4(\omega - 1).$$

is strictly decreasing on $(0, 2)$ since

$$d'(\omega) = 2(\omega\mu^2 - 2) < 2(\omega - 2) < 0.$$

Moreover $d(0) = 4 > 0$ and $d(2) = 4\mu^2 - 4 < 0$. As a function of ω , $\lambda(\mu)$ changes from real to complex at

$$\omega = \tilde{\omega}(\mu) := \frac{2}{1 + \sqrt{1 - \mu^2}}. \quad (9.43)$$

In the complex case we find

$$|\lambda(\mu)| = \frac{1}{4} \left((\omega\mu)^2 + 4(\omega - 1) - (\omega\mu)^2 \right) = \omega - 1, \quad \tilde{\omega}(\mu) < \omega < 2.$$

In the real case both roots of (9.42) are positive and the larger one is

$$\lambda(\mu) = \frac{1}{4} \left(\omega\mu + \sqrt{(\omega\mu)^2 - 4(\omega - 1)} \right)^2, \quad 0 < \omega \leq \tilde{\omega}(\mu). \quad (9.44)$$

Both $\lambda(\mu)$ and $\tilde{\omega}(\mu)$ are strictly increasing as functions of μ . It follows that $|\lambda(\mu)|$ is maximized for $\mu = \rho(\mathbf{G}_J) =: \beta$ and for this value of μ we obtain (9.28) for $0 < \omega \leq \tilde{\omega}(\beta) = \omega^*$.

Evidently $\rho(\mathbf{G}_\omega) = \omega - 1$ is strictly increasing in $\omega^* < \omega < 2$. Equation (9.30) will follow if we can show that $\rho(\mathbf{G}_\omega)$ is strictly decreasing in $0 < \omega < \omega^*$. By differentiation

$$\frac{d}{d\omega} \left(\omega\beta + \sqrt{(\omega\beta)^2 - 4(\omega - 1)} \right) = \frac{\beta\sqrt{(\omega\beta)^2 - 4(\omega - 1)} + \omega\beta^2 - 2}{\sqrt{(\omega\beta)^2 - 4(\omega - 1)}}.$$

Since $\beta^2(\omega^2\beta^2 - 4\omega + 4) < (2 - \omega\beta^2)^2$ the numerator is negative and the strict decrease of $\rho(\mathbf{G}_\omega)$ in $0 < \omega < \omega^*$ follows.

9.7 Review Questions

9.7.1 Consider a matrix $A \in \mathbb{C}^{n \times n}$ with nonzero diagonal elements.

- Define the J and GS method in component form,
- Do they always converge?
- Give a necessary and sufficient condition that $A^n \rightarrow \mathbf{0}$.
- Is there a matrix norm $\| \cdot \|$ consistent on $\mathbb{C}^{n \times n}$ such that $\|A\| < \rho(A)$?

9.7.2 What is a Neumann series? when does it converge?

9.7.3 How do we define convergence of a fixed point iteration $x_{k+1} = Gx_k + c$?
When does it converge?

9.7.4 Define Richardson's method.

Chapter 10

The Conjugate Gradient Method

The **conjugate gradient method** was introduced in [10] as a direct method for solving a symmetric positive definite linear system, or equivalently minimizing a quadratic function. Today its main use is as an iterative method for solving large sparse linear systems and we focus on this here. On a test problem we show that it performs as well as the SOR method with optimal acceleration parameter, and we do not have to estimate any such parameter. We also consider the **preconditioned conjugate gradient method** which can be used to speed up convergence of the conjugate gradient method.

The conjugate gradient method can also be used for minimization and we first discuss a related minimization method known as steepest descent.

Throughout this chapter $\mathbf{A} \in \mathbb{R}^{n \times n}$ will be a symmetric positive definite matrix. Thus, $\mathbf{A}^T = \mathbf{A}$ and $\mathbf{y}^T \mathbf{A} \mathbf{y} > 0$ for all nonzero $\mathbf{y} \in \mathbb{R}^n$. We recall that \mathbf{A} has positive eigenvalues and that the spectral (2-norm) condition number of \mathbf{A} is given by $\kappa := \frac{M}{m}$, where M and m are the largest and smallest eigenvalue of \mathbf{A} .

10.1 Quadratic Minimization

We start by discussing some aspect of quadratic minimization and its relation to solving linear systems.

Consider for $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$, the quadratic function $Q : \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$Q(\mathbf{y}) := \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}. \quad (10.1)$$

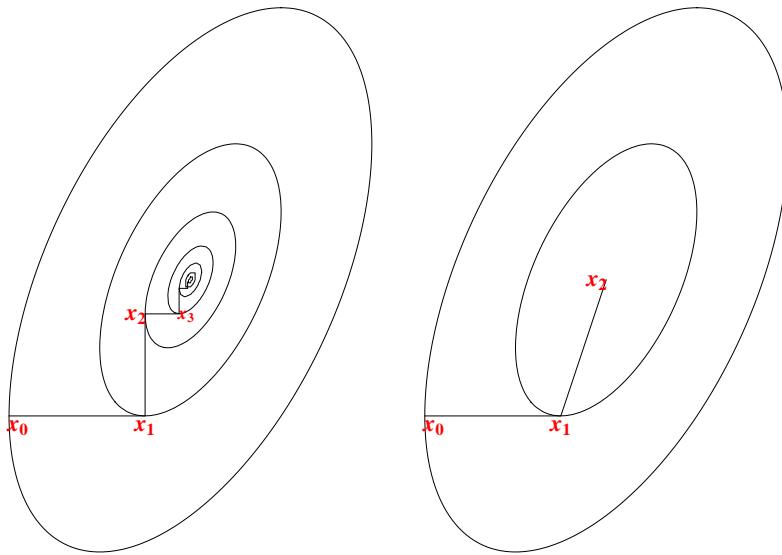


Figure 10.1: Level curves for $Q(x, y)$ given by (10.2). Also shown is a steepest descent iteration (left) and a conjugate gradient iteration (right) to find the minimum of Q . (cf Examples 10.3, 10.11)

As an example, some level curves of

$$Q(x, y) := \frac{1}{2} \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = x^2 - xy + y^2 \quad (10.2)$$

are shown in Figure 10.1. The level curves are ellipses and the graph of Q is a paraboloid (cf. Exercise 10.1).

Exercise 10.1 (Paraboloid)

Let $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T$ be the spectral decomposition of \mathbf{A} , i.e., \mathbf{U} is orthonormal and $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is diagonal. Define new variables $\mathbf{v} = [v_1, \dots, v_n]^T := \mathbf{U}^T \mathbf{y}$, and set $\mathbf{c} := \mathbf{U}^T \mathbf{b} = [c_1, \dots, c_n]^T$. Show that

$$Q(\mathbf{y}) = \frac{1}{2} \sum_{j=1}^n \lambda_j v_j^2 - \sum_{j=1}^n c_j v_j.$$

Lemma 10.2 (Quadratic function)

A vector $\mathbf{x} \in \mathbb{R}^n$ minimizes Q if and only if $\mathbf{Ax} = \mathbf{b}$. Moreover, the residual

$\mathbf{r}(\mathbf{y}) = \mathbf{b} - \mathbf{A}\mathbf{y}$ at any $\mathbf{y} \in \mathbb{R}^n$ is equal to the negative gradient, i.e., $\mathbf{r}(\mathbf{y}) = -\nabla Q(\mathbf{y})$, where $\nabla := \left[\frac{\partial}{\partial y_1}, \dots, \frac{\partial}{\partial y_n} \right]^T$.

Proof. For any $\mathbf{y}, \mathbf{h} \in \mathbb{R}^n$ and $\varepsilon \in \mathbb{R}$

$$Q(\mathbf{y} + \varepsilon\mathbf{h}) = Q(\mathbf{y}) - \varepsilon\mathbf{h}^T \mathbf{r}(\mathbf{y}) + \frac{1}{2}\varepsilon^2 \mathbf{h}^T \mathbf{A}\mathbf{h}, \text{ where } \mathbf{r}(\mathbf{y}) := \mathbf{b} - \mathbf{A}\mathbf{y}. \quad (10.3)$$

If $\mathbf{y} = \mathbf{x}$, $\varepsilon = 1$, and $\mathbf{A}\mathbf{x} = \mathbf{b}$ then (10.3) simplifies to $Q(\mathbf{x} + \mathbf{h}) = Q(\mathbf{x}) + \frac{1}{2}\mathbf{h}^T \mathbf{A}\mathbf{h}$, and since \mathbf{A} is symmetric positive definite $Q(\mathbf{x} + \mathbf{h}) > Q(\mathbf{x})$ for all nonzero $\mathbf{h} \in \mathbb{R}^n$. It follows that \mathbf{x} is the unique minimum of Q . Conversely, if $\mathbf{A}\mathbf{x} \neq \mathbf{b}$, say $\mathbf{r}(\mathbf{x}) > 0$, and \mathbf{h} is nonzero then $Q(\mathbf{x} + \varepsilon\mathbf{h}) < Q(\mathbf{x})$ for $\varepsilon > 0$ sufficiently small. Thus \mathbf{x} does not minimize Q . By (10.3) for $\mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} \frac{\partial}{\partial y_i} Q(\mathbf{y}) &:= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (Q(\mathbf{y} + \varepsilon\mathbf{e}_i) - Q(\mathbf{y})) \\ &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \left(-\varepsilon \mathbf{e}_i^T \mathbf{r}(\mathbf{y}) + \frac{1}{2}\varepsilon^2 \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i \right) = -\mathbf{e}_i^T \mathbf{r}(\mathbf{y}), \quad i = 1, \dots, n, \end{aligned}$$

showing that $\mathbf{r}(\mathbf{y}) = -\nabla Q(\mathbf{y})$. \square

A general class of minimization algorithms for Q is given as follows:

1. Choose $\mathbf{x}_0 \in \mathbb{R}^n$.
2. For $k = 0, 1, 2, \dots$

Choose a “search direction” \mathbf{p}_k ,
 Choose a “step length” α_k ,
 Compute $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$.

We would like to generate a sequence $\{\mathbf{x}_k\}$ that converges quickly to the minimum \mathbf{x} of Q .

For a fixed direction \mathbf{p}_k we say that α_k is the **optimal the step length** if $Q(\mathbf{x}_{k+1})$ is as small as possible, i.e.

$$Q(\mathbf{x}_{k+1}) = Q(\mathbf{x}_k + \alpha_k \mathbf{p}_k) = \min_{\alpha \in \mathbb{R}} Q(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

By (10.3) we have $Q(\mathbf{x}_k + \alpha \mathbf{p}_k) = Q(\mathbf{x}_k) - \alpha \mathbf{p}_k^T \mathbf{r}_k + \frac{1}{2}\alpha^2 \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$, where $\mathbf{r}_k := \mathbf{b} - \mathbf{A}\mathbf{x}_k$. Since $\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k \geq 0$ we find a minimum α_k by solving $\frac{\partial}{\partial \alpha} Q(\mathbf{x}_k + \alpha \mathbf{p}_k) = 0$. It follows that the optimal α_k is uniquely given by

$$\alpha_k := \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}. \quad (10.5)$$

10.2 Steepest Descent

In the method of **Steepest Descent**, also known as the **Gradient Method** we choose $\mathbf{p}_k = \mathbf{r}_k$, the negative gradient, and the optimal α_k . Starting from \mathbf{x}_0 and $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ we compute for $k = 0, 1, 2 \dots$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \left(\frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k} \right) \mathbf{r}_k.$$

Computationally a step in the steepest descent iteration can proceed as follows

$$\begin{aligned} \mathbf{t}_k &= \mathbf{A}\mathbf{r}_k, \\ \alpha_k &= (\mathbf{r}_k^T \mathbf{r}_k) / (\mathbf{r}_k^T \mathbf{t}_k), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{r}_k, \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{t}_k. \end{aligned} \tag{10.6}$$

Here, and in general, the following updating of the residual is used:

$$\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{A}\mathbf{x}_{k+1} = \mathbf{b} - \mathbf{A}(\mathbf{x}_k + \alpha_k \mathbf{p}_k) = \mathbf{r}_k - \alpha_k \mathbf{t}_k, \quad \mathbf{t}_k := \mathbf{A}\mathbf{p}_k. \tag{10.7}$$

Example 10.3 (Steepest descent iteration)

Suppose $Q(x, y)$ is given by (10.2). Starting with $\mathbf{x}_0 = [-1, -1/2]^T$ and $\mathbf{r}_0 = -\mathbf{A}\mathbf{x}_0 = [3/2, 0]^T$ we find

$$\begin{aligned} \mathbf{t}_0 &= 3 \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}, \quad \mathbf{x}_1 = -4^{-1} \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \mathbf{r}_1 = 3 * 4^{-1} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \mathbf{t}_1 &= 3 * 4^{-1} \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad \mathbf{x}_2 = -4^{-1} \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}, \quad \mathbf{r}_2 = 3 * 4^{-1} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}, \end{aligned}$$

and in general for $k \geq 1$

$$\begin{aligned} \mathbf{t}_{2k-2} &= 3 * 4^{1-k} \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}, \quad \mathbf{x}_{2k-1} = -4^{-k} \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \mathbf{r}_{2k-1} = 3 * 4^{-k} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \mathbf{t}_{2k-1} &= 3 * 4^{-k} \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad \mathbf{x}_{2k} = -4^{-k} \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}, \quad \mathbf{r}_{2k} = 3 * 4^{-k} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}, \end{aligned}$$

and that $\alpha_k = 1/2$ for all k . See the left part of Figure 10.1. Since $\|\mathbf{r}_{j+1}\|_2 / \|\mathbf{r}_j\| = 1/2$ for all j the convergence is not too impressive.

Exercise 10.4 (Steepest descent iteration)

Verify the numbers in Example 10.3.

10.2.1 Convergence Analysis for Steepest Descent

We want to show that the method always converges and give an upper bound for the rate of convergence. For this the following inequality will be used.

Theorem 10.5 (Kantorovich inequality)

For any symmetric positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$

$$1 \leq \frac{(\mathbf{y}^T \mathbf{A} \mathbf{y})(\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y})}{(\mathbf{y}^T \mathbf{y})^2} \leq \frac{(M+m)^2}{4Mm} \quad \mathbf{y} \neq \mathbf{0}, \quad \mathbf{y} \in \mathbb{R}^n, \quad (10.8)$$

where M and m are the largest and smallest eigenvalue of \mathbf{A} , respectively.

Proof. For $j = 1, \dots, n$ let $(\lambda_j, \mathbf{u}_j)$ be orthonormal eigenpairs of \mathbf{A} and $\mathbf{y} \in \mathbb{R}^n$. By Theorem 0.66 $(\lambda_j^{-1}, \mathbf{u}_j)$ are eigenpairs for \mathbf{A}^{-1} . Let $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{u}_j$ be the eigenvector expansion of \mathbf{y} . By orthonormality, (cf. (6.7))

$$a := \frac{\mathbf{y}^T \mathbf{A} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \sum_{i=1}^n t_i \lambda_i, \quad b := \frac{\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \sum_{i=1}^n \frac{t_i}{\lambda_i}, \quad (10.9)$$

where

$$t_i = \frac{c_i^2}{\sum_{j=1}^n c_j^2} \geq 0, \quad i = 1, \dots, n \text{ and } \sum_{i=1}^n t_i = 1. \quad (10.10)$$

Thus a and b are **convex combinations** of the eigenvalues of \mathbf{A} and \mathbf{A}^{-1} , respectively. Let c be a positive constant to be chosen later. By the geometric/arithmetic mean inequality (8.26) and (10.9)

$$\sqrt{ab} = \sqrt{(ac)(b/c)} \leq (ac + b/c)/2 = \frac{1}{2} \sum_{i=1}^n t_i (\lambda_i c + 1/(\lambda_i c)) = \frac{1}{2} \sum_{i=1}^n t_i f(\lambda_i c),$$

where $f : [mc, Mc] \rightarrow \mathbb{R}$ is given by $f(x) := x + 1/x$. By (10.10)

$$\sqrt{ab} \leq \frac{1}{2} \max_{mc \leq x \leq Mc} f(x).$$

Since $f \in C^2$ and f'' is positive it follows from Lemma 8.38 that f is a convex function. But a convex function takes its maximum at one of the endpoints of the range (cf. Exercise 10.6) and we obtain

$$\sqrt{ab} \leq \frac{1}{2} \max\{f(mc), f(Mc)\}. \quad (10.11)$$

Choosing $c := 1/\sqrt{mM}$ we find $f(mc) = f(Mc) = \sqrt{\frac{M}{m}} + \sqrt{\frac{m}{M}} = \frac{M+m}{\sqrt{mM}}$. By (10.11) we obtain

$$\frac{(\mathbf{y}^T \mathbf{A} \mathbf{y})(\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y})}{(\mathbf{y}^T \mathbf{y})^2} = ab \leq \frac{(M+m)^2}{4Mm},$$

the upper bound in (10.8). For the lower bound we use Cauchy-Swarz inequality as follows

$$1 = \left(\sum_{i=1}^n t_i \right)^2 = \left(\sum_{i=1}^n (t_i \lambda_i)^{1/2} (t_i / \lambda_i)^{1/2} \right)^2 \leq \left(\sum_{i=1}^n t_i \lambda_i \right) \left(\sum_{i=1}^n t_i / \lambda_i \right) = ab.$$

□

Exercise 10.6 (Maximum of a convex function)

Show that if $f : [a, b] \rightarrow \mathbb{R}$ is convex then $\max_{a \leq x \leq b} f(x) \leq \max\{f(a), f(b)\}$.

The convergence analysis for the steepest descent method is in terms of a special inner product. We define the \mathbf{A} -inner product and the corresponding \mathbf{A} -norm by

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{A} \mathbf{y}, \quad \|\mathbf{y}\|_{\mathbf{A}} := \sqrt{\langle \mathbf{y}, \mathbf{y} \rangle}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad (10.12)$$

Exercise 10.7 (The \mathbf{A} -inner product)

Show that if \mathbf{A} is symmetric positive definite then the \mathbf{A} -inner product is an inner product.

The following theorem gives upper bounds for the \mathbf{A} -norm of the error in steepest descent .

Theorem 10.8 (Error bound for steepest descent)

For the \mathbf{A} -norms of the errors in the steepest descent method (10.6) the following upper bounds hold

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k < e^{-\frac{2}{\kappa} k}, \quad , k \geq 0, \quad (10.13)$$

where $\kappa = \text{cond}_2(\mathbf{A}) := M/m$ is the spectral condition number of \mathbf{A} , and M and m are the largest and smallest eigenvalue of \mathbf{A} , respectively.

Proof. Let $\boldsymbol{\epsilon}_j := \mathbf{x} - \mathbf{x}_j$, $j = 0, 1, \dots$, where $\mathbf{A}\mathbf{x} = \mathbf{b}$. It is enough to show that

$$\frac{\|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2}{\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2} \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^2, \quad k = 0, 1, 2, \dots, \quad (10.14)$$

for then $\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}} \leq \left(\frac{\kappa - 1}{\kappa + 1} \right) \|\boldsymbol{\epsilon}_{k-1}\|_{\mathbf{A}} \leq \dots \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k \|\boldsymbol{\epsilon}_0\|_{\mathbf{A}}$. It follows from (10.6) that

$$\boldsymbol{\epsilon}_{k+1} = \boldsymbol{\epsilon}_k - \alpha_k \mathbf{r}_k, \quad \alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}.$$

We find

$$\begin{aligned}\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 &= \boldsymbol{\epsilon}_k^T \mathbf{A} \boldsymbol{\epsilon}_k = \mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k, \\ \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2 &= (\boldsymbol{\epsilon}_k - \alpha_k \mathbf{r}_k)^T \mathbf{A} (\boldsymbol{\epsilon}_k - \alpha_k \mathbf{r}_k) \\ &= \boldsymbol{\epsilon}_k^T \mathbf{A} \boldsymbol{\epsilon}_k - 2\alpha_k \mathbf{r}_k^T \mathbf{A} \boldsymbol{\epsilon}_k + \alpha_k^2 \mathbf{r}_k^T \mathbf{A} \mathbf{r}_k = \|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}.\end{aligned}$$

Combining these and using Kantorovich inequality

$$\frac{\|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2}{\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2} = 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)} \leq 1 - \frac{4mM}{(m+M)^2} = \left(\frac{M-m}{M+m}\right)^2 = \left(\frac{\kappa-1}{\kappa+1}\right)^2$$

and (10.14) is proved.

The inequality

$$\frac{x-1}{x+1} < e^{-2/x} \quad \text{for } x > 1 \quad (10.15)$$

follows from the familiar series expansion of the exponential function. Indeed, with $y = 1/x$, using $2^k/k! = 2$, $k = 1, 2$, and $2^k/k! < 2$ for $k > 2$, we find

$$e^{2/x} = e^{2y} = \sum_{k=0}^{\infty} \frac{(2y)^k}{k!} < 1 + 2 \sum_{k=1}^{\infty} y^k = \frac{1+y}{1-y} = \frac{x+1}{x-1}$$

and (10.15) follows. \square

In general the first upper bound in Theorem 10.8 is quite sharp. In fact for the iteration in Example 10.3 the first inequality in (10.13) is an equality (cf. Exercise 10.9). The second inequality is sharp when κ is large.

Exercise 10.9 (A test for the error bound)

Show that in Example 10.3 that $\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} = \left(\frac{\kappa-1}{\kappa+1}\right)^k = 2^{-k}$ for $k \geq 0$.

Theorem 10.8 implies

1. Since $\frac{\kappa-1}{\kappa+1} < 1$ the steepest descent method always converges for a symmetric positive definite matrix..
2. The upper bound for the rate of convergence depends on the condition number κ . The convergence can be slow when $\frac{\kappa-1}{\kappa+1}$ is close to one, and this happens even for a moderately ill-conditioned \mathbf{A} .

Exercise 10.10 (Orthogonality in steepest descent)

Show that $\mathbf{r}_k^T \mathbf{r}_{k+1} = 0$ in the method of steepest descent. Does this mean that all the residuals are orthogonal?

10.3 The Conjugate Gradient Method

The conjugate gradient method is of the form (10.4), where all residuals are orthogonal, while the search directions are \mathbf{A} -orthogonal. In symbols $\mathbf{r}_i^T \mathbf{r}_j = \mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0$ for $i \neq j$. Moreover, every step length is optimal.

For the derivation we choose a starting vector $\mathbf{x}_0 \in \mathbb{R}^n$. If $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0 = \mathbf{0}$ then \mathbf{x}_0 is the exact solution and we are finished, otherwise we initially make a steepest descent step with $\mathbf{p}_0 = \mathbf{r}_0$. Since $\mathbf{p}_0^T \mathbf{A} \mathbf{p}_0$ is nonzero, $\mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{p}_0$, and $\mathbf{r}_1 = \mathbf{r}_0 - \alpha_0 \mathbf{A} \mathbf{p}_0$ are well defined. The choice of α_0 ensures that \mathbf{r}_1 and \mathbf{r}_0 are orthogonal. Indeed, $\mathbf{r}_1^T \mathbf{r}_0 = (\mathbf{r}_0 - \alpha_0 \mathbf{A} \mathbf{p}_0)^T \mathbf{r}_0 = 0$ since $\mathbf{p}_0 = \mathbf{r}_0$.

For the general case we define for $j \geq 0$

$$\mathbf{p}_j := \mathbf{r}_j - \sum_{i=0}^{j-1} \left(\frac{\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} \right) \mathbf{p}_i, \quad (10.16)$$

$$\mathbf{x}_{j+1} := \mathbf{x}_j + \alpha_j \mathbf{p}_j \quad \alpha_j := \frac{\mathbf{r}_j^T \mathbf{r}_j}{\mathbf{p}_j^T \mathbf{A} \mathbf{p}_j}, \quad (10.17)$$

$$\mathbf{r}_{j+1} = \mathbf{r}_j - \alpha_j \mathbf{A} \mathbf{p}_j. \quad (10.18)$$

The last equation follows from (10.7). We show in Exercise 10.16 that the step length α_j is optimal for all j .

Using induction, suppose for some $k \geq 0$ that $\mathbf{r}_j \neq 0$, and that $\mathbf{r}_i^T \mathbf{r}_j = 0$, for $i \neq j$ and $i, j \leq k$. we want to show that $\mathbf{r}_i^T \mathbf{r}_j = 0$, for $i \neq j$ and $i, j \leq k+1$. For this it is enough to show that $\mathbf{r}_{k+1}^T \mathbf{r}_j = 0$, for $j < k$. We showed this for $k = 0$ above.

Observe that \mathbf{p}_j is computed by the Gram-Schmidt orthogonalization process applied to the linearly independent residuals $\mathbf{r}_0, \dots, \mathbf{r}_j$ using the \mathbf{A} -inner product. It follows from Theorem 0.38 that \mathbf{p}_j is nonzero and $\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0$ for $i \neq j$ and $i, j \leq k$.

Consider next the orthogonality of the residuals. We have

$$\begin{aligned} \mathbf{r}_{k+1}^T \mathbf{r}_j &\stackrel{(10.18)}{=} (\mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k)^T \mathbf{r}_j \\ &\stackrel{(10.16)}{=} \mathbf{r}_k^T \mathbf{r}_j - \alpha_k \mathbf{p}_k^T \mathbf{A} (\mathbf{p}_j + \sum_{i=0}^{j-1} \left(\frac{\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} \right) \mathbf{p}_i) \\ &\stackrel{\mathbf{A}\text{-ort}}{=} \mathbf{r}_k^T \mathbf{r}_j - \alpha_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_j, \quad j = 0, 1, \dots, k. \end{aligned}$$

That the final expression is zero follows by orthogonality and \mathbf{A} -orthogonality for $j < k$ and by the definition of α_k for $j = k$.

The expression (10.16) for \mathbf{p}_k can be greatly simplified. All terms except the last one vanish since by orthogonality of the residuals

$$\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i \stackrel{(10.18)}{=} \mathbf{r}_j^T \left(\frac{\mathbf{r}_i - \mathbf{r}_{i+1}}{\alpha_i} \right) = 0, \quad i = 0, 1, \dots, j-2.$$

For the last term with $k = j - 1$

$$\beta_k := -\frac{\mathbf{r}_{k+1}^T \mathbf{A} \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \stackrel{(10.18)}{=} \frac{\mathbf{r}_{k+1}^T (\mathbf{r}_{k+1} - \mathbf{r}_k)}{\alpha_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \stackrel{(10.17)}{=} \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}. \quad (10.19)$$

To summarize, in the **conjugate gradient method** we start with $\mathbf{p}_0 = \mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ and then generate a sequence of vectors $\{\mathbf{x}_k\}$ as follows:

For $k = 0, 1, 2, \dots$

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}, \quad (10.20)$$

$$\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (10.21)$$

$$\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k, \quad \beta_k := \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}. \quad (10.22)$$

The conjugate gradient method is also a direct method. Since $\dim \mathbb{R}^n = n$ the $n + 1$ residuals $\mathbf{r}_0, \dots, \mathbf{r}_n$ cannot all be nonzero and for orthogonal residuals we find the exact solution in at most n iterations.

For computation we use the following formulas for $k = 0, 1, 2, \dots$

$$\begin{aligned} \mathbf{t}_k &= \mathbf{A} \mathbf{p}_k, \\ \alpha_k &= (\mathbf{r}_k^T \mathbf{r}_k) / (\mathbf{r}_k^T \mathbf{t}_k), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{p}_k, \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{t}_k, \\ \beta_k &= (\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}) / (\mathbf{r}_k^T \mathbf{r}_k), \\ \mathbf{p}_{k+1} &:= \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k. \end{aligned} \quad (10.23)$$

Example 10.11 (Conjugate gradient iteration)

Consider (10.23) applied to the positive definite linear system $\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$.

Starting as in Example 10.3 with $\mathbf{x}_0 = \begin{bmatrix} -1 \\ -1/2 \end{bmatrix}$ we find $\mathbf{p}_0 = \mathbf{r}_0 = \begin{bmatrix} 3/2 \\ 0 \end{bmatrix}$ and then

$$\begin{aligned} \mathbf{t}_0 &= \begin{bmatrix} 3 \\ -3/2 \end{bmatrix}, \quad \alpha_0 = 1/2, \quad \mathbf{x}_1 = \begin{bmatrix} -1/4 \\ -1/2 \end{bmatrix}, \quad \mathbf{r}_1 = \begin{bmatrix} 0 \\ 3/4 \end{bmatrix}, \quad \beta_0 = 1/4, \quad \mathbf{p}_1 = \begin{bmatrix} 3/8 \\ 3/4 \end{bmatrix}, \\ \mathbf{t}_1 &= \begin{bmatrix} 0 \\ 9/8 \end{bmatrix}, \quad \alpha_1 = 2/3, \quad \mathbf{x}_2 = \mathbf{0}, \quad \mathbf{r}_2 = \mathbf{0}. \end{aligned}$$

Thus \mathbf{x}_2 is the exact solution, see the right part in Figure 10.1.

Exercise 10.12 (Conjugate gradient iteration, II)

Do one iteration with the conjugate gradient method when $\mathbf{x}_0 = \mathbf{0}$. (Answer: $\mathbf{x}_1 = \left(\frac{\mathbf{b}^T \mathbf{b}}{\mathbf{b}^T \mathbf{A} \mathbf{b}} \right) \mathbf{b}$.)

Exercise 10.13 (Conjugate gradient iteration, III)

Do two conjugate gradient iterations for the system

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

starting with $\mathbf{x}_0 = \mathbf{0}$.

10.3.1 The Best Approximation Property

In this section we show a best approximation property which will be used for the convergence analysis.

The iterates in the conjugate gradient method are \mathbf{A} -orthogonal projections into certain subspaces of \mathbb{R}^n called **Krylov spaces**.

They are defined by $\mathbb{W}_0 = \{\mathbf{0}\}$ and

$$\mathbb{W}_k = \text{span}(\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \mathbf{A}^2\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0), \quad k = 1, 2, 3, \dots.$$

The Krylov spaces are nested subspaces

$$\mathbb{W}_0 \subset \mathbb{W}_1 \subset \mathbb{W}_2 \subset \dots \subset \mathbb{W}_n \subset \mathbb{R}^n$$

with $\dim(\mathbb{W}_k) \leq k$ for all $k \geq 0$. Moreover, If $\mathbf{v} \in \mathbb{W}_k$ then $\mathbf{Av} \in \mathbb{W}_{k+1}$.

Lemma 10.14 (Krylov space)

We have

$$\mathbf{x}_k - \mathbf{x}_0 \in \mathbb{W}_k, \quad \mathbf{r}_k, \mathbf{p}_k \in \mathbb{W}_{k+1}, \quad k = 0, 1, \dots, \quad (10.24)$$

and

$$\mathbf{r}_k^T \mathbf{w} = \mathbf{p}_k^T \mathbf{A} \mathbf{w} = 0, \quad \mathbf{w} \in \mathbb{W}_k. \quad (10.25)$$

Proof. Since $\mathbf{p}_0 = \mathbf{r}_0$ (10.24) clearly holds for $k = 0$. Suppose it holds for some $k \geq 0$. Then by (10.23), $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k \in \mathbb{W}_{k+2}$ and $\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k \in \mathbb{W}_{k+2}$ and $\mathbf{x}_{k+1} - \mathbf{x}_0 \stackrel{(10.17)}{=} \mathbf{x}_k - \mathbf{x}_0 + \alpha_k \mathbf{p}_k \in \mathbb{W}_{k+1}$. Thus (10.24) follows by induction. Since any $\mathbf{w} \in \mathbb{W}_k$ is a linear combination of $\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}$ and also $\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}\}$, (10.25) follows. \square

Theorem 10.15 (Best approximation property)

Suppose $\mathbf{Ax} = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive definite and $\{\mathbf{x}_k\}$ is generated by the conjugate gradient method (10.23). Then

$$\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}} = \min_{\mathbf{w} \in \mathbb{W}_k} \|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}. \quad (10.26)$$

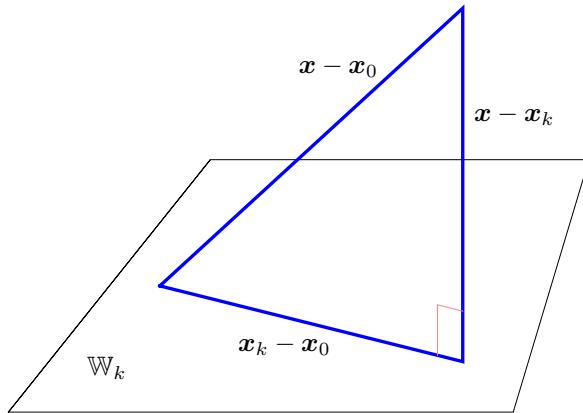


Figure 10.2: The orthogonal projection of $\mathbf{x} - \mathbf{x}_0$ into \mathbb{W}_k .

Proof. Fix k , let $\mathbf{w} \in \mathbb{W}_k$ and $\mathbf{u} := \mathbf{x}_k - \mathbf{x}_0 - \mathbf{w}$. By (10.24) $\mathbf{u} \in \mathbb{W}_k$ and then (10.25) implies that $\mathbf{r}_k^T \mathbf{u} = 0$. Since $(\mathbf{x} - \mathbf{x}_k)^T \mathbf{A} \mathbf{u} = \mathbf{r}_k^T \mathbf{u}$ we find

$$\begin{aligned}\|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}^2 &= (\mathbf{x} - \mathbf{x}_k + \mathbf{u})^T \mathbf{A} (\mathbf{x} - \mathbf{x}_k + \mathbf{u}) \\ &= (\mathbf{x} - \mathbf{x}_k) \mathbf{A} (\mathbf{x} - \mathbf{x}_k) + 2\mathbf{r}_k^T \mathbf{u} + \mathbf{u}^T \mathbf{A} \mathbf{u} \\ &= \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2 + \|\mathbf{u}\|_{\mathbf{A}}^2 \geq \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2.\end{aligned}$$

Taking square roots the result follows. \square

If $\mathbf{x}_0 = \mathbf{0}$ then (10.26) says that \mathbf{x}_k is the element in \mathbb{W}_k that is closest to the solution \mathbf{x} in the \mathbf{A} -norm. More generally, if $\mathbf{x}_0 \neq \mathbf{0}$ then $\mathbf{x} - \mathbf{x}_k = (\mathbf{x} - \mathbf{x}_0) - (\mathbf{x}_k - \mathbf{x}_0)$ and $\mathbf{x}_k - \mathbf{x}_0$ is the element in \mathbb{W}_k that is closest to $\mathbf{x} - \mathbf{x}_0$ in the \mathbf{A} -norm. This is the orthogonal projection of $\mathbf{x} - \mathbf{x}_0$ into \mathbb{W}_k , see Figure 10.2.

Exercise 10.16 (The cg step length is optimal)

Show that the step length α_k in the conjugate gradient method is optimal.

Exercise 10.17 (Starting value in cg)

Show that the conjugate gradient method (10.23) for $\mathbf{Ax} = \mathbf{b}$ starting with \mathbf{x}_0 is the same as applying the method to the system $\mathbf{Ay} = \mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}_0$ starting with $\mathbf{y}_0 = \mathbf{0}$.¹²

¹²Hint: The conjugate gradient method for $\mathbf{Ay} = \mathbf{r}_0$ can be written $\mathbf{y}_{k+1} := \mathbf{y}_k + \gamma_k \mathbf{q}_k$,

10.4 The Conjugate Gradient Algorithm

In this section we give numerical examples and discuss implementation.

The formulas in (10.23) form a basis for an algorithm.

Algorithm 10.18 (Conjugate Gradient Iteration)

The symmetric positive definite linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ is solved by the conjugate gradient method. \mathbf{x} is a starting vector for the iteration. The iteration is stopped when $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2 \leq \text{tol}$ or $k > \text{itmax}$. K is the number of iterations used.

```

1 function [x,K]=cg(A,b,x,tol,itmax)
2 r=b-A*x; p=r; rho=r'*r;
3 rho0=rho;
4 for k=0:itmax
5   if sqrt(rho/rho0)<= tol
6     K=k; return
7   end
8   t=A*p; a=rho/(p'*t);
9   x=x+a*p; r=r-a*t;
10  rhos=rho; rho=r'*r;
11  p=r+(rho/rhos)*p;
12 end
13 K=itmax+1;
```

The work involved in each iteration is

1. one matrix times vector ($\mathbf{t} = \mathbf{Ap}$),
2. two inner products ($(\mathbf{p}^T \mathbf{t}$ and $\mathbf{r}^T \mathbf{r}$),
3. three vector-plus-scalar-times-vector ($\mathbf{x} = \mathbf{x} + a\mathbf{p}$, $\mathbf{r} = \mathbf{r} - a\mathbf{t}$ and $\mathbf{p} = \mathbf{r} + (\rho/\rho_{\text{hos}})\mathbf{p}$),

The dominating part is the computation of $\mathbf{t} = \mathbf{Ap}$.

10.4.1 Numerical Example

We test the conjugate gradient method on the example used in Chapter 9. For a similar test for the steepest descent method see Exercise 10.24. The matrix is given by the Kronecker sum $\mathbf{T}_2 := \mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1$ where $\mathbf{T}_1 = \text{tridiag}_m(a, d, a)$. We recall that this matrix is symmetric positive definite if $d > 0$ and $d \geq 2|a|$. We set $h = 1/(m+1)$ and $\mathbf{f} = [1, \dots, 1]^T \in \mathbb{R}^n$.

We consider two problems.

$\gamma_k := \frac{\mathbf{s}_k^T \mathbf{s}_k}{\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k}$, $\mathbf{s}_{k+1} := \mathbf{s}_k - \gamma_k \mathbf{A} \mathbf{q}_k$, $\mathbf{q}_{k+1} := \mathbf{s}_{k+1} + \delta_k \mathbf{q}_k$, $\delta_k := \frac{\mathbf{s}_{k+1}^T \mathbf{s}_{k+1}}{\mathbf{s}_k^T \mathbf{s}_k}$. Show that $\mathbf{y}_k = \mathbf{x}_k - \mathbf{x}_0$, $\mathbf{s}_k = \mathbf{r}_k$, and $\mathbf{q}_k = \mathbf{p}_k$, for $k = 0, 1, 2, \dots$

1. $a = 1/9, d = 5/18$, the Averaging matrix.
2. $a = -1, d = 2$, the Poisson matrix.

10.4.2 Implementation Issues

Note that for our test problems \mathbf{T}_2 only has $O(5n)$ nonzero elements. Therefore, taking advantage of the sparseness of \mathbf{T}_2 we can compute \mathbf{t} in Algorithm 10.18 in $O(n)$ arithmetic operations. With such an implementation the total number of arithmetic operations in one iteration is $O(n)$. We also note that it is not necessary to store the matrix \mathbf{T}_2 .

To use the Conjugate Gradient Algorithm on the test matrix for large n it is advantageous to use a matrix equation formulation. We define matrices $\mathbf{V}, \mathbf{R}, \mathbf{P}, \mathbf{B}, \mathbf{T} \in \mathbb{R}^{m \times m}$ by $\mathbf{x} = \text{vec}(\mathbf{V})$, $\mathbf{r} = \text{vec}(\mathbf{R})$, $\mathbf{p} = \text{vec}(\mathbf{P})$, $\mathbf{t} = \text{vec}(\mathbf{T})$, and $h^2\mathbf{f} = \text{vec}(\mathbf{B})$. Then $\mathbf{T}_2\mathbf{x} = h^2\mathbf{f} \iff \mathbf{T}_1\mathbf{V} + \mathbf{V}\mathbf{T}_1 = \mathbf{B}$, and $\mathbf{t} = \mathbf{T}_2\mathbf{p} \iff \mathbf{T} = \mathbf{T}_1\mathbf{P} + \mathbf{P}\mathbf{T}_1$.

This leads to the following algorithm for testing the conjugate gradient algorithm. For simplicity we start with $\mathbf{x}_0 = \mathbf{0}$ and use an \mathbf{f} with all elements equal to one.

Algorithm 10.19 (Testing Conjugate Gradient)

$$\mathbf{A} = \text{tridiag}_m(a, d, a) \otimes \mathbf{I}_m + \mathbf{I}_m \otimes \text{tridiag}_m(a, d, a) \in \mathbb{R}^{(m^2) \times (m^2)}$$

```

1 function [V,K]=cgtest(m,a,d,tol,itmax)
2 R=ones(m)/(m+1)^2; rho=sum(sum(R.*R)); rho0=rho; P=R;
3 V=zeros(m,m); T1=sparse(tridiagonal(a,d,a,m));
4 for k=1:itmax
5   if sqrt(rho/rho0)<= tol
6     K=k; return
7   end
8   T=T1*P+P*T1;
9   a=rho/sum(sum(P.*T)); V=V+a*P; R=R-a*T;
10  rhos=rho; rho=sum(sum(R.*R)); P=R+(rho/rhos)*P;
11 end
12 K=itmax+1;
```

For both the averaging- and Poisson matrix we use $tol = 10^{-8}$.

For the averaging matrix we obtain the values in Table 10.20.

The convergence is quite rapid. It appears that the number of iterations can be bounded independently of n , and therefore we solve the problem in $O(n)$ operations. This is the best we can do for a problem with n unknowns.

Consider next the Poisson problem. In Table 10.21 we list K , the required number of iterations, and K/\sqrt{n} .

The results show that K is much smaller than n and appears to be proportional to \sqrt{n} . This is the same speed as for SOR and we don't have to estimate

n	2 500	10 000	40 000	1 000 000	4 000 000
K	19	18	18	16	15

Table 10.20: The number of iterations K for the averaging problem on a $\sqrt{n} \times \sqrt{n}$ grid for various n

n	2 500	10 000	40 000	160 000
K	94	188	370	735
K/\sqrt{n}	1.88	1.88	1.85	1.84

Table 10.21: The number of iterations K for the Poisson problem on a $\sqrt{n} \times \sqrt{n}$ grid for various n

any acceleration parameter.

10.4.3 The Spectral Condition Numbers

Using the best approximation property we will show in Section 10.5 the following upper bound for the \mathbf{A} -norm of the error in the conjugate gradient method.

Theorem 10.22 (Error bounds for cg)

We have

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k < 2e^{-\frac{2}{\sqrt{\kappa}}k}, \quad k \geq 0, \quad (10.27)$$

where $\kappa = M/m$, and where M and m are the largest and smallest eigenvalue of \mathbf{A} , respectively.

Thus the upper bound for the rate of convergence is determined by the square root of the spectral condition number. This is much better than the estimate (10.13) for the steepest descent method. Especially for problems with large condition numbers.

So what is the spectral condition number of \mathbf{T}_2 ? The eigenvalues were given in (4.22) as

$$\lambda_{j,k} = 2d + 2a \cos(j\pi h) + 2a \cos(k\pi h), \quad j, k = 1, \dots, m. \quad (10.28)$$

For the averaging problem it follows that the largest and smallest eigenvalue of \mathbf{T}_2 are $M = \frac{5}{9} + \frac{4}{9} \cos(\pi h)$ and $m = \frac{5}{9} - \frac{4}{9} \cos(\pi h)$. Thus

$$\kappa_A = \frac{M}{m} = \frac{5 + 4 \cos(\pi h)}{5 - 4 \cos(\pi h)} \leq 9,$$

and the condition number is independent of n . This verifies what we observed in Table 10.20. The number of iterations can be bounded independently of the size n of the problem.

The eigenvalues for the Poisson problem can also be found from (10.28). We find $M = 4(1 + \cos(\pi h))$ and $m = 4(1 - \cos(\pi h))$ and then

$$\kappa_P = \frac{M}{m} = \frac{1 + \cos(\pi h)}{1 - \cos(\pi h)} = 1/\tan(\pi h^2/2).$$

Thus $\kappa_P \approx 4(m+1)^2/\pi^2 = O(n)$ (see also Exercise 8.35) and we solve the discrete Poisson problem in $O(n^{3/2})$ arithmetic operations. This is the same as for the SOR method and for the fast method without the FFT. In comparison the Cholesky Algorithm requires $O(n^2)$ arithmetic operations both for the averaging and the Poisson problem.

Exercise 10.23 (Krylov space and cg iterations)

Consider the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ where

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 4 \\ 0 \\ 0 \end{bmatrix}.$$

- a) Determine the vectors defining the Krylov spaces for $k \leq 3$ taking as initial approximation $\mathbf{x} = \mathbf{0}$. Answer: $[\mathbf{b}, \mathbf{Ab}, \mathbf{A}^2\mathbf{b}] = \begin{bmatrix} 4 & 8 & 20 \\ 0 & -4 & -16 \\ 0 & 0 & 4 \end{bmatrix}$.
- b) Carry out three CG-iterations on $\mathbf{Ax} = \mathbf{b}$. Answer:

$$[\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3] = \begin{bmatrix} 0 & 2 & 8/3 & 3 \\ 0 & 0 & 4/3 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

$$[\mathbf{r}_0, \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3] = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 4/3 & 0 \end{bmatrix},$$

$$[\mathbf{Ap}_0, \mathbf{Ap}_1, \mathbf{Ap}_2] = \begin{bmatrix} 8 & 0 & 0 \\ -4 & 3 & 0 \\ 0 & -2 & 16/9 \end{bmatrix},$$

$$[\mathbf{p}_0, \mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3] = \begin{bmatrix} 4 & 1 & 4/9 & 0 \\ 0 & 2 & 8/9 & 0 \\ 0 & 0 & 12/9 & 0 \end{bmatrix},$$

c) Verify that

- $\dim(\mathbb{W}_k) = k$ for $k = 0, 1, 2, 3$.
- \mathbf{x}_3 is the exact solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$.
- $\mathbf{r}_0, \dots, \mathbf{r}_{k-1}$ is an orthogonal basis for \mathbb{W}_k for $k = 1, 2, 3$.
- $\mathbf{p}_0, \dots, \mathbf{p}_{k-1}$ is an \mathbf{A} -orthogonal basis for \mathbb{W}_k for $k = 1, 2, 3$.
- $\{|r_k\epsilon|\}$ is monotonically decreasing.
- $\{|x_k - x\epsilon|\}$ is monotonically decreasing.

Exercise 10.24 (Program code for testing steepest descent)

Write a function $K=sdtest(m,a,d,tol,itmax)$ to test the Steepest descent method on the matrix \mathbf{T}_2 . Make the analogues of Table 10.20 and Table 10.21. For Table 10.21 it is enough to test for say $n = 100, 400, 1600, 2500$, and tabulate K/n instead of K/\sqrt{n} in the last row. Conclude that the upper bound (10.13) is realistic. Compare also with the number of iterations for the J and GS method in Table 9.1.

Exercise 10.25 (Compare Richardson and steepest descent)

Go back and study the Richardson method (9.25) where a constant α is used. Suppose we use the α^* in (9.26). What seems to be the main difference between this method and the steepest descent method?

Exercise 10.26 (Using cg to solve normal equations)

Consider solving the least squares problem by using the conjugate gradient method on the normal equations $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$. Explain why only the following modifications in Algorithm 10.18 are necessary

1. $r=\mathbf{A}'(\mathbf{b}-\mathbf{A}^*\mathbf{x})$; $p=r$;
2. $a=rho/(t'*t)$;
3. $r=r-a^*\mathbf{A}'*t$;

Note that the condition number of the normal equations is $\text{cond}_2(\mathbf{A})^2$, the square of the condition number of \mathbf{A} .

10.5 Proof of Convergence

We prove Theorem 10.22.

By Theorem 10.15 \mathbf{x}_k is the best approximation to the solution \mathbf{x} in the \mathbf{A} -norm. We convert this best approximation property into a minimization problem involving polynomials. In the following we let Π_k denote the class of univariate polynomials of degree $\leq k$ with real coefficients.

Lemma 10.27 (Krylov space and polynomials)

Suppose $\mathbf{A}\mathbf{x} = \mathbf{b}$ where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive definite with orthonormal eigenpairs $(\lambda_j, \mathbf{u}_j)$, $j = 1, 2, \dots, n$, and let $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ with $\mathbf{x}_0 \in \mathbb{R}^n$. If $\mathbf{w} = \sum_{j=0}^{k-1} a_j \mathbf{A}^j \mathbf{r}_0 \in \mathbb{W}_k$ for some $k \geq 0$ then

$$\mathbf{w} = P(\mathbf{A})\mathbf{r}_0, \quad P(\mathbf{A}) = a_0 I + a_1 \mathbf{A} + a_2 \mathbf{A}^2 + \cdots + a_{k-1} \mathbf{A}^{k-1}.$$

where $P(\mathbf{A})$ is a matrix polynomial corresponding to the ordinary polynomial $P(t) = a_0 + a_1 t + \cdots + a_{k-1} t^{k-1}$ of degree $\leq k-1$. Moreover,

$$\|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 = \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j} Q(\lambda_j)^2, \quad Q(t) := 1 - tP(t), \quad (10.29)$$

where $\mathbf{b} = \sum_{j=1}^n \sigma_j \mathbf{u}_j$.

Proof. Clearly $\mathbf{w} = P(\mathbf{A})\mathbf{r}_0$. We find

$$\|\mathbf{x} - P(\mathbf{A})\mathbf{b}\|_{\mathbf{A}}^2 = \mathbf{c}^T \mathbf{A}^{-1} \mathbf{c}, \quad \mathbf{c} = Q(\mathbf{A})\mathbf{b}, \quad Q(\mathbf{A}) = I - \mathbf{A}P(\mathbf{A}). \quad (10.30)$$

Using the eigenvector expansion for \mathbf{b} and (0.66) we obtain

$$\mathbf{c} = \sum_{j=1}^n \sigma_j Q(\lambda_j) \mathbf{u}_j, \quad \mathbf{A}^{-1} \mathbf{c} = \sum_{i=1}^n \sigma_i \frac{Q(\lambda_i)}{\lambda_i} \mathbf{u}_i. \quad (10.31)$$

Combining (10.30),(10.31), and using orthonormality we find

$$\|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 = \mathbf{c}^T \mathbf{A}^{-1} \mathbf{c} = \left(\sum_{j=1}^n \sigma_j Q(\lambda_j) \mathbf{u}_j \right)^T \left(\sum_{i=1}^n \sigma_i \frac{Q(\lambda_i)}{\lambda_i} \mathbf{u}_i \right) = \sum_{j=1}^n \sigma_j^2 \frac{Q(\lambda_j)^2}{\lambda_j}.$$

□

We will use the following Theorem to estimate the rate of convergence.

Theorem 10.28 (cg and best polynomial approximation)

Suppose $[a, b]$ with $0 < a < b$ is an interval containing all the eigenvalues of \mathbf{A} . Then for all $Q \in \Pi_k$ with $Q(0) = 1$ we have

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq \max_{a \leq x \leq b} |Q(x)|.$$

Proof. We find $\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}^2 = \mathbf{r}_0^T \mathbf{A}^{-1} \mathbf{r}_0 = \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j}$. Therefore, by the best approximation property and (10.29), for any $\mathbf{w} \in \mathbb{W}_k$

$$\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2 \leq \|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 \leq \max_{a \leq x \leq b} |Q(x)|^2 \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j} = \max_{a \leq x \leq b} |Q(x)|^2 \|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}^2$$

and the result follows by taking square roots. \square

In the next section we use properties of the Chebyshev polynomials to show that

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_A}{\|\mathbf{x} - \mathbf{x}_0\|_A} \leq \min_{\substack{Q \in \Pi_k \\ Q(0)=1}} \max_{m \leq x \leq M} |Q(x)| = \frac{2}{\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^k + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k}, \quad (10.32)$$

where $\kappa = M/m$. But this implies the first inequality in (10.27). The second inequality follows from (10.15)

10.5.1 Chebyshev Polynomials

Suppose $a < b$, $c \notin [a, b]$ and $k \in \mathbb{N}$. Consider the set \mathcal{S}_k of all polynomials Q of degree $\leq k$ such that $Q(c) = 1$. For any continuous function f on $[a, b]$ we define

$$\|f\|_\infty = \max_{a \leq x \leq b} |f(x)|.$$

We want to find a polynomial $Q^* \in \mathcal{S}_k$ such that (cf. Corollary 10.28, where $c = 0 < a < b$)

$$\|Q^*\|_\infty = \min_{Q \in \mathcal{S}_k} \|Q\|_\infty.$$

We will show that Q^* is a suitably shifted and normalized version of the **Chebyshev polynomial**. The Chebyshev polynomial T_n of degree n can be defined recursively by

$$T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t), \quad n \geq 1, \quad t \in \mathbb{R},$$

starting with $T_0(t) = 1$ and $T_1(t) = t$. Thus $T_2(t) = 2t^2 - 1$, $T_3(t) = 4t^3 - 3t$ etc. In general T_n is a polynomial of degree n .

There are some convenient closed form expressions for T_n .

Lemma 10.29 (Closed forms of Chebyshev polynomials)

For $n \geq 0$

1. $T_n(t) = \cos(n \arccos t)$ for $t \in [-1, 1]$,
2. $T_n(t) = \frac{1}{2} [(t + \sqrt{t^2 - 1})^n + (t - \sqrt{t^2 - 1})^{-n}]$ for $|t| \geq 1$.

Proof. 1. With $P_n(t) = \cos(n \arccos t)$ we have $P_n(t) = \cos n\phi$, where $t = \cos \phi$. Therefore

$$P_{n+1}(t) + P_{n-1}(t) = \cos(n+1)\phi + \cos(n-1)\phi = 2 \cos \phi \cos n\phi = 2tP_n(t)$$

and it follows that P_n satisfies the same recurrence relation as T_n . Since $P_0 = T_0$ and $P_1 = T_1$ we have $P_n = T_n$ for all $n \geq 0$.

2. Fix t with $|t| \geq 1$ and let $x_n := T_n(t)$ for $n \geq 0$. The recurrence relation for the Chebyshev polynomials can then be written

$$x_{n+1} - 2tx_n + x_{n-1} = 0 \text{ for } n \geq 1, \text{ with } x_0 = 1, x_1 = t. \quad (10.33)$$

To solve this difference equation we insert $x_n = z^n$ into (10.33) and obtain $z^{n+1} - 2tz^n + z^{n-1} = 0$ or $z^2 - 2tz + 1 = 0$. The roots of this equation are

$$z_1 = t + \sqrt{t^2 - 1}, \quad z_2 = t - \sqrt{t^2 - 1} = (t + \sqrt{t^2 - 1})^{-1}.$$

Now z_1^n , z_2^n and more generally $c_1 z_1^n + c_2 z_2^n$ are solutions of (10.33) for any constants c_1 and c_2 . We find these constants from the initial conditions $x_0 = c_1 + c_2 = 1$ and $x_1 = c_1 z_1 + c_2 z_2 = t$. Since $z_1 + z_2 = 2t$ the solution is $c_1 = c_2 = \frac{1}{2}$. \square

We show that the unique solution to our minimization problem is

$$Q^*(x) = \frac{T_k(u(x))}{T_k(u(c))}, \quad u(x) = \frac{b+a-2x}{b-a}. \quad (10.34)$$

Clearly $Q^* \in \mathcal{X}_k$.

Theorem 10.30 (A minimal norm problem)

Suppose $a < b$, $c \notin [a, b]$ and $k \in \mathbb{N}$. If $Q \in \mathcal{S}_k$ and $Q \neq Q^*$ then $\|Q\|_\infty > \|Q^*\|_\infty$.

Proof. Recall that a nonzero polynomial p of degree k can have at most k zeros. If $p(z) = p'(z) = 0$, we say that p has a double zero at z . Counting such a zero as two zeros it is still true that a nonzero polynomial of degree k has at most k zeros.

$|Q^*|$ takes on its maximum $1/|T_k(u(c))|$ at the $k+1$ points μ_0, \dots, μ_k in $[a, b]$ such that $u(\mu_i) = \cos(i\pi/k)$ for $i = 0, 1, \dots, k$. Suppose $Q \in S_k$ and that $\|Q\| \leq \|Q^*\|$. We have to show that $Q \equiv Q^*$. Let $f \equiv Q - Q^*$. We want to show that f has at least k zeros in $[a, b]$. Since f is a polynomial of degree $\leq k$ and $f(c) = 0$, this means that $f \equiv 0$ or equivalently $Q \equiv Q^*$.

Consider $I_j = [\mu_{j-1}, \mu_j]$ for a fixed j . Let

$$\sigma_j = f(\mu_{j-1})f(\mu_j).$$

We have $\sigma_j \leq 0$. For if say $Q^*(\mu_j) > 0$ then

$$Q(\mu_j) \leq \|Q\|_\infty \leq \|Q^*\|_\infty = Q^*(\mu_j)$$

so that $f(\mu_j) \leq 0$. Moreover,

$$-Q(\mu_{j-1}) \leq \|Q\|_\infty \leq \|Q^*\|_\infty = -Q^*(\mu_{j-1}).$$

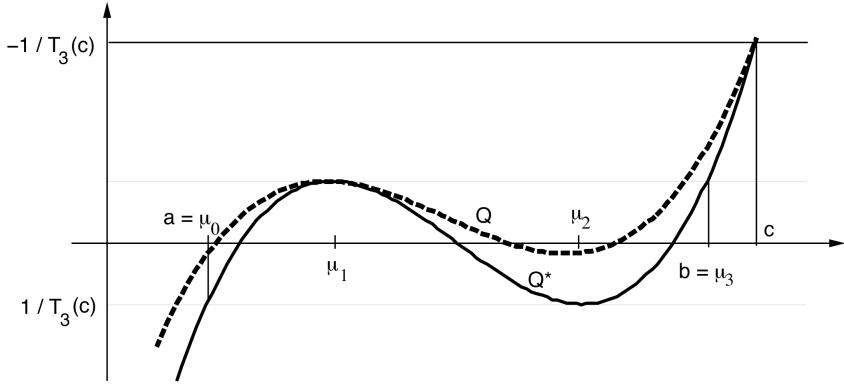


Figure 10.3: This is an illustration of the proof of Theorem 10.30 for $k = 3$. $f \equiv Q - Q^*$ has a double zero at μ_1 and one zero between μ_2 and μ_3 .

Thus $f(\mu_{j-1}) \geq 0$ and it follows that $\sigma_j \leq 0$. Similarly, $\sigma_j \leq 0$ if $Q^*(\mu_j) < 0$.

If $\sigma_j < 0$, f must have a zero in I_j since it is continuous. Suppose $\sigma_j = 0$. Then $f(\mu_{j-1}) = 0$ or $f(\mu_j) = 0$. If $f(\mu_j) = 0$ then $Q(\mu_j) = Q^*(\mu_j)$. But then μ_j is a maximum or minimum both for Q and Q^* . If $\mu_j \in (a, b)$ then $Q'(\mu_j) = Q^{*\prime}(\mu_j) = 0$. Thus $f(\mu_j) = f'(\mu_j) = 0$, and f has a double zero at μ_j . We can count this as one zero for I_j and one for I_{j+1} . If $\mu_j = b$, we still have a zero in I_j . Similarly, if $f(\mu_{j-1}) = 0$, a double zero of f at μ_{j-1} appears if $\mu_{j-1} \in (a, b)$. We count this as one zero for I_{j-1} and one for I_j .

In this way we associate one zero of f for each of the k intervals I_j , $j = 1, 2, \dots, k$. We conclude that f has at least k zeros in $[a, b]$. \square

Exercise 10.31 (An explicit formula for the Chebyshev polynomial)
Show that

$$T_n(t) = \cosh(n \operatorname{arccosh} t) \text{ for } |t| \geq 1,$$

where $\operatorname{arccosh}$ is the inverse function of $\cosh x := (e^x + e^{-x})/2$.

Theorem 10.30 with $a = m$, $b = M$, and $c = 0$ implies that the minimizing polynomial in (10.32) is given by

$$Q^*(x) = T_k \left(\frac{M+m-2x}{M-m} \right) / T_k \left(\frac{M+m}{M-m} \right), \quad (10.35)$$

where m and M , the smallest and largest eigenvalue of \mathbf{A} . By Lemma 10.29

$$\max_{m \leq x \leq M} \left| T_k \left(\frac{M+m-2x}{M-m} \right) \right| = \max_{-1 \leq t \leq 1} |T_k(t)| = 1. \quad (10.36)$$

Moreover with $t = (M + m)/(M - m)$ we have

$$t + \sqrt{t^2 - 1} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}, \quad \kappa = M/m.$$

Thus again by Lemma 10.29 we find

$$T_k \left(\frac{M+m}{M-m} \right) = T_k \left(\frac{\kappa+1}{\kappa-1} \right) = \frac{1}{2} \left[\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1} \right)^k + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^k \right] \quad (10.37)$$

and (10.32) and the first inequality in follows.

10.5.2 Monotonicity of the error

We end this chapter by showing that the Euclidian norm of the error is strictly decreasing.

Theorem 10.32 (The error in cg is strictly decreasing)

Let in the conjugate gradient method m be the smallest integer such that $\mathbf{r}_{m+1} = \mathbf{0}$. For $k \leq m$ we have $\|\boldsymbol{\epsilon}_{k+1}\|_2 < \|\boldsymbol{\epsilon}_k\|_2$. More precisely,

$$\|\boldsymbol{\epsilon}_k\|_2^2 - \|\boldsymbol{\epsilon}_{k+1}\|_2^2 = \frac{\|\mathbf{p}_k\|_2^2}{\|\mathbf{p}_k\|_{\mathbf{A}}^2} (\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 + \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2)$$

where $\boldsymbol{\epsilon}_j = \mathbf{x} - \mathbf{x}_j$ and $\mathbf{A}\mathbf{x} = b$.

Proof. For $j \leq m$

$$\boldsymbol{\epsilon}_j = \mathbf{x}_{m+1} - \mathbf{x}_j = \mathbf{x}_m - \mathbf{x}_j + \alpha_m \mathbf{p}_m = \mathbf{x}_{m-1} - \mathbf{x}_j + \alpha_{m-1} \mathbf{p}_{m-1} + \alpha_m \mathbf{p}_m = \dots$$

so that

$$\boldsymbol{\epsilon}_j = \sum_{i=j}^m \alpha_i \mathbf{p}_i, \quad \alpha_i = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}. \quad (10.38)$$

By (10.38) and \mathbf{A} -orthogonality

$$\|\boldsymbol{\epsilon}_j\|_{\mathbf{A}}^2 = \boldsymbol{\epsilon}_j^T \mathbf{A} \boldsymbol{\epsilon}_j = \sum_{i=j}^m \alpha_i^2 \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i = \sum_{i=j}^m \frac{(\mathbf{r}_i^T \mathbf{r}_i)^2}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}. \quad (10.39)$$

By (10.22) and Lemma 10.14

$$\mathbf{p}_i^T \mathbf{p}_k = (\mathbf{r}_i + \beta_{i-1} \mathbf{p}_{i-1})^T \mathbf{p}_k = \beta_{i-1} \mathbf{p}_{i-1}^T \mathbf{p}_k = \dots = \beta_{i-1} \cdots \beta_k (\mathbf{p}_k^T \mathbf{p}_k),$$

and since $\beta_{i-1} \cdots \beta_k = (\mathbf{r}_i^T \mathbf{r}_i) / (\mathbf{r}_k^T \mathbf{r}_k)$ we obtain

$$\mathbf{p}_i^T \mathbf{p}_k = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_k^T \mathbf{r}_k} \mathbf{p}_k^T \mathbf{p}_k, \quad i \geq k. \quad (10.40)$$

Since

$$\|\boldsymbol{\epsilon}_k\|_2^2 = \|\boldsymbol{\epsilon}_{k+1} + \mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 = \|\boldsymbol{\epsilon}_{k+1} + \alpha_k \mathbf{p}_k\|_2^2,$$

we obtain

$$\begin{aligned} \|\boldsymbol{\epsilon}_k\|_2^2 - \|\boldsymbol{\epsilon}_{k+1}\|_2^2 &= \alpha_k (\mathbf{p}_k^T \boldsymbol{\epsilon}_{k+1} + \alpha_k \mathbf{p}_k^T \mathbf{p}_k) \\ &\stackrel{(10.38)}{=} \alpha_k \left(2 \sum_{i=k+1}^m \alpha_i \mathbf{p}_i^T \mathbf{p}_k + \alpha_k \mathbf{p}_k^T \mathbf{p}_k \right) = \left(\sum_{i=k}^m + \sum_{i=k+1}^m \right) \alpha_k \alpha_i \mathbf{p}_i^T \mathbf{p}_k \\ &\stackrel{(10.40)}{=} \left(\sum_{i=k}^m + \sum_{i=k+1}^m \right) \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_k^T \mathbf{r}_k} \mathbf{p}_k^T \mathbf{p}_k \\ &\stackrel{(10.39)}{=} \frac{\|\mathbf{p}_k\|_2^2}{\|\mathbf{p}_k\|_{\mathbf{A}}^2} (\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 + \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2). \end{aligned}$$

and the Theorem is proved. \square

10.6 Preconditioning

For problems $\mathbf{A}\mathbf{x} = \mathbf{b}$ of size n , where both n and $\text{cond}_2(\mathbf{A})$ are large, it is often possible to improve the performance of the conjugate gradient method by using a technique known as **preconditioning**. Instead of $\mathbf{A}\mathbf{x} = \mathbf{b}$ we consider an equivalent system $\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b}$, where \mathbf{B} is nonsingular and $\text{cond}_2(\mathbf{B}\mathbf{A})$ is smaller than $\text{cond}_2(\mathbf{A})$. The matrix \mathbf{B} will in many cases be the inverse of another matrix, $\mathbf{B} = \mathbf{M}^{-1}$. We cannot use CG on $\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b}$ directly since $\mathbf{B}\mathbf{A}$ in general is not symmetric even if both \mathbf{A} and \mathbf{B} are. But if \mathbf{B} (and hence \mathbf{M}) is symmetric positive definite then we can apply CG to a symmetrized system and then transform the recurrence formulas to an iterative method for the original system $\mathbf{A}\mathbf{x} = \mathbf{b}$. This iterative method is known as the **preconditioned conjugate gradient method**. We shall see that the convergence properties of this method is determined by the eigenvalues of $\mathbf{B}\mathbf{A}$.

Suppose \mathbf{B} is symmetric positive definite. By Theorem 3.31 there is a non-singular matrix \mathbf{C} such that $\mathbf{B} = \mathbf{C}^T \mathbf{C}$. (\mathbf{C} is only needed for the derivation and will not appear in the final formulas). Now

$$\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b} \Leftrightarrow \mathbf{C}^T (\mathbf{C}\mathbf{A}\mathbf{C}^T) \mathbf{C}^{-T} \mathbf{x} = \mathbf{C}^T \mathbf{C}\mathbf{b} \Leftrightarrow (\mathbf{C}\mathbf{A}\mathbf{C}^T)\mathbf{y} = \mathbf{C}\mathbf{b}, \quad \& \quad \mathbf{x} = \mathbf{C}^T \mathbf{y}.$$

We have 3 linear systems

$$\mathbf{A}\mathbf{x} = \mathbf{b} \tag{10.41}$$

$$\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b} \tag{10.42}$$

$$(\mathbf{C}\mathbf{A}\mathbf{C}^T)\mathbf{y} = \mathbf{C}\mathbf{b}, \quad \& \quad \mathbf{x} = \mathbf{C}^T \mathbf{y}. \tag{10.43}$$

Note that (10.41) and (10.43) are symmetric positive definite linear systems. In addition to being symmetric positive definite the matrix \mathbf{CAC}^T is similar to \mathbf{BA} . Indeed,

$$\mathbf{C}^T(\mathbf{CAC}^T)\mathbf{C}^{-T} = \mathbf{BA}.$$

Thus \mathbf{CAC}^T and \mathbf{BA} have the same eigenvalues. Therefore if we apply the conjugate gradient method to (10.43) then the rate of convergence will be determined by the eigenvalues of \mathbf{BA} .

We apply the conjugate gradient method to $(\mathbf{CAC}^T)\mathbf{y} = \mathbf{Cb}$. Denoting the search direction by \mathbf{q}_k and the residual by $\mathbf{z}_k = \mathbf{Cb} - \mathbf{CAC}^T\mathbf{y}_k$ we obtain the following from (10.20), (10.21), and (10.22).

$$\begin{aligned}\mathbf{y}_{k+1} &= \mathbf{y}_k + \alpha_k \mathbf{q}_k, & \alpha_k &= \mathbf{z}_k^T \mathbf{z}_k / \mathbf{q}_k^T (\mathbf{CAC}^T) \mathbf{q}_k, \\ \mathbf{z}_{k+1} &= \mathbf{z}_k - \alpha_k (\mathbf{CAC}^T) \mathbf{q}_k, \\ \mathbf{q}_{k+1} &= \mathbf{z}_{k+1} + \beta_k \mathbf{q}_k, & \beta_k &= \mathbf{z}_{k+1}^T \mathbf{z}_{k+1} / \mathbf{z}_k^T \mathbf{z}_k.\end{aligned}$$

With

$$\mathbf{x}_k := \mathbf{C}^T \mathbf{y}_k, \quad \mathbf{p}_k := \mathbf{C}^T \mathbf{q}_k, \quad \mathbf{s}_k := \mathbf{C}^T \mathbf{z}_k, \quad \mathbf{r}_k := \mathbf{C}^{-1} \mathbf{z}_k \quad (10.44)$$

this can be transformed into

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \alpha_k = \frac{\mathbf{s}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}, \quad (10.45)$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (10.46)$$

$$\mathbf{s}_{k+1} = \mathbf{s}_k - \alpha_k \mathbf{B} \mathbf{A} \mathbf{p}_k, \quad (10.47)$$

$$\mathbf{p}_{k+1} = \mathbf{s}_{k+1} + \beta_k \mathbf{p}_k, \quad \beta_k = \frac{\mathbf{s}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{s}_k^T \mathbf{r}_k}. \quad (10.48)$$

Here \mathbf{x}_k will be an approximation to the solution \mathbf{x} of $\mathbf{Ax} = \mathbf{b}$, $\mathbf{r}_k = \mathbf{b} - \mathbf{Ax}_k$ is the residual in the original system, and $\mathbf{s}_k = \mathbf{Bb} - \mathbf{BAx}_k$ is the residual in the preconditioned system. This follows since by (10.44)

$$\mathbf{r}_k = \mathbf{C}^{-1} \mathbf{z}_k = \mathbf{b} - \mathbf{C}^{-1} \mathbf{CAC}^T \mathbf{y}_k = \mathbf{b} - \mathbf{Ax}_k$$

and $\mathbf{s}_k = \mathbf{C}^T \mathbf{z}_k = \mathbf{C}^T \mathbf{Cr}_k = \mathbf{Br}_k$. We start with $\mathbf{r}_0 = \mathbf{b} - \mathbf{Ax}_0$, $\mathbf{p}_0 = \mathbf{s}_0 = \mathbf{Br}_0$ and obtain the following preconditioned conjugate gradient algorithm for determining approximations \mathbf{x}_k to the solution of a symmetric positive definite system $\mathbf{Ax} = \mathbf{b}$.

Algorithm 10.33 (Preconditioned conjugate gradient)

The symmetric positive definite linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ is solved by the preconditioned conjugate gradient method on the system $\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b}$, where \mathbf{B} is symmetric positive definite. \mathbf{x} is a starting vector for the iteration. The iteration is stopped when $\|\mathbf{r}_k\|_2/\|\mathbf{r}_0\|_2 \leq \text{tol}$ or $k > \text{itmax}$. K is the number of iterations used.

```

1 function [x,K]=pcg(A,B,b,x,tol ,itmax)
2 r=b-A*x; p=B*r; s=p; rho=s'*r;
3 rho0=rho;
4 for k=0:itmax
5   if sqrt(rho/rho0)<= tol
6     K=k; return
7   end
8   t=A*p; a=rho/(p'*t);
9   x=x+a*p; r=r-a*t;
10  w=B*t; s=s-a*w;
11  rhos=rho; rho=s'*r;
12  p=r+(rho/rhos)*p;
13 end
14 K=itmax+1;
```

This algorithm is quite similar to Algorithm 10.18. It differs in the calculation of ρ . The main additional work is contained in $w = B * t$. We'll discuss this further in connection with an example. There the inverse of \mathbf{B} is known and we have to solve a linear system to find w .

We have the following convergence result for this algorithm.

Theorem 10.34 (Error bound preconditioned cg)

Suppose we apply a symmetric positive definite preconditioner \mathbf{B} to the symmetric positive definite system $\mathbf{A}\mathbf{x} = \mathbf{b}$. Then the quantities \mathbf{x}_k computed in Algorithm 10.33 satisfy the following bound:

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where $\kappa = \lambda_{\max}/\lambda_{\min}$ is the ratio of the largest and smallest eigenvalue of $\mathbf{B}\mathbf{A}$.

Proof. Since Algorithm 10.33 is equivalent to solving (10.43) by the conjugate gradient method Theorem 10.22 implies that

$$\frac{\|\mathbf{y} - \mathbf{y}_k\|_{\mathbf{CAC}^T}}{\|\mathbf{y} - \mathbf{y}_0\|_{\mathbf{CAC}^T}} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where \mathbf{y}_k is the conjugate gradient approximation to the solution \mathbf{y} of (10.43) and κ is the ratio of the largest and smallest eigenvalue of \mathbf{CAC}^T . Since \mathbf{BA} and

$\mathbf{C}\mathbf{A}\mathbf{C}^T$ are similar this is the same as the κ in the theorem. By (10.44) we have

$$\begin{aligned}\|\mathbf{y} - \mathbf{y}_k\|_{\mathbf{C}\mathbf{A}\mathbf{C}^T}^2 &= (\mathbf{y} - \mathbf{y}_k)^T (\mathbf{C}\mathbf{A}\mathbf{C}^T)(\mathbf{y} - \mathbf{y}_k) \\ &= (\mathbf{C}^T(\mathbf{y} - \mathbf{y}_k))^T \mathbf{A}(\mathbf{C}^T(\mathbf{y} - \mathbf{y}_k)) = \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2\end{aligned}$$

and the proof is complete. \square

We conclude that \mathbf{B} should satisfy the following requirements for a problem of size n :

1. The eigenvalues of $\mathbf{B}\mathbf{A}$ should be located in a narrow interval. Preferably one should be able to bound the length of the interval independently of n .
2. The evaluation of $\mathbf{B}\mathbf{x}$ for a given vector \mathbf{x} should not be expensive in storage and arithmetic operations, ideally $O(n)$ for both.

10.7 Preconditioning Example

Throughout this section we use the same grid and notation as in Section 2.4. Let $h = 1/(m+1)$.

We recall the Poisson problem

$$-\nabla^2 u = -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f(x, y) \quad \text{for } (x, y) \in \Omega = (0, 1)^2 \quad (10.49)$$

$$u = 0 \text{ on } \partial\Omega,$$

where f is a given function, Ω is the unit square in the plane, and $\partial\Omega$ is the boundary of Ω . For numerical solution we have the **discrete Poisson problem** which can either be written as a matrix equation

$$\begin{aligned}h^2 f_{j,k} &= 4v_{j,k} - v_{j-1,k} - v_{j+1,k} - v_{j,k-1} - v_{j,k+1}, \quad j, k = 1, \dots, m \\ v_{0,k} &= v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0, \quad j, k = 0, 1, \dots, m+1,\end{aligned}$$

or as a system $\mathbf{A}_p \mathbf{x} = \mathbf{b}$, where $\mathbf{x} = \text{vec}(v_{i,j})$, $\mathbf{b} = h^2 \text{vec}(f_{i,j})$ and the elements $a_{i,j}$ of \mathbf{A}_p are given by

$$\begin{aligned}a_{ii} &= 4, & i &= 1, \dots, n \\ a_{i+1,i} = a_{i,i+1} &= -1, & i &= 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m \\ a_{i+m,i} = a_{i,i+m} &= -1, & i &= 1, \dots, n-m \\ a_{ij} &= 0, & & \text{otherwise.}\end{aligned}$$

10.7.1 A Banded Matrix

Consider the problem

$$\begin{aligned} -\frac{\partial}{\partial x} \left(c(x, y) \frac{\partial u}{\partial x} \right) - \frac{\partial}{\partial y} \left(c(x, y) \frac{\partial u}{\partial y} \right) &= f(x, y) & (x, y) \in \Omega = (0, 1)^2 \\ u(x, y) &= 0 & (x, y) \in \partial\Omega. \end{aligned} \quad (10.50)$$

Here Ω is the open unit square while $\partial\Omega$ is the boundary of Ω . The functions f and c are given and we seek a function $u = u(x, y)$ such that (10.50) holds. We assume that c and f are defined and continuous on Ω and that $c(x, y) > 0$ for all $(x, y) \in \Omega$. The problem (10.50) reduces to the Poisson problem in the special case where $c(x, y) = 1$ for $(x, y) \in \Omega$.

As for the Poisson problem we solve (10.50) numerically on a grid of points

$$\{(jh, kh) : j, k = 0, 1, \dots, m+1\}, \quad \text{where } h = 1/(m+1),$$

and where m is a positive integer. Let (x, y) be one of the interior grid points. For univariate functions f, g we use the central difference approximations

$$\begin{aligned} \frac{\partial}{\partial t} \left(f(t) \frac{\partial}{\partial t} g(t) \right) &\approx \left(f(t + \frac{h}{2}) \frac{\partial}{\partial t} g(t + h/2) - f(t - \frac{h}{2}) \frac{\partial}{\partial t} g(t - h/2) \right) / h \\ &\approx \left(f(t + \frac{h}{2})(g(t + h) - g(t)) - f(t - \frac{h}{2})(g(t) - g(t - h)) \right) / h^2 \end{aligned}$$

to obtain

$$\frac{\partial}{\partial x} \left(c \frac{\partial u}{\partial x} \right)_{j,k} \approx \frac{c_{j+\frac{1}{2},k}(v_{j+1,k} - v_{j,k}) - c_{j-\frac{1}{2},k}(v_{j,k} - v_{j-1,k})}{h^2}$$

and

$$\frac{\partial}{\partial y} \left(c \frac{\partial u}{\partial y} \right)_{j,k} \approx \frac{c_{j,k+\frac{1}{2}}(v_{j,k+1} - v_{j,k}) - c_{j,k-\frac{1}{2}}(v_{j,k} - v_{j,k-1})}{h^2},$$

where $c_{p,q} = c(ph, qh)$ and $v_{j,k} \approx u(jh, kh)$. With these approximations the discrete analog of (10.50) turns out to be

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= h^2 f_{j,k} & j, k = 1, \dots, m \\ v_{j,k} &= 0 & j = 0, m+1 \text{ all } k \text{ or } k = 0, m+1 \text{ all } j, \end{aligned} \quad (10.51)$$

where

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= (c_{j,k-\frac{1}{2}} + c_{j-\frac{1}{2},k} + c_{j+\frac{1}{2},k} + c_{j,k+\frac{1}{2}}) v_{j,k} \\ &\quad - c_{j,k-\frac{1}{2}} v_{j,k-1} - c_{j-\frac{1}{2},k} v_{j-1,k} - c_{j+\frac{1}{2},k} v_{j+1,k} - c_{j,k+\frac{1}{2}} v_{j,k+1} \end{aligned} \quad (10.52)$$

and $f_{j,k} = f(jh, kh)$.

As before we let $\mathbf{V} = (v_{j,k}) \in \mathbb{R}^{m \times m}$ and $\mathbf{F} = (f_{j,k}) \in \mathbb{R}^{m \times m}$. The corresponding linear system can be written $\mathbf{A}\mathbf{x} = \mathbf{b}$ where $\mathbf{x} = \text{vec}(\mathbf{V})$, $\mathbf{b} = h^2 \text{vec}(\mathbf{F})$, and the n -by- n coefficient matrix \mathbf{A} is given by

$$\begin{aligned} a_{i,i} &= c_{j_i, k_i - \frac{1}{2}} + c_{j_i - \frac{1}{2}, k_i} + c_{j_i + \frac{1}{2}, k_i} + c_{j_i, k_i + \frac{1}{2}}, & i = 1, 2, \dots, n \\ a_{i+1,i} &= a_{i,i+1} = -c_{j_i + \frac{1}{2}, k_i}, & i \bmod m \neq 0 \\ a_{i+m,i} &= a_{i,i+m} = -c_{j_i, k_i + \frac{1}{2}}, & i = 1, 2, \dots, n-m \\ a_{i,j} &= 0 & \text{otherwise,} \end{aligned} \tag{10.53}$$

where (j_i, k_i) with $1 \leq j_i, k_i \leq m$ is determined uniquely from the equation $i = j_i + (k_i - 1)m$ for $i = 1, \dots, n$. When $c(x, y) = 1$ for all $(x, y) \in \Omega$ then we recover the Poisson matrix.

In general we cannot write \mathbf{A} as a matrix equation of the form (4.15). But we can show that \mathbf{A} is symmetric and it is positive definite as long as the function c is positive on Ω . Recall that a matrix \mathbf{A} is positive definite if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \neq 0$.

Theorem 10.1 (Positive definite matrix)

If $c(x, y) > 0$ for $(x, y) \in \Omega$ then the matrix \mathbf{A} given by (10.53) is symmetric positive definite.

Proof.

To each $x \in \mathbb{R}^n$ there corresponds a matrix $\mathbf{V} \in \mathbb{R}^{m \times m}$ such that $x = \text{vec}(\mathbf{V})$. We claim that

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{j=1}^m \sum_{k=0}^m c_{j, k + \frac{1}{2}} (v_{j, k+1} - v_{j, k})^2 + \sum_{k=1}^m \sum_{j=0}^m c_{j + \frac{1}{2}, k} (v_{j+1, k} - v_{j, k})^2, \tag{10.54}$$

where $v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0$ for $j, k = 0, 1, \dots, m+1$. Since $c_{j + \frac{1}{2}, k}$ and $c_{j, k + \frac{1}{2}}$ correspond to values of c in Ω for the values of j, k in the sums it follows that they are positive and from (10.54) we see that $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for all $x \in \mathbb{R}^n$. Moreover if $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$ then all quadratic factors are zero and $v_{j, k+1} = v_{j, k}$ for $k = 0, 1, \dots, m$ and $j = 1, \dots, m$. Now $v_{j,0} = v_{j,m+1} = 0$ implies that $\mathbf{V} = \mathbf{0}$ and hence $x = 0$. Thus \mathbf{A} is symmetric positive definite.

It remains to prove (10.54). From the connection between (10.52) and (10.53)

we have

$$\begin{aligned}
 \mathbf{x}^T \mathbf{A} \mathbf{x} &= \sum_{j=1}^m \sum_{k=1}^m -(\mathbf{P}_h v)_{j,k} v_{j,k} \\
 &= \sum_{j=1}^m \sum_{k=1}^m \left(c_{j,k-\frac{1}{2}} v_{j,k}^2 + c_{j-\frac{1}{2},k} v_{j,k}^2 + c_{j+\frac{1}{2},k} v_{j,k}^2 + c_{j,k+\frac{1}{2}} v_{j,k}^2 \right. \\
 &\quad \left. - c_{j,k-\frac{1}{2}} v_{j,k-1} v_{j,k} - c_{j,k+\frac{1}{2}} v_{j,k} v_{j,k+1} \right. \\
 &\quad \left. - c_{j-\frac{1}{2},k} v_{j-1,k} v_{j,k} - c_{j+\frac{1}{2},k} v_{j,k} v_{j+1,k} \right).
 \end{aligned}$$

Using the homogenous boundary conditions we have

$$\begin{aligned}
 \sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k}^2 &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1}^2, \\
 \sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k-1} v_{j,k} &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1} v_{j,k}, \\
 \sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j,k}^2 &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k}^2, \\
 \sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j-,k} v_{j,k} &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k} v_{j,k}.
 \end{aligned}$$

It follows that

$$\begin{aligned}
 \mathbf{x}^T \mathbf{A} \mathbf{x} &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} (v_{j,k}^2 + v_{j,k+1}^2 - 2v_{j,k} v_{j,k+1}) \\
 &\quad + \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} (v_{j,k}^2 + v_{j+1,k}^2 - 2v_{j,k} v_{j+1,k})
 \end{aligned}$$

and (10.54) follows. \square \square

10.7.2 Applying Preconditioning

Consider solving $\mathbf{Ax} = \mathbf{b}$, where \mathbf{A} is given by (10.53) and $\mathbf{b} \in \mathbb{R}^n$. Since \mathbf{A} is positive definite it is nonsingular and the system has a unique solution $\mathbf{x} \in \mathbb{R}^n$. Moreover we can use either Cholesky factorization or the block tridiagonal solver

n	2500	10000	22500	40000	62500
K	222	472	728	986	1246
K/\sqrt{n}	4.44	4.72	4.85	4.93	4.98
K_{pre}	22	23	23	23	23

Table 10.2: The number of iterations K (no preconditioning) and K_{pre} (with preconditioning) for the problem (10.50) using the discrete Poisson problem as a preconditioner.

to find \mathbf{x} . Since the bandwidth of \mathbf{A} is $m = \sqrt{n}$ both of these methods require $O(n^2)$ arithmetic operations for large n .

If we choose $c(x, y) \equiv 1$ in (10.50), we get the Poisson problem (10.49). With this in mind, we may think of the coefficient matrix \mathbf{A}_p arising from the discretization of the Poisson problem as an approximation to the matrix (10.53). This suggests using $\mathbf{B} = \mathbf{A}_p^{-1}$, the inverse of the discrete Poisson matrix as a preconditioner for the system (10.51).

Consider Algorithm 10.33. With this preconditioner the calculation $\mathbf{w} = \mathbf{B}\mathbf{t}$ takes the form $\mathbf{A}_p\mathbf{w}_k = \mathbf{t}_k$.

In Section 5.2 we developed a Simple fast Poisson Solver, Cf. Algorithm 5.1. This method can be utilized to solve $\mathbf{A}_p\mathbf{w} = \mathbf{t}$.

Consider the specific problem where

$$c(x, y) = e^{-x+y} \text{ and } f(x, y) = 1.$$

We have used Algorithm 10.18 (conjugate gradient without preconditioning), and Algorithm 10.33 (conjugate gradient with preconditioning) to solve the problem (10.50). We used $\mathbf{x}_0 = 0$ and $\epsilon = 10^{-8}$. The results are shown in Table 10.2.

Without preconditioning the number of iterations still seems to be more or less proportional to \sqrt{n} although the convergence is slower than for the constant coefficient problem. Using preconditioning speeds up the convergence considerably. The number of iterations appears to be bounded independently of n . This illustrates that preconditioning is needed when solving nontrivial problems.

Using a preconditioner increases the work in each iteration. For the present example the number of arithmetic operations in each iteration changes from $O(n)$ without preconditioning to $O(n^{3/2})$ or $O(n \log_2 n)$ with preconditioning. This is not a large increase and both the number of iterations and the computing time is reduced significantly.

Let us finally show that the number $\kappa = \lambda_{max}/\lambda_{min}$ which determines the rate of convergence for the preconditioned conjugate gradient method applied to (10.50) can be bounded independently of n .

Theorem 10.3 (Eigenvalues of preconditioned matrix)

Suppose $0 < c_0 \leq c(x, y) \leq c_1$ for all $(x, y) \in [0, 1]^2$. For the eigenvalues of the matrix $\mathbf{B}\mathbf{A} = \mathbf{A}_p^{-1}\mathbf{A}$ just described we have

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} \leq \frac{c_1}{c_0}.$$

Proof.

Suppose $\mathbf{A}_p^{-1}\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ for some $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$. Then $\mathbf{A}\mathbf{x} = \lambda\mathbf{A}_p\mathbf{x}$. Multiplying this by \mathbf{x}^T and solving for λ we find

$$\lambda = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{A}_p \mathbf{x}}.$$

We computed $\mathbf{x}^T \mathbf{A} \mathbf{x}$ in (10.54) and we obtain $\mathbf{x}^T \mathbf{A}_p \mathbf{x}$ by setting all the c 's there equal to one

$$\mathbf{x}^T \mathbf{A}_p \mathbf{x} = \sum_{i=1}^m \sum_{j=0}^m (v_{i,j+1} - v_{i,j})^2 + \sum_{j=1}^m \sum_{i=0}^m (v_{i+1,j} - v_{i,j})^2.$$

Thus $\mathbf{x}^T \mathbf{A}_p \mathbf{x} > 0$ and bounding all the c 's in (10.54) from below by c_0 and above by c_1 we find

$$c_0(\mathbf{x}^T \mathbf{A}_p \mathbf{x}) \leq \mathbf{x}^T \mathbf{A} \mathbf{x} \leq c_1(\mathbf{x}^T \mathbf{A}_p \mathbf{x})$$

which implies that $c_0 \leq \lambda \leq c_1$ for all eigenvalues λ of $\mathbf{B}\mathbf{A} = \mathbf{A}_p^{-1}\mathbf{A}$. \square

Using $c(x, y) = e^{-x+y}$ as above, we find $c_0 = e^{-2}$ and $c_1 = 1$. Thus $\kappa \leq e^2 \approx 7.4$, a quite acceptable matrix condition which explains the convergence results from our numerical experiment.

10.8 Review Questions

10.8.1 Does the steepest descent and conjugate gradient method always converge?

10.8.2 What kind of orthogonalities occur in the conjugate gradient method?

10.8.3 What is a Krylow space?

10.8.4 What is a convex function?

10.8.5 How do SOR and conjugate gradient compare?

Part IV

Orthonormal Transformations and Least Squares

Chapter 11

Orthonormal and Unitary Transformations

Row operations are used in Gaussian elimination to reduce a matrix to triangular form. This can also be described as transformations by elementary lower triangular matrices (cf. (1.14)). These are not the only kind of transformations that can be used for such a task. In this chapter we study how transformations by orthonormal and unitary matrices can be used to reduce a square matrix to upper triangular form and more generally a rectangular matrix to upper triangular (also called upper trapezoidal) form. This lead to a decomposition of the matrix known as a QR decomposition and a reduced form which we refer to as a QR factorization. The QR decomposition and factorization will be used in later chapters to solve least squares- and eigenvalue problems.

It cannot be repeated too often that orthonormal transformations have the advantage that they preserve the Euclidian norm of a vector, and the spectral norm and Frobenius norm of a matrix, see Lemma 7.22 and Theorem 8.19. This means that when an orthonormal transformation is applied to an inaccurate vector or matrix then the error will not grow. Thus in general an orthonormal transformation is numerically stable.

11.1 The Householder Transformation

Definition 11.1 (Householder Transformation)

A matrix $H \in \mathbb{C}^{n \times n}$ of the form

$$H := I - uu^*, \text{ where } u \in \mathbb{C}^n \text{ and } u^*u = 2$$

is called a **Householder transformation**. The name **elementary reflector** is also used.

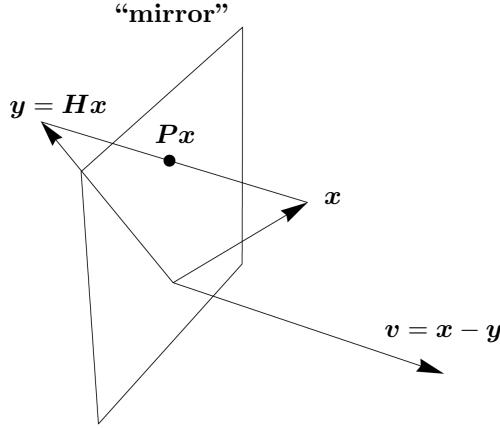


Figure 11.1: The Householder transformation in Exercise 11.2

In the real case and for $n = 2$ we find $\mathbf{H} = \begin{bmatrix} 1-u_1^2 & -u_1 u_2 \\ -u_2 u_1 & 1-u_2^2 \end{bmatrix}$. A Householder transformation is Hermitian and unitary. Indeed, $\mathbf{H}^* = (\mathbf{I} - \mathbf{u}\mathbf{u}^*)^* = \mathbf{H}$ and

$$\mathbf{H}^* \mathbf{H} = \mathbf{H}^2 = (\mathbf{I} - \mathbf{u}\mathbf{u}^*)(\mathbf{I} - \mathbf{u}\mathbf{u}^*) = \mathbf{I} - 2\mathbf{u}\mathbf{u}^* + \mathbf{u}(\mathbf{u}^*\mathbf{u})\mathbf{u}^* = \mathbf{I}.$$

In the real case \mathbf{H} is symmetric and orthonormal.

There are several ways to represent a Householder transformation. Householder used $\mathbf{I} - 2\mathbf{u}\mathbf{u}^*$, where $\mathbf{u}^*\mathbf{u} = 1$. For any nonzero $\mathbf{v} \in \mathbb{R}^n$ the matrix

$$\mathbf{H} := \mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^*}{\mathbf{v}^*\mathbf{v}} \quad (11.1)$$

is a Householder transformation. Indeed, $\mathbf{H} = \mathbf{I} - \mathbf{u}\mathbf{u}^*$, where $\mathbf{u} := \sqrt{2} \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$. Moreover, if $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ with $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$ and $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$ then $\mathbf{H}\mathbf{x} = \mathbf{y}$ (Cf. Exercise 11.2).

Exercise 11.2 (Reflector)

Suppose $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ with $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$ and $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$.

- (a) Show that $\mathbf{H}\mathbf{x} := (\mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{y}$. ¹³
- (b) Show that $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$ implies that $\mathbf{x} - \mathbf{y}$ is orthogonal to $\mathbf{x} + \mathbf{y}$ and conclude that $\mathbf{Px} := \mathbf{x} - \frac{\mathbf{v}^T\mathbf{x}}{\mathbf{v}^T\mathbf{v}}\mathbf{v}$ is the orthogonal projection of \mathbf{x} into the subspace $\text{span}(\mathbf{x} + \mathbf{y})$. The vector \mathbf{y} is the reflected image of \mathbf{x} , where the subspace $\{\mathbf{w} \in \mathbb{R}^n : \mathbf{w}^T\mathbf{v} = 0\}$ acts as a "mirror". See Figure 11.1.

¹³Hint: Show first that $\mathbf{v}^T\mathbf{v} = 2\mathbf{v}^T\mathbf{x}$

- (c) Determine the matrices \mathbf{H}, \mathbf{P} when $\mathbf{x} := [1, 0, 1]^T$ and $\mathbf{y} := [-1, 0, 1]^T$. Check that \mathbf{Px} is the projection into the "mirror" $\{(x, y, z) \in \mathbb{R}^3 : x = 0\}$ and $\mathbf{Hx} = \mathbf{y}$.

A main use of Householder transformations is to produce zeros in vectors.

Theorem 11.3 (Zeros in vectors)

Suppose $\mathbf{x} \in \mathbb{C}^n$ is nonzero and define $\rho \in \mathbb{C}$ and $\mathbf{z}, \mathbf{u} \in \mathbb{C}^n$ by

$$\rho := \begin{cases} x_1/|x_1|, & \text{if } x_1 \neq 0, \\ 1, & \text{otherwise.} \end{cases}, \quad \mathbf{x} =: \rho \|\mathbf{x}\|_2 \mathbf{z}, \quad \mathbf{u} := \frac{\mathbf{z} + \mathbf{e}_1}{\sqrt{1 + z_1}}. \quad (11.2)$$

Then $\mathbf{u}^* \mathbf{u} = 2$ and

$$\mathbf{Hx} := (\mathbf{I} - \mathbf{u}\mathbf{u}^*)\mathbf{x} = a\mathbf{e}_1, \quad a := -\rho \|\mathbf{x}\|_2. \quad (11.3)$$

Proof. Since $|\rho| = 1$ and $\|\mathbf{x}\|_2 = |\rho| \|\mathbf{x}\|_2 \|\mathbf{z}\|_2$ it follows that $\|\mathbf{z}\|_2 = 1$. Moreover, $z_1 = |x_1|/\|\mathbf{x}\|_2$ is real so that $\mathbf{u}^* \mathbf{u} = \frac{(\mathbf{z} + \mathbf{e}_1)^*(\mathbf{z} + \mathbf{e}_1)}{1 + z_1} = \frac{2 + 2z_1}{1 + z_1} = 2$. Finally,

$$\begin{aligned} \mathbf{Hx} &= \mathbf{x} - (\mathbf{u}^* \mathbf{x})\mathbf{u} = \rho \|\mathbf{x}\|_2 (\mathbf{z} - (\mathbf{u}^* \mathbf{z})\mathbf{u}) = \rho \|\mathbf{x}\|_2 (\mathbf{z} - \frac{(\mathbf{z}^* + \mathbf{e}_1^*)\mathbf{z}}{1 + z_1} (\mathbf{z} + \mathbf{e}_1)) \\ &= \rho \|\mathbf{x}\|_2 (\mathbf{z} - (\mathbf{z} + \mathbf{e}_1)) = -\rho \|\mathbf{x}\|_2 \mathbf{e}_1 = a\mathbf{e}_1. \end{aligned}$$

□

The formulas in Theorem 11.3 are implemented in the following algorithm adapted from [22]. To any given $\mathbf{x} \in \mathbb{C}^n$ a number a and a vector \mathbf{u} with $\mathbf{u}^* \mathbf{u} = 2$ is computed so that $(\mathbf{I} - \mathbf{u}\mathbf{u}^*)\mathbf{x} = a\mathbf{e}_1$.

Algorithm 11.4 (Generate a Householder transformation)

```

1 function [ u , a ]=housegen( x )
2 a=norm( x );
3 if a==0
4     u=x; u( 1 )=sqrt( 2 ); return;
5 end
6 if x( 1 )== 0
7     r=1;
8 else
9     r=x( 1 )/abs( x( 1 ) );
10 end
11 u=conj( r )*x/a;
12 u( 1 )=u( 1 )+1;
13 u=u/sqrt( u( 1 ) );
14 a=-r*a;
15 end

```

Note that

- If $\mathbf{x} = \mathbf{0}$ then any \mathbf{u} with $\|\mathbf{u}\|_2 = \sqrt{2}$ can be used in the Householder transformation. In the algorithm we use $\mathbf{u} = \sqrt{2}\mathbf{e}_1$ in this case.
- In Theorem 11.3 the first component of \mathbf{z} is $z_1 = |x_1|/\|\mathbf{x}\|_2 \geq 0$. Since $\|\mathbf{z}\|_2 = 1$ we have $1 \leq 1 + z_1 \leq 2$. It follows that \mathbf{u} is well defined and we avoid cancelation error when computing $1 + z_1$.

Exercise 11.5 (What does algorithm housegen do when $\mathbf{x} = \mathbf{e}_1$?)

Determine \mathbf{H} in Algorithm 11.4 when $\mathbf{x} = \mathbf{e}_1$.

Householder transformations can also be used to zero out only the lower part of a vector. Suppose $\mathbf{y} \in \mathbb{C}^k$, $\mathbf{z} \in \mathbb{C}^{n-k}$. Let $\hat{\mathbf{u}}$ and a be the output of Algorithm 11.4 called with $\mathbf{x} = \mathbf{z}$, i.e., $[\hat{\mathbf{u}}, a] = \text{housegen}(\mathbf{z})$ and set $\mathbf{u}^T = [\mathbf{0}^T, \hat{\mathbf{u}}^T] \in \mathbb{R}^n$. Then

$$\mathbf{H} = \mathbf{I} - \mathbf{u}\mathbf{u}^* = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} - \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{u}} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \hat{\mathbf{u}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{H}} \end{bmatrix},$$

where $\hat{\mathbf{H}} = \mathbf{I} - \hat{\mathbf{u}}\hat{\mathbf{u}}^*$. Since $\mathbf{u}^*\mathbf{u} = \hat{\mathbf{u}}^*\hat{\mathbf{u}} = 2$ we see that \mathbf{H} and $\hat{\mathbf{H}}$ are Householder transformations.

Exercise 11.6 (Examples of Householder transformations)

If $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ with $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$ and $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$ then it follows from Exercise 11.2 that $(\mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{y}$. Use this to construct a Householder transformation \mathbf{H} such that $\mathbf{H}\mathbf{x} = \mathbf{y}$ in the following cases.

a) $\mathbf{x} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}.$

b) $\mathbf{x} = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 0 \\ 3 \\ 0 \end{bmatrix}.$

Exercise 11.7 (2×2 Householder transformation)

Show that a real 2×2 Householder transformation can be written in the form

$$\mathbf{H} = \begin{bmatrix} -\cos \phi & \sin \phi \\ \sin \phi & \cos \phi \end{bmatrix}.$$

Find $\mathbf{H}\mathbf{x}$ if $\mathbf{x} = [\cos \phi, \sin \phi]^T$.

11.2 Householder Triangulation

We say that a matrix $\mathbf{R} \in \mathbb{C}^{m \times n}$ is **upper trapezoidal**, if $r_{i,j} = 0$ for $j < i$ and $i = 1, 2, \dots, m$. Here are three upper trapezoidal matrices corresponding to $m < n$, $m = n$, and $m > n$.

$$\begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix}, \quad \begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix}, \quad \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{bmatrix}.$$

In this section we consider a method for bringing a matrix to upper trapezoidal form using Householder transformations. We treat the cases $m > n$ and $m \leq n$ separately and consider first $m > n$. We describe how to find a sequence $\mathbf{H}_1, \dots, \mathbf{H}_n$ of Householder transformations such that

$$\mathbf{A}_{n+1} := \mathbf{H}_n \mathbf{H}_{n-1} \cdots \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{R},$$

and where \mathbf{R}_1 is upper triangular. We define

$$\mathbf{A}_1 := \mathbf{A}, \quad \mathbf{A}_{k+1} = \mathbf{H}_k \mathbf{A}_k, \quad k = 1, 2, \dots, n.$$

Suppose \mathbf{A}_k is upper triangular in its first $k - 1$ columns (which is true for $k = 1$)

$$\mathbf{A}_k = \left[\begin{array}{ccc|cccccc} a_{1,1}^1 & \cdots & a_{1,k-1}^1 & a_{1,k}^1 & \cdots & a_{1,j}^1 & \cdots & a_{1,n}^1 \\ \ddots & & \vdots & \vdots & & \vdots & & \vdots \\ & a_{k-1,k-1}^{k-1} & & a_{k-1,k}^{k-1} & \cdots & a_{k-1,j}^{k-1} & \cdots & a_{k-1,n}^{k-1} \\ \hline & & a_{k,k}^k & \cdots & a_{k,j}^k & \cdots & a_{k,n}^k & \\ & & \vdots & & \vdots & & \vdots & \\ & a_{i,k}^k & \cdots & a_{i,j}^k & \cdots & a_{i,n}^k & & \\ & \vdots & & \vdots & & \vdots & & \\ & a_{m,k}^k & \cdots & a_{m,j}^k & \cdots & a_{m,n}^k & & \end{array} \right] \quad (11.4)$$

$$= \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \mathbf{D}_k \end{bmatrix}.$$

Let $\hat{\mathbf{H}}_k := \mathbf{I} - \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^*$ be a Householder transformation that maps the first column $[a_{k,k}^k, \dots, a_{m,k}^k]^T$ of \mathbf{D}_k to a multiple of \mathbf{e}_1 , $\hat{\mathbf{H}}_k(\mathbf{D}_k \mathbf{e}_1) = a_k \mathbf{e}_1$. Using Algorithm 11.4 we have $[\hat{\mathbf{u}}_k, a_k] = \text{housegen}(\mathbf{D}_k \mathbf{e}_1)$. Set $\mathbf{H}_k := \begin{bmatrix} \mathbf{I}_{k-1} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{H}}_k \end{bmatrix}$. Then

$$\mathbf{A}_{k+1} := \mathbf{H}_k \mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \hat{\mathbf{H}}_k \mathbf{D}_k \end{bmatrix} = \begin{bmatrix} \mathbf{B}_{k+1} & \mathbf{C}_{k+1} \\ \mathbf{0} & \mathbf{D}_{k+1} \end{bmatrix},$$

where $\mathbf{B}_{k+1} \in \mathbb{C}^{k \times k}$ is upper triangular and $\mathbf{D}_{k+1} \in \mathbb{C}^{(m-k) \times (n-k)}$. Thus \mathbf{A}_{k+1} is upper triangular in its first k columns and the reduction has been carried one step further. At the end $\mathbf{R} := \mathbf{A}_{n+1} = [\mathbf{R}_1 \mathbf{0}]$, where \mathbf{R}_1 is upper triangular.

The process can also be applied to $\mathbf{A} \in \mathbb{C}^{m \times n}$ if $m \leq n$. In this case $m - 1$ Householder transformations will suffice and $\mathbf{H}_{m-1} \cdots \mathbf{H}_1 \mathbf{A}$ is upper trapezoidal.

In an algorithm we can store most of the vectors $\hat{\mathbf{u}}_k = [u_{kk}, \dots, u_{mk}]^T$ and \mathbf{A}_k in \mathbf{A} . However, the elements $u_{k,k}$ and $a_k = r_{k,k}$ have to compete for the diagonal in \mathbf{A} . For $m = 4$ and $n = 3$ the two possibilities look as follows:

$$\mathbf{A} = \begin{bmatrix} u_{11} & r_{12} & r_{13} \\ u_{21} & u_{22} & r_{23} \\ u_{31} & u_{32} & u_{33} \\ u_{41} & u_{42} & u_{43} \end{bmatrix} \quad \text{or} \quad \mathbf{A} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ u_{21} & r_{22} & r_{23} \\ u_{31} & u_{32} & r_{33} \\ u_{41} & u_{42} & u_{43} \end{bmatrix}.$$

Whatever alternative is chosen the loser, if needed, has to be stored in a separate vector. In the following algorithm we store a_k in \mathbf{A} . We also apply the Householder transformations to a second matrix \mathbf{B} . The algorithm can then be used to solve linear systems with several right hand sides.

Algorithm 11.8 (Householder Triangulation)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$, $\mathbf{B} \in \mathbb{C}^{m \times r}$ and let $s := \min(n, m - 1)$. The algorithm uses `housegen` to compute Householder transformations $\mathbf{H}_1, \dots, \mathbf{H}_s$ such that $\mathbf{R} = \mathbf{H}_s \cdots \mathbf{H}_1 \mathbf{A}$ is upper trapezoidal and $\mathbf{C} = \mathbf{H}_s \cdots \mathbf{H}_1 \mathbf{B}$. If \mathbf{B} is the empty matrix then \mathbf{C} is the empty matrix with m rows and 0 columns.

```

1 function [R,C] = housetriang(A,B)
2 [m,n]=size(A); r=size(B,2); A=[A,B];
3 for k=1:min(n,m-1)
4   [v,A(k:k)] = housegen(A(k:m,k));
5   C=A(k:m,k+1:n+r); A(k:m,k+1:n+r)=C-v*(v'*C);
6 end
7 R=triu(A(:,1:n)); C=A(:,n+1:n+r);

```

Here $v = \hat{\mathbf{u}}_k$ and we have used $\hat{\mathbf{H}}_k \mathbf{C} = (\mathbf{I} - vv^*)\mathbf{C} = \mathbf{C} - v(v^*\mathbf{C})$ for the update. The Matlab command `triu` extracts the upper triangular part of \mathbf{A} putting zeros in rows $n + 1, \dots, m$.

The algorithm can be used to solve linear systems and least squares problem.

11.2.1 Solving Linear Systems using Unitary Transformations

Consider now the linear system $\mathbf{Ax} = \mathbf{b}$, where \mathbf{A} is square. Using Algorithm 11.8 we obtain an upper triangular system $\mathbf{Rx} = \mathbf{c}$ that is nonsingular if \mathbf{A} is non-singular. Thus, it can be solved by back substitution and we have a method for solving linear systems that is an alternative to Gaussian elimination. The two

methods are similar since they both reduce \mathbf{A} to upper triangular form using certain transformations.

Which method is better? Here is a short discussion

- Advantages with Householder:
 - Always works for nonsingular systems.
 - Row interchanges are not necessary, but see [4].
 - Numerically stable.
- Advantages with Gauss
 - Half the number of arithmetic operations compared to Householder.
 - Row interchanges are often not necessary.
 - Usually stable (but no guarantee).

Linear systems can be constructed where Gaussian elimination will fail numerically even if row interchanges are used, see [30]. On the other hand the transformations used in Householder triangulation are unitary so the method is quite stable. So why is Gaussian elimination more popular than Householder triangulation? One reason is that the number of arithmetic operations in (11.5) when $m = n$ is $4n^3/3 = 2G_n$, which is twice the number for Gaussian elimination. We show this below. Numerical stability can be a problem with Gaussian elimination, but years and years of experience shows that it works well for most practical problems and pivoting is often not necessary. Also Gaussian elimination often wins for banded and sparse problems.

11.2.2 The number of Arithmetic Operations

The bulk of the work in Algorithm 11.8 is the computation of $\mathbf{C} - \mathbf{v} * (\mathbf{v}^T * \mathbf{C})$ for each k . In the real case it can be determined from the following lemma.

Lemma 11.9 (Updating a Householder transformation)

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{u} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^n$. The computation of $\mathbf{A} - \mathbf{u}(\mathbf{u}^T \mathbf{A})$ and $\mathbf{A} - (\mathbf{A}\mathbf{v})\mathbf{v}^T$ both cost approximately $4mn$ arithmetic operations.

Proof. It costs $2mn$ arithmetic operations to compute $\mathbf{w}^T := \mathbf{u}^T \mathbf{A}$, mn arithmetic operations to compute $\mathbf{W} = \mathbf{u}\mathbf{w}^T$ and mn arithmetic operations for the final subtraction $\mathbf{A} - \mathbf{W}$, a total of $4mn$ arithmetic operations. Taking the transpose we obtain the same count for $\mathbf{A} - (\mathbf{A}\mathbf{v})\mathbf{v}^T$. \square

Since in Algorithm 11.8, $\mathbf{C} \in \mathbb{C}^{(m-k+1) \times (n+r-k)}$ and $m \geq n$ the cost of computing the update $\mathbf{C} - \mathbf{v} * (\mathbf{v}^T * \mathbf{C})$ is $4(m-k)(n+r-k)$ arithmetic operations.

This implies that the work in Algorithm 11.8 can be estimated as

$$\int_0^n 4(m-k)(n+r-k)dk = 2m(n+r)^2 - \frac{2}{3}(n+r)^3. \quad (11.5)$$

For $m = n$ and $r = 0$ this gives $4n^3/3$ for the number of arithmetic operations to bring a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ to upper triangular form using Householder transformations.

11.3 The QR Decomposition and QR Factorization

Gaussian elimination without row interchanges gives a LU-factorization $\mathbf{A} = \mathbf{L}\mathbf{U}$ of $\mathbf{A} \in \mathbb{R}^{n \times n}$. Consider Householder triangulation of \mathbf{A} . Applying Algorithm 11.8 gives $\mathbf{R} = \mathbf{H}_{n-1} \cdots \mathbf{H}_1 \mathbf{A}$ implying the factorization $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where $\mathbf{Q} = \mathbf{H}_1 \cdots \mathbf{H}_{n-1}$ is orthonormal and \mathbf{R} is upper triangular. This is known as a QR-factorization of \mathbf{A} .

For a rectangular matrix we define the following.

Definition 11.10 (QR decomposition)

Let $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m, n \in \mathbb{N}$. We say that $\mathbf{A} = \mathbf{Q}\mathbf{R}$ is a **QR decomposition** of \mathbf{A} if $\mathbf{Q} \in \mathbb{C}^{m \times m}$ is square and unitary and \mathbf{R} is upper trapezoidal. If $m \geq n$ then \mathbf{R} takes the form

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix}$$

where $\mathbf{R}_1 \in \mathbb{C}^{n \times n}$ is upper triangular and $\mathbf{0}_{m-n,n} \in \mathbb{C}^{m-n,n}$ is the zero matrix. For $m \geq n$ we call $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ a **QR factorization** of \mathbf{A} if $\mathbf{Q}_1 \in \mathbb{C}^{m \times n}$ has orthonormal columns and $\mathbf{R}_1 \in \mathbb{C}^{n \times n}$ is upper triangular.

A QR factorization is obtained from a QR decomposition $\mathbf{A} = \mathbf{Q}\mathbf{R}$ by simply using the first n columns of \mathbf{Q} and the first n rows of \mathbf{R} . Indeed, if we partition \mathbf{Q} as $[\mathbf{Q}_1, \mathbf{Q}_2]$ and $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$, where $\mathbf{Q}_1 \in \mathbb{R}^{m \times n}$ and $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ then $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ is a QR factorization of \mathbf{A} . On the other hand a QR factorization $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ of \mathbf{A} can be turned into a QR decomposition by extending the set of columns $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ of \mathbf{Q}_1 into an orthonormal basis $\{\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{q}_{n+1}, \dots, \mathbf{q}_m\}$ for \mathbb{R}^m and adding $m - n$ rows of zeros to \mathbf{R}_1 . We then obtain the QR decomposition $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_m]$ and $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$.

Example 11.11 (QR decomposition and factorization)

An example of a QR decomposition is

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & -1 & -1 \\ 1 & 1 & 1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix} \times \begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \\ 0 & 0 & 0 \end{bmatrix} = \mathbf{Q}\mathbf{R},$$

while a QR factorization $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ is obtained by dropping the last column of \mathbf{Q} and the last row of \mathbf{R} , so that

$$\mathbf{A} = \frac{1}{2} \begin{bmatrix} 1 & 1 & -1 \\ 1 & 1 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1.$$

Consider existence and uniqueness.

Theorem 11.12 (Existence of QR decomposition)

Any matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m, n \in \mathbb{N}$ has a QR decomposition. If $m \geq n$ and \mathbf{A} is real then the QR factorization is unique if \mathbf{A} has linearly independent columns and \mathbf{R} has positive diagonal elements.

Proof. The function `housegen(x)` returns a Householder transformation for any $\mathbf{x} \in \mathbb{C}^n$. Thus with $\mathbf{B} = \mathbf{I}$ in Algorithm 11.8 we obtain a QR decomposition $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where $\mathbf{Q} = \mathbf{C}^* = \mathbf{H}_1 \cdots \mathbf{H}_s$, is unitary. Thus a QR factorization always exists.

For uniqueness, if $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ is a QR factorization of \mathbf{A} and \mathbf{R}_1 has positive diagonal elements then $\mathbf{A}^T \mathbf{A} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{R}_1^T \mathbf{R}_1$ is the Cholesky factorization of $\mathbf{A}^T \mathbf{A}$. Since the Cholesky factorization is unique it follows that \mathbf{R}_1 is unique and since necessarily $\mathbf{Q}_1 = \mathbf{A}\mathbf{R}_1^{-1}$, it must also be unique. \square

Example 11.13 (QR decomposition and factorization)

Consider finding the QR decomposition and factorization of the matrix $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$ using the method of the uniqueness proof of Theorem 11.12. We find $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix}$. The Cholesky factorization of $\mathbf{B} = \mathbf{R}^T \mathbf{R}$ is given by $\mathbf{R} = \frac{1}{\sqrt{5}} \begin{bmatrix} 5 & -4 \\ 0 & 3 \end{bmatrix}$. Now $\mathbf{R}^{-1} = \frac{1}{3\sqrt{5}} \begin{bmatrix} 3 & 4 \\ 0 & 5 \end{bmatrix}$ so $\mathbf{Q} = \mathbf{A}\mathbf{R}^{-1} = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}$. Since \mathbf{A} is square $\mathbf{A} = \mathbf{Q}\mathbf{R}$ is both the QR decomposition and QR factorization of \mathbf{A} .

The QR factorization can be used to prove a classical determinant inequality.

Theorem 11.14 (Hadamard's Inequality)

For any $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{n \times n}$ we have

$$|\det(\mathbf{A})| \leq \prod_{j=1}^n \|\mathbf{a}_j\|_2. \quad (11.6)$$

Equality holds if and only if \mathbf{A} has a zero column or the columns of \mathbf{A} are orthogonal.

Proof. Let $\mathbf{A} = \mathbf{QR}$ be a QR factorization of \mathbf{A} . Since

$$1 = \det(\mathbf{I}) = \det(\mathbf{Q}^*\mathbf{Q}) = \det(\mathbf{Q}^*)\det(\mathbf{Q}) = \det(\mathbf{Q})^*\det(\mathbf{Q}) = |\det(\mathbf{Q})|^2$$

we have $|\det(\mathbf{Q})| = 1$. Let $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]$. Then $(\mathbf{A}^*\mathbf{A})_{jj} = \|\mathbf{a}_j\|_2^2 = (\mathbf{R}^*\mathbf{R})_{jj} = \|\mathbf{r}_j\|_2^2$, and

$$|\det(\mathbf{A})| = |\det(\mathbf{QR})| = |\det(\mathbf{R})| = \prod_{j=1}^n |r_{jj}| \leq \prod_{j=1}^n \|\mathbf{r}_j\|_2 = \prod_{j=1}^n \|\mathbf{a}_j\|_2.$$

The inequality is proved. If equality holds then either $\det(\mathbf{A}) = 0$ and \mathbf{A} has a zero column, or $\det(\mathbf{A}) \neq 0$ and $r_{jj} = \|\mathbf{r}_j\|_2$ for $j = 1, \dots, n$. This happens if and only if \mathbf{R} is diagonal. But then $\mathbf{A}^*\mathbf{A} = \mathbf{R}^*\mathbf{R}$ is diagonal, which means that the columns of \mathbf{A} are orthogonal. \square

Exercise 11.15 (QR decomposition)

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{Q} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} 2 & 2 \\ 0 & 2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Show that \mathbf{Q} is orthonormal and that \mathbf{QR} is a QR decomposition of \mathbf{A} . Find a QR factorization of \mathbf{A} .

Exercise 11.16 (Householder triangulation)

a) Let

$$\mathbf{A} := \begin{bmatrix} 1 & 0 & 1 \\ -2 & -1 & 0 \\ 2 & 2 & 1 \end{bmatrix}.$$

Find Householder transformations $\mathbf{H}_1, \mathbf{H}_2 \in \mathbb{R}^{3 \times 3}$ such that $\mathbf{H}_2\mathbf{H}_1\mathbf{A}$ is upper triangular.

b) Find the QR factorization of \mathbf{A} where \mathbf{R} has positive diagonal elements.

11.3.1 QR and Gram-Schmidt

The Gram-Schmidt orthogonalization of the columns of \mathbf{A} can be used to find the QR factorization of \mathbf{A} .

Theorem 11.17 (QR and Gram-Schmidt)

Suppose $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m \times n}$ has rank n and define

$$\mathbf{v}_1 = \mathbf{a}_1, \quad \mathbf{v}_j = \mathbf{a}_j - \sum_{i=1}^{j-1} \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \mathbf{v}_i, \quad \text{for } j = 2, \dots, n. \quad (11.7)$$

Let

$$\mathbf{Q}_1 := [\mathbf{q}_1, \dots, \mathbf{q}_n], \quad \mathbf{q}_j = \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|_2}, \quad j = 1, \dots, n,$$

$$\mathbf{R}_1 := \begin{bmatrix} \|\mathbf{v}_1\|_2 & \mathbf{a}_2^T \mathbf{q}_1 & \mathbf{a}_3^T \mathbf{q}_1 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_1 & \mathbf{a}_n^T \mathbf{q}_1 \\ 0 & \|\mathbf{v}_2\|_2 & \mathbf{a}_3^T \mathbf{q}_2 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_2 & \mathbf{a}_n^T \mathbf{q}_2 \\ 0 & 0 & \|\mathbf{v}_3\|_2 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_3 & \mathbf{a}_n^T \mathbf{q}_3 \\ \ddots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \ddots & \ddots & \|\mathbf{v}_{n-1}\|_2 & \mathbf{a}_n^T \mathbf{q}_{n-1} \\ 0 & 0 & 0 & \|\mathbf{v}_n\|_2 & \mathbf{a}_n^T \mathbf{q}_n \end{bmatrix}. \quad (11.8)$$

Then $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ is the unique QR factorization of \mathbf{A} .

Proof. Let \mathbf{Q}_1 and \mathbf{R}_1 be given by (11.8). The matrix \mathbf{Q}_1 is well defined and has orthonormal columns, since $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ is an orthonormal basis for $\text{span}(\mathbf{A})$ by Theorem 0.38. By (11.7)

$$\mathbf{a}_j = \mathbf{v}_j + \sum_{i=1}^{j-1} \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \mathbf{v}_i = r_{jj} \mathbf{q}_j + \sum_{i=1}^{j-1} \mathbf{q}_i r_{ij} = \mathbf{Q}_1 \mathbf{R}_1 \mathbf{e}_j, \quad j = 1, \dots, n.$$

Clearly \mathbf{R}_1 has positive diagonal elements and the factorization is unique. \square

Example 11.18 (QR using Gram-Schmidt)

Consider finding the QR decomposition and factorization of the matrix $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = [\mathbf{a}_1, \mathbf{a}_2]$ using Gram-Schmidt. Using (11.7) we find $\mathbf{v}_1 = \mathbf{a}_1$ and $\mathbf{v}_2 = \mathbf{a}_2 - \frac{\mathbf{a}_2^T \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1 = \frac{3}{5} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$. Thus $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2]$, where $\mathbf{q}_1 = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 \\ -1 \end{bmatrix}$ and $\mathbf{q}_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$. By (11.8) we find

$$\mathbf{R}_1 = \mathbf{R} = \begin{bmatrix} \|\mathbf{v}_1\|_2 & \mathbf{a}_2^T \mathbf{q}_1 \\ 0 & \|\mathbf{v}_2\|_2 \end{bmatrix} = \frac{1}{\sqrt{5}} \begin{bmatrix} 5 & -4 \\ 0 & 3 \end{bmatrix}$$

and this agrees with what we found in Example 11.13.

Exercise 11.19 (QR using Gram-Schmidt, II)

Construct \mathbf{Q}_1 and \mathbf{R}_1 in Example 11.11 using Gram-Schmidt orthogonalization.

The Gram-Schmidt orthogonalization process should not be used to compute the QR factorization numerically. The columns of \mathbf{Q}_1 computed in floating point arithmetic using Gram-Schmidt orthogonalization will often be far from orthogonal. There is a modified version of Gram-Schmidt which behaves better numerically, but this will not be considered here, see [2]. Instead we consider Householder transformations.

11.4 Givens Rotations

In some applications, the matrix we want to triangulate has a special structure. Suppose for example that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is square and upper Hessenberg as illustrated by a **Wilkinson diagram** for $n = 4$

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix}.$$

Only one element in each column needs to be annihilated and a full Householder transformation will be inefficient. In this case we can use a simpler transformation.

Definition 11.20 (Givens rotation, plane rotation)

A **plane rotation** (also called a **Given's rotation**) is a matrix $\mathbf{P} \in \mathbb{R}^{2,2}$ of the form

$$\mathbf{P} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \text{ where } c^2 + s^2 = 1.$$

A plane rotation is orthonormal and there is a unique angle $\theta \in [0, 2\pi)$ such that $c = \cos \theta$ and $s = \sin \theta$. Moreover, the identity matrix is a plane rotation corresponding to $\theta = 0$.

Exercise 11.21 (Plane rotation)

Show that if $\mathbf{x} = \begin{bmatrix} r \cos \alpha \\ r \sin \alpha \end{bmatrix}$ then $\mathbf{P}\mathbf{x} = \begin{bmatrix} r \cos(\alpha - \theta) \\ r \sin(\alpha - \theta) \end{bmatrix}$. Thus \mathbf{P} rotates a vector \mathbf{x} in the plane an angle θ clockwise. See Figure 11.2.

Suppose

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \neq \mathbf{0}, \quad c := \frac{x_1}{r}, \quad s := \frac{x_2}{r}, \quad r := \|\mathbf{x}\|_2.$$

Then

$$\mathbf{P}\mathbf{x} = \frac{1}{r} \begin{bmatrix} x_1 & x_2 \\ -x_2 & x_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{1}{r} \begin{bmatrix} x_1^2 + x_2^2 \\ 0 \end{bmatrix} = \begin{bmatrix} r \\ 0 \end{bmatrix},$$

and we have introduced a zero in \mathbf{x} . We can take $\mathbf{P} = \mathbf{I}$ when $\mathbf{x} = \mathbf{0}$.

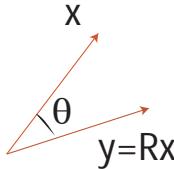


Figure 11.2: A plane rotation.

For an n -vector $\mathbf{x} \in \mathbb{R}^n$ and $1 \leq i < j \leq n$ we define a **rotation in the i,j -plane** as a matrix $\mathbf{P}_{ij} = (p_{kl}) \in \mathbb{R}^{n \times n}$ by $p_{kl} = \delta_{kl}$ except for positions ii, jj, ij, ji , which are given by

$$\begin{bmatrix} p_{ii} & p_{ij} \\ p_{ji} & p_{jj} \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \text{ where } c^2 + s^2 = 1.$$

Thus, for $n = 4$,

$$\mathbf{P}_{1,2} = \begin{bmatrix} c & s & 0 & 0 \\ -s & c & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{P}_{13} = \begin{bmatrix} c & 0 & s & 0 \\ 0 & 1 & 0 & 0 \\ -s & 0 & c & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{P}_{23} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & s & c & 0 \\ 0 & -s & c & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Premultiplying a matrix by a rotation in the i,j -plane changes only rows i and j of the matrix, while postmultiplying the matrix by such a rotation only changes column i and j . In particular, if $\mathbf{B} = \mathbf{P}_{ij}\mathbf{A}$ and $\mathbf{C} = \mathbf{A}\mathbf{P}_{ij}$ then $\mathbf{B}(k, :) = \mathbf{A}(k, :)$, $\mathbf{C}(:, k) = \mathbf{A}(:, k)$ for all $k \neq i, j$ and

$$\begin{bmatrix} \mathbf{B}(i, :) \\ \mathbf{B}(j, :) \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \mathbf{A}(i, :) \\ \mathbf{A}(j, :) \end{bmatrix}, \quad [\mathbf{C}(:, i) \ \mathbf{C}(:, j)] = [\mathbf{A}(:, i) \ \mathbf{A}(:, j)] \begin{bmatrix} c & s \\ -s & c \end{bmatrix}. \quad (11.9)$$

An upper Hessenberg matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be transformed to upper triangular form using rotations $\mathbf{P}_{i,i+1}$ for $i = 1, \dots, n-1$. For $n = 4$ the process can be illustrated as follows.

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{12}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ \mathbf{0} & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{23}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{34}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & 0 & r_{33} & r_{34} \\ 0 & 0 & 0 & r_{44} \end{bmatrix}.$$

For an algorithm see Exercise 11.22.

Exercise 11.22 (Solving upper Hessenberg system using roations)

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be upper Hessenberg and nonsingular, and let $\mathbf{b} \in \mathbb{R}^n$. The following algorithm solves the linear system $\mathbf{Ax} = \mathbf{b}$ using rotations $\mathbf{P}_{k,k+1}$ for $k = 1, \dots, n-1$. Determine the number of arithmetic operations of this algorithm.

Algorithm 11.23 (Upper Hessenberg linear system)

Suppose $A \in \mathbb{R}^{n \times n}$ is nonsingular and upper Hessenberg and that $b \in \mathbb{R}^n$. This algorithm uses Given's rotations to solve the linear system $Ax = b$. It uses Algorithm 1.13.

```

1 function x=rothestri(A,b)
2 n=length(A); A=[A b];
3 for k=1:n-1
4     r=norm([A(k,k) ,A(k+1,k)]) ;
5     if r>0
6         c=A(k,k)/r; s=A(k+1,k)/r ;
7         A([k k+1],k+1:n+1)=[c s;-s c]*A([k k+1],k+1:n+1);
8     end
9     A(k,k)=r; A(k+1,k)=0;
10 end
11 x=backsolve(A(:,1:n),A(:,n+1));

```

11.5 Review Questions

11.5.1 What is a Householder transformation?

11.5.2 Why are they good for numerical work?

11.5.3 What are the main differences between solving a linear system by Gaussian elimination and Householder transformations?

11.5.4 What are the differences between a QR decomposition and a QR factorization?

11.5.5 Does any matrix have a QR decomposition?

11.5.6 What is a Givens transformation?

Chapter 12

Least Squares

In this chapter $m, n \in \mathbb{N}$, $\mathbf{A} \in \mathbb{C}^{m \times n}$ and $\mathbf{b} \in \mathbb{C}^m$, unless otherwise specified. Consider the linear system $\mathbf{Ax} = \mathbf{b}$ of m equations in n unknowns. It is over-determined, if $m > n$, square, if $m = n$, and underdetermined, if $m < n$. In either case the system can only be solved approximately if $\mathbf{b} \notin \text{span}(\mathbf{A})$. One way to solve $\mathbf{Ax} = \mathbf{b}$ approximately is to select a vector norm $\|\cdot\|$ and look for $\mathbf{x} \in \mathbb{C}^n$ which minimizes $\|\mathbf{Ax} - \mathbf{b}\|$. The choice $\|\cdot\| = \|\cdot\|_2$, the Euclidean norm, is particularly convenient since it leads to a linear system. Only this norm is considered here.

Definition 12.1 (Least squares problem)

To find $\mathbf{x} \in \mathbb{C}^n$ that minimizes $E : \mathbb{C}^n \rightarrow \mathbb{R}$ given by

$$E(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|_2^2,$$

is called a **least squares problem**. A minimizer \mathbf{x} is called a **least squares solution**.

Since the square root function is monotone, minimizing $E(\mathbf{x})$ or $\sqrt{E(\mathbf{x})}$ is equivalent.

12.1 Existence, Uniqueness, and Characterization

We first consider existence, uniqueness and characterization of solutions.

Theorem 12.2 (Existence)

The least squares problem always has a solution.

Proof. We use the orthogonal column space decomposition $\mathbb{C}^m = \text{span}(\mathbf{A}) \overset{\perp}{\oplus} \ker(\mathbf{A}^*)$, in Theorem 0.43 and write $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$, where $\mathbf{b}_1 \in \text{span}(\mathbf{A})$ and $\mathbf{b}_2 \in \ker(\mathbf{A}^*)$ are the orthogonal projections into $\text{span}(\mathbf{A})$ and $\ker(\mathbf{A}^*)$, respectively. Observe that $\mathbf{b}_1 - \mathbf{Ax} \in \text{span}(\mathbf{A})$ for any $\mathbf{x} \in \mathbb{C}^n$ and $(\mathbf{b}_1 - \mathbf{Ax})^* \mathbf{b}_2 = 0$. By Pythagoras

$$\|\mathbf{b} - \mathbf{Ax}\|_2^2 = \|(\mathbf{b}_1 - \mathbf{Ax}) + \mathbf{b}_2\|_2^2 = \|\mathbf{b}_1 - \mathbf{Ax}\|_2^2 + \|\mathbf{b}_2\|_2^2 \geq \|\mathbf{b}_2\|_2^2.$$

We obtain the minimum value $\|\mathbf{b}_2\|_2$ of $\|\mathbf{b} - \mathbf{Ax}\|_2^2$ for any \mathbf{x} such that $\mathbf{Ax} = \mathbf{b}_1$. Since $\mathbf{b}_1 \in \text{span}(\mathbf{A})$ we can always find such an \mathbf{x} and existence follows. \square

Theorem 12.3 (Uniqueness)

The least squares solution is unique if and only if \mathbf{A} has linearly independent columns.

Proof. By what we just showed any solution \mathbf{x} of the least squares problem satisfies $\mathbf{Ax} = \mathbf{b}_1$. There is a unique such \mathbf{x} if and only if $\text{rank}(\mathbf{A}) = n$. \square

Theorem 12.4 (Characterization)

$\mathbf{x} \in \mathbb{C}^n$ is a least squares solution if and only if

$$\mathbf{A}^* \mathbf{Ax} = \mathbf{A}^* \mathbf{b}. \quad (12.1)$$

Proof. Suppose $\mathbf{A}^* \mathbf{Ax} = \mathbf{A}^* \mathbf{b}$. Then $\mathbf{A}^*(\mathbf{b} - \mathbf{Ax}) = \mathbf{0}$ or $\mathbf{A}^*(\mathbf{b}_1 + \mathbf{b}_2 - \mathbf{Ax}) = \mathbf{0}$. Since $\mathbf{A}^* \mathbf{b}_2 = \mathbf{0}$ it follows that $\mathbf{A}^*(\mathbf{b}_1 - \mathbf{Ax}) = \mathbf{0}$. But then $\mathbf{b}_1 - \mathbf{Ax} \in \text{span}(\mathbf{A}) \cap \ker(\mathbf{A}^*)$ which implies that $\mathbf{b}_1 - \mathbf{Ax} = \mathbf{0}$ and \mathbf{x} is a least squares solution. Conversely, if \mathbf{x} is a least squares solution then $\mathbf{b}_1 - \mathbf{Ax} = \mathbf{0}$ so that $\mathbf{A}^*(\mathbf{b}_1 - \mathbf{Ax}) = \mathbf{0}$, and therefore $\mathbf{A}^*(\mathbf{b} - \mathbf{Ax}) = \mathbf{0}$. Thus \mathbf{x} satisfies (12.1). \square

Note that only square or overdetermined systems can have unique solutions. The linear system (12.1) is called the **normal equations**. Since a least squares solution exists the normal equations has at least one solution for any \mathbf{b} .

One way to solve the least squares problem is to write E as a quadratic function and set partial derivatives equal to zero. If \mathbf{A} and \mathbf{b} have real components we find

$$E(\mathbf{x}) := (\mathbf{Ax} - \mathbf{b})^* (\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^* \mathbf{Bx} - 2\mathbf{c}^* \mathbf{x} + \beta,$$

where

$$\mathbf{B} = \mathbf{A}^* \mathbf{A}, \quad \mathbf{c} = \mathbf{A}^* \mathbf{b}, \quad \beta = \mathbf{b}^* \mathbf{b}.$$

Taking partial derivatives and using Lemma 10.2 we have

$$\nabla E(\mathbf{x}) := \left[\frac{\partial E}{\partial x_1}, \dots, \frac{\partial E}{\partial x_n} \right]^T = 2(\mathbf{Bx} - \mathbf{c}) = \mathbf{0},$$

and again we see that any least squares solution must satisfy the normal equations. The characterization theorem 12.4 also states the converse. Any solution of the normal equations is a least squares solution.

Example 12.5 (Average)

Consider the least squares problem defined by

$$\begin{aligned} x_1 &= 1 \\ x_1 &= 1, \quad \mathbf{A} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{x} = [x_1], \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}, \\ x_1 &= 2 \end{aligned}$$

We find

$$\|\mathbf{Ax} - \mathbf{b}\|_2^2 = (x_1 - 1)^2 + (x_1 - 1)^2 + (x_1 - 2)^2 = 3x_1^2 - 8x_1 + 6.$$

Setting the first derivative with respect to x_1 equal to zero we obtain $6x_1 - 8 = 0$ or $x_1 = 4/3$, the average of b_1, b_2, b_3 . The second derivative is positive and $x_1 = 4/3$ is a global minimum. The normal equation is $3x_1 = 4$.

Example 12.6 (Input/output model)

Suppose we have a simple input/output model. To every input $\mathbf{u} \in \mathbb{R}^n$ we obtain an output $y \in \mathbb{R}$. Assuming we have a linear relation

$$y = \mathbf{u}^T \mathbf{x} = \sum_{i=1}^n u_i x_i,$$

between \mathbf{u} and y , how can we determine \mathbf{x} ?

Performing $m \geq n$ experiments we obtain a table of values

$$\begin{array}{c|c|c|c|c} \mathbf{u} & \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_m \\ \hline y & y_1 & y_2 & \cdots & y_m \end{array} .$$

We would like to find \mathbf{x} such that

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_m^T \end{bmatrix} \mathbf{x} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \mathbf{b}.$$

We can estimate \mathbf{x} by solving the least squares problem $\min \|\mathbf{Ax} - \mathbf{b}\|_2^2$.

12.2 Some Curve Fitting examples

Given

- size: $1 \leq n \leq m$,
- sites: $\mathcal{S} := \{t_1, t_2, \dots, t_m\} \subset [a, b]$, $\mathbf{t} := [t_1, \dots, t_m]^T \in \mathbb{R}^m$,
- y -values: $\mathbf{y} = [y_1, y_2, \dots, y_m]^T \in \mathbb{R}^m$,
- positive weights: $\mathbf{W} := \text{diag}(w_1, w_2, \dots, w_m) \in \mathbb{R}^{m \times m}$, $w_k > 0$, $k = 1, \dots, m$,
- functions: $\phi_j : [a, b] \rightarrow \mathbb{R}$, $j = 1, \dots, n$.

Find a function (curve fit) $p : [a, b] \rightarrow \mathbb{R}$ given by $p := \sum_{j=1}^n x_j \phi_j$ such that $p(t_k) \approx y_k$ for $k = 1, \dots, m$.

An example is shown in Figure 12.1. Here $\phi_1(t) = 1$ and $\phi_2(t) = t$ and p is a straight line (linear regression).

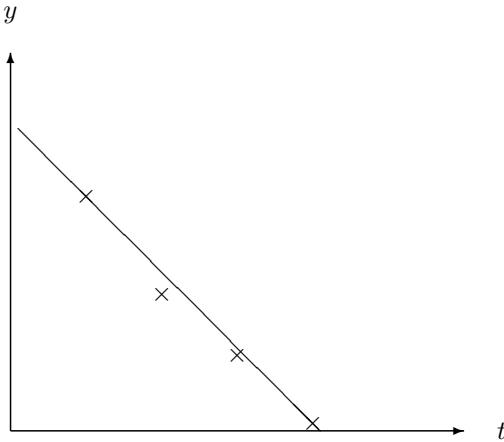


Figure 12.1: A least squares fit to data.

The curve fitting problem can be defined from the overdetermined linear system $\mathbf{Ax} = \mathbf{b}$, where

$$\mathbf{A} := \mathbf{W}^{1/2} \begin{bmatrix} \phi_1(t_1) & \cdots & \phi_n(t_1) \\ \vdots & & \vdots \\ \phi_1(t_m) & \cdots & \phi_n(t_m) \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad \mathbf{b} := \mathbf{W}^{1/2} \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m. \quad (12.2)$$

Then we find $\mathbf{x} \in \mathbb{R}^n$ as a solution of the corresponding least squares problem given by

$$E(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \sum_{k=1}^m w_k \left(\sum_{j=1}^n x_j \phi_j(t_k) - y_k \right)^2. \quad (12.3)$$

Typical examples of functions ϕ_j are polynomials, trigonometric functions, exponential functions, or splines. The numbers w_k are called weights. If y_k is an accurate observation, we can choose a large weight w_k . This will force $p(t_k) - y_k$ to be small. Similarly, a small w_k will allow $p(t_k) - y_k$ to be large. If an estimate for the standard deviation δy_k in y_k is known for each k , we can choose $w_k = 1/(\delta y_k)^2$, $k = 1, 2, \dots, m$.

In many cases the normal equations have a unique solution.

Lemma 12.7 (Curve fitting)

Let \mathbf{A} be given by (12.2). The matrix $\mathbf{B} := \mathbf{A}^T \mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive semi-definite. If $\{\phi_1, \dots, \phi_n\}$ is linearly independent on \mathcal{S} , i.e.,

$$p(t_k) := \sum_{j=1}^n x_j \phi_j(t_k) = 0, \quad k = 1, \dots, m \Rightarrow x_1 = \dots = x_n = 0 \quad (12.4)$$

then \mathbf{B} is symmetric positive definite.

Proof. \mathbf{B} is symmetric positive semidefinite by Corollary 3.27. By the same corollary \mathbf{A} is positive definite if \mathbf{A} has linearly independent columns, i.e., $\mathbf{Ax} = 0$ implies $\mathbf{x} = \mathbf{0}$. But if $(\mathbf{Ax})_k = \sqrt{w_k} \sum_{j=1}^n x_j \phi_j(t_k) = 0$, $k = 1, \dots, m$ then (12.4) implies that $x_j = 0$, $j = 1, \dots, n$. \square

Example 12.8 (Straight line fit)

Consider $n = 2$, $w_i = 1$, $i = 1, \dots, m$, $\phi_1(t) = 1$, and $\phi_2(t) = t$. The normal equations can be written

$$\begin{bmatrix} m & \sum t_k \\ \sum t_k & \sum t_k^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \sum y_k \\ \sum t_k y_k \end{bmatrix}. \quad (12.5)$$

Here k ranges from 1 to m in the sums. Recall that a nonzero polynomial of degree at most n has at most n roots. Therefore, by the Lemma 12.7, this 2×2 system is symmetric positive definite if the t 's contain at least 2 distinct points. With the data

t	1.0	2.0	3.0	4.0
y	3.1	1.8	1.0	0.1

we try a least squares fit of the form

$$p(t) = x_1 + x_2 t.$$

We can find x_1 and x_2 by solving the linear system (12.5). In this case we obtain

$$\begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 10.1 \end{bmatrix}. \quad (12.6)$$

The solution is $x_1 = 3.95$ and $x_2 = -0.98$. The data and the polynomial $p(t)$ are shown in Figure 12.1.

Example 12.9 (Ill conditioning and the Hilbert matrix)

The normal equations can be ill-conditioned. Consider the curve fitting problem using the polynomials $\phi_j(t) := t^{j-1}$, for $j = 1, \dots, n$, equidistant sites $t_k = (k-1)/(m-1)$, and $w_k = 1$ for $k = 1, \dots, m$. The normal equations are $\mathbf{B}_n \mathbf{x} = \mathbf{c}_n$, where for $n = 3$

$$\mathbf{B}_3 \mathbf{x} := \begin{bmatrix} m & \sum t_k & \sum t_k^2 \\ \sum t_k & \sum t_k^2 & \sum t_k^3 \\ \sum t_k^2 & \sum t_k^3 & \sum t_k^4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \sum y_k \\ \sum t_k y_k \\ \sum t_k^2 y_k \end{bmatrix}.$$

\mathbf{B}_n is symmetric positive definite if at least n of the t 's are distinct. However \mathbf{B}_n is extremely ill-conditioned even for moderate n . Indeed, $\frac{1}{m} \mathbf{B}_n \approx \mathbf{H}_n$, where $\mathbf{H}_n \in \mathbb{R}^{n \times n}$ is the **Hilbert Matrix** with i, j element $1/(i+j-1)$. Thus for $n = 3$

$$\mathbf{H}_3 = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix}.$$

The elements of $\frac{1}{m} \mathbf{B}_n$ are related to Riemann sums approximations to the elements of \mathbf{H}_n . In fact,

$$\frac{1}{m} b_{i,j} = \frac{1}{m} \sum_{k=1}^m t_k^{i+j-2} = \frac{1}{m} \sum_{k=1}^m \left(\frac{k-1}{m-1} \right)^{i+j-2} \approx \int_0^1 x^{i+j-2} dx = \frac{1}{i+j-1} = h_{i,j}.$$

The elements of \mathbf{H}^{-1} are determined in Exercise 0.63. We find $K_1(\mathbf{H}_6) \approx 3 \cdot 10^7$. It appears that $\frac{1}{m} \mathbf{B}_n$ and hence \mathbf{B}_n is ill-conditioned for moderate n at least if m is large. The cure for this problem is to use a different basis for polynomials. Orthogonal polynomials is an excellent choice. Another possibility is to use the basis $(t - \tilde{t})^{j-1}$, $j = 1, \dots, n$, for a suitable \tilde{t} , see Exercise 12.11.

Exercise 12.10 (Straight line fit (linear regression))

Suppose $(t_i, y_i)_{i=1}^m$ are m points in the plane. We consider the over-determined systems

$$\begin{array}{lll} \text{(i)} & x_1 = y_1 & \text{(ii)} & x_1 + t_1 x_2 = y_1 \\ & x_1 = y_2 & & x_1 + t_2 x_2 = y_2 \\ & \vdots & & \vdots \\ & x_1 = y_m & & x_1 + t_m x_2 = y_m \end{array}$$

- a) Find the normal equations for (i) and the least squares solution.
- b) Find the normal equations for (ii) and give a geometric interpretation of the least squares solution.

Exercise 12.11 (Straight line fit using shifted power form)

Related to (ii) in Exercise 12.10 we have the overdetermined system

$$(iii) \quad x_1 + (t_i - \hat{t})x_2 = y_i, \quad i = 1, 2, \dots, m,$$

where $\hat{t} = (t_1 + \dots + t_m)/m$.

- a) Find the normal equations for (iii) and give a geometric interpretation of the least squares solution.
- b) Fit a straight line to the points (t_i, y_i) : (998.5, 1), (999.5, 1.9), (1000.5, 3.1) and (1001.5, 3.5) using a). Draw a sketch of the solution.

Exercise 12.12 (Fitting a circle to points)

In this problem we derive an algorithm to fit a circle $(t - c_1)^2 + (y - c_2)^2 = r^2$ to $m \geq 3$ given points $(t_i, y_i)_{i=1}^m$ in the (t, y) -plane. We obtain the overdetermined system

$$(t_i - c_1)^2 + (y_i - c_2)^2 = r^2, \quad i = 1, \dots, m, \quad (12.7)$$

of m equations in the three unknowns c_1, c_2 and r . This system is nonlinear, but it can be solved from the linear system

$$t_i x_1 + y_i x_2 + x_3 = t_i^2 + y_i^2, \quad i = 1, \dots, m, \quad (12.8)$$

and then setting $c_1 = x_1/2$, $c_2 = x_2/2$ and $r^2 = c_1^2 + c_2^2 + x_3$.

- a) Derive (12.8) from (12.7). Explain how we can find c_1, c_2, r once $[x_1, x_2, x_3]$ is determined.
- b) Formulate (12.8) as a linear least squares problem for suitable \mathbf{A} and \mathbf{b} .
- c) Does the matrix \mathbf{A} in b) have linearly independent columns?
- d) Use (12.8) to find the circle passing through the three points (1, 4), (3, 2), (1, 0).

12.3 The Least Squares Problem and the Singular Value Decomposition

Suppose m, n and $\mathbf{A} \in \mathbb{C}^{m \times n}$ has rank r . Recall that the singular value decomposition $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$ and the reduced form $\mathbf{A} = \mathbf{U}_1\Sigma_1\mathbf{V}_1^*$, called the singular value factorization, are related as follows

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^* = [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} = \mathbf{U}_1\Sigma_1\mathbf{V}_1^*. \quad (12.9)$$

Here $\Sigma_1 := \text{diag}(\sigma_1, \dots, \sigma_r)$ with $\sigma_1 > \sigma_2 > \dots > \sigma_r > 0$, and $\mathbf{U} \in \mathbb{C}^{m \times m}$, $\mathbf{V} \in \mathbb{C}^{n \times n}$ are unitary matrices.

12.3.1 Orthogonal Projections

It follows from Theorem 7.15 that \mathbf{U}_1 is an orthonormal basis for $\text{span}(\mathbf{A})$ and \mathbf{U}_2 is an orthonormal basis for $\ker(\mathbf{A}^*)$. Formulas for orthogonal projections follow from the singular value factorization.

Theorem 12.13 (Orthogonal projections)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$, $\mathbf{b} \in \mathbb{C}^m$ and let $\mathbf{A} = \mathbf{U}_1\Sigma_1\mathbf{V}_1^*$ be the singular value factorization of \mathbf{A} . Then

$$\mathbf{b}_1 := \mathbf{A}\mathbf{A}^\dagger\mathbf{b}, \quad \mathbf{A}^\dagger := \mathbf{V}_1\Sigma_1^{-1}\mathbf{U}_1^*, \quad (12.10)$$

is the orthogonal projection of \mathbf{b} into $\text{span}(\mathbf{A})$, and

$$\mathbf{b}_2 := (\mathbf{I} - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b} \quad (12.11)$$

is the orthogonal projection of \mathbf{b} into the orthogonal complement $\ker(\mathbf{A}^*)$ of $\text{span}(\mathbf{A})$. Moreover,

$$\mathbf{A}^\dagger\mathbf{b} = \mathbf{A}^\dagger\mathbf{b}_1. \quad (12.12)$$

Proof. Using (12.9)

$$\mathbf{b} = \mathbf{U}\mathbf{U}^*\mathbf{b} = [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{bmatrix} \mathbf{U}_1^* \\ \mathbf{U}_2^* \end{bmatrix} \mathbf{b} = \mathbf{U}_1\mathbf{U}_1^*\mathbf{b} + \mathbf{U}_2\mathbf{U}_2^*\mathbf{b} =: \mathbf{b}_1 + \mathbf{b}_2.$$

$\mathbf{b}_1 = \mathbf{U}_1(\mathbf{U}_1^*\mathbf{b})$ is in $\text{span}(\mathbf{A})$ since \mathbf{U}_1 is an orthonormal basis for $\text{span}(\mathbf{A})$. Similarly, \mathbf{b}_2 is in $\ker(\mathbf{A}^*)$ since \mathbf{U}_2 is an orthonormal basis for $\ker(\mathbf{A}^*)$. Moreover, $\mathbf{A}\mathbf{A}^\dagger\mathbf{b} = \mathbf{U}_1\mathbf{U}_1^*\mathbf{b}$, and then $\mathbf{b}_2 = \mathbf{b} - \mathbf{b}_1 = (\mathbf{I} - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b}$. Since $\mathbf{A}^\dagger\mathbf{A}\mathbf{A}^\dagger = \mathbf{A}^\dagger$ we find $\mathbf{A}^\dagger\mathbf{b}_1 = \mathbf{A}^\dagger\mathbf{A}\mathbf{A}^\dagger\mathbf{b} = \mathbf{A}^\dagger\mathbf{b}$ and (12.12) follows. \square

Example 12.14 (Orthogonal projections)

The singular value decomposition of $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$ is $\mathbf{A} = \mathbf{I}_3\mathbf{A}\mathbf{I}_2$. Thus $\mathbf{U}_1 =$

$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$ and $\mathbf{U}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$. Moreover $\mathbf{A}^\dagger = \mathbf{I}_2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{I}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$. If $\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$, then $\mathbf{b}_1 = \mathbf{A}\mathbf{A}^\dagger\mathbf{b} = \mathbf{U}_1\mathbf{U}_1^T\mathbf{b} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ 0 \end{bmatrix}$ and $\mathbf{b}_2 = (\mathbf{I}_3 - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b} = \mathbf{U}_2\mathbf{U}_2^T\mathbf{b} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ b_3 \end{bmatrix}$.

Theorem 12.15 (General least squares solution)

If \mathbf{x} is a least square solution then $\mathbf{x} = \mathbf{A}^\dagger\mathbf{b} + \mathbf{z}$ for some $\mathbf{z} \in \ker(\mathbf{A})$. Conversely, $\mathbf{A}^\dagger\mathbf{b} + \mathbf{z}$ is a least square solution for any $\mathbf{z} \in \ker(\mathbf{A})$.

Proof. If \mathbf{x} is a least square solution then $\mathbf{Ax} = \mathbf{b}_1$. Define $\mathbf{z} := \mathbf{x} - \mathbf{A}^\dagger\mathbf{b}$. Then $\mathbf{Az} = \mathbf{Ax} - \mathbf{AA}^\dagger\mathbf{b} = \mathbf{b}_1 - \mathbf{b}_1 = \mathbf{0}$ and $\mathbf{z} \in \ker(\mathbf{A})$. Conversely, if $\mathbf{x} = \mathbf{A}^\dagger\mathbf{b} + \mathbf{z}$ with $\mathbf{z} \in \ker(\mathbf{A})$ then $\mathbf{Ax} = \mathbf{AA}^\dagger\mathbf{b} + \mathbf{Az} = \mathbf{b}_1$. \square

12.3.2 The Generalized Inverse

Consider the matrix

$$\mathbf{A}^\dagger := \mathbf{V}_1 \boldsymbol{\Sigma}_1^{-1} \mathbf{U}_1^*. \quad (12.13)$$

If \mathbf{A} is square and nonsingular then $\mathbf{A}^\dagger \mathbf{A} = \mathbf{AA}^\dagger = \mathbf{I}$ and \mathbf{A}^\dagger is the usual inverse of \mathbf{A} . It could be a problem that our definition of \mathbf{A}^\dagger depends on the particular singular value factorization. However, we show in Exercises 12.16, 12.17 that it is independent of the choice of singular value factorization. We call \mathbf{A}^\dagger the **generalized inverse** or **pseudo inverse** of \mathbf{A} . Any matrix has a generalized inverse, and so \mathbf{A}^\dagger is a generalization of the usual inverse.

We show in Exercise 12.19 that if \mathbf{A} has linearly independent columns then

$$\mathbf{A}^\dagger = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*. \quad (12.14)$$

Exercise 12.16 (The generalized inverse)

Show that $\mathbf{B} := \mathbf{A}^\dagger$ satisfies (1) $\mathbf{ABA} = \mathbf{A}$, (2) $\mathbf{BAB} = \mathbf{B}$, (3) $(\mathbf{BA})^* = \mathbf{BA}$, and (4) $(\mathbf{AB})^* = \mathbf{AB}$.

Exercise 12.17 (Uniqueness of generalized inverse)

Given $\mathbf{A} \in \mathbb{C}^{m \times n}$, and suppose $\mathbf{B}, \mathbf{C} \in \mathbb{C}^{n \times m}$ satisfy

$$\begin{array}{lll} \mathbf{ABA} = \mathbf{A} & (1) & \mathbf{ACA} = \mathbf{A}, \\ \mathbf{BAB} = \mathbf{B} & (2) & \mathbf{CAC} = \mathbf{C}, \\ (\mathbf{AB})^* = \mathbf{AB} & (3) & (\mathbf{AC})^* = \mathbf{AC}, \\ (\mathbf{BA})^* = \mathbf{BA} & (4) & (\mathbf{CA})^* = \mathbf{CA}. \end{array}$$

Verify the following proof that $\mathbf{B} = \mathbf{C}$.

$$\begin{aligned}\mathbf{B} &= (\mathbf{B}\mathbf{A})\mathbf{B} = (\mathbf{A}^*)\mathbf{B}^*\mathbf{B} = (\mathbf{A}^*\mathbf{C}^*)\mathbf{A}^*\mathbf{B}^*\mathbf{B} = \mathbf{C}\mathbf{A}(\mathbf{A}^*\mathbf{B}^*)\mathbf{B} \\ &= \mathbf{C}\mathbf{A}(\mathbf{B}\mathbf{A}) = (\mathbf{C})\mathbf{A}\mathbf{B} = \mathbf{C}(\mathbf{A}\mathbf{C})\mathbf{A}\mathbf{B} = \mathbf{C}\mathbf{C}^*\mathbf{A}^*(\mathbf{A}\mathbf{B}) \\ &= \mathbf{C}\mathbf{C}^*(\mathbf{A}^*\mathbf{B}^*\mathbf{A}^*) = \mathbf{C}(\mathbf{C}^*\mathbf{A}^*) = \mathbf{C}\mathbf{A}\mathbf{C} = \mathbf{C}.\end{aligned}$$

Exercise 12.18 (Verify that a matrix is a generalized inverse)

Show that the matrices $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}$ and $\mathbf{B} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$ satisfy the axioms in Exercise 12.16. Thus we can conclude that $\mathbf{B} = \mathbf{A}^\dagger$ without computing the singular value decomposition of \mathbf{A} .

Exercise 12.19 (Linearly independent columns and generalized inverse)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ has linearly independent columns. Show that $\mathbf{A}^* \mathbf{A}$ is nonsingular and $\mathbf{A}^\dagger = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*$. If \mathbf{A} has linearly independent rows, then show that $\mathbf{A} \mathbf{A}^*$ is nonsingular and $\mathbf{A}^\dagger = \mathbf{A}^* (\mathbf{A} \mathbf{A}^*)^{-1}$.

Exercise 12.20 (The generalized inverse of a vector)

Show that $\mathbf{u}^\dagger = (\mathbf{u}^* \mathbf{u})^{-1} \mathbf{u}^*$ if $\mathbf{u} \in \mathbb{C}^{n,1}$ is nonzero.

Exercise 12.21 (The generalized inverse of an outer product)

If $\mathbf{A} = \mathbf{u}\mathbf{v}^*$ where $\mathbf{u} \in \mathbb{C}^m$, $\mathbf{v} \in \mathbb{C}^n$ are nonzero, show that

$$\mathbf{A}^\dagger = \frac{1}{\alpha} \mathbf{A}^*, \quad \alpha = \|\mathbf{u}\|_2^2 \|\mathbf{v}\|_2^2.$$

Exercise 12.22 (The generalized inverse of a diagonal matrix)

Show that $\text{diag}(\lambda_1, \dots, \lambda_n)^\dagger = \text{diag}(\lambda_1^\dagger, \dots, \lambda_n^\dagger)$ where

$$\lambda_i^\dagger = \begin{cases} 1/\lambda_i, & \lambda_i \neq 0 \\ 0 & \lambda_i = 0. \end{cases}$$

Exercise 12.23 (Properties of the generalized inverse)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$. Show that

- a) $(\mathbf{A}^*)^\dagger = (\mathbf{A}^\dagger)^*$.
- b) $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$.
- c) $(\alpha \mathbf{A})^\dagger = \frac{1}{\alpha} \mathbf{A}^\dagger$, $\alpha \neq 0$.

Exercise 12.24 (The generalized inverse of a product)

Suppose $k, m, n \in \mathbb{N}$, $\mathbf{A} \in \mathbb{C}^{m \times n}$, $\mathbf{B} \in \mathbb{C}^{n \times k}$. Suppose \mathbf{A} has linearly independent columns and \mathbf{B} has linearly independent rows.

- Show that $(\mathbf{AB})^\dagger = \mathbf{B}^\dagger \mathbf{A}^\dagger$. Hint: Let $\mathbf{E} = \mathbf{AF}$, $\mathbf{F} = \mathbf{B}^\dagger \mathbf{A}^\dagger$. Show by using $\mathbf{A}^\dagger \mathbf{A} = \mathbf{BB}^\dagger = \mathbf{I}$ that \mathbf{F} is the generalized inverse of \mathbf{E} .
- Find $\mathbf{A} \in \mathbb{R}^{1,2}$, $\mathbf{B} \in \mathbb{R}^{2,1}$ such that $(\mathbf{AB})^\dagger \neq \mathbf{B}^\dagger \mathbf{A}^\dagger$.

Exercise 12.25 (The generalized inverse of the conjugate transpose)

Show that $\mathbf{A}^* = \mathbf{A}^\dagger$ if and only if all singular values of \mathbf{A} are either zero or one.

Exercise 12.26 (Linearly independent columns)

Show that if \mathbf{A} has rank n then $\mathbf{A}(\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* \mathbf{b}$ is the projection of \mathbf{b} into $\text{span}(\mathbf{A})$. (Cf. Exercise 12.19.)

Exercise 12.27 (Analaysis of the general linear system)

Consider the linear system $\mathbf{Ax} = \mathbf{b}$ where $\mathbf{A} \in \mathbb{C}^{n \times n}$ has rank $r > 0$ and $\mathbf{b} \in \mathbb{C}^n$. Let

$$\mathbf{U}^* \mathbf{A} \mathbf{V} = \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

represent the singular value decomposition of \mathbf{A} .

- Let $\mathbf{c} = [c_1, \dots, c_n]^T = \mathbf{U}^* \mathbf{b}$ and $\mathbf{y} = [y_1, \dots, y_n]^T = \mathbf{V}^* \mathbf{x}$. Show that $\mathbf{Ax} = \mathbf{b}$ if and only if
- $$\begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{y} = \mathbf{c}.$$
- Show that $\mathbf{Ax} = \mathbf{b}$ has a solution \mathbf{x} if and only if $c_{r+1} = \dots = c_n = 0$.
 - Deduce that a linear system $\mathbf{Ax} = \mathbf{b}$ has either no solution, one solution or infinitely many solutions.

Exercise 12.28 (Fredholm's alternative)

For any $\mathbf{A} \in \mathbb{C}^{m \times n}$, $\mathbf{b} \in \mathbb{C}^n$ show that one and only one of the following systems has a solution

$$(1) \quad \mathbf{Ax} = \mathbf{b}, \quad (2) \quad \mathbf{A}^* \mathbf{y} = \mathbf{0}, \quad \mathbf{y}^* \mathbf{b} \neq 0.$$

In other words either $\mathbf{b} \in \text{span}(\mathbf{A})$, or we can find $\mathbf{y} \in \ker(\mathbf{A}^*)$ such that $\mathbf{y}^* \mathbf{b} \neq 0$. This is called **Fredholm's alternative**.

12.4 Numerical Solution

We assume that $m \geq n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Any least squares solution is a solution of the normal equations $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$. Two examples illustrate this.

Example 12.29 (Unique solution of normal equations)

Consider the least squares problem with

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

This is the matrix in Example 11.11. We find

$$\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 4 & 4 & 6 \\ 4 & 20 & 26 \\ 6 & 26 & 70 \end{bmatrix}, \quad \mathbf{c} := \mathbf{A}^T \mathbf{b} = \begin{bmatrix} 4 \\ 4 \\ 6 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

The solution is unique since $\text{rank}(\mathbf{A}) = 3$ and therefore \mathbf{B} is symmetric positive definite.

Example 12.30 (Nonunique solution of normal equations)

Consider the least squares problem with

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

We find

$$\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}, \quad \mathbf{c} := \mathbf{A}^T \mathbf{b} = \begin{bmatrix} b_1 + b_2 \\ b_1 + b_2 \end{bmatrix}.$$

\mathbf{B} is symmetric positive semidefinite, but not positive definite. Any $\mathbf{x} = [x_1, x_2]^T \in \mathbb{R}^2$ with $x_1 + x_2 = (b_1 + b_2)/2$ is a solution.

Numerical methods can be based on normal equations, QR factorization, or Singular Value Factorization. We discuss each of these approaches in turn.

12.4.1 Normal Equations

Suppose \mathbf{A} has linearly independent columns. The coefficient matrix $\mathbf{B} := \mathbf{A}^T \mathbf{A}$ in the normal equations is symmetric positive definite, and we can solve these equations using the Cholesky factorization of \mathbf{B} . Consider forming the normal equations. We can use either a column oriented- or a row oriented approach.

$$\begin{aligned} 1. \text{ inner product: } (\mathbf{A}^T \mathbf{A})_{i,j} &= \sum_{k=1}^m a_{k,i} a_{k,j}, \quad i, j = 1, \dots, n, \\ (\mathbf{A}^T \mathbf{b})_i &= \sum_{k=1}^m a_{k,i} b_k, \quad i = 1, \dots, n, \end{aligned}$$

$$2. \text{ outer product: } \mathbf{A}^T \mathbf{A} = \sum_{k=1}^m \begin{bmatrix} a_{k1} \\ \vdots \\ a_{kn} \end{bmatrix} [a_{k1} \ \cdots \ a_{kn}], \quad \mathbf{A}^T \mathbf{b} = \sum_{k=1}^m \begin{bmatrix} a_{k1} \\ \vdots \\ a_{kn} \end{bmatrix} b_k.$$

The outer product form is suitable for large problems since it uses only one pass through the data importing one row of \mathbf{A} at a time from some separate storage.

Consider the number of operations to find the least squares solution. We need $2m$ arithmetic operations for each inner product. Since \mathbf{B} is symmetric we only need to compute $n(n+1)/2$ such inner products. It follows that \mathbf{B} can be computed in approximately mn^2 arithmetic operations. In conclusion the number of operations are mn^2 to find \mathbf{B} , $2mn$ to find $\mathbf{A}^T \mathbf{b}$, $n^3/3$ to find \mathbf{R} , n^2 to solve $\mathbf{R}^T \mathbf{y} = \mathbf{c}$ and n^2 to solve $\mathbf{R}\mathbf{x} = \mathbf{y}$. If $m \approx n$ it takes $\frac{4}{3}n^3 = 2G_n$ arithmetic operations. If m is much bigger than n the number of operations is approximately mn^2 , the work to compute \mathbf{B} .

A problem with the normal equations approach is that the linear system can be poorly conditioned. In fact the 2-norm condition number of $\mathbf{B} := \mathbf{A}^T \mathbf{A}$ is the square of the condition number of \mathbf{A} . This follows, since the eigenvalues of \mathbf{B} are the square of the singular values of \mathbf{A} so that

$$K_2(\mathbf{B}) = \frac{\sigma_1^2}{\sigma_n^2} = \left(\frac{\sigma_1}{\sigma_n} \right)^2 = K_2(\mathbf{A})^2.$$

If \mathbf{A} is ill-conditioned, this could make the normal equations approach problematic. One difficulty which can be encountered is that the computed $\mathbf{A}^T \mathbf{A}$ might not be positive definite. See Problem 12.39 for an example.

12.4.2 QR Factorization

Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ has rank n and let $\mathbf{b} \in \mathbb{R}^m$. The QR factorization can be used to solve the least squares problem. Suppose $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$ is a QR factorization of \mathbf{A} . Since \mathbf{Q}_1 has orthonormal columns we find

$$\mathbf{A}^T \mathbf{A} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{R}_1^T \mathbf{R}_1, \quad \mathbf{A}^T \mathbf{b} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{b}.$$

Since \mathbf{A} has rank n the matrix \mathbf{R}_1^T is nonsingular and can be canceled. Thus

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b} \implies \mathbf{R}_1 \mathbf{x} = \mathbf{c}_1, \quad \mathbf{c}_1 := \mathbf{Q}_1^T \mathbf{b}.$$

We can use Householder transformations or Givens rotations to find \mathbf{R}_1 and \mathbf{c}_1 . Consider using the Householder triangulation algorithm Algorithm 11.8. We find $\mathbf{R} = \mathbf{Q}^T \mathbf{A}$ and $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$, where $\mathbf{A} = \mathbf{Q}\mathbf{R}$ is the QR decomposition of \mathbf{A} . The

matrices \mathbf{R}_1 and \mathbf{c}_1 are located in the first n rows of \mathbf{R} and \mathbf{c} . Using also Algorithm 1.13 we have the following method to solve the full rank least squares problem.

1. $[\mathbf{R}, \mathbf{c}] = \text{housetriang}(\mathbf{A}, \mathbf{b})$.
2. $\mathbf{x} = \text{rbacksolve}(\mathbf{R}(1:n, 1:n), \mathbf{c}(1:n), n)$.

Example 12.31 (Solution using QR factorization)

Consider the least squares problem with

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} \quad \text{and } \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

This is the matrix in Example 11.11. The least squares solution \mathbf{x} is found by solving the system

$$\begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

and we find $\mathbf{x} = [1, 0, 0]^T$.

Using Householder triangulation is a useful alternative to normal equations for solving full rank least squares problems. It can even be extended to rank deficient problems, see [2]. The 2 norm condition number for the system $\mathbf{R}_1 \mathbf{x} = \mathbf{c}_1$ is $K_2(\mathbf{R}_1) = K_2(\mathbf{Q}_1 \mathbf{R}_1) = K_2(\mathbf{A})$, and as discussed in the previous section this is the square root of $K_2(\mathbf{A}^T \mathbf{A})$, the condition number for the normal equations. Thus if \mathbf{A} is mildly ill-conditioned the normal equations can be quite ill-conditioned and solving the normal equations can give inaccurate results. On the other hand Algorithm 11.8 is quite stable.

But using Householder transformations requires more work. The leading term in the number of arithmetic operations in Algorithm 11.8 is approximately $2mn^2 - 2n^3/3$, (cf. (11.5) while the number of arithmetic operations needed to form the normal equations, taking advantage of symmetry is approximately mn^2 . Thus for m much larger than n using Householder triangulation requires twice as many arithmetic operations as the an approach based on the normal equations. Also, Householder triangulation have problems taking advantage of the structure in sparse problems.

12.4.3 Singular Value Factorization

This method can be used even if \mathbf{A} does not have full rank. It requires knowledge of the generalized inverse of \mathbf{A} . By Theorem 12.15

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$$

is a least squares solution for any $\mathbf{z} \in \ker(\mathbf{A})$.

Example 12.32 (Solution using singular value factorization)

The generalized inverse of $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}$ is $\mathbf{A}^\dagger = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$. (cf. Exercise 12.18). Moreover, $[-1, 1]^T$ is a basis for $\ker(\mathbf{A})$. If $\mathbf{b} = [b_1, b_2, b_3]^T$, then for any $\mathbf{z} \in \mathbb{R}$ the vector

$$\mathbf{x} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + z \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} \frac{b_1+b_2}{4} - z \\ \frac{b_1+b_2}{4} + z \end{bmatrix} \quad (12.15)$$

is a solution of $\min \|\mathbf{Ax} - \mathbf{b}\|_2$ and this gives all solutions.

If \mathbf{A} has linearly independent columns then $\ker(\mathbf{A}) = \{\mathbf{0}\}$ and $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$ is the unique solution.

When $\text{rank}(\mathbf{A})$ is less than the number of columns of \mathbf{A} then $\ker(\mathbf{A}) \neq \{\mathbf{0}\}$, and we have a choice of \mathbf{z} . One possible choice is $\mathbf{z} = \mathbf{0}$ giving the solution $\mathbf{A}^\dagger \mathbf{b}$.

Theorem 12.33 (Minimal solution)

The least squares solution with minimal Euclidian norm is $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$ corresponding to $\mathbf{z} = \mathbf{0}$.

Proof. Suppose $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$, with $\mathbf{z} \in \ker(\mathbf{A})$. Recall that if the right singular vectors of \mathbf{A} are partitioned as $[\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_n] = [\mathbf{V}_1, \mathbf{V}_2]$, then \mathbf{V}_2 is a basis for $\ker(\mathbf{A})$. Moreover, $\mathbf{V}_2^* \mathbf{V}_1 = \mathbf{0}$ since \mathbf{V} has orthonormal columns. If $\mathbf{A}^\dagger = \mathbf{V}_1 \Sigma^{-1} \mathbf{U}_1^*$ and $\mathbf{z} \in \ker(\mathbf{A})$ then $\mathbf{z} = \mathbf{V}_2 \mathbf{y}$ for some $\mathbf{y} \in \mathbb{C}^{n-r}$ and we obtain

$$\mathbf{z}^* \mathbf{A}^\dagger \mathbf{b} = \mathbf{y}^* \mathbf{V}_2^* \mathbf{V}_1 \Sigma^{-1} \mathbf{U}_1^* \mathbf{b} = \mathbf{0}.$$

Thus \mathbf{z} and $\mathbf{A}^\dagger \mathbf{b}$ are orthogonal so that by Pythagoras $\|\mathbf{x}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{b} + \mathbf{z}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{b}\|_2^2 + \|\mathbf{z}\|_2^2 \geq \|\mathbf{A}^\dagger \mathbf{b}\|_2^2$ with equality for $\mathbf{z} = \mathbf{0}$. \square

12.5 Perturbation Theory for Least Squares

In this section we consider what effect small changes in the data \mathbf{A}, \mathbf{b} have on the solution \mathbf{x} of the least squares problem $\min \|\mathbf{Ax} - \mathbf{b}\|_2$.

If \mathbf{A} has linearly independent columns then we can write the least squares solution \mathbf{x} (the solution of $\mathbf{A}^* \mathbf{A} \mathbf{x} = \mathbf{A}^* \mathbf{b}$) as

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}, \quad \mathbf{A}^\dagger := (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*.$$

12.5.1 Perturbing the Right Hand Side

Let us now consider the effect of a perturbation in \mathbf{b} on \mathbf{x} .

Theorem 12.34 (Perturbing the Right Hand Side)

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ has linearly independent columns, and let $\mathbf{b}, \mathbf{e} \in \mathbb{C}^m$. Let $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ be the solutions of $\min \|\mathbf{Ax} - \mathbf{b}\|_2$ and $\min \|\mathbf{Ay} - \mathbf{b} - \mathbf{e}\|_2$. Finally, let $\mathbf{b}_1, \mathbf{e}_1$ be the projections of \mathbf{b} and \mathbf{e} into $\text{span}(\mathbf{A})$. If $\mathbf{b}_1 \neq \mathbf{0}$, we have for any operator norm

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|. \quad (12.16)$$

Proof. (12.12) implies that $\mathbf{A}^\dagger \mathbf{b} = \mathbf{A}^\dagger \mathbf{b}_1$ and $\mathbf{A}^\dagger \mathbf{e} = \mathbf{A}^\dagger \mathbf{e}_1$. Subtracting $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}_1$ from $\mathbf{y} = \mathbf{A}^\dagger \mathbf{b}_1 + \mathbf{A}^\dagger \mathbf{e}_1$ we have $\mathbf{y} - \mathbf{x} = \mathbf{A}^\dagger \mathbf{e}_1$. Thus $\|\mathbf{y} - \mathbf{x}\| = \|\mathbf{A}^\dagger \mathbf{e}_1\| \leq \|\mathbf{A}^\dagger\| \|\mathbf{e}_1\|$. Moreover, $\|\mathbf{b}_1\| = \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$. Therefore $\|\mathbf{y} - \mathbf{x}\| / \|\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{A}^\dagger\| \|\mathbf{e}_1\| / \|\mathbf{b}_1\|$ proving the rightmost inequality. From $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{e}_1$ and $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}_1$ we obtain the leftmost inequality. \square

(12.16) is analogous to the bound (8.12) for linear systems. We see that the number $K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|$ generalizes the condition number $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ for a square matrix. The main difference between (12.16) and (8.12) is however that $\|\mathbf{e}\| / \|\mathbf{b}\|$ in (8.12) has been replaced by $\|\mathbf{e}_1\| / \|\mathbf{b}_1\|$, the projections of \mathbf{e} and \mathbf{b} on $\text{span}(\mathbf{A})$. If \mathbf{b} lies almost entirely in $\ker(\mathbf{A}^*)$, i.e. $\|\mathbf{b}\| / \|\mathbf{b}_1\|$ is large, $\|\mathbf{e}_1\| / \|\mathbf{b}_1\|$ can be much larger than $\|\mathbf{e}\| / \|\mathbf{b}\|$. This is illustrated in Figure 12.2. If \mathbf{b} is almost orthogonal to $\text{span}(\mathbf{A})$, $\|\mathbf{e}_1\| / \|\mathbf{b}_1\|$ will normally be much larger than $\|\mathbf{e}\| / \|\mathbf{b}\|$. Note that $\|\mathbf{e}_1\| / \|\mathbf{b}_1\|$ is also present in the lower bound.

Example 12.35 (Perturbing the Right Hand Side)

Suppose

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 10^{-4} \\ 0 \\ 1 \end{bmatrix}, \quad \mathbf{e} = \begin{bmatrix} 10^{-6} \\ 0 \\ 0 \end{bmatrix}.$$

For this example we can compute $K(\mathbf{A})$ by finding \mathbf{A}^\dagger explicitly. Indeed,

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad (\mathbf{A}^T \mathbf{A})^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}, \quad \mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

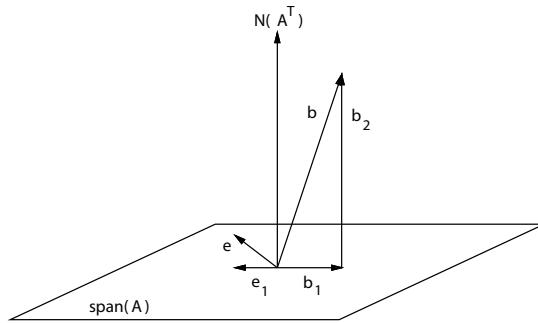


Figure 12.2: Graphical interpretation of the bounds in Theorem 12.34.

Thus $K_\infty(\mathbf{A}) = \|\mathbf{A}\|_\infty \|\mathbf{A}^\dagger\|_\infty = 2 \cdot 2 = 4$ is quite small.

Consider now the projections \mathbf{b}_1 and \mathbf{e}_1 . We find $\mathbf{A}\mathbf{A}^\dagger = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$. Hence

$$\mathbf{b}_1 = \mathbf{A}\mathbf{A}^\dagger \mathbf{b} = [10^{-4}, 0, 0]^T, \quad \text{and} \quad \mathbf{e}_1 = \mathbf{A}\mathbf{A}^\dagger \mathbf{e} = [10^{-6}, 0, 0]^T.$$

Thus $\|\mathbf{e}_1\|_\infty / \|\mathbf{b}_1\|_\infty = 10^{-2}$ and (12.16) takes the form

$$\frac{1}{4} 10^{-2} \leq \frac{\|\mathbf{y} - \mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \leq 4 \cdot 10^{-2}.$$

To verify the bounds we compute the solutions as $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} = [10^{-4}, 0]^T$ and $\mathbf{y} = \mathbf{A}^\dagger (\mathbf{b} + \mathbf{e}) = [10^{-4} + 10^{-6}, 0]^T$. Hence

$$\frac{\|\mathbf{x} - \mathbf{y}\|_\infty}{\|\mathbf{x}\|_\infty} = \frac{10^{-6}}{10^{-4}} = 10^{-2},$$

Exercise 12.36 (Condition number)

Let

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

a) Determine the projections \mathbf{b}_1 and \mathbf{b}_2 of \mathbf{b} on $\text{span}(\mathbf{A})$ and $\ker(\mathbf{A}^T)$.

b) Compute $K(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^\dagger\|_2$.

For each \mathbf{A} we can find \mathbf{b} and \mathbf{e} so that we have equality in the upper bound in (12.16). The lower bound is best possible in a similar way.

Exercise 12.37 (Equality in perturbation bound)

- a) Let $\mathbf{A} \in \mathbb{C}^{m \times n}$. Show that we have equality to the right in (12.16) if $\mathbf{b} = \mathbf{A}\mathbf{y}_A$, $\mathbf{e}_1 = \mathbf{y}_{A^\dagger}$ where $\|\mathbf{A}\mathbf{y}_A\| = \|\mathbf{A}\|$, $\|\mathbf{A}^\dagger \mathbf{y}_{A^\dagger}\| = \|\mathbf{A}^\dagger\|$.
- b) Show that we have equality to the left if we switch \mathbf{b} and \mathbf{e} in a).
- c) Let \mathbf{A} be as in Example 12.35. Find extremal \mathbf{b} and \mathbf{e} when the l_∞ norm is used.

12.5.2 Perturbing the Matrix

The analysis of the effects of a perturbation \mathbf{E} in \mathbf{A} is quite difficult. The following result is stated without proof, see [16, p. 51]. For other estimates see [2] and [24].

Theorem 12.38 (Perturbing the Matrix)

Suppose $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{m \times n}$, $m > n$, where \mathbf{A} has linearly independent columns and $\alpha := 1 - \|\mathbf{E}\|_2 \|\mathbf{A}^\dagger\|_2 > 0$. Then $\mathbf{A} + \mathbf{E}$ has linearly independent columns. Let $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2 \in \mathbb{C}^m$ where \mathbf{b}_1 and \mathbf{b}_2 are the projections on $\text{span}(\mathbf{A})$ and $\ker(\mathbf{A}^*)$ respectively. Suppose $\mathbf{b}_1 \neq \mathbf{0}$. Let \mathbf{x} and \mathbf{y} be the solutions of $\min \|\mathbf{Ax} - \mathbf{b}\|_2$ and $\min \|(\mathbf{A} + \mathbf{E})\mathbf{y} - \mathbf{b}\|_2$. Then

$$\rho = \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{\alpha} K(1 + \beta K) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2}, \quad \beta = \frac{\|\mathbf{b}_2\|_2}{\|\mathbf{b}_1\|_2}, \quad K = \|\mathbf{A}\|_2 \|\mathbf{A}^\dagger\|_2. \quad (12.17)$$

(12.17) says that the relative error in \mathbf{y} as an approximation to \mathbf{x} can be at most $K(1 + \beta K)/\alpha$ times as large as the size $\|\mathbf{E}\|_2/\|\mathbf{A}\|_2$ of the relative perturbation in \mathbf{A} . β will be small if \mathbf{b} lies almost entirely in $\text{span}(\mathbf{A})$, and we have approximately $\rho \leq \frac{1}{\alpha} K \|\mathbf{E}\|_2 / \|\mathbf{A}\|_2$. This corresponds to the estimate (8.17) for linear systems. If β is not small, the term $\frac{1}{\alpha} K^2 \beta \|\mathbf{E}\|_2 / \|\mathbf{A}\|_2$ will dominate. In other words, the condition number is roughly $K(\mathbf{A})$ if β is small and $K(\mathbf{A})^2 \beta$ if β is not small. Note that β is large if \mathbf{b} is almost orthogonal to $\text{span}(\mathbf{A})$ and that $\mathbf{b}_2 = \mathbf{b} - \mathbf{Ax}$ is the residual of \mathbf{x} .

Exercise 12.39 (Problem using normal equations)

Consider the least squares problems where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1+\epsilon \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 2 \\ 3 \\ 2 \end{bmatrix}, \quad \epsilon \in \mathbb{R}.$$

- a) Find the normal equations and the exact least squares solution.

- b) Suppose ϵ is small and we replace the $(2, 2)$ entry $3+2\epsilon+\epsilon^2$ in $\mathbf{A}^T \mathbf{A}$ by $3+2\epsilon$. (This will be done in a computer if $\epsilon < \sqrt{u}$, u being the round-off unit). For example, if $u = 10^{-16}$ then $\sqrt{u} = 10^{-8}$. Solve $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$ for \mathbf{x} and compare with the \mathbf{x} found in a). (We will get a much more accurate result using the QR factorization or the singular value decomposition on this problem).

12.6 Perturbation Theory for Singular Values

In this section we consider what effect a small change in the matrix \mathbf{A} has on the singular values.

We recall the Hoffman-Wielandt Theorem for singular values, Theorem 7.28. If $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ are rectangular matrices with singular values $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ and $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$, then

$$\sum_{j=1}^n |\alpha_j - \beta_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2.$$

This shows that the singular values of a matrix are well conditioned. Changing the Frobenius norm of a matrix by small amount only changes the singular values by a small amount.

Using the 2-norm we have a similar result involving only one singular value.

Theorem 12.40 (Perturbation of singular values)

Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ be rectangular matrices with singular values $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ and $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$. Then

$$|\alpha_j - \beta_j| \leq \|\mathbf{A} - \mathbf{B}\|_2, \text{ for } j = 1, 2, \dots, n. \quad (12.18)$$

Proof. Fix j and let \mathcal{S} be the $n - j + 1$ dimensional subspace for which the minimum in Theorem 7.27 is obtained for \mathbf{A} . Then

$$\alpha_j = \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{B} + (\mathbf{A} - \mathbf{B}))\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Bx}\|_2}{\|\mathbf{x}\|_2} + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \beta_j + \|\mathbf{A} - \mathbf{B}\|_2.$$

By symmetry we obtain $\beta_j \leq \alpha_j + \|\mathbf{A} - \mathbf{B}\|_2$ and the proof is complete. \square

The following result is an analogue of Theorem 8.33.

Theorem 12.41 (Generalized inverse when perturbing the matrix)

Let $\mathbf{A}, \mathbf{E} \in \mathbb{R}^{m \times n}$ have singular values $\alpha_1 \geq \dots \geq \alpha_n$ and $\epsilon_1 \geq \dots \geq \epsilon_n$. If $\|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2 < 1$ then

1. $\text{rank}(\mathbf{A} + \mathbf{E}) \geq \text{rank}(\mathbf{A})$,

$$2. \quad \|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2} = \frac{1}{\alpha_r - \epsilon_1},$$

where r is the rank of \mathbf{A} .

Proof. Suppose \mathbf{A} has rank r and let $\mathbf{B} := \mathbf{A} + \mathbf{E}$ have singular values $\beta_1 \geq \dots \geq \beta_n$. In terms of singular values the inequality $\|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2 < 1$ can be written $\epsilon_1/\alpha_r < 1$ or $\alpha_r > \epsilon_1$. By Theorem 12.40 we have $\alpha_r - \beta_r \leq \epsilon_1$, which implies $\beta_r \geq \alpha_r - \epsilon_1 > 0$, and this shows that $\text{rank}(\mathbf{A} + \mathbf{E}) > r$. To prove 2., the inequality $\beta_r \geq \alpha_r - \epsilon_1$ implies that

$$\|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{1}{\beta_r} \leq \frac{1}{\alpha_r - \epsilon_1} = \frac{1/\alpha_r}{1 - \epsilon_1/\alpha_r} = \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2}.$$

□

12.7 Review Questions

12.7.1 Do the normal equations always have a solution?

12.7.2 When is the least squares solution unique?

12.7.3 Express the general least squares solution in terms of the generalized inverse.

12.7.4 Consider perturbing the right-hand side in a linear equation and a least squares problem. What is the main difference in the perturbation inequalities?

12.7.5 Why does one often prefer using QR factorization instead of normal equations for solving least squares problems?

12.7.6 What is an orthogonal sum?

12.7.7 How is an orthogonal projection defined?

Part V

Eigenvalues and Eigenvectors

Chapter 13

Numerical Eigenvalue Problems

13.1 Eigenpairs

Eigenpairs have applications in quantum mechanics, differential equations, elasticity in mechanics, etc, etc. Typical computational problems involve

- Finding one or a few of the eigenvalues.
- Finding one or a few of the eigenpairs.
- Finding all eigenvalues.
- Finding all eigenpairs.

In this and the next chapter we consider numerical methods for finding one or more of the eigenvalues and eigenvectors of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$. Maybe the first method which comes to mind is to form the characteristic polynomial $\pi_{\mathbf{A}}$ of \mathbf{A} , and then use a polynomial root finder, like Newton's method to determine one or several of the eigenvalues.

It turns out that this is not suitable as an all purpose method. One reason is that a small change in one of the coefficients of $\pi_{\mathbf{A}}(\lambda)$ can lead to a large change in the roots of the polynomial. For example, if $\pi_{\mathbf{A}}(\lambda :) = \lambda^{16}$ and $q(\lambda) = \lambda^{16} - 10^{-16}$ then the roots of $\pi_{\mathbf{A}}$ are all equal to zero, while the roots of q are $\lambda_j = 10^{-1} e^{2\pi i j / 16}$, $j = 1, \dots, 16$. The roots of q have absolute value 0.1 and a perturbation in one of the polynomial coefficients of magnitude 10^{-16} has led to an error in the roots of approximately 0.1. The situation can be somewhat remedied by representing the polynomials using a different basis.

We will see that for many matrices the eigenvalues are less sensitive to perturbations in the elements of the matrix. In this text we will only consider methods

which work directly with the matrix.

13.2 Perturbation of Eigenvalues

In this section we study the following problem. Given matrices $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$, where we think of \mathbf{E} as a perturbation of \mathbf{A} . By how much do the eigenvalues of \mathbf{A} and $\mathbf{A} + \mathbf{E}$ differ? Not surprisingly this problem is more complicated than the corresponding problem for linear systems.

We illustrate this by considering two examples. Suppose $\mathbf{A}_0 := \mathbf{0}$ is the zero matrix. If $\lambda \in \sigma(\mathbf{A}_0 + \mathbf{E}) = \sigma(\mathbf{E})$, then $|\lambda| \leq \|\mathbf{E}\|_\infty$ by Theorem 9.6, and any zero eigenvalue of \mathbf{A}_0 is perturbed by at most $\|\mathbf{E}\|_\infty$. On the other hand consider for $\epsilon > 0$ the matrices

$$\mathbf{A}_1 := \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix}, \quad \mathbf{E} := \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \epsilon & 0 & 0 & \cdots & 0 & 0 \end{bmatrix} = \epsilon \mathbf{e}_n \mathbf{e}_1^T.$$

The characteristic polynomial of $\mathbf{A}_1 + \mathbf{E}$ is $\pi(\lambda) := (-1)^n(\lambda^n - \epsilon)$, and the zero eigenvalues of \mathbf{A}_1 are perturbed by the amount $|\lambda| = \|\mathbf{E}\|_\infty^{1/n}$. Thus, for $n = 16$, a perturbation of say $\epsilon = 10^{-16}$ gives a change in eigenvalue of 0.1.

The following theorem shows that a dependence $\|\mathbf{E}\|_\infty^{1/n}$ is the worst that can happen.

Theorem 13.1 (Elsner's Theorem)

Suppose $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$. To every $\mu \in \sigma(\mathbf{A} + \mathbf{E})$ there is a $\lambda \in \sigma(\mathbf{A})$ such that

$$|\mu - \lambda| \leq K \|\mathbf{E}\|_2^{1/n}, \quad K = (\|\mathbf{A}\|_2 + \|\mathbf{A} + \mathbf{E}\|_2)^{1-1/n}. \quad (13.1)$$

Proof. Suppose \mathbf{A} has eigenvalues $\lambda_1, \dots, \lambda_n$ and let λ_1 be one which is closest to μ . Let \mathbf{u}_1 with $\|\mathbf{u}_1\|_2 = 1$ be an eigenvector corresponding to μ , and extend \mathbf{u}_1 to an orthonormal basis $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ of \mathbb{C}^n . Note that

$$\begin{aligned} \|(\mu \mathbf{I} - \mathbf{A})\mathbf{u}_1\|_2 &= \|(\mathbf{A} + \mathbf{E})\mathbf{u}_1 - \mathbf{A}\mathbf{u}_1\|_2 = \|\mathbf{E}\mathbf{u}_1\|_2 \leq \|\mathbf{E}\|_2, \\ \prod_{j=2}^n \|(\mu \mathbf{I} - \mathbf{A})\mathbf{u}_j\|_2 &\leq \prod_{j=2}^n (|\mu| + \|\mathbf{A}\mathbf{u}_j\|_2) \leq ((\|\mathbf{A} + \mathbf{E}\|_2 + \|\mathbf{A}\|_2)^{n-1}). \end{aligned}$$

Using this and Hadamard's inequality (11.6) we find

$$\begin{aligned} |\mu - \lambda_1|^n &\leq \prod_{j=1}^n |\mu - \lambda_j| = |\det(\mu\mathbf{I} - \mathbf{A})| = |\det((\mu\mathbf{I} - \mathbf{A})[\mathbf{u}_1, \dots, \mathbf{u}_n])| \\ &\leq \|(\mu\mathbf{I} - \mathbf{A})\mathbf{u}_1\|_2 \prod_{j=2}^n \|(\mu\mathbf{I} - \mathbf{A})\mathbf{u}_j\|_2 \leq \|\mathbf{E}\|_2 (\|(\mathbf{A} + \mathbf{E})\|_2 + \|\mathbf{A}\|_2)^{n-1}. \end{aligned}$$

The result follows by taking n th roots in this inequality. \square

It follows from this theorem that the eigenvalues depend continuously on the elements of the matrix. The factor $\|\mathbf{E}\|_2^{1/n}$ shows that this dependence is almost, but not quite, differentiable. As an example, the eigenvalues of the matrix $\begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix}$ are $1 \pm \sqrt{\epsilon}$ and this expression is not differentiable at $\epsilon = 0$.

Recall that a matrix is nondefective if the eigenvectors form a basis for \mathbb{C}^n . For nondefective matrices we can get rid of the annoying exponent $1/n$ in $\|\mathbf{E}\|_2$. The following theorem is proved in Section 13.5. For a more general discussion see [24].

Theorem 13.2 (Linearly independent eigenvectors)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ has linearly independent eigenvectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ be the eigenvector matrix. Suppose $\mathbf{E} \in \mathbb{C}^{n \times n}$ and let μ be an eigenvalue of $\mathbf{A} + \mathbf{E}$. Then we can find an eigenvalue λ of \mathbf{A} such that

$$|\lambda - \mu| \leq K_p(\mathbf{X})\|\mathbf{E}\|_p, \quad 1 \leq p \leq \infty, \text{ where } K_p(\mathbf{X}) := \|\mathbf{X}\|_p \|\mathbf{X}^{-1}\|_p. \quad (13.2)$$

The equation (13.2) shows that for a nondefective matrix the absolute error can be magnified by at most $K_p(\mathbf{X})$, the condition number of the eigenvector matrix with respect to inversion. If $K_p(\mathbf{X})$ is small then a small perturbation changes the eigenvalues by small amounts.

Even if we get rid of the factor $1/n$, the equation (13.2) illustrates that it can be difficult or sometimes impossible to compute accurate eigenvalues and eigenvectors of matrices with almost linearly dependent eigenvectors. On the other hand the eigenvalue problem for normal matrices is better conditioned. Indeed, if \mathbf{A} is normal then it has a set of orthonormal eigenvectors and the eigenvector matrix is unitary. If we restrict attention to the 2-norm then $K_2(\mathbf{X}) = 1$ and (13.2) implies the following result.

Theorem 13.3 (Perturbations, normal matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is normal and let μ be an eigenvalue of $\mathbf{A} + \mathbf{E}$ for some $\mathbf{E} \in \mathbb{C}^{n \times n}$. Then we can find an eigenvalue λ of \mathbf{A} such that $|\lambda - \mu| \leq \|\mathbf{E}\|_2$.

For an even stronger result for Hermitian matrices see Corollary 6.46. We conclude that the situation for the absolute error in an eigenvalue of a Hermitian matrix is quite satisfactory. Small perturbations in the elements are not magnified in the eigenvalues.

13.2.1 Gerschgorin's Theorem

The following theorem is useful for locating eigenvalues of an arbitrary square matrix.

Theorem 13.4 (Gerschgorin's Circle Theorem)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$. Define for $i = 1, 2, \dots, n$

$$R_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|,$$

$$C_j = \{z \in \mathbb{C} : |z - a_{jj}| \leq c_j\}, \quad c_j := \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|.$$

Then any eigenvalue of \mathbf{A} lies in $R \cap C$ where $R = R_1 \cup R_2 \cup \dots \cup R_n$ and $C = C_1 \cup C_2 \cup \dots \cup C_n$.

Proof. Suppose (λ, \mathbf{x}) is an eigenpair for \mathbf{A} . We claim that $\lambda \in R_i$, where i is such that $|x_i| = \|\mathbf{x}\|_\infty$. Indeed, $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ implies that $\sum_j a_{ij}x_j = \lambda x_i$ or $(\lambda - a_{ii})x_i = \sum_{j \neq i} a_{ij}x_j$. Dividing by x_i and taking absolute values we find

$$|\lambda - a_{ii}| = \left| \sum_{j \neq i} a_{ij}x_j/x_i \right| \leq \sum_{j \neq i} |a_{ij}| |x_j/x_i| \leq r_i$$

since $|x_j/x_i| \leq 1$ for all j . Thus $\lambda \in R_i$.

Since λ is also an eigenvalue of \mathbf{A}^T , it must be in one of the row disks of \mathbf{A}^T . But these are the column disks C_j of \mathbf{A} . Hence $\lambda \in C_j$ for some j . \square

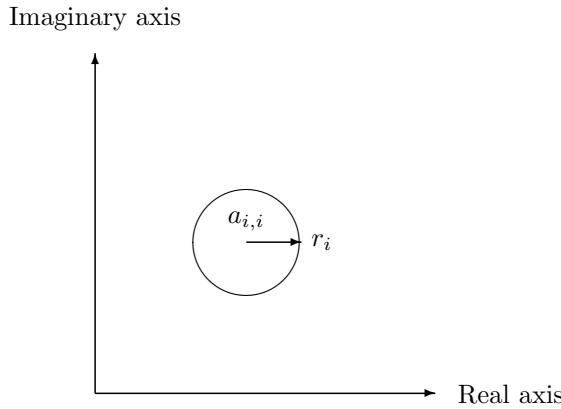
The set R_i is a subset of the complex plane consisting of all points inside a circle with center at a_{ii} and radius r_i , c.f. Figure 13.1. R_i is called a (Gerschgorin) row disk.

An eigenvalue λ lies in the union of the row disks R_1, \dots, R_n and also in the union of the column disks C_1, \dots, C_n . If \mathbf{A} is Hermitian then $R_i = C_i$ for $i = 1, 2, \dots, n$. Moreover, in this case the eigenvalues of \mathbf{A} are real, and the Gerschgorin disks can be taken to be intervals on the real line.

Example 13.5 (Gerschgorin)

Let $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ be the second derivative matrix. Since \mathbf{A} is Hermitian we have $R_i = C_i$ for all i and the eigenvalues are real. We find

$$R_1 = R_m = \{z \in \mathbb{R} : |z - 2| \leq 1\}, \quad \text{and } R_i = \{z \in \mathbb{R} : |z - 2| \leq 2\}, \quad i = 2, 3, \dots, m-1.$$

Figure 13.1: The Gershgorin disk R_i .

We conclude that $\lambda \in [0, 4]$ for any eigenvalue λ of \mathbf{T} . To check this, we recall that by Lemma 4.11 the eigenvalues of \mathbf{T} are given by

$$\lambda_j = 4 \left[\sin \frac{j\pi}{2(m+1)} \right]^2, \quad j = 1, 2, \dots, m.$$

When m is large the smallest eigenvalue $4 \left[\sin \frac{\pi}{2(m+1)} \right]^2$ is very close to zero and the largest eigenvalue $4 \left[\sin \frac{m\pi}{2(m+1)} \right]^2$ is very close to 4. Thus Gershgorin's theorem gives a remarkably good estimate for large m .

Sometimes some of the Gershgorin disks are distinct and we have

Corollary 13.6 (Disjoint Gershgorin disks)

If p of the Gershgorin row disks are disjoint from the others, the union of these disks contains precisely p eigenvalues. The same result holds for the column disks.

Proof. Consider a family of matrices

$$\mathbf{A}(t) := \mathbf{D} + t(\mathbf{A} - \mathbf{D}), \quad \mathbf{D} := \text{diag}(a_{11}, \dots, a_{nn}), \quad t \in [0, 1].$$

We have $\mathbf{A}(0) = \mathbf{D}$ and $\mathbf{A}(1) = \mathbf{A}$. As a function of t , every eigenvalue of $\mathbf{A}(t)$ is a continuous function of t . This follows from Theorem 13.1, see Exercise 13.7. The row disks $R_i(t)$ of $\mathbf{A}(t)$ have radius proportional to t , indeed

$$R_i(t) = \{z \in \mathbb{C} : |z - a_{ii}| \leq t r_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|.$$

Clearly $0 \leq t_1 < t_2 \leq 1$ implies $R_i(t_1) \subset R_i(t_2)$ and $R_i(1)$ is a row disk of \mathbf{A} for all i . Suppose $\bigcup_{k=1}^p R_{i_k}(1)$ are disjoint from the other disks of \mathbf{A} and set $R^p(t) := \bigcup_{k=1}^p R_{i_k}(t)$ for $t \in [0, 1]$. Now $R^p(0)$ contains only the p eigenvalues $a_{i_1, i_1}, \dots, a_{i_p, i_p}$ of $\mathbf{A}(0) = \mathbf{D}$. As t increases from zero to one the set $R^p(t)$ is disjoint from the other row disks of \mathbf{A} and by the continuity of the eigenvalues cannot loose or gain eigenvalues. It follows that $R^p(1)$ must contain p eigenvalues of \mathbf{A} . \square

Exercise 13.7 (Continuity of eigenvalues)

Suppose $t_1, t_2 \in [0, 1]$ and that μ is an eigenvalue of $\mathbf{A}(t_2)$. Show, using Theorem 13.1 with $\mathbf{A} = \mathbf{A}(t_1)$ and $\mathbf{E} = \mathbf{A}(t_2) - \mathbf{A}(t_1)$, that $\mathbf{A}(t_1)$ has an eigenvalue λ such that

$$|\lambda - \mu| \leq C(t_2 - t_1)^{1/n}, \text{ where } C \leq 2(\|\mathbf{D}\|_2 + \|\mathbf{A} - \mathbf{D}\|_2).$$

Thus, as a function of t , every eigenvalue of $\mathbf{A}(t)$ is a continuous function of t .

Example 13.8 Consider the matrix $\mathbf{A} = \begin{bmatrix} 1 & \epsilon_1 & \epsilon_2 \\ \epsilon_3 & 2 & \epsilon_4 \\ \epsilon_5 & \epsilon_6 & 3 \end{bmatrix}$, where $|\epsilon_i| \leq 10^{-15}$ all i . By Corollary 13.6 the eigenvalues $\lambda_1, \lambda_2, \lambda_3$ of \mathbf{A} are distinct and satisfy $|\lambda_j - j| \leq 2 \times 10^{-15}$ for $j = 1, 2, 3$.

Exercise 13.9 (Nonsingularity using Gershgorin)

Consider the matrix

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 0 & 0 \\ 1 & 4 & 1 & 0 \\ 0 & 1 & 4 & 1 \\ 0 & 0 & 1 & 4 \end{pmatrix}.$$

Show using Gershgorin's theorem that \mathbf{A} is nonsingular.

Exercise 13.10 (Gershgorin, strictly diagonally dominant matrix)

Show using Gershgorin's theorem that a strictly diagonally dominant matrix \mathbf{A} ($|a_{i,i}| > \sum_{j \neq i} |a_{i,j}|$ for all i) is nonsingular.

13.3 Unitary Similarity Transformation of a Matrix into Upper Hessenberg Form

Before attempting to find eigenvalues and eigenvectors of a matrix (exceptions are made for certain sparse matrices), it is often advantageous to reduce it by similarity transformations to a simpler form. Orthogonal or unitary similarity

transformations are particularly important since they are insensitive to noise in the elements of the matrix. In this section we show how this reduction can be carried out.

Recall that a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is upper Hessenberg if $a_{i,j} = 0$ for $j = 1, 2, \dots, i-2$, $i = 3, 4, \dots, n$. We will reduce $\mathbf{A} \in \mathbb{R}^{n \times n}$ to upper Hessenberg form by unitary similarity transformations. Let $\mathbf{A}_1 = \mathbf{A}$ and define $\mathbf{A}_{k+1} = \mathbf{H}_k \mathbf{A}_k \mathbf{H}_k$ for $k = 1, 2, \dots, n-2$. Here \mathbf{H}_k is a Householder transformation chosen to introduce zeros in the elements of column k of \mathbf{A}_k under the subdiagonal. The final matrix \mathbf{A}_{n-1} will be upper Hessenberg.

If $\mathbf{A}_1 = \mathbf{A}$ is symmetric, the matrix \mathbf{A}_{n-1} will be symmetric and tridiagonal. For if $\mathbf{A}_k^T = \mathbf{A}_k$ then

$$\mathbf{A}_{k+1}^T = (\mathbf{H}_k \mathbf{A}_k \mathbf{H}_k)^T = \mathbf{H}_k \mathbf{A}_k^T \mathbf{H}_k = \mathbf{A}_{k+1}.$$

Since \mathbf{A}_{n-1} is upper Hessenberg and symmetric, it must be tridiagonal.

To describe the reduction to upper Hessenberg or tridiagonal form in more detail we partition \mathbf{A}_k as follows

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{D}_k & \mathbf{E}_k \end{bmatrix}.$$

Suppose $\mathbf{B}_k \in \mathbb{R}^{k,k}$ is upper Hessenberg, and the first $k-1$ columns of $\mathbf{D}_k \in \mathbb{R}^{n-k,k}$ are zero, i.e. $\mathbf{D}_k = [\mathbf{0}, \mathbf{0}, \dots, \mathbf{0}, \mathbf{d}_k]$. Let $\mathbf{V}_k = \mathbf{I} - \mathbf{v}_k \mathbf{v}_k^T \in \mathbb{R}^{n-k, n-k}$ be a Householder transformation such that $\mathbf{V}_k \mathbf{d}_k = \alpha_k \mathbf{e}_1$, where $\alpha_k^2 = \mathbf{d}_k^T \mathbf{d}_k$. Define

$$\mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

The matrix \mathbf{H}_k is a Householder transformation, and we find

$$\begin{aligned} \mathbf{A}_{k+1} &= \mathbf{H}_k \mathbf{A}_k \mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{D}_k & \mathbf{E}_k \end{bmatrix} \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \mathbf{V}_k \\ \mathbf{V}_k \mathbf{D}_k & \mathbf{V}_k \mathbf{E}_k \mathbf{V}_k \end{bmatrix}. \end{aligned}$$

Now $\mathbf{V}_k \mathbf{D}_k = [\mathbf{V}_k \mathbf{0}, \dots, \mathbf{V}_k \mathbf{0}, \mathbf{V}_k \mathbf{d}_k] = (\mathbf{0}, \dots, \mathbf{0}, \alpha_k \mathbf{e}_1)$. Moreover, the matrix \mathbf{B}_k is not affected by the \mathbf{H}_k transformation. Therefore the upper left $(k+1) \times (k+1)$ corner of \mathbf{A}_{k+1} is upper Hessenberg and the reduction is carried one step further. The reduction stops with \mathbf{A}_{n-1} which is upper Hessenberg.

To find \mathbf{A}_{k+1} we use Algorithm 11.4 to find \mathbf{v}_k and α_k . We store \mathbf{v}_k in the k th column of a matrix \mathbf{L} as $\mathbf{L}(k+1 : n, k) = \mathbf{v}_k$. This leads to the following algorithm.

Algorithm 13.11 (Householder reduction to Hessenberg form) This algorithm uses Householder similarity transformations to reduce a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ to upper Hessenberg form. The reduced matrix \mathbf{B} is tridiagonal if \mathbf{A} is symmetric. Details of the transformations are stored in a lower triangular matrix \mathbf{L} . The elements of \mathbf{L} can be used to assemble an orthonormal matrix \mathbf{Q} such that $\mathbf{B} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}$. Algorithm 11.4 is used in each step of the reduction.

```

1 function [L,B] = hesshousegen(A)
2 n=length(A); L=zeros(n,n); B=A;
3 for k=1:n-2
4     [v,B(k+1,k)]=housegen(B(k+1:n,k));
5     L(k+1:n,k)=v; B(k+2:n,k)=zeros(n-k-1,1);
6     C=B(k+1:n,k+1:n); B(k+1:n,k+1:n)=C-v*(v'*C);
7     C=B(1:n,k+1:n); B(1:n,k+1:n)=C-(C*v)*v';
8 end
```

Exercise 13.12 (Number of arithmetic operations)

Show that the number of arithmetic operations for Algorithm 13.11 is $\frac{10}{3}n^3 = 5G_n$.

We can use the output of Algorithm 13.11 to assemble the matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ such that \mathbf{Q} is orthonormal and $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$ is upper Hessenberg. We need to compute the product $\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_{n-2}$, where $\mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{v}_k \mathbf{v}_k^T \end{bmatrix}$ and $\mathbf{v}_k \in \mathbb{R}^{n-k}$. Since $\mathbf{v}_1 \in \mathbb{R}^{n-1}$ and $\mathbf{v}_{n-2} \in \mathbb{R}^2$ it is most economical to assemble the product from right to left. We compute

$$\mathbf{Q}_{n-1} = \mathbf{I} \text{ and } \mathbf{Q}_k = \mathbf{H}_k \mathbf{Q}_{k+1} \text{ for } k = n-2, n-3, \dots, 1.$$

Suppose \mathbf{Q}_{k+1} has the form $\begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_{k+1} \end{bmatrix}$, where $\mathbf{U}_{k+1} \in \mathbb{R}^{n-k, n-k}$. Then

$$\mathbf{Q}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{v}_k \mathbf{v}_k^T \end{bmatrix} * \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_k - \mathbf{v}_k (\mathbf{v}_k^T \mathbf{U}_k) \end{bmatrix}.$$

This leads to the following algorithm.

Algorithm 13.13 (Assemble Householder transformations)

Suppose $[L, B] = \text{hesshousegen}(A)$ is the output of Algorithm 13.11. This algorithm assembles an orthonormal matrix \mathbf{Q} from the columns of \mathbf{L} such that $\mathbf{B} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}$ is upper Hessenberg.

```

1 function Q = accumulateQ(L)
2 n=length(L); Q=eye(n);
3 for k=n-2:-1:1
4     v=L(k+1:n,k); C=Q(k+1:n,k+1:n);
5     Q(k+1:n,k+1:n)=C-v*(v'*C);
6 end
```

Exercise 13.14 (Number of arithmetic operations)

Show that the number of arithmetic operations required by Algorithm 13.13 is $\frac{4}{3}n^3 = 2G_n$.

Exercise 13.15 (Tridiagonalize a symmetric matrix)

If \mathbf{A} is symmetric we can modify Algorithm 13.11 as follows. To find \mathbf{A}_{k+1} from \mathbf{A}_k we have to compute $\mathbf{V}_k \mathbf{E}_k \mathbf{V}_k$ where \mathbf{E}_k is symmetric. Dropping subscripts we have to compute a product of the form $\mathbf{G} = (\mathbf{I} - \mathbf{v}\mathbf{v}^T)\mathbf{E}(\mathbf{I} - \mathbf{v}\mathbf{v}^T)$. Let $\mathbf{w} := \mathbf{E}\mathbf{v}$, $\beta := \frac{1}{2}\mathbf{v}^T\mathbf{w}$ and $\mathbf{z} := \mathbf{w} - \beta\mathbf{v}$. Show that $\mathbf{G} = \mathbf{E} - \mathbf{v}\mathbf{z}^T - \mathbf{z}\mathbf{v}^T$. Since \mathbf{G} is symmetric, only the sub- or superdiagonal elements of \mathbf{G} need to be computed. Computing \mathbf{G} in this way, it can be shown that we need $O(4n^3/3)$ operations to tridiagonalize a symmetric matrix by orthonormal similarity transformations. This is less than half the work to reduce a nonsymmetric matrix to upper Hessenberg form. We refer to [23] for a detailed algorithm.

13.4 Computing a Selected Eigenvalue of a Symmetric Matrix

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be symmetric with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. In this section we consider a method to compute an approximation to the m th eigenvalue λ_m for some $1 \leq m \leq n$. Using Householder similarity transformations as outlined in the previous section we can assume that \mathbf{A} is symmetric and tridiagonal.

$$\mathbf{A} = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}. \quad (13.3)$$

Suppose one of the off-diagonal elements is equal to zero, say $c_i = 0$. We then have $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}$, where

$$\mathbf{A}_1 = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{i-2} & d_{i-1} & c_{i-1} \\ & & & c_{i-1} & d_i \end{bmatrix} \text{ and } \mathbf{A}_2 = \begin{bmatrix} d_{i+1} & c_{i+1} & & & \\ c_{i+1} & d_{i+2} & c_{i+2} & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}.$$

Thus \mathbf{A} is block diagonal and each diagonal block is tridiagonal. By 6. of Theorem 0.66 we can split the eigenvalue problem into two smaller problems involving \mathbf{A}_1 and \mathbf{A}_2 . We assume that this reduction has been carried out so that \mathbf{A} is irreducible, i. e., $c_i \neq 0$ for $i = 1, \dots, n-1$.

We first show that irreducibility implies that the eigenvalues are distinct.

Lemma 13.16 (Distinct eigenvalues of a tridiagonal matrix)

An irreducible, tridiagonal and symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has n real and distinct eigenvalues.

Proof. Let \mathbf{A} be given by (13.3). By Theorem 6.39 the eigenvalues are real. Define for $x \in \mathbb{R}$ the polynomial $p_k(x) := \det(x\mathbf{I}_k - \mathbf{A}_k)$ for $k = 1, \dots, n$, where \mathbf{A}_k is the upper left $k \times k$ corner of \mathbf{A} (the leading principal submatrix of order k). The eigenvalues of \mathbf{A} are the roots of the polynomial p_n . Using the last column to expand for $k \geq 2$ the determinant $p_{k+1}(x)$ we find

$$p_{k+1}(x) = (x - d_{k+1})p_k(x) - c_k^2 p_{k-1}(x). \quad (13.4)$$

Since $p_1(x) = x - d_1$ and $p_2(x) = (x - d_2)(x - d_1) - c_1^2$ this also holds for $k = 0, 1$ if we define $p_{-1}(x) = 0$ and $p_0(x) = 1$. For M sufficiently large we have

$$p_2(-M) > 0, \quad p_2(d_1) < 0, \quad p_2(+M) > 0.$$

Since p_2 is continuous there are $y_1 \in (-M, d_1)$ and $y_2 \in (d_1, M)$ such that $p_2(y_1) = p_2(y_2) = 0$. It follows that the root d_1 of p_1 separates the roots of p_2 , so y_1 and y_2 must be distinct. Consider next

$$p_3(x) = (x - d_3)p_2(x) - c_2^2 p_1(x) = (x - d_3)(x - y_1)(x - y_2) - c_2^2(x - d_1).$$

Since $y_1 < d_1 < y_2$ we have for M sufficiently large

$$p_3(-M) < 0, \quad p_3(y_1) > 0, \quad p_3(y_2) < 0, \quad p_3(+M) > 0.$$

Thus the roots x_1, x_2, x_3 of p_3 are separated by the roots y_1, y_2 of p_2 . In the general case suppose for $k \geq 2$ that the roots z_1, \dots, z_{k-1} of p_{k-1} separate the roots y_1, \dots, y_k of p_k . Choose M so that $y_0 := -M < y_1, y_{k+1} := M > y_k$. Then

$$y_0 < y_1 < z_1 < y_2 < z_2 \cdots < z_{k-1} < y_k < y_{k+1}.$$

We claim that for M sufficiently large

$$p_{k+1}(y_j) = (-1)^{k+1-j}|p_{k+1}(y_j)| \neq 0, \text{ for } j = 0, 1, \dots, k+1.$$

This holds for $j = 0, k+1$, and for $j = 1, \dots, k$ since

$$p_{k+1}(y_j) = -c_k^2 p_{k-1}(y_j) = -c_k^2(y_j - z_1) \cdots (y_j - z_{k-1}).$$

It follows that the roots x_1, \dots, x_{k+1} are separated by the roots y_1, \dots, y_k of p_k and by induction the roots of p_n (the eigenvalues of \mathbf{A}) are distinct. \square

13.4.1 The Inertia Theorem

We say that two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ are **congruent** if $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$ for some nonsingular matrix $\mathbf{E} \in \mathbb{C}^{n \times n}$. By Theorem 6.37 a Hermitian matrix \mathbf{A} is both congruent and similar to a diagonal matrix \mathbf{D} , $\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{D}$ where \mathbf{U} is unitary. The eigenvalues of \mathbf{A} are the diagonal elements of \mathbf{D} . Let $\pi(\mathbf{A})$, $\zeta(\mathbf{A})$ and $v(\mathbf{A})$ denote the number of positive, zero and negative eigenvalues of \mathbf{A} . If \mathbf{A} is Hermitian then all eigenvalues are real and $\pi(\mathbf{A}) + \zeta(\mathbf{A}) + v(\mathbf{A}) = n$.

Theorem 13.17 (Sylvester's Inertia Theorem)

If $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ are Hermitian and congruent then $\pi(\mathbf{A}) = \pi(\mathbf{B})$, $\zeta(\mathbf{A}) = \zeta(\mathbf{B})$ and $v(\mathbf{A}) = v(\mathbf{B})$.

Proof. Suppose $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$, where \mathbf{E} is nonsingular. Assume first that \mathbf{A} and \mathbf{B} are diagonal matrices. Suppose $\pi(\mathbf{A}) = k$ and $\pi(\mathbf{B}) = m < k$. We shall show that this leads to a contradiction. Let \mathbf{E}_1 be the upper left $m \times k$ corner of \mathbf{E} . Since $m < k$, we can find a nonzero \mathbf{x} such that $\mathbf{E}_1 \mathbf{x} = \mathbf{0}$ (cf. Lemma 0.44). Let $\mathbf{y}^T = [\mathbf{x}^T, \mathbf{0}^T] \in \mathbb{C}^n$, and $\mathbf{z} = [z_1, \dots, z_n]^T = \mathbf{E} \mathbf{y}$. Then $z_i = 0$ for $i = 1, 2, \dots, m$. If \mathbf{A} has positive eigenvalues $\lambda_1, \dots, \lambda_k$ and \mathbf{B} has eigenvalues μ_1, \dots, μ_n , where $\mu_i \leq 0$ for $i \geq m+1$ then

$$\mathbf{y}^* \mathbf{A} \mathbf{y} = \sum_{i=1}^n \lambda_i |y_i|^2 = \sum_{i=1}^k \lambda_i |x_i|^2 > 0.$$

But

$$\mathbf{y}^* \mathbf{A} \mathbf{y} = \mathbf{y}^* \mathbf{E}^* \mathbf{B} \mathbf{E} \mathbf{y} = \mathbf{z}^* \mathbf{B} \mathbf{z} = \sum_{i=m+1}^n \mu_i |z_i|^2 \leq 0,$$

a contradiction.

We conclude that $\pi(\mathbf{A}) = \pi(\mathbf{B})$ if \mathbf{A} and \mathbf{B} are diagonal. Moreover, $v(\mathbf{A}) = \pi(-\mathbf{A}) = \pi(-\mathbf{B}) = v(\mathbf{B})$ and $\zeta(\mathbf{A}) = n - \pi(\mathbf{A}) - v(\mathbf{A}) = n - \pi(\mathbf{B}) - v(\mathbf{B}) = \zeta(\mathbf{B})$. This completes the proof for diagonal matrices.

Let in the general case \mathbf{U}_1 and \mathbf{U}_2 be unitary matrices such that $\mathbf{U}_1^* \mathbf{A} \mathbf{U}_1 = \mathbf{D}_1$ and $\mathbf{U}_2^* \mathbf{B} \mathbf{U}_2 = \mathbf{D}_2$ where \mathbf{D}_1 and \mathbf{D}_2 are diagonal matrices. Since $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$, we find $\mathbf{D}_1 = \mathbf{F}^* \mathbf{D}_2 \mathbf{F}$ where $\mathbf{F} = \mathbf{U}_2^* \mathbf{E} \mathbf{U}_1$ is nonsingular. Thus \mathbf{D}_1 and \mathbf{D}_2 are congruent diagonal matrices. But since \mathbf{A} and \mathbf{D}_1 , \mathbf{B} and \mathbf{D}_2 have the same eigenvalues, we find $\pi(\mathbf{A}) = \pi(\mathbf{D}_1) = \pi(\mathbf{D}_2) = \pi(\mathbf{B})$. Similar results hold for ζ and v . \square

Corollary 13.18 (Counting eigenvalues using the LDLT factorization)

Suppose $\mathbf{A} = \text{tridiag}(c_i, d_i, c_i) \in \mathbb{R}^{n \times n}$ is symmetric and that $\alpha \in \mathbb{R}$ is such that $\mathbf{A} - \alpha \mathbf{I}$ has an symmetric LU factorization, i.e. $\mathbf{A} - \alpha \mathbf{I} = \mathbf{L} \mathbf{D} \mathbf{L}^T$ where \mathbf{L} is unit

lower triangular and \mathbf{D} is diagonal. Then the number of eigenvalues of \mathbf{A} strictly less than α equals the number of negative diagonal elements in \mathbf{D} . The diagonal elements $d_1(\alpha), \dots, d_n(\alpha)$ in \mathbf{D} can be computed recursively as follows

$$d_1(\alpha) = d_1 - \alpha, \quad d_k(\alpha) = d_k - \alpha - c_{k-1}^2/d_{k-1}(\alpha), \quad k = 2, 3, \dots, n. \quad (13.5)$$

Proof. Since the diagonal elements in \mathbf{R} in an LU factorization equal the diagonal elements in \mathbf{D} in an \mathbf{LDL}^T factorization we see that the formulas in (13.5) follows immediately from (2.4). Since \mathbf{L} is nonsingular, $\mathbf{A} - \alpha\mathbf{I}$ and \mathbf{D} are congruent. By the previous theorem $v(\mathbf{A} - \alpha\mathbf{I}) = v(\mathbf{D})$, the number of negative diagonal elements in \mathbf{D} . If $\mathbf{Ax} = \lambda\mathbf{x}$ then $(\mathbf{A} - \alpha\mathbf{I})\mathbf{x} = (\lambda - \alpha)\mathbf{x}$, and $\lambda - \alpha$ is an eigenvalue of $\mathbf{A} - \alpha\mathbf{I}$. But then $v(\mathbf{A} - \alpha\mathbf{I})$ equals the number of eigenvalues of \mathbf{A} which are less than α . \square

Exercise 13.19 (Counting eigenvalues)

Consider the matrix in Exercise 13.9. Determine the number of eigenvalues greater than 4.5.

Exercise 13.20 (Overflow in LDLT factorization)

Let for $n \in \mathbb{N}$

$$\mathbf{A}_n = \begin{bmatrix} 10 & 1 & 0 & \cdots & 0 \\ 1 & 10 & 1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & 1 & 10 & 1 \\ 0 & \cdots & 0 & 1 & 10 \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

- a) Let d_k be the diagonal elements of \mathbf{D} in a symmetric factorization of \mathbf{A}_n . Show that $5 + \sqrt{24} < d_k \leq 10$, $k = 1, 2, \dots, n$.
- b) Show that $D_n := \det(\mathbf{A}_n) > (5 + \sqrt{24})^n$. Give $n_0 \in \mathbb{N}$ such that your computer gives an overflow when D_{n_0} is computed in floating point arithmetic.

Exercise 13.21 (Simultaneous diagonalization)

(Simultaneous diagonalization of two symmetric matrices by a congruence transformation). Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ where $\mathbf{A}^T = \mathbf{A}$ and \mathbf{B} is symmetric positive definite. Let $\mathbf{B} = \mathbf{U}^T \mathbf{D} \mathbf{U}$ where \mathbf{U} is orthonormal and $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$. Let $\hat{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{U} \mathbf{A} \mathbf{U}^T \mathbf{D}^{-1/2}$ where

$$\mathbf{D}^{-1/2} := \text{diag}(d_1^{-1/2}, \dots, d_n^{-1/2}).$$

a) Show that $\hat{\mathbf{A}}$ is symmetric.

Let $\hat{\mathbf{A}} = \hat{\mathbf{U}}^T \hat{\mathbf{D}} \hat{\mathbf{U}}$ where $\hat{\mathbf{U}}$ is orthonormal and $\hat{\mathbf{D}}$ is diagonal. Set $\mathbf{E} = \mathbf{U}^T \mathbf{D}^{-1/2} \hat{\mathbf{U}}^T$.

b) Show that \mathbf{E} is nonsingular and that $\mathbf{E}^T \mathbf{A} \mathbf{E} = \hat{\mathbf{D}}$, $\mathbf{E}^T \mathbf{B} \mathbf{E} = \mathbf{I}$.

For a more general result see Theorem 10.1 in [15].

13.4.2 Approximating λ_m

Corollary 13.18 can be used to determine the m th eigenvalue of \mathbf{A} , where $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. For this we use interval bisection. Using Gershgorin's theorem we first find an interval $[a, b]$, such that $[a, b]$ contains the eigenvalues of \mathbf{A} . Let for $x \in [a, b]$

$$\rho(x) := \#\{k : d_k(x) < 0 \text{ for } k = 1, \dots, n\}$$

be the number of eigenvalues of \mathbf{A} which are strictly less than x . Clearly $\rho(a) = 0$, $\rho(b) = n$ and $\rho(e) - \rho(d)$ is the number of eigenvalues in $[d, e]$. Let $c = (a + b)/2$ and $k := \rho(c)$. If $k \geq m$ then $\lambda_m \leq c$ and $\lambda_m \in [a, c]$, while if $k < m$ then $\lambda_m \geq c$ and $\lambda_m \in [c, b]$. Continuing with the interval containing λ_m we generate a sequence $\{[a_j, b_j]\}$ of intervals, each containing λ_m and $b_j - a_j = 2^{-j}(b - a)$.

As it stands this method will fail if in (13.5) one of the $d_k(\alpha)$ is zero. One possibility is to replace such a $d_k(\alpha)$ by a suitable small number, say $\delta_k = \pm|c_k|\epsilon_M$, where the negative sign is used if $c_k < 0$, and ϵ_M is the Machine epsilon, typically 2×10^{-16} for Matlab. This replacement is done if $|d_k(\alpha)| < |\delta_k|$.

Exercise 13.22 (Program code for one eigenvalue)

Suppose $\mathbf{A} = \text{tridiag}(\mathbf{c}, \mathbf{d}, \mathbf{c})$ is symmetric and tridiagonal with elements d_1, \dots, d_n on the diagonal and c_1, \dots, c_{n-1} on the neighboring subdiagonals. Let $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ be the eigenvalues of \mathbf{A} . We shall write a program to compute one eigenvalue λ_m for a given m using bisection and the method outlined in Section 13.4.2.

- a) Write a function `k=count(c, d, x)` which for given x counts the number of eigenvalues of \mathbf{A} strictly less than x . Use the replacement described above if one of the $d_j(x)$ is close to zero.
- b) Write a function `lambda=findeigv(c, d, m)` which first estimates an interval $[a, b]$ containing all eigenvalues of \mathbf{A} and then generates a sequence $\{[a_k, b_k]\}$ of intervals each containing λ_m . Iterate until $b_k - a_k \leq (b - a)\epsilon_M$, where ϵ_M is Matlab's machine epsilon `eps`. Typically $\epsilon_M \approx 2.22 \times 10^{-16}$.
- c) Test the program on $\mathbf{T} := \text{tridiag}(-1, 2, -1)$ of size 100. Compare the exact value of λ_5 with your result and the result obtained by using Matlab's built-in function `eig`.

Exercise 13.23 (Determinant of upper Hessenberg matrix)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ is upper Hessenberg and $x \in \mathbb{C}$. We will study two algorithms to compute $f(x) = \det(\mathbf{A} - x\mathbf{I})$.

- a) Show that Gaussian elimination without pivoting requires $O(n^2)$ arithmetic operations.
- b) Show that the number of arithmetic operations is the same if partial pivoting is used.
- c) Estimate the number of arithmetic operations if Given's rotations are used.
- d) Compare the two methods discussing advantages and disadvantages.

13.5 Perturbation Proofs

We first show that the p -norm of a diagonal matrix is equal to its spectral radius.

Lemma 13.24 (p -norm of a diagonal matrix)

If $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is a diagonal matrix then $\|\mathbf{A}\|_p = \rho(\mathbf{A})$ for $1 \leq p \leq \infty$.

Proof. For $p = \infty$ the proof is left as an exercise. For any $\mathbf{x} \in \mathbb{C}^n$ and $p < \infty$ we have

$$\|\mathbf{Ax}\|_p = \|[\lambda_1 x_1, \dots, \lambda_n x_n]^T\|_p = \left(\sum_{j=1}^n |\lambda_j|^p |x_j|^p \right)^{1/p} \leq \rho(\mathbf{A}) \|\mathbf{x}\|_p.$$

Thus $\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p} \leq \rho(\mathbf{A})$. But from Theorem 9.6 we have $\rho(\mathbf{A}) \leq \|\mathbf{A}\|_p$ and the proof is complete. \square

Exercise 13.25 (∞ -norm of a diagonal matrix)

Give a direct proof that $\|\mathbf{A}\|_\infty = \rho(\mathbf{A})$ if \mathbf{A} is diagonal.

Suppose now (μ, \mathbf{x}) is an approximation to an eigenpair of a matrix \mathbf{A} . One way to check the accuracy is to compute the residual $\mathbf{r} := \mathbf{Ax} - \mu\mathbf{x}$. For an exact eigenpair the residual is zero and we could hope that a small residual implies an accurate eigenpair.

Theorem 13.26 (Absolute errors)

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ has linearly independent eigenvectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ be the eigenvector matrix. To any $\mu \in \mathbb{C}$ and $\mathbf{x} \in \mathbb{C}^n$ with $\|\mathbf{x}\|_p = 1$ we can find an eigenvalue λ of \mathbf{A} such that

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{r}\|_p, \quad 1 \leq p \leq \infty, \tag{13.6}$$

where $\mathbf{r} := \mathbf{Ax} - \mu\mathbf{x}$ and $K_p(\mathbf{X}) := \|\mathbf{X}\|_p \|\mathbf{X}^{-1}\|_p$. If for some $\mathbf{E} \in \mathbb{C}^{n \times n}$ it holds that (μ, \mathbf{x}) is an eigenpair for $\mathbf{A} + \mathbf{E}$, then we can find an eigenvalue λ of \mathbf{A} such that

$$|\lambda - \mu| \leq K_p(\mathbf{X})\|\mathbf{E}\|_p, \quad 1 \leq p \leq \infty, \quad (13.7)$$

Proof. If $\mu \in \sigma(\mathbf{A})$ then we can take $\lambda = \mu$ and (13.6), (13.7) hold trivially. So assume $\mu \notin \sigma(\mathbf{A})$. Since \mathbf{A} is nondefective it can be diagonalized, we have $\mathbf{A} = \mathbf{XD}^{-1}\mathbf{X}^{-1}$, where $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ and $(\lambda_j, \mathbf{x}_j)$ are the eigenpairs of \mathbf{A} for $j = 1, \dots, n$. Define $\mathbf{D}_1 := \mathbf{D} - \mu\mathbf{I}$. Then $\mathbf{D}_1^{-1} = \text{diag}((\lambda_1 - \mu)^{-1}, \dots, (\lambda_n - \mu)^{-1})$ exists and

$$\mathbf{XD}_1^{-1}\mathbf{X}^{-1}\mathbf{r} = (\mathbf{X}(\mathbf{D} - \mu\mathbf{I})\mathbf{X}^{-1})^{-1}\mathbf{r} = (\mathbf{A} - \mu\mathbf{I})^{-1}(\mathbf{A} - \mu\mathbf{I})\mathbf{x} = \mathbf{x}.$$

Using this and Lemma 13.24 we obtain

$$1 = \|\mathbf{x}\|_p = \|\mathbf{XD}_1^{-1}\mathbf{X}^{-1}\mathbf{r}\|_p \leq \|\mathbf{D}_1^{-1}\|_p K_p(\mathbf{X})\|\mathbf{r}\|_p = \frac{K_p(\mathbf{X})\|\mathbf{r}\|_p}{\min_j |\lambda_j - \mu|}.$$

But then (13.6) follows. If $(\mathbf{A} + \mathbf{E})\mathbf{x} = \mu\mathbf{x}$ then $\mathbf{0} = \mathbf{Ax} - \mu\mathbf{x} + \mathbf{Ex} = \mathbf{r} + \mathbf{Ex}$. But then $\|\mathbf{r}\|_p = \|\mathbf{Ex}\|_p \leq \|\mathbf{E}\|_p$. Inserting this in (13.6) proves (13.7). \square

For the accuracy of an eigenvalue of small magnitude we are interested in the size of the relative error.

Theorem 13.27 (Relative errors)

Suppose in Theorem 13.26 that $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular. To any $\mu \in \mathbb{C}$ and $\mathbf{x} \in \mathbb{C}^n$ with $\|\mathbf{x}\|_p = 1$, we can find an eigenvalue λ of \mathbf{A} such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{r}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (13.8)$$

where $\mathbf{r} := \mathbf{Ax} - \mu\mathbf{x}$. If for some $\mathbf{E} \in \mathbb{C}^{n \times n}$ it holds that (μ, \mathbf{x}) is an eigenpair for $\mathbf{A} + \mathbf{E}$, then we can find an eigenvalue λ of \mathbf{A} such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})\|\mathbf{A}^{-1}\mathbf{E}\|_p \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{E}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (13.9)$$

Proof. Applying Theorem 9.6 to \mathbf{A}^{-1} we have for any $\lambda \in \sigma(\mathbf{A})$

$$\frac{1}{\lambda} \leq \|\mathbf{A}^{-1}\|_p = \frac{K_p(\mathbf{A})}{\|\mathbf{A}\|_p}$$

and (13.8) follows from (13.6). To prove (13.9) we define the matrices $\mathbf{B} := \mu\mathbf{A}^{-1}$ and $\mathbf{F} := -\mathbf{A}^{-1}\mathbf{E}$. If (λ_j, \mathbf{x}) are the eigenpairs for \mathbf{A} then $(\frac{\mu}{\lambda_j}, \mathbf{x})$ are the eigenpairs for \mathbf{B} for $j = 1, \dots, n$. Since (μ, \mathbf{x}) is an eigenpair for $\mathbf{A} + \mathbf{E}$ we find

$$(\mathbf{B} + \mathbf{F} - \mathbf{I})\mathbf{x} = (\mu\mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{E} - \mathbf{I})\mathbf{x} = \mathbf{A}^{-1}(\mu\mathbf{I} - (\mathbf{E} + \mathbf{A}))\mathbf{x} = \mathbf{0}.$$

Thus $(1, \mathbf{x})$ is an eigenpair for $\mathbf{B} + \mathbf{F}$. Applying Theorem 13.26 to this eigenvalue we can find $\lambda \in \sigma(\mathbf{A})$ such that $|\frac{\mu}{\lambda} - 1| \leq K_p(\mathbf{X})\|\mathbf{F}\|_p = K_p(\mathbf{X})\|\mathbf{A}^{-1}\mathbf{E}\|_p$ which proves the first estimate in (13.9). The second inequality in (13.9) follows from the submultiplicativity of the p -norm. \square

13.6 Review Questions

13.6.1 Suppose $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$. To every $\mu \in \sigma(\mathbf{A} + \mathbf{E})$ there is a $\lambda \in \sigma(\mathbf{A})$ which is in some sense close to μ .

- What is the general result (Elsner's theorem)?
- what if \mathbf{A} is non defective?
- what if \mathbf{A} is normal?
- what if \mathbf{A} is Hermitian?

13.6.2 Can Gershgorin's theorem be used to check if a matrix is nonsingular?

13.6.3 How many arithmetic operation does it take to reduce a matrix by similarity transformations to upper Hessenberg form by Householder transformations?

13.6.4 Give a condition ensuring that a tridiagonal symmetric matrix has real and distinct eigenvalues:

13.6.5 What is the content of Sylvester's inertia theorem?

13.6.6 Give an application of this theorem.

Chapter 14

The QR Algorithm

The QR algorithm is a method to find all eigenvalues and eigenvectors of a matrix. It is related to a simpler method called the power method and we start studying this method and its variants.

14.1 The Power Method

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ have eigenpairs $(\lambda_j, \mathbf{v}_j)$, $j = 1, \dots, n$. Given $\mathbf{z}_0 \in \mathbb{C}^n$ we assume that

- (i) $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$,
 - (ii) $\mathbf{z}_0^T \mathbf{v}_1 \neq 0$
 - (iii) \mathbf{A} has linearly independent eigenvectors.
- (14.1)

The first assumption means that \mathbf{A} has a dominant eigenvalue λ_1 of algebraic multiplicity one. The second assumption says that \mathbf{z}_0 has a component in the direction \mathbf{v}_1 . The third assumption is not necessary, but is included in order to simplify the analysis.

The **power method** is a technique to compute the dominant eigenvector \mathbf{v}_1 of \mathbf{A} . As a by product we can also find the corresponding eigenvalue. We define a sequence $\{\mathbf{z}_k\}$ of vectors in \mathbb{C}^n by

$$\mathbf{z}_k := \mathbf{A}^k \mathbf{z}_0 = \mathbf{A} \mathbf{z}_{k-1}, \quad k = 1, 2, \dots \quad (14.2)$$

To see what happens let $\mathbf{z}_0 = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_n \mathbf{v}_n$, where by assumption (ii) of (14.1) we have $c_1 \neq 0$. Since $\mathbf{A}^k \mathbf{v}_j = \lambda_j^k \mathbf{v}_j$ for all j we see that

$$\mathbf{z}_k = c_1 \lambda_1^k \mathbf{v}_1 + c_2 \lambda_2^k \mathbf{v}_2 + \dots + c_n \lambda_n^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (14.3)$$

Dividing by λ_1^k we find

$$\frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1 + c_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \cdots + c_n \left(\frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (14.4)$$

Assumption (i) of (14.1) implies that $(\lambda_j/\lambda_1)^k \rightarrow 0$ as $k \rightarrow \infty$ for all $j \geq 2$ and we obtain

$$\lim_{k \rightarrow \infty} \frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1, \quad (14.5)$$

the dominant eigenvector of \mathbf{A} . It can be shown that this also holds for defective matrices as long as (i) and (ii) of (14.1) hold, see for example page 58 of [23].

In practice we need to scale the iterates \mathbf{z}_k somehow and we normally do not know λ_1 . Instead we choose a norm on \mathbb{C}^n , set $\mathbf{x}_0 = \mathbf{z}_0/\|\mathbf{z}_0\|$ and generate for $k = 1, 2, \dots$ unit vectors as follows:

$$\begin{aligned} (i) \quad & \mathbf{y}_k = \mathbf{A}\mathbf{x}_{k-1} \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k/\|\mathbf{y}_k\|. \end{aligned} \quad (14.6)$$

Lemma 14.1 (Convergence of the power method)

Suppose (14.1) holds. Then

$$\lim_{k \rightarrow \infty} \left(\frac{|\lambda_1|}{\lambda_1} \right)^k \mathbf{x}_k = \frac{c_1}{|c_1|} \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}.$$

In particular, if $\lambda_1 > 0$ and $c_1 > 0$ then the sequence $\{\mathbf{x}_k\}$ will converge to the eigenvector $\mathbf{u}_1 := \mathbf{v}_1/\|\mathbf{v}_1\|$ of unit length.

Proof. By induction on k it follows that $\mathbf{x}_k = \mathbf{z}_k/\|\mathbf{z}_k\|$ for all $k \geq 0$, where $\mathbf{z}_k = \mathbf{A}^k \mathbf{z}_0$. Indeed, this holds for $k = 1$, and if it holds for $k - 1$ then $\mathbf{y}_k = \mathbf{A}\mathbf{x}_{k-1} = \mathbf{A}\mathbf{z}_{k-1}/\|\mathbf{z}_{k-1}\| = \mathbf{z}_k/\|\mathbf{z}_{k-1}\|$ and $\mathbf{x}_k = (\mathbf{z}_k/\|\mathbf{z}_{k-1}\|)(\|\mathbf{z}_{k-1}\|/\|\mathbf{z}_k\|) = \mathbf{z}_k/\|\mathbf{z}_k\|$. But then

$$\mathbf{x}_k = \frac{\mathbf{z}_k}{\|\mathbf{z}_k\|} = \frac{c_1 \lambda_1^k}{|c_1 \lambda_1^k|} \frac{\mathbf{v}_1 + \frac{c_2}{c_1} \left(\frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \cdots + \frac{c_n}{c_1} \left(\frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n}{\|\mathbf{v}_1 + \frac{c_2}{c_1} \left(\frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \cdots + \frac{c_n}{c_1} \left(\frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n\|}, \quad k = 0, 1, 2, \dots,$$

and this implies the lemma. \square

Suppose we know an approximate eigenvector \mathbf{u} of \mathbf{A} , but not the corresponding eigenvalue μ . One way of estimating μ is to minimize the Euclidian norm of the residual $r(\lambda) := \mathbf{A}\mathbf{u} - \lambda\mathbf{u}$.

Theorem 14.2 (The Rayleigh quotient minimizes the residual)

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mathbf{u} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$, and let $\rho : \mathbb{C} \rightarrow \mathbb{R}$ be given by $\rho(\lambda) = \|\mathbf{A}\mathbf{u} - \lambda\mathbf{u}\|_2$. Then ρ is minimized when $\lambda := \frac{\mathbf{u}^* \mathbf{A} \mathbf{u}}{\mathbf{u}^* \mathbf{u}}$, the Rayleigh quotient for \mathbf{A} .

Proof. Assume $\mathbf{u}^*\mathbf{u} = 1$ and extend \mathbf{u} to an orthonormal basis $\{\mathbf{u}, \mathbf{U}\}$ for \mathbb{C}^n . Then $\mathbf{U}^*\mathbf{u} = \mathbf{0}$ and

$$\begin{bmatrix} \mathbf{u}^* \\ \mathbf{U}^* \end{bmatrix} (\mathbf{A}\mathbf{u} - \lambda\mathbf{u}) = \begin{bmatrix} \mathbf{u}^*\mathbf{A}\mathbf{u} - \lambda\mathbf{u}^*\mathbf{u} \\ \mathbf{U}^*\mathbf{A}\mathbf{u} - \lambda\mathbf{U}^*\mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{u}^*\mathbf{A}\mathbf{u} - \lambda \\ \mathbf{U}^*\mathbf{A}\mathbf{u} \end{bmatrix}.$$

By unitary invariance of the Euclidian norm

$$\rho(\lambda)^2 = |\mathbf{u}^*\mathbf{A}\mathbf{u} - \lambda|^2 + \|\mathbf{U}^*\mathbf{A}\mathbf{u}\|_2^2,$$

and ρ has a global minimum at $\lambda = \mathbf{u}^*\mathbf{A}\mathbf{u}$. \square

Exercise 14.3 (Orthogonal vectors)

Show that \mathbf{u} and $\mathbf{A}\mathbf{u} - \lambda\mathbf{u}$ are orthogonal when $\lambda = \frac{\mathbf{u}^*\mathbf{A}\mathbf{u}}{\mathbf{u}^*\mathbf{u}}$.

Using Rayleigh quotients we can incorporate the calculation of the eigenvalue into the power iteration. We can then compute the residual and stop the iteration when the residual is sufficiently small. But what is sufficiently small? Recall that if \mathbf{A} is nonsingular and nondefective with eigenvector matrix \mathbf{X} and (μ, \mathbf{u}) is an approximate eigenpair with $\|\mathbf{u}\|_2 = 1$, then by (13.8) we can find an eigenvalue λ of \mathbf{A} such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_2(\mathbf{X})K_2(\mathbf{A}) \frac{\|\mathbf{A}\mathbf{u} - \mu\mathbf{u}\|_2}{\|\mathbf{A}\|_2}.$$

Thus if the relative residual is small and both \mathbf{A} and \mathbf{X} are well conditioned then the relative error in the eigenvalue will be small.

This discussion leads to the power method with Rayleigh quotient computation. Given $\mathbf{A} \in \mathbb{C}^{n \times n}$, a starting vector $\mathbf{z} \in \mathbb{C}^n$, a maximum number K of iterations, and a convergence tolerance tol . The power method combined with a Rayleigh quotient estimate for the eigenvalue is used to compute a dominant eigenpair (l, \mathbf{x}) of \mathbf{A} with $\|\mathbf{x}\|_2 = 1$. The integer it returns the number of iterations needed in order for $\|\mathbf{Ax} - l\mathbf{x}\|_2 / \|\mathbf{A}\|_F < tol$. If no such eigenpair is found in K iterations the value $it = K + 1$ is returned.

Algorithm 14.4 (The Power Method)

```

1 function [l,x,it]=powerit(A,z,K,tol)
2 af=norm(A,'fro'); x=z/norm(z);
3 for k=1:K
4     y=A*x; l=x'*y;
5     if norm(y-l*x)/af<tol
6         it=k; x=y/norm(y); return
7     end
8     x=y/norm(y);
9 end
10 it=K+1;
```

Example 14.5 (Power method)

We try powerit on the three matrices

$$\mathbf{A}_1 := \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 1.7 & -0.4 \\ 0.15 & 2.2 \end{bmatrix}, \text{ and } \mathbf{A}_3 = \begin{bmatrix} 1 & 2 \\ -3 & 4 \end{bmatrix}.$$

In each case we start with the random vector $\mathbf{z} = [0.6602, 0.3420]$ and $\text{tol} = 10^{-6}$. For \mathbf{A}_1 we get convergence in 7 iterations, for \mathbf{A}_2 it takes 174 iterations, and for \mathbf{A}_3 we do not get convergence.

The matrix \mathbf{A}_3 does not have a dominant eigenvalue since the two eigenvalues are complex conjugate of each other. Thus the basic condition (i) of (14.1) is not satisfied and the power method diverges. The enormous difference in the rate of convergence for \mathbf{A}_1 and \mathbf{A}_2 can be explained by looking at (14.4). The rate of convergence depends on the ratio $\frac{|\lambda_2|}{|\lambda_1|}$. If this ratio is small then the convergence is fast, while it can be quite slow if the ratio is close to one. The eigenvalues of \mathbf{A}_1 are $\lambda_1 = 5.3723$ and $\lambda_2 = -0.3723$ giving a quite small ratio of 0.07 and the convergence is fast. On the other hand the eigenvalues of \mathbf{A}_2 are $\lambda_1 = 2$ and $\lambda_2 = 1.9$ and the corresponding ratio is 0.95 resulting in slow convergence.

A variant of the power method is the **shifted power method**. In this method we choose a number s and apply the power method to the matrix $\mathbf{A} - s\mathbf{I}$. The number s is called a shift since it shifts an eigenvalue λ of \mathbf{A} to $\lambda - s$ of $\mathbf{A} - s\mathbf{I}$. Sometimes the convergence can be faster if the shift is chosen intelligently. For example, if we apply the shifted power method to \mathbf{A}_2 in Example 14.5 with shift 1.8, then with the same starting vector and tol as above, we get convergence in 17 iterations instead of 174 for the unshifted algorithm.

14.1.1 The Inverse Power Method

Another variant of the power method with Rayleigh quotient is the **inverse power method**. This method can be used to determine any eigenpair (λ, \mathbf{x}) of \mathbf{A} as long as λ has algebraic multiplicity one. In the inverse power method we apply the power method to the inverse matrix $(\mathbf{A} - s\mathbf{I})^{-1}$, where s is a shift. If \mathbf{A} has eigenvalues $\lambda_1, \dots, \lambda_n$ in no particular order then $(\mathbf{A} - s\mathbf{I})^{-1}$ has eigenvalues

$$\mu_1(s) = (\lambda_1 - s)^{-1}, \mu_2(s) = (\lambda_2 - s)^{-1}, \dots, \mu_n(s) = (\lambda_n - s)^{-1}.$$

Suppose λ_1 is a simple eigenvalue of \mathbf{A} . Then $\lim_{s \rightarrow \lambda_1} |\mu_1(s)| = \infty$, while $\lim_{s \rightarrow \lambda_1} \mu_j(s) = (\lambda_j - \lambda_1)^{-1} < \infty$ for $j = 2, \dots, n$. Hence, by choosing s sufficiently close to λ_1 the inverse power method will converge to that eigenvalue.

For the inverse power method (14.6) is replaced by

$$(i) \quad (\mathbf{A} - s\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1} \quad (14.7)$$

$$(ii) \quad \mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|.$$

Note that we solve the linear system rather than computing the inverse matrix. Normally the PLU factorization of $\mathbf{A} - s\mathbf{I}$ is precomputed in order to speed up the computation.

A variant of the inverse power method is known simply as **Rayleigh quotient iteration**. In this method we change the shift from iteration to iteration, using the previous Rayleigh quotient s_{k-1} as the current shift. In each iteration we need to compute the following quantities

- (i) $(\mathbf{A} - s_{k-1}\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1}$,
- (ii) $\mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|$,
- (iii) $s_k = \mathbf{x}_k^* \mathbf{A} \mathbf{x}_k$,
- (iv) $\mathbf{r}_k = \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k$.

We can avoid the calculation of $\mathbf{A} \mathbf{x}_k$ in (iii) and (iv). Let

$$\rho_k := \frac{\mathbf{y}_k^* \mathbf{x}_{k-1}}{\mathbf{y}_k^* \mathbf{y}_k}, \quad \mathbf{w}_k := \frac{\mathbf{x}_{k-1}}{\|\mathbf{y}_k\|_2}.$$

Then

$$\begin{aligned} s_k &= \frac{\mathbf{y}_k^* \mathbf{A} \mathbf{y}_k}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^* (\mathbf{A} - s_{k-1}\mathbf{I})\mathbf{y}_k}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^* \mathbf{x}_{k-1}}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \rho_k, \\ \mathbf{r}_k &= \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k = \frac{\mathbf{A} \mathbf{y}_k - (s_{k-1} + \rho_k) \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \frac{\mathbf{x}_{k-1} - \rho_k \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \mathbf{w}_k - \rho_k \mathbf{x}_k. \end{aligned}$$

Another problem is that the linear system in i) becomes closer and closer to singular as s_k converges to the eigenvalue. Thus the system becomes more and more ill-conditioned and we can expect large errors in the computed \mathbf{y}_k . This is indeed true, but we are lucky. Most of the error occurs in the direction of the eigenvector and this error disappears when we normalize \mathbf{y}_k in ii). Miraculously, the normalized eigenvector will be quite accurate.

Given an approximation (s, \mathbf{x}) to an eigenpair (λ, \mathbf{v}) of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$. The following algorithm computes a hopefully better approximation to (λ, \mathbf{v}) by doing one Rayleigh quotient iteration. The length nr of the new residual is also returned

Algorithm 14.6 (Rayleigh quotient iteration)

```

1 function [x,s,nr]=rayleight(A,x,s)
2 n=length(x);
3 y=(A-s*eye(n,n))\x;
4 yn=norm(y);
5 w=x/yn;
6 x=y/yn;
7 rho=x'*w;
8 s=s+rho;
9 nr=norm(w-rho*x);

```

k	1	2	3	4	5
$\ r\ _2$	1.0e+000	7.7e-002	1.6e-004	8.2e-010	2.0e-020
$ s - \lambda_1 $	3.7e-001	-1.2e-002	-2.9e-005	-1.4e-010	-2.2e-016

Table 14.8: Quadratic convergence of Rayleigh quotient iteration.

Since the shift changes from iteration to iteration the computation of y in `rayleighit` will require $O(n^3)$ arithmetic operations for a full matrix. For such a matrix it might pay to reduce it to an upper Hessenberg form or tridiagonal form before starting the iteration. However, if we have a good approximation to an eigenpair then only a few iterations are necessary to obtain close to machine accuracy.

If Rayleigh quotient iteration converges the convergence will be quadratic and sometimes even cubic. We illustrate this with an example.

Example 14.7 (Rayleigh quotient iteration)

The smallest eigenvalue of the matrix $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ is $\lambda_1 = (5 - \sqrt{33})/2 \approx -0.37$. Starting with $x = [1, 1]^T$ and $s = 0$ `rayleighit` converges to this eigenvalue and corresponding eigenvector. In Table 14.8 we show the rate of convergence by iterating `rayleighit` 5 times. The errors are approximately squared in each iteration indicating quadratic convergence.

14.2 The basic QR Algorithm

The QR algorithm is an iterative method to compute all eigenvalues and eigenvectors of a matrix $A \in \mathbb{C}^{n \times n}$. The matrix is reduced to triangular form by a sequence of unitary similarity transformations computed from the QR factorization of A . Recall that for a square matrix the QR factorization and the QR decomposition are the same. If $A = QR$ is a QR factorization then $Q \in \mathbb{C}^{n \times n}$ is unitary, $Q^*Q = I$ and $R \in \mathbb{C}^{n \times n}$ is upper triangular.

The basic QR algorithm takes the following form:

$$\boxed{\begin{aligned} A_1 &= A \\ \text{for } k &= 1, 2, \dots \\ Q_k R_k &= A_k \quad (\text{QR factorization of } A_k) \\ A_{k+1} &= R_k Q_k. \\ \text{end} \end{aligned}} \tag{14.8}$$

The determination of the QR factorization of A_k and the computation of $R_k Q_k$ is called a QR step. It is not at all clear that a QR step does anything

useful. At this point, since $\mathbf{R}_k = \mathbf{Q}_k^* \mathbf{A}_k \mathbf{Q}_k$ we find

$$\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k = \mathbf{Q}_k^* \mathbf{A}_k \mathbf{Q}_k, \quad (14.9)$$

so \mathbf{A}_{k+1} is unitary similar to \mathbf{A}_k . By induction \mathbf{A}_{k+1} is unitary similar to \mathbf{A} . Thus, each \mathbf{A}_k has the same eigenvalues as \mathbf{A} . We shall see that the basic QR algorithm is related to the power method.

Here are two examples to illustrate what happens.

Example 14.9 (QR iteration; real eigenvalues)

We start with

$$\mathbf{A}_1 = \mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \left(\frac{1}{\sqrt{5}} \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \right) * \left(\frac{1}{\sqrt{5}} \begin{bmatrix} 5 & 4 \\ 0 & 3 \end{bmatrix} \right) = \mathbf{Q}_1 \mathbf{R}_1$$

and obtain

$$\mathbf{A}_2 = \mathbf{R}_1 \mathbf{Q}_1 = \frac{1}{5} \begin{bmatrix} 5 & 4 \\ 0 & 3 \end{bmatrix} * \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 14 & 3 \\ 3 & 6 \end{bmatrix} = \begin{bmatrix} 2.8 & 0.6 \\ 0.6 & 1.2 \end{bmatrix}.$$

Continuing we find

$$\mathbf{A}_4 \approx \begin{bmatrix} 2.997 & -0.074 \\ -0.074 & 1.0027 \end{bmatrix}, \quad \mathbf{A}_{10} \approx \begin{bmatrix} 3.0000 & -0.0001 \\ -0.0001 & 1.0000 \end{bmatrix}$$

\mathbf{A}_{10} is almost diagonal and contains approximations to the eigenvalues $\lambda_1 = 3$ and $\lambda_2 = 1$ on the diagonal.

Example 14.10 (QR iteration; complex eigenvalues)

Applying the QR iteration (14.8) to the matrix

$$\mathbf{A}_1 = \mathbf{A} = \begin{bmatrix} 0.9501 & 0.8913 & 0.8214 & 0.9218 \\ 0.2311 & 0.7621 & 0.4447 & 0.7382 \\ 0.6068 & 0.4565 & 0.6154 & 0.1763 \\ 0.4860 & 0.0185 & 0.7919 & 0.4057 \end{bmatrix}$$

we obtain

$$\mathbf{A}_{14} = \left[\begin{array}{c|ccc} 2.323 & 0.047223 & -0.39232 & -0.65056 \\ \hline -2.1e-10 & 0.13029 & 0.36125 & 0.15946 \\ -4.1e-10 & -0.58622 & 0.052576 & -0.25774 \\ \hline 1.2e-14 & 3.3e-05 & -1.1e-05 & 0.22746 \end{array} \right].$$

This matrix is almost quasi-triangular and estimates for the eigenvalues $\lambda_1, \dots, \lambda_4$ of \mathbf{A} can now easily be determined from the diagonal blocks of \mathbf{A}_{14} . The 1×1 blocks

give us two real eigenvalues $\lambda_1 \approx 2.323$ and $\lambda_4 \approx 0.2275$. The middle 2×2 block has complex eigenvalues resulting in $\lambda_2 \approx 0.0914 + 0.4586i$ and $\lambda_3 \approx 0.0914 - 0.4586i$. From Gershgorin's circle theorem 13.4 and Corollary 13.6 it follows that the approximations to the real eigenvalues are quite accurate. We would also expect the complex eigenvalues to have small absolute errors.

These two examples illustrate what happens in general. The sequence $(\mathbf{A}_k)_k$ converges to the triangular Schur form (Cf. Theorem 6.29) if all the eigenvalues are real or the quasi-triangular Schur form (Cf. Definition 6.33) if some of the eigenvalues are complex.

14.2.1 The Relation to the Power Method

Let us show that the basic QR algorithm is related to the power method. We obtain the QR factorization of the powers \mathbf{A}^k as follows:

Theorem 14.11 (QR and power)

For $k = 1, 2, 3, \dots$, the QR factorization of \mathbf{A}^k is $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$, where

$$\tilde{\mathbf{Q}}_k := \mathbf{Q}_1 \cdots \mathbf{Q}_k \text{ and } \tilde{\mathbf{R}}_k := \mathbf{R}_k \cdots \mathbf{R}_1, \quad (14.10)$$

and $\mathbf{Q}_1, \dots, \mathbf{Q}_k, \mathbf{R}_1, \dots, \mathbf{R}_k$ are the matrices generated by the basic QR algorithm (14.8).

Proof. By (14.9)

$$\mathbf{A}_k = \mathbf{Q}_{k-1}^* \mathbf{A}_{k-1} \mathbf{Q}_{k-1} = \mathbf{Q}_{k-1}^* \mathbf{Q}_{k-2}^* \mathbf{A}_{k-2} \mathbf{Q}_{k-2} \mathbf{Q}_{k-1} = \cdots = \tilde{\mathbf{Q}}_{k-1}^* \mathbf{A} \tilde{\mathbf{Q}}_{k-1}. \quad (14.11)$$

The proof is by induction on k . Clearly $\tilde{\mathbf{Q}}_1 \tilde{\mathbf{R}}_1 = \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{A}_1$. Suppose $\tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{R}}_{k-1} = \mathbf{A}^{k-1}$ for some $k \geq 2$. Since $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k$ and using (14.11)

$$\tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k = \tilde{\mathbf{Q}}_{k-1} (\mathbf{Q}_k \mathbf{R}_k) \tilde{\mathbf{R}}_{k-1} = \tilde{\mathbf{Q}}_{k-1} \mathbf{A}_k \tilde{\mathbf{R}}_{k-1} = (\tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{Q}}_{k-1}^*) \mathbf{A} \tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{R}}_{k-1} = \mathbf{A}^k.$$

□

Since $\tilde{\mathbf{R}}_k$ is upper triangular, its first column is a multiple of \mathbf{e}_1 so that

$$\mathbf{A}^k \mathbf{e}_1 = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k \mathbf{e}_1 = \tilde{r}_{11}^{(k)} \tilde{\mathbf{Q}}_k \mathbf{e}_1 \text{ or } \tilde{\mathbf{q}}_1^{(k)} := \tilde{\mathbf{Q}}_k \mathbf{e}_1 = \frac{1}{\tilde{r}_{11}^{(k)}} \mathbf{A}^k \mathbf{e}_1.$$

Since $\|\tilde{\mathbf{q}}_1^{(k)}\|_2 = 1$ the first column of $\tilde{\mathbf{Q}}_k$ is the result of applying the normalized power iteration (14.6) to the starting vector $\mathbf{x}_0 = \mathbf{e}_1$. If this iteration converges we conclude that the first column of $\tilde{\mathbf{Q}}_k$ must converge to a dominant eigenvector of \mathbf{A} . It can be shown that the first column of \mathbf{A}_k must then converge to $\lambda_1 \mathbf{e}_1$, where λ_1 is a dominant eigenvalue of \mathbf{A} . This is clearly what happens in Examples 14.9 and 14.10. Indeed, what is observed in practice is that the sequence $(\tilde{\mathbf{Q}}_k^* \mathbf{A} \tilde{\mathbf{Q}}_k)_k$ converges to a (quasi-triangular) Schur form of \mathbf{A} .

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{12}^*} \begin{bmatrix} x & x & x & x \\ \mathbf{x} & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{23}^*} \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & \mathbf{x} & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{34}^*} \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & \mathbf{x} & x \end{bmatrix}.$$

Figure 14.1: Post multiplication in a QR step.

14.2.2 Invariance of the Hessenberg Form

One QR step requires $O(n^3)$ arithmetic operations for a matrix \mathbf{A} of order n . By an initial reduction of \mathbf{A} to upper Hessenberg form \mathbf{H}_1 using Algorithm 13.11, the cost of a QR step can be reduced to $O(n^2)$. Consider a QR step on \mathbf{H}_1 . We first determine plane rotations $\mathbf{P}_{i,i+1}$, $i = 1, \dots, n-1$ so that $\mathbf{P}_{n-1,n} \cdots \mathbf{P}_{1,2} \mathbf{H}_1 = \mathbf{R}_1$ is upper triangular. The details were described in Section 11.4. Thus $\mathbf{H}_1 = \mathbf{Q}_1 \mathbf{R}_1$, where $\mathbf{Q}_1 = \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$ is a QR factorization of \mathbf{H}_1 . To finish the QR step we compute $\mathbf{R}_1 \mathbf{Q}_1 = \mathbf{R}_1 \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$. This postmultiplication step is illustrated by the Wilkinson diagram in Figure 14.1.

The postmultiplication by $\mathbf{P}_{i,i+1}$ introduces a nonzero in position $(i+1, i)$ leaving the other elements marked by a zero in Figure 14.1 unchanged. Thus the final matrix $\mathbf{R} \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$ is upper Hessenberg and a QR step leaves the Hessenberg form invariant.

In conclusion, to compute \mathbf{A}_{k+1} from \mathbf{A}_k requires $O(n^2)$ arithmetic operations if \mathbf{A}_k is upper Hessenberg and $O(n)$ arithmetic operations if \mathbf{A}_k is tridiagonal.

14.2.3 Deflation

If a subdiagonal element $a_{i+1,i}$ of an upper Hessenberg matrix \mathbf{A} is equal to zero, then the eigenvalues of \mathbf{A} are the union of the eigenvalues of the two smaller matrices $A(1:i, 1:i)$ and $A(i+1:n, i+1:n)$. Thus if during the iteration the $(i+1, i)$ element of \mathbf{A}_k is sufficiently small then we can continue the iteration on the two smaller submatrices separately.

To see what effect this can have on the eigenvalues of \mathbf{A} suppose $|a_{i+1,i}^{(k)}| \leq \epsilon$. Let $\hat{\mathbf{A}}_k := \mathbf{A}_k - a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T$ be the matrix obtained from \mathbf{A}_k by setting the $(i+1, i)$ element equal to zero. Since $\mathbf{A}_k = \tilde{\mathbf{Q}}_{k-1}^* \mathbf{A} \tilde{\mathbf{Q}}_{k-1}$ we have

$$\hat{\mathbf{A}}_k = \tilde{\mathbf{Q}}_{k-1}^* (\mathbf{A} + \mathbf{E}) \tilde{\mathbf{Q}}_{k-1}, \quad \mathbf{E} = \tilde{\mathbf{Q}}_{k-1} (a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T) \tilde{\mathbf{Q}}_{k-1}^*.$$

Since $\tilde{\mathbf{Q}}_{k-1}$ is unitary, $\|\mathbf{E}\|_F = \|a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T\|_F = |a_{i+1,i}^{(k)}| \leq \epsilon$ and setting $a_{i+1,i}^{(k)} = 0$ amounts to a perturbation in the original \mathbf{A} of at most ϵ . For how to chose ϵ see the discussion on page 94–95 in [23].

This deflation occurs often in practice and can with a proper implementation reduce the computation time considerably. It should be noted that to find the

eigenvectors of the original matrix one has to continue with some care, see [23].

14.3 The Shifted QR Algorithms

Like in the inverse power method it is possible to speed up the convergence by introducing shifts. The **explicitly shifted QR algorithm** works as follows:

```

 $\mathbf{A}_1 = \mathbf{A}$ 
for  $k = 1, 2, \dots$ 
    Choose a shift  $s_k$ 
     $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k - s_k \mathbf{I}$       (QR factorization of  $\mathbf{A}_k - s_k \mathbf{I}$ )
     $\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I}$ .
end

```

Since $\mathbf{R}_k = \mathbf{Q}_k^* (\mathbf{A}_k - s_k \mathbf{I})$ we find

$$\mathbf{A}_{k+1} = \mathbf{Q}_k^* (\mathbf{A}_k - s_k \mathbf{I}) \mathbf{Q}_k + s_k \mathbf{I} = \mathbf{Q}_k^* \mathbf{A}_k \mathbf{Q}_k$$

and \mathbf{A}_{k+1} and \mathbf{A}_k are unitary similar.

The shifted QR algorithm is related to the power method with shift, cf. Theorem 14.11 and also the inverse power method. In fact the last column of \mathbf{Q}_k is the result of one iteration of the inverse power method to \mathbf{A}^* with shift s_k . Indeed, since $\mathbf{A} - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k$ we have $(\mathbf{A} - s_k \mathbf{I})^* = \mathbf{R}_k^* \mathbf{Q}_k^*$ and $(\mathbf{A} - s_k \mathbf{I})^* \mathbf{Q}_k = \mathbf{R}_k^*$. Thus, since \mathbf{R}_k^* is lower triangular with n, n element $\bar{r}_{nn}^{(k)}$ we find $(\mathbf{A} - s_k \mathbf{I})^* \mathbf{Q}_k \mathbf{e}_n = \mathbf{R}_k^* \mathbf{e}_n = \bar{r}_{nn}^{(k)} \mathbf{e}_n$ from which the conclusion follows.

The shift $s_k := \mathbf{e}_n^T \mathbf{A}_k \mathbf{e}_n$ is called the **Rayleigh quotient shift**, while the eigenvalue of the lower right 2×2 corner of \mathbf{A}_k closest to the n, n element of \mathbf{A}_k is called the **Wilkinson shift**. This shift can be used to find complex eigenvalues of a real matrix. The convergence is very fast and at least quadratic both for the Rayleigh quotient shift and the Wilkinson shift.

By doing two QR iterations at a time it is possible to find both real and complex eigenvalues without using complex arithmetic. The corresponding algorithm is called the **implicitly shifted QR algorithm**

After having computed the eigenvalues we can compute the eigenvectors in steps. First we find the eigenvectors of the triangular or quasi-triangular matrix. We then compute the eigenvectors of the upper Hessenberg matrix and finally we get the eigenvectors of \mathbf{A} .

Practical experience indicates that only $O(n)$ iterations are needed to find all eigenvalues of \mathbf{A} . Thus both the explicit- and implicit shift QR algorithms are normally $O(n^3)$ algorithms.

For further remarks and detailed algorithms see [23].

14.4 A Convergence Theorem

There is no theorem which proves convergence of the QR algorithm in general. The following theorem shows convergence of the basic QR algorithm under somewhat restrictive assumptions.

Theorem 14.12 (Convergence of basis QR)

Suppose in the basic QR algorithm (14.8) that

1. $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be diagonalized, $\mathbf{X}^{-1}\mathbf{AX} = \mathbf{\Lambda} := \text{diag}(\lambda_1, \dots, \lambda_n)$.
2. The eigenvalues $\lambda_1, \dots, \lambda_n$ are real with $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n| > 0$.
3. The inverse of the eigenvector matrix has an LU factorization $\mathbf{X}^{-1} = \mathbf{LR}$.

Let $\tilde{\mathbf{Q}}_k = \mathbf{Q}_1 \dots \mathbf{Q}_k$ for $k \geq 1$. Then there is a diagonal matrix \mathbf{D}_k with diagonal elements ± 1 such that $\tilde{\mathbf{Q}}_k \mathbf{D}_k \rightarrow \mathbf{Q}$, where $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$ is triangular and \mathbf{Q} is the Q-factor in the QR factorization of the eigenvector matrix \mathbf{X} .

Proof. In this proof we assume that every QR factorization has an \mathbf{R} with positive diagonal elements so that the factorization is unique. Let $\mathbf{X} = \mathbf{QR}$ be the QR factorization of \mathbf{X} . We observe that $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$ is upper triangular. For since $\mathbf{X}^{-1} \mathbf{AX} = \mathbf{\Lambda}$ we have $\mathbf{R}^{-1} \mathbf{Q}^T \mathbf{A} \mathbf{QR} = \mathbf{\Lambda}$ so that $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{R} \mathbf{\Lambda} \mathbf{R}^{-1}$ is upper triangular. Since $\mathbf{A}_{k+1} = \tilde{\mathbf{Q}}_k^T \mathbf{A} \tilde{\mathbf{Q}}_k$, it is enough to show that $\tilde{\mathbf{Q}}_k \mathbf{D}_k \rightarrow \mathbf{Q}$ for some diagonal matrix \mathbf{D}_k with diagonal elements ± 1 .

We define the nonsingular matrices

$$\mathbf{F}_k := \mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \mathbf{R}^{-1} = \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k, \quad \mathbf{G}_k := \hat{\mathbf{R}}_k \mathbf{R} \mathbf{\Lambda}^k \mathbf{R}, \quad \mathbf{D}_k := \text{diag}\left(\frac{\delta_1}{|\delta_1|}, \dots, \frac{\delta_n}{|\delta_n|}\right),$$

where $\delta_1, \dots, \delta_n$ are the diagonal elements in the upper triangular matrix \mathbf{G}_k and $\mathbf{F}_k = \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k$ is the QR factorization of \mathbf{F}_k . Then

$$\begin{aligned} \mathbf{A}^k &= \mathbf{X} \mathbf{\Lambda}^k \mathbf{X}^{-1} = \mathbf{Q} \mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{R} = \mathbf{Q} (\mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \mathbf{R}^{-1}) (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) \\ &= \mathbf{Q} \mathbf{F}_k (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) = \mathbf{Q} \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) = (\mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}) (\mathbf{D}_k \mathbf{G}_k), \end{aligned}$$

and this is the QR factorization of \mathbf{A}^k . Indeed, $\mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}$ is a product of orthonormal matrices and therefore orthonormal. Moreover $\mathbf{D}_k \mathbf{G}_k$ is a product of upper triangular matrices and therefore upper triangular. Note that \mathbf{D}_k is chosen so that this matrix has positive diagonal elements. By Theorem 14.11 $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$ is also the QR factorization of \mathbf{A}^k , and we must have $\tilde{\mathbf{Q}}_k = \mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}$ or $\tilde{\mathbf{Q}}_k \mathbf{D}_k = \mathbf{Q} \hat{\mathbf{Q}}_k$. The theorem will follow if we can show that $\hat{\mathbf{Q}}_k \rightarrow \mathbf{I}$.

The matrix $\Lambda^k L \Lambda^{-k}$ is lower triangular with elements $(\frac{\lambda_i}{\lambda_j})^k l_{ij}$ on and under the diagonal. Thus for $n = 3$

$$\Lambda^k L \Lambda^{-k} = \begin{bmatrix} 1 & 0 & 0 \\ (\frac{\lambda_2}{\lambda_1})^k l_{21} & 1 & 0 \\ (\frac{\lambda_3}{\lambda_1})^k l_{31} & (\frac{\lambda_3}{\lambda_2})^k l_{32} & 1 \end{bmatrix}.$$

By Assumption 2. it follows that $\Lambda^k L \Lambda^{-k} \rightarrow I$, and hence $F_k \rightarrow I$. Since $\hat{R}_k^T \hat{R}_k$ is the Cholesky factorization of $F_k^T F_k$ it follows that $\hat{R}_k^T \hat{R}_k \rightarrow I$. By the continuity of the Cholesky factorization it holds $\hat{R}_k \rightarrow I$ and hence $\hat{R}_k^{-1} \rightarrow I$. But then $\hat{Q}_k = F_k \hat{R}_k^{-1} \rightarrow I$. \square

Exercise 14.13 (QR convergence detail)

Use Theorem 8.33 to show that $\hat{R}_k \rightarrow I$ implies $\hat{R}_k^{-1} \rightarrow I$.

14.5 Review Questions

14.5.1 What is the main use of the power method?

14.5.2 Can the QR method be used to find all eigenvectors of a matrix?

14.5.3 Can the power method be used to find an eigenvalue?

14.5.4 Do the power method converge to an eigenvector corresponding to a complex eigenvalue?

14.5.5 What is the inverse power method?

14.5.6 Give a relation between the QR algorithm and the power method.

14.5.7 How can we make the basic QR algorithm converge faster?

Part VI

Appendix

Appendix A

Determinants

The first systematic treatment of determinants was given by Cauchy in 1812. He adopted the word “determinant” which was introduced by Gauss in 1801. The first use of determinants was made by Leibniz in 1693 in a letter to De L’Hôpital. By the beginning of the 20th century the theory of determinants filled four volumes of almost 2000 pages (Muir, 1906–1923. Historic references can be found in this work). The main use of determinants in this text will be to study the characteristic polynomial of a matrix.

In this section we give the elementary properties of determinants that we need.

A.1 Permutations

For $n \in \mathbb{N}$, let $N_n = \{1, 2, \dots, n\}$. A *permutation* is a function $\sigma : N_n \rightarrow N_n$ which is one-to-one and onto. That is, $\{\sigma(1), \sigma(2), \dots, \sigma(n)\}$ is a rearrangement of $\{1, 2, \dots, n\}$. If $n = 2$, there are two permutations $\{1, 2\}$ and $\{2, 1\}$, while for $n = 3$ we have six permutations $\{1, 2, 3\}$, $\{1, 3, 2\}$, $\{2, 1, 3\}$, $\{2, 3, 1\}$, $\{3, 1, 2\}$ and $\{3, 2, 1\}$. We denote the set of all permutations on N_n by S_n . There are $n!$ elements in S_n .

If σ, τ are two permutations in S_n , we can define their product $\sigma\tau$ as

$$\sigma\tau = \{\sigma(\tau(1)), \sigma(\tau(2)), \dots, \sigma(\tau(n))\}.$$

For example if $\sigma = \{1, 3, 2\}$ and $\tau = \{3, 2, 1\}$, then $\sigma\tau = \{\sigma(3), \sigma(2), \sigma(1)\} = \{2, 3, 1\}$, while $\tau\sigma = \{\tau(1), \tau(3), \tau(2)\} = \{3, 1, 2\}$. Thus in general $\sigma\tau \neq \tau\sigma$. It is easily shown that the product of two permutations σ, τ is a permutation, i.e. $\sigma\tau : N_n \rightarrow N_n$ is one-to-one and onto.

The permutation $\epsilon = \{1, 2, \dots, n\}$ is called the *identity permutation* in S_n . We have $\epsilon\sigma = \sigma\epsilon = \sigma$ for all $\sigma \in S_n$.

Since each $\sigma \in S_n$ is one-to-one and onto, it has a unique inverse σ^{-1} . To define $\sigma^{-1}(j)$ for $j \in N_n$, we find the unique i such that $\sigma(i) = j$. Then $\sigma^{-1}(j) = i$. We have $\sigma^{-1}\sigma = \sigma\sigma^{-1} = \epsilon$. As an example, if $\sigma = \{2, 3, 1\}$ then $\sigma^{-1} = \{3, 1, 2\}$, and $\sigma^{-1}\sigma = \sigma\sigma^{-1} = \{1, 2, 3\} = \epsilon$.

With each $\sigma \in S_n$ we can associate a + or - sign. We define

$$\text{sign}(\sigma) = \frac{g(\sigma)}{|g(\sigma)|},$$

where

$$g(\sigma) = \prod_{i=2}^n (\sigma(i) - \sigma(1))(\sigma(i) - \sigma(2)) \cdots (\sigma(i) - \sigma(i-1)).$$

For example if $\epsilon = \{1, 2, 3, 4\}$ and $\sigma = \{4, 3, 1, 2\}$, then

$$\begin{aligned} g(\epsilon) &= (2-1)(3-1)(3-2)(4-1)(4-2)(4-3) = 1! \cdot 2! \cdot 3! > 0, \\ g(\sigma) &= (3-4)(1-4)(1-3)(2-4)(2-3)(2-1) \\ &= (-1)(-3)(-2)(-2)(-1) \cdot 1 = -1! \cdot 2! \cdot 3! < 0. \end{aligned}$$

Thus $\text{sign}(\epsilon) = +1$ and $\text{sign}(\sigma) = -1$.

$g(\sigma)$ contains one positive factor $(2-1)$ and five negative ones. The negative factors are called *inversions*. The number of inversions equals the number of times a bigger integer precedes a smaller one in σ . That is, in $\{4, 3, 1, 2\}$ 4 precedes 3, 1 and 2 (three inversions corresponding to the negative factors $(3-4)$, $(1-4)$ and $(2-4)$ in $g(\sigma)$), and 3 precedes 1 and 2 ($(1-3)$ and $(2-3)$ in $g(\sigma)$). This makes it possible to compute $\text{sign}(\sigma)$ without actually writing down $g(\sigma)$.

In general, the sign function has the following properties

1. $\text{sign}(\epsilon) = 1$.
2. $\text{sign}(\sigma\tau) = \text{sign}(\sigma)\text{sign}(\tau)$ for $\sigma, \tau \in S_n$.
3. $\text{sign}(\sigma^{-1}) = \text{sign}(\sigma)$ for $\sigma \in S_n$.

Since all factors in $g(\epsilon)$ are positive, we have $g(\epsilon) = |g(\epsilon)|$ and $\text{sign}(\epsilon) = 1$. This proves 1. To prove 2 we first note that for any S_n

$$\text{sign}(\sigma) = \frac{g(\sigma)}{g(\epsilon)}.$$

Since $g(\sigma)$ and $g(\epsilon)$ contain the same factors apart from signs and $g(\epsilon) > 0$, we have $|g(\sigma)| = g(\epsilon)$. Now

$$\text{sign}(\sigma\tau) = \frac{g(\sigma\tau)}{g(\epsilon)} = \frac{g(\sigma\tau)}{g(\tau)} \frac{g(\tau)}{g(\epsilon)} = \frac{g(\sigma\tau)}{g(\tau)} \text{sign}(\tau).$$

We have to show that $g(\sigma\tau)/g(\tau) = g(\sigma)/g(\epsilon)$. We write $g(\sigma)/g(\epsilon)$ in the form

$$\frac{g(\sigma)}{g(\epsilon)} = \prod_{i=2}^n \prod_{j=1}^{i-1} r_\sigma(i, j), \quad r_\sigma(i, j) = \frac{\sigma(i) - \sigma(j)}{i - j}.$$

Now

$$\frac{g(\sigma\tau)}{g(\tau)} = \frac{\prod_{i=2}^n (\sigma(\tau(i)) - \sigma(\tau(1))) \cdots (\sigma(\tau(i)) - \sigma(\tau(i-1)))}{\prod_{i=2}^n (\tau(i) - \tau(1)) \cdots (\tau(i) - \tau(i-1))} = \prod_{i=2}^n \prod_{j=1}^{i-1} r_\sigma(\tau(i), \tau(j)).$$

τ is a permutation so $g(\sigma)/g(\epsilon)$ and $g(\sigma\tau)/g(\tau)$ contain the same factors. Moreover, the sign of the factors are the same since $r(i, j) = r(j, i)$ for all $i \neq j$. Thus $g(\sigma)/g(\epsilon) = g(\sigma\tau)/g(\tau)$, and 2 is proved. Finally, 3 follows from 1 and 2; $1 = \text{sign}(\epsilon) = \text{sign}(\sigma\sigma^{-1}) = \text{sign}(\sigma)\text{sign}(\sigma^{-1})$ so that σ and σ^{-1} have the same sign.

Example A.1 (Properties of permutations)

It can be shown that $\rho(\sigma\tau) = (\rho\sigma)\tau$ for $\rho, \sigma, \tau \in S_n$, i.e. multiplication of permutations is associative. (In fact, we have

1. Multiplication is associative.
2. There exists an identity permutation ϵ .
3. Every permutation has an inverse.

Thus the set S_n of permutations is a group with respect to multiplication. S_n is called the symmetric group of degree n .

A.2 Basic Properties of Determinants

For any $A \in \mathbb{C}^{n \times n}$ the determinant of A is defined the number

$$\det(A) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma(1),1} a_{\sigma(2),2} \cdots a_{\sigma(n),n}. \quad (\text{A.1})$$

This sum ranges of all $n!$ permutations of $\{1, 2, \dots, n\}$. We also denote the determinant by (Cayley, 1841)

$$\left| \begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right|.$$

From the definition we have

$$\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{21}a_{12}.$$

The first term on the right corresponds to the identity permutation ϵ given by $\epsilon(i) = i$, $i = 1, 2$. The second term comes from the permutation $\sigma = \{2, 1\}$. For $n = 3$

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11}a_{22}a_{33} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} - a_{31}a_{22}a_{13}.$$

The following is a list of properties of determinants.

1. **Triangular matrix** The determinant of a triangular matrix is the product of the diagonal elements. $\det(A) = a_{11}a_{22} \cdots a_{nn}$. In particular $\det(I) = 1$.
2. **Transpose** $\det(A^T) = \det(A)$.
3. **Homogeneity** For any $\beta_i \in \mathbb{C}$, $i = 1, 2, \dots, n$, we have

$$\det([\beta_1 \mathbf{a}_1, \beta_2 \mathbf{a}_2, \dots, \beta_n \mathbf{a}_n]) = \beta_1 \beta_2 \cdots \beta_n \det([\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]).$$

4. **Permutation of columns** If $\tau \in S_n$ then

$$\det(\mathbf{B}) := \det([\mathbf{a}_{\tau(1)}, \mathbf{a}_{\tau(2)}, \dots, \mathbf{a}_{\tau(n)}]) = \text{sign}(\tau) \det([\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]).$$

5. **Additivity**

$$\begin{aligned} \det([\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \mathbf{a}_k + \mathbf{a}'_k, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n]) \\ = \det([\mathbf{a}_1, \dots, \mathbf{a}_n]) + \det([\mathbf{a}_1, \dots, \mathbf{a}'_k, \dots, \mathbf{a}_n]). \end{aligned}$$

6. **Singular matrix** $\det(A) = 0$ if and only if A is singular.
7. **Product rule** If $A, B \in \mathbb{C}^{n \times n}$ then $\det(AB) = \det(A)\det(B)$.
8. **Block triangular** If A is block triangular with diagonal blocks B and C then $\det(A) = \det(B)\det(C)$.

Proof.

1. If $\sigma \neq \epsilon$, we can find distinct integers i and j such that $\sigma(i) > i$ and $\sigma(j) < j$. But then $a_{\sigma(i),i} = 0$ if \mathbf{A} is upper triangular and $a_{\sigma(j),j} = 0$ if \mathbf{A} is lower triangular. Hence

$$\det(\mathbf{A}) = \text{sign}(\epsilon) a_{\epsilon(1),1} a_{\epsilon(2),2} \cdots a_{\epsilon(n),n} = a_{1,1} a_{2,2} \cdots a_{n,n}.$$

Since the identity matrix is triangular with all diagonal elements equal to one, we have that $\det(\mathbf{I}) = 1$.

2. By definition of \mathbf{A}^T and the det-function

$$\det(\mathbf{A}^T) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{1,\sigma(1)} a_{2,\sigma(2)} \cdots a_{n,\sigma(n)}.$$

Consider an element $a_{i,\sigma(i)}$. If $\sigma(i) = j$ then

$$a_{i,\sigma(i)} = a_{\sigma^{-1}(j),j}.$$

Since $\sigma(1), \sigma(2), \dots, \sigma(n)$ ranges through $\{1, 2, \dots, n\}$, we obtain

$$\begin{aligned} \det(\mathbf{A}^T) &= \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \sum_{\sigma \in S_n} \text{sign}(\sigma^{-1}) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \sum_{\sigma^{-1} \in S_n} \text{sign}(\sigma^{-1}) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \det(\mathbf{A}). \end{aligned}$$

3. This follows immediately from the definition of $\det[(\beta_1 \mathbf{a}_1, \beta_2 \mathbf{a}_2, \dots, \beta_n \mathbf{a}_n)]$.

4. We have

$$\det(\mathbf{B}) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma(1),\tau(1)} a_{\sigma(2),\tau(2)} \cdots a_{\sigma(n),\tau(n)}.$$

Fix i in $\{1, 2, \dots, n\}$. Let $k = \sigma(i)$ and $m = \tau(i)$. Then $\tau^{-1}(m) = i$ and $\sigma(\tau^{-1}(m)) = k$. Hence

$$a_{\sigma(i),\tau(i)} = a_{k,m} = a_{\sigma\tau^{-1}(m),m}.$$

Moreover, $\text{sign}(\sigma) = \text{sign}(\tau)\text{sign}(\sigma\tau^{-1})$. Thus

$$\det(\mathbf{B}) = \text{sign}(\tau) \sum_{\sigma \in S_n} \text{sign}(\sigma\tau^{-1}) a_{\sigma\tau^{-1}(1),1} a_{\sigma\tau^{-1}(2),2} \cdots a_{\sigma\tau^{-1}(n),n}.$$

But as σ ranges over S_n , $\sigma\tau^{-1}$ also ranges over S_n . Hence

$$\det(\mathbf{B}) = \text{sign}(\tau) \det[(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)].$$

5. This follows at once from the definition.

6. We observe that the determinant of a matrix is equal to the product of the eigenvalues and that a matrix is singular if and only if zero is an eigenvalue (cf. Theorems 0.68, 0.69). But then the result follows.

7. To better understand the general proof, we do the 2×2 case first. Let $\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2)$, $\mathbf{B} = (\mathbf{b}_1, \mathbf{b}_2)$. Then

$$\mathbf{AB} = (\mathbf{Ab}_1, \mathbf{Ab}_2) = (b_{1,1}\mathbf{a}_1 + b_{2,1}\mathbf{a}_2, b_{1,2}\mathbf{a}_1 + b_{2,2}\mathbf{a}_2).$$

Using the additivity, we obtain

$$\begin{aligned}\det(\mathbf{AB}) &= \det(b_{1,1}\mathbf{a}_1, b_{1,2}\mathbf{a}_1) + \det(b_{2,1}\mathbf{a}_2, b_{1,2}\mathbf{a}_1) \\ &\quad + \det(b_{1,1}\mathbf{a}_1, b_{2,2}\mathbf{a}_2) + \det(b_{2,1}\mathbf{a}_2, b_{2,2}\mathbf{a}_2).\end{aligned}$$

Next we have by homogeneity

$$\begin{aligned}\det(\mathbf{AB}) &= b_{1,1}b_{1,2}\det(\mathbf{a}_1, \mathbf{a}_1) + b_{2,1}b_{1,2}\det(\mathbf{a}_2, \mathbf{a}_1) \\ &\quad + b_{1,1}b_{2,2}\det(\mathbf{a}_1, \mathbf{a}_2) + b_{2,1}b_{2,2}\det(\mathbf{a}_2, \mathbf{a}_2).\end{aligned}$$

Property 6 implies that $\det(\mathbf{a}_1, \mathbf{a}_1) = \det(\mathbf{a}_2, \mathbf{a}_2) = 0$. Using Property 4, we obtain $\det(\mathbf{a}_2, \mathbf{a}_1) = -\det(\mathbf{a}_1, \mathbf{a}_2)$ and

$$\det(\mathbf{AB}) = (b_{1,1}b_{2,2} - b_{2,1}b_{1,2})\det(\mathbf{a}_1, \mathbf{a}_2) = \det(\mathbf{B})\det(\mathbf{A}).$$

The proof for $n > 2$ follows the $n = 2$ case step by step. Let $\mathbf{C} = (\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n) = \mathbf{AB}$. Then

$$\mathbf{c}_i = \mathbf{Ab}_i = b_{1,i}\mathbf{a}_1 + b_{2,i}\mathbf{a}_2 + \cdots + b_{n,i}\mathbf{a}_n, \quad i = 1, 2, \dots, n.$$

Using the additivity, we obtain

$$\det(\mathbf{AB}) = \sum_{i_1=1}^n \sum_{i_2=1}^n \cdots \sum_{i_n=1}^n \det[(b_{i_1,1}\mathbf{a}_{i_1}, b_{i_2,2}\mathbf{a}_{i_2}, \dots, b_{i_n,n}\mathbf{a}_{i_n})].$$

Next we have by homogeneity

$$\det(\mathbf{AB}) = \sum_{i_1=1}^n \sum_{i_2=1}^n \cdots \sum_{i_n=1}^n b_{i_1,1}b_{i_2,2} \cdots b_{i_n,n} \det[(\mathbf{a}_{i_1}, \mathbf{a}_{i_2}, \dots, \mathbf{a}_{i_n})].$$

Property 6 implies that $\det[(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_n})] = 0$ if any two of the indices i_1, \dots, i_n are equal. Therefore we only get a contribution to the sum whenever i_1, \dots, i_n is a permutation of $\{1, 2, \dots, n\}$. Thus

$$\det(\mathbf{AB}) = \sum_{\sigma \in S_n} b_{\sigma(1),1} \cdots b_{\sigma(n),n} \det[(\mathbf{a}_{\sigma(1)}, \dots, \mathbf{a}_{\sigma(n)})].$$

By Property 4 we obtain

$$\det(\mathbf{AB}) = \sum_{\sigma \in S_n} \text{sign}(\tau) b_{\sigma(1),1} \cdots b_{\sigma(n),n} \det[(\mathbf{a}_1, \dots, \mathbf{a}_n)].$$

According to the definition of $\det(\mathbf{B})$ this is equal to $\det(\mathbf{B})\det(\mathbf{A})$.

8. Suppose \mathbf{A} is block upper triangular. Let

$$S_{n,k} = \{\sigma \in S_n : \sigma(i) \leq k \text{ if } i \leq k, \text{ and } \sigma(i) \geq k+1 \text{ if } i \geq k+1\}.$$

We claim that $a_{\sigma(1),1} \cdots a_{\sigma(n),n} = 0$ if $\sigma \notin S_{n,k}$, because if $\sigma(i) > k$ for some $i \leq k$ then $a_{\sigma(i),i} = 0$ since it lies in the zero part of \mathbf{A} . If $\sigma(i) \leq k$ for some $i \geq k+1$, we must have $\sigma(j) > k$ for some $j \leq k$ to make “room” for $\sigma(i)$, and $a_{\sigma(j),j} = 0$. It follows that

$$\det(\mathbf{A}) = \sum_{\sigma \in S_{n,k}} \text{sign}(\sigma) a_{\sigma(1),1} \cdots a_{\sigma(n),n}.$$

Define

$$\rho(i) = \begin{cases} \sigma(i) & i = 1, \dots, k \\ i & i = k+1, \dots, n, \end{cases} \quad \tau(i) = \begin{cases} i & i = 1, \dots, k \\ \sigma(i) & i = k+1, \dots, n. \end{cases}$$

If $\sigma \in S_{n,k}$, ρ and τ will be permutations. Moreover, $\sigma = \rho\tau$. Define $\hat{\rho}$ and $\hat{\tau}$ in S_k and S_{n-k} respectively by $\hat{\rho}(i) = \rho(i)$, $i = 1, \dots, k$, and $\hat{\tau}(i) = \tau(i+k)-k$ for $i = 1, \dots, n-k$. As σ ranges over $S_{n,k}$, $\hat{\rho}$ and $\hat{\tau}$ will take on all values in S_k and S_{n-k} respectively. Since $\text{sign}(\hat{\rho}) = \text{sign}(\rho)$ and $\text{sign}(\hat{\tau}) = \text{sign}(\tau)$, we find

$$\text{sign}(\sigma) = \text{sign}(\rho)\text{sign}(\tau) = \text{sign}(\hat{\rho})\text{sign}(\hat{\tau}).$$

Then

$$\begin{aligned} \det(\mathbf{A}) &= \sum_{\hat{\rho} \in S_k} \sum_{\hat{\tau} \in S_{n-k}} \text{sign}(\hat{\rho})\text{sign}(\hat{\tau}) b_{\hat{\rho}(1),1} \cdots b_{\hat{\rho}(k),k} d_{\hat{\tau}(1),1} \cdots d_{\hat{\tau}(n-k),n-k} \\ &= \det(\mathbf{B}) \det(\mathbf{D}). \end{aligned}$$

□

A.3 The Adjoint Matrix and Cofactor Expansion

We start with a useful formula for the solution of a linear system.

Let $\mathbf{A}_j(\mathbf{b})$ denote the matrix obtained from \mathbf{A} by replacing the j th column of \mathbf{A} by \mathbf{b} . For example,

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \quad \mathbf{A}_1(\mathbf{b}) = \begin{bmatrix} 3 & 2 \\ 6 & 1 \end{bmatrix}, \quad \mathbf{A}_2(\mathbf{b}) = \begin{bmatrix} 1 & 3 \\ 2 & 6 \end{bmatrix}, \\ \mathbf{I} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{I}_1(\mathbf{x}) = \begin{bmatrix} x_1 & 0 \\ x_2 & 1 \end{bmatrix}, \quad \mathbf{I}_2(\mathbf{x}) = \begin{bmatrix} 1 & x_1 \\ 0 & x_2 \end{bmatrix}. \end{aligned}$$

Theorem A.2 (Cramer's rule (1750))

Suppose $\mathbf{A} \in \mathbb{C}^{n \times n}$ with $\det(\mathbf{A}) \neq 0$ and $\mathbf{b} \in \mathbb{C}^n$. Let $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ be the unique solution of $\mathbf{Ax} = \mathbf{b}$. Then

$$x_j = \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n.$$

Proof. Since $1 = \det(\mathbf{I}) = \det(\mathbf{AA}^{-1}) = \det(\mathbf{A})\det(\mathbf{A}^{-1})$ we have $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$. Then

$$\begin{aligned} \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})} &= \det(\mathbf{A}^{-1}\mathbf{A}_j(\mathbf{b})) \\ &= \det([\mathbf{A}^{-1}\mathbf{a}_1, \dots, \mathbf{A}^{-1}\mathbf{a}_{j-1}, \mathbf{A}^{-1}\mathbf{b}, \mathbf{A}^{-1}\mathbf{a}_{j+1}, \dots, \mathbf{A}^{-1}\mathbf{a}_n]) \\ &= \det([\mathbf{e}_1, \dots, \mathbf{e}_{j-1}, \mathbf{x}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]) = x_j, \end{aligned}$$

where we used Property 8 for the last equality. \square

Let $\mathbf{A}_{i,j}$ denote the submatrix of \mathbf{A} obtained by deleting the i th row and j th column of \mathbf{A} . For example,

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \quad \mathbf{A}_{1,1} = \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}, \quad \mathbf{A}_{1,2} = \begin{bmatrix} 4 & 6 \\ 7 & 9 \end{bmatrix}, \\ \mathbf{A}_{2,1} &= \begin{bmatrix} 2 & 3 \\ 8 & 9 \end{bmatrix}, \quad \mathbf{A}_{2,2} = \begin{bmatrix} 1 & 3 \\ 7 & 9 \end{bmatrix}, \quad \text{etc.} \end{aligned}$$

Definition A.3 (Cofactor and Adjoint)

For $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $1 \leq i, j \leq n$ the determinant $\det(\mathbf{A}_{ij})$ is called the **cofactor** of a_{ij} . The matrix $\text{adj}(\mathbf{A}) \in \mathbb{C}^{n \times n}$ with elements $(-1)^{i+j} \det(\mathbf{A}_{j,i})$ is called the **adjoint** of \mathbf{A} .

Theorem A.4 (The inverse as an adjoint)

If $\mathbf{A} \in \mathbb{C}^{n \times n}$ is nonsingular then

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A}).$$

Proof. Let $\mathbf{A}^{-1} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, where $\mathbf{x}_j = [x_{1j}, \dots, x_{nj}]^T$. The equation $\mathbf{AA}^{-1} = \mathbf{I}$ implies that $\mathbf{Ax}_j = \mathbf{e}_j$ for $j = 1, \dots, n$ and by Cramer's rule

$$x_{ij} = \frac{\det(\mathbf{A}_i(\mathbf{e}_j))}{\det(\mathbf{A})} = (-1)^{i+j} \frac{\det(\mathbf{A}_{ji})}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n.$$

For the last equality we first interchange the first and i th column of $\mathbf{A}_i(\mathbf{e}_j)$. By Property 4 it follows that $\det(\mathbf{A}_i(\mathbf{e}_j)) = (-1)^{i-1} \det([\mathbf{e}_j, \mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n])$. We then interchange row j and row 1. Using Property 8 we obtain

$$\det(\mathbf{A}_i(\mathbf{e}_j)) = (-1)^{i+j-2} \det(\mathbf{A}_{ji}) = (-1)^{i+j} \det(\mathbf{A}_{ji}).$$

□

Corollary A.5 (The adjoint and the inverse)

For any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have

$$\mathbf{A} \text{ adj}(\mathbf{A}) = \text{adj}(\mathbf{A})\mathbf{A} = \det(\mathbf{A})\mathbf{I}. \quad (\text{A.2})$$

Proof. If \mathbf{A} is nonsingular then (A.2) follows from Theorem A.4. We simply multiply by \mathbf{A} from the left and from the right. Suppose next that \mathbf{A} is singular with m zero eigenvalues $\lambda_1, \dots, \lambda_m$ and nonzero eigenvalues $\lambda_{m+1}, \dots, \lambda_n$. We define $\epsilon_0 := \min_{m+1 \leq j \leq n} |\lambda_j|$. For any $\epsilon \in (0, \epsilon_0)$ the matrix $\mathbf{A} + \epsilon\mathbf{I}$ has nonzero eigenvalues $\epsilon, \dots, \epsilon, \lambda_{m+1} + \epsilon, \dots, \lambda_n + \epsilon$ and hence is nonsingular. By what we have proved

$$(\mathbf{A} + \epsilon\mathbf{I}) \text{ adj}(\mathbf{A} + \epsilon\mathbf{I}) = \text{adj}(\mathbf{A} + \epsilon\mathbf{I})(\mathbf{A} + \epsilon\mathbf{I}) = \det(\mathbf{A} + \epsilon\mathbf{I})\mathbf{I}. \quad (\text{A.3})$$

Since the elements in $\mathbf{A} + \epsilon\mathbf{I}$ and $\text{adj}(\mathbf{A} + \epsilon\mathbf{I})$ depend continuously on ϵ we can take limits in (A.3) to obtain (A.2). □

Corollary A.6 (Cofactor expansion)

For any $\mathbf{A} \in \mathbb{C}^{n \times n}$ we have

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } i = 1, \dots, n, \quad (\text{A.4})$$

$$\det(\mathbf{A}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } j = 1, \dots, n. \quad (\text{A.5})$$

Proof. By (A.2) we have $\mathbf{A} \text{ adj}(\mathbf{A}) = \det(\mathbf{A})\mathbf{I}$. But then $\det(\mathbf{A}) = \mathbf{e}_i^T \mathbf{A} \text{adj}(\mathbf{A}) \mathbf{e}_i = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij})$ which is (A.4). Applying this row expansion to \mathbf{A}^T we find $\det(\mathbf{A}^T) = \sum_{j=1}^n (-1)^{i+j} a_{ji} \det(\mathbf{A}_{ji})$. Switching the roles of i and j proves (A.5). □

A.4 Computing Determinants

A determinant of an n -by- n matrix computed from the definition can contain up to $n!$ terms and we need other methods to compute determinants.

A matrix can be reduced to upper triangular form using elementary row operations. We can then use Property 1. to compute the determinant. The elementary operations using either rows or columns are

1. Interchanging two rows(columns).
2. Multiply a row(column) by a scalar α .
3. Add a constant multiple of one row(column) to another row(column).

Let \mathbf{B} be the result of performing an elementary operation on \mathbf{A} . For the three elementary operations the numbers $\det(\mathbf{A})$ and $\det(\mathbf{B})$ are related as follows.

1. $\det(\mathbf{B}) = -\det(\mathbf{A})$ (from Property 4.)
2. $\det(\mathbf{B}) = \alpha \det(\mathbf{A})$ (from Property 3.)
3. $\det(\mathbf{B}) = \det(\mathbf{A})$ (from Properties 5., 7.)

It follows from Property 2. that it is enough to show this for column operations. The proof of 1. and 2. are immediate. For 3. suppose we add α times column k to column i for some $k \neq i$. Then using Properties 5. and 7. we find

$$\begin{aligned} \det(\mathbf{B}) &= \det([a_1, \dots, a_{i-1}, a_i + \alpha a_k, a_{i+1}, \dots, a_n]) \\ &\stackrel{5}{=} \det(\mathbf{A}) + \det([a_1, \dots, a_{i-1}, \alpha a_k, a_{i+1}, \dots, a_n]) \stackrel{7}{=} \det(\mathbf{A}) \end{aligned}$$

A.5 Some Useful Determinant Formulas

Suppose $\mathbf{A} \in \mathbb{C}^{m \times n}$ and suppose for an integer $r \leq \min\{m, n\}$ that $\mathbf{i} = \{i_1, \dots, i_r\}$ and $\mathbf{j} = \{j_1, \dots, j_r\}$ are integers with $1 \leq i_1 < i_2 < \dots < i_r \leq m$ and $1 \leq j_1 < j_2 < \dots < j_r$. We let

$$\mathbf{A}(\mathbf{i}, \mathbf{j}) = \begin{bmatrix} a_{i_1, j_1} & \cdots & a_{i_1, j_r} \\ \vdots & & \vdots \\ a_{i_r, j_1} & \cdots & a_{i_r, j_r} \end{bmatrix}$$

be the submatrix of \mathbf{A} consisting of rows i_1, \dots, i_r and columns j_1, \dots, j_r . The following formula bears a strong resemblance to the formula for matrix multiplication.

Theorem A.7 (Cauchy-Binet formula)

Let $\mathbf{A} \in \mathbb{C}^{m \times p}$, $\mathbf{B} \in \mathbb{C}^{p \times n}$ and $\mathbf{C} = \mathbf{AB}$. Suppose $1 \leq r \leq \min\{m, n, p\}$ and let

$\mathbf{i} = \{i_1, \dots, i_r\}$ and $\mathbf{j} = \{j_1, \dots, j_r\}$ be integers with $1 \leq i_1 < i_2 < \dots < i_r \leq m$ and $1 \leq j_1 < j_2 < \dots < j_r \leq n$. Then

$$\det(\mathbf{C}(\mathbf{i}, \mathbf{j})) = \sum_{\mathbf{k}} \det(\mathbf{A}(\mathbf{i}, \mathbf{k})) \det(\mathbf{B}(\mathbf{k}, \mathbf{j})), \quad (\text{A.6})$$

where we sum over all $\mathbf{k} = \{k_1, \dots, k_r\}$ with $1 \leq k_1 < k_2 < \dots < k_r \leq p$.

Appendix B

Computer Arithmetic

B.1 Absolute and Relative Errors

Suppose a and b are real or complex scalars. If b is an approximation to a then there are different ways of measuring the error in b .

Definition B.1 (Absolute Error)

The absolute error in b as an approximation to a is the number $\epsilon := |a - b|$. The number $e := b - a$ is called the error in b as an approximation to a . This is what we have to add to a to get b .

Note that the absolute error is symmetric in a and b , so that ϵ is also the absolute error in a as an approximation to b

Definition B.2 (Relative Error) If $a \neq 0$ then the relative error in b as an approximation to a is defined by

$$\rho = \rho_b := \frac{|b - a|}{|a|}.$$

We say that a and b agree to approximately $-\log_{10} \rho$ digits.

As an example, if $a := 31415.9265$ and $b := 31415.8951$, then $\rho = 0.999493 * 10^{-6}$ and a and b agree to approximately 6 digits.

We have $b = a(1 + r)$ for some r if and only if $\rho = |r|$.

We can also consider the relative error $\rho_a := |a - b|/|b|$ in a as an approximation to b .

Lemma B.3 (Relative errors)

If $a, b \neq 0$ and $\rho_b < 1$ then $\rho_a \leq \rho_b/(1 - \rho_b)$.

Proof. Since $|a|\rho_b = |b - a| \geq |a| - |b|$ we obtain $|b| \geq |a| - |a - b| = (1 - \rho_b)|a|$. Then

$$\rho_a = \frac{|b - a|}{|b|} \leq \frac{|b - a|}{(1 - \rho_b)|a|} = \frac{\rho_b}{1 - \rho_b}.$$

□

If ρ_b is small then ρ_a is small and it does not matter whether we choose ρ_a or ρ_b to discuss relative error.

B.2 Floating Point Numbers

We shall assume that the reader is familiar with different number systems (binary, octal, decimal, hexadecimal) and how to convert from one number system to another. We use $(x)_\beta$ to indicate a number written to the base β . If no parenthesis and subscript are used, the base 10 is understood. For instance,

$$\begin{aligned}(100)_2 &= 4, \\ (0.1)_2 &= 0.5, \\ 0.1 &= (0.1)_{10} = (0.0001100110011001\dots)_2.\end{aligned}$$

In general,

$$x = (c_m c_{m-1} \dots c_0.d_1 d_2 \dots d_n)_\beta$$

means

$$x = \sum_{i=0}^m c_i \beta^i + \sum_{i=1}^n d_i \beta^{-i}, \quad 0 \leq c_i, d_i \leq \beta - 1.$$

We can move the decimal point by adding an exponent:

$$y = x \cdot \beta^e,$$

for example

$$(0.1)_{10} = (1.100110011001\dots)_2 \cdot 2^{-4}.$$

We turn now to a description of the floating-point numbers. We will only describe a **standard system**, namely the binary IEEE floating-point standard. Although it is not used by all systems, it has been widely adopted and is used in Matlab. For a more complete introduction to the subject see [11],[22].

We denote the real numbers which are represented in our computer by \mathcal{F} . The set \mathcal{F} are characterized by three integers t , and \underline{e}, \bar{e} . We define

$$\epsilon_M := 2^{-t}, \quad \text{machine epsilon,} \tag{B.1}$$

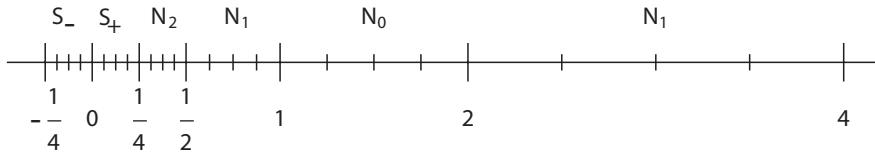


Figure B.1: Distribution of some positive floating-point numbers

and

$$\mathcal{F} := \{0\} \cup \mathcal{S} \cup \mathcal{N}, \text{ where}$$

$$\begin{aligned} \mathcal{N} &:= \mathcal{N}_+ \cup \mathcal{N}_-, \quad \mathcal{N}_+ := \cup_{e=\underline{e}}^{\bar{e}} \mathcal{N}_e, \quad \mathcal{N}_- := -\mathcal{N}_+, \\ \mathcal{N}_e &:= \{(1.d_1d_2 \cdots d_t)_2\} * 2^e = \{1, 1 + \epsilon_M, 1 + 2\epsilon_M, \dots, 2 - \epsilon_M\} * 2^e, \\ \mathcal{S} &:= \mathcal{S}_+ \cup \mathcal{S}_-, \quad \mathcal{S}_+ := \{\epsilon_M, 2\epsilon_M, 3\epsilon_M, \dots, 1 - \epsilon_M\} * 2^{\underline{e}}, \quad \mathcal{S}_- := -\mathcal{S}_+. \end{aligned} \quad (\text{B.2})$$

Example B.4 (Floating numbers)

Suppose $t := 2$, $\bar{e} = 3$ and $\underline{e} := -2$. Then $\epsilon_M = 1/4$ and we find

$$\begin{aligned} \mathcal{N}_{-2} &= \left\{ \frac{1}{4}, \frac{5}{16}, \frac{3}{8}, \frac{7}{16} \right\}, \quad \mathcal{N}_{-1} = \left\{ \frac{1}{2}, \frac{5}{8}, \frac{3}{4}, \frac{7}{8} \right\}, \quad \mathcal{N}_0 = \left\{ 1, \frac{5}{4}, \frac{3}{2}, \frac{3}{4}, \frac{7}{4} \right\}, \\ \mathcal{N}_1 &= \left\{ 2, \frac{5}{2}, 3, \frac{7}{2} \right\}, \quad \mathcal{N}_2 = \{4, 5, 6, 7\}, \quad \mathcal{N}_3 = \{8, 10, 12, 14\}, \\ \mathcal{S}_+ &= \left\{ \frac{1}{16}, \frac{1}{8}, \frac{3}{16} \right\}, \quad \mathcal{S}_- = \left\{ -\frac{3}{16}, -\frac{1}{8}, -\frac{1}{16} \right\}. \end{aligned}$$

The position of some of these sets on the real line is shown in Figure B.1

1. The elements of \mathcal{N} are called **normalized (floating-point) numbers**. They consists of three parts, the sign +1 or -1, the **mantissa** $(1.d_1d_2 \cdots d_t)_2$, and the **exponent part** 2^e .
2. the elements in \mathcal{N}_+ has the sign +1 indicated by the bit $\sigma = 0$ and the elements in \mathcal{N}_- has the sign bit $\sigma = 1$. Thus the sign of a number is $(-1)^\sigma$. The standard system has two zeros +0 and -0.
3. The mantissa is a number between 1 and 2. It consists of $t + 1$ binary digits.
4. The number e in the exponent part is restricted to the range $\underline{e} \leq e \leq \bar{e}$.
5. The positive normalized numbers are located in the interval $[r_m, r_M]$, where

$$r_m := 2^{\underline{e}}, \quad r_M := (2 - \epsilon_M) * 2^{\bar{e}}. \quad (\text{B.3})$$

6. The elements in \mathcal{S} are called **subnormal** or **denormalized**. As for normalized numbers they consists of three parts, but the mantissa is less than one in size. The main use of subnormal numbers is to soften the effect of underflow. If a number is in the range $(0, (1 - \epsilon_M/2) * 2^{\underline{e}})$, then it is rounded to the nearest subnormal number or to zero.
7. Two additional symbols "Inf" and "NaN" are used for special purposes.
8. The symbol **Inf** is used to represent numbers outside the interval $[-r_M, r_M]$ (**overflow**), and results of arithmetic operations of the form $x/0$, where $x \in \mathcal{N}$. Inf has a sign, +Inf and -Inf.
9. The symbol **NaN** stands for "not a number". a NaN results from illegal operations of the form $0/0, 0 * \text{Inf}, \text{Inf}/\text{Inf}, \text{Inf} - \text{Inf}$ and so on.
10. The choices of t , \bar{e} , and \underline{e} are to some extent determined by the architecture of the computer. A floating-point number, say x , occupies $n := 1 + \tau + t$ bits, where 1 bit is used for the sign, τ bits for the exponent, and t bits for the fractional part of the mantissa.

τ	t	
σ	exp	frac

Here $\sigma = 0$ if $x > 0$ and $\sigma = 1$ if $x < 0$, and $\text{exp} \in \{0, 1, 2, 3, \dots, 2^\tau - 1\}$ is an integer. The integer frac is the fractional part $d_1 d_2 \cdots d_t$ of the mantissa. The value of a normalized number in the standard system is

$$x = (-1)^\sigma * (1.\text{frac})_2 * 2^{\text{exp}-b}, \text{ where } b := 2^{\tau-1} - 1. \quad (\text{B.4})$$

The integer b is called the **bias**.

11. To explain the choice of b we note that the extreme values $\text{exp} = 0$ and $\text{exp} = 2^\tau - 1$ are used for special purposes. The value $\text{exp} = 0$ is used for the number zero and the subnormal numbers, while $\text{exp} = 2^\tau - 1$ is used for Inf and NaN. Since $2b = 2^\tau - 2$, the remaining numbers of exp , i.e., $\text{exp} \in \{1, 2, \dots, 2^\tau - 2\}$ correspond to e in the set $\{1 - b, 2 - b, \dots, b\}$. Thus in a standard system we have

$$\underline{e} = 1 - b, \quad \bar{e} = b := 2^{\tau-1} - 1. \quad (\text{B.5})$$

12. The most common choices of τ and t are shown in the following table

precision	τ	t	b	$\epsilon_M = 2^{-t}$	$r_m = 2^{1-b}$	r_M
half	5	10	15	9.8×10^{-4}	6.1×10^{-5}	6.6×10^4
single	8	23	127	1.2×10^{-7}	1.2×10^{-38}	3.4×10^{38}
double	11	52	1023	2.2×10^{-16}	2.2×10^{-308}	1.8×10^{308}
quad	15	112	16383	1.9×10^{-34}	3.4×10^{-4932}	1.2×10^{4932}

Here b is given by (B.5) and r_M by (B.3). The various lines correspond to a normalized number occupying **half** a word of 32 bits, one word (**single precision**), two words (**double precision**), and 4 words (**quad precision**).

B.3 Rounding and Arithmetic Operations

The standard system is a closed system. Every $x \in \mathbb{R}$ has a representation as either a floating-point number, or Inf or NaN, and every arithmetic operation produces a result. We denote the computer representation of a real number x by $\text{fl}(x)$.

B.3.1 Rounding

To represent a real number x there are three cases.

$$\text{fl}(x) = \begin{cases} \text{Inf}, & \text{if } x > r_M, \\ -\text{Inf}, & \text{if } x < -r_M, \\ \text{round to zero}, & \text{otherwise.} \end{cases}$$

To represent a real number with $|x| \leq r_M$ the system chooses a machine number $\text{fl}(x)$ closest to x . This is known as **rounding**. When x is midway between two numbers in \mathcal{F} we can either choose the one of larger magnitude (**round away from zero**), or pick the one with a zero last bit (**round to zero**). The standard system uses round to zero. As an example, if $x = 1 + \epsilon_M/2$, then x is midway between 1 and $1 + \epsilon_M$. Therefore $\text{fl}(x) = 1 + \epsilon_M$ if round away from zero is used, while $\text{fl}(x) = 1$ if x is rounded to zero. This is because the machine representation of 1 has `frac` = 0.

The following lemma gives a bound for the relative error in rounding.

Theorem B.5 (Relative error in rounding)

If $r_m \leq |x| \leq r_M$ then

$$\text{fl}(x) = x(1 + \delta), \quad |\delta| \leq u_M := \frac{1}{2}\epsilon_M = 2^{-t-1}.$$

Proof. Suppose $2^e < x < 2^{e+1}$. Then $\text{fl}(x) \in \{1, 1 + \epsilon_M, 1 + 2\epsilon_M, \dots, 2 - \epsilon_M\} * 2^e$. These numbers are uniformly spaced with spacing $\epsilon_M * 2^e$ and therefore $|\text{fl}(x) - x| \leq \frac{1}{2}\epsilon_M 2^e \leq \frac{1}{2}\epsilon_M * |x|$. The proof for a negative x is similar. \square

The number u_M is called the **rounding unit**.

B.3.2 Arithmetic Operations

Suppose $x, y \in \mathcal{N}$. In a standard system we have

$$\text{fl}(x \circ y) = (x \circ y)(1 + \delta), \quad |\delta| \leq u_M, \quad \circ \in \{+, -, *, /, \sqrt{\cdot}\}, \quad (\text{B.6})$$

where u_M is the rounding unit of the system. This means that the computed value is as good as the rounded exact answer. This is usually achieved by using one or several extra digits known as **guard digits** in the calculation.

B.4 Backward Rounding-Error Analysis

The computed sum of two numbers $\alpha_1, \alpha_2 \in \mathcal{N}$ satisfy $\text{fl}(\alpha_1 \circ \alpha_2) = (\alpha_1 + \alpha_2)(1 + \delta)$, where $|\delta| \leq u_M$, the rounding unit. If we write this as $\text{fl}(\alpha_1 \circ \alpha_2) = \tilde{\alpha}_1 + \tilde{\alpha}_2$, where $\tilde{\alpha}_i := \alpha_i(1 + \delta)$ for $i = 1, 2$, we see that the computed sum is the exact sum of two numbers which approximate the exact summands with small relative error, $|\delta| \leq u_M$. The error in the addition has been boomeranged back on the data α_1, α_2 , and in this context we call δ the **backward error**. A similar interpretation is valid for the other arithmetic operations $-$, $*$, $/$, $\sqrt{\cdot}$, and we assume it also holds for the elementary functions \sin, \cos, \exp, \log and so on.

Suppose more generally we want to compute the value of an expression $\phi(\alpha_1, \dots, \alpha_n)$. Here $\alpha_1, \dots, \alpha_n \in \mathcal{N}$ are given data, and we are using the arithmetic operations, and implementations of the standard elementary functions, in the computation. A **backward error analysis** consists of showing that the computed result is obtained as the exact result of using data $\beta := [\beta_1, \dots, \beta_n]^T$ instead of $\alpha := [\alpha_1, \dots, \alpha_n]$. In symbols

$$\tilde{\phi}(\alpha_1, \dots, \alpha_n) = \phi(\beta_1, \dots, \beta_n).$$

If we can show that the relative error in β as an approximation to α is $O(u_M)$ either componentwise or norm-wise in some norm, then we say that the algorithm to compute $\phi(\alpha_1, \dots, \alpha_n)$ is **backward stable**. Normally the constant K in the $O(u_M)$ term will grow with n . Typically $K = p(n)$ for some polynomial p is acceptable, while an exponential growth of K can be problematic.

B.4.1 Computing a Sum

We illustrate this discussion by computing the backward error in the sum of n numbers $s := \alpha_1 + \dots + \alpha_n$, where $\alpha_i \in \mathcal{N}$ for all i . We have the following

algorithm.

```

 $s_1 := \alpha_1$ 
for  $k = 2 : n$ 
     $s_k := \text{fl}(s_{k-1} + \alpha_k)$ 
end
 $\tilde{s} := s_n$ 

```

Using a standard system we obtain for $n = 3$

$$\begin{aligned}s_2 &= \text{fl}(\alpha_1 + \alpha_2) = \alpha_1(1 + \delta_2) + \alpha_2(1 + \delta_2), \\ s_3 &= \text{fl}(s_2 + \alpha_3) = s_2(1 + \delta_3) + \alpha_3(1 + \delta_3) = \alpha_1(1 + \eta_1) + \alpha_2(1 + \eta_2) + \alpha_3(1 + \eta_3), \\ \eta_1 &= \eta_2 = (1 + \delta_2)(1 + \delta_3), \quad \eta_3 = (1 + \delta_3), \quad |\delta_i| \leq u_M.\end{aligned}$$

In general, with $\delta_1 := 0$,

$$\tilde{s} = \sum_{i=1}^n \alpha_i(1 + \eta_i). \quad \eta_i = (1 + \delta_i) \dots (1 + \delta_n), \quad |\delta_i| \leq u_M, \quad i = 1, \dots, n. \quad (\text{B.7})$$

With $\phi(\alpha_1, \dots, \alpha_n) := \alpha_1 + \dots + \alpha_n$ this shows that

$$\tilde{s} = \tilde{\phi}(\alpha_1, \dots, \alpha_n) = \phi(\beta_1, \dots, \beta_n), \quad \beta_i = \alpha_i(1 + \eta_i). \quad (\text{B.8})$$

The following lemma gives a convenient bound on the η factors.

Lemma B.6 (Bound on factors)

Suppose for integers k, m with $0 \leq m \leq k$ and $k \geq 1$ that

$$1 + \eta_k := \frac{(1 + \delta_1) \dots (1 + \delta_m)}{(1 + \delta_{m+1}) \dots (1 + \delta_k)}, \quad |\delta_j| \leq u_M, \quad j = 1, \dots, k.$$

If $ku_M \leq \frac{1}{11}$ then

$$|\eta_k| \leq ku'_M, \quad \text{where } u'_M := 1.1u_M. \quad (\text{B.9})$$

Proof. We first show that

$$ku_M \leq \alpha < 1 \implies |\eta_k| \leq k \frac{u_M}{1 - \alpha}. \quad (\text{B.10})$$

For convenience we use $u := u_M$ in the proof. Since $u < 1$ we have $1/(1 - u) = 1 + u + u^2/(1 - u) > 1 + u$ and we obtain

$$(1 - u)^k \leq \frac{(1 - u)^m}{(1 + u)^{k-m}} \leq 1 + \eta_k \leq \frac{(1 + u)^m}{(1 - u)^{k-m}} \leq (1 - u)^{-k}.$$

The proof of (B.10) will be complete if we can show that

$$1 - ku \leq (1 - u)^k, \quad (1 - u)^{-k} \leq 1 + ku'.$$

The first inequality is an easy induction on k . If it holds for k , then

$$(1 - u)^{k+1} = (1 - u)^k(1 - u) \geq (1 - ku)(1 - u) = 1 - (k + 1)u + ku^2 \geq 1 - (k + 1)u.$$

The second inequality is a consequence of the first,

$$(1 - u)^{-k} \leq (1 - ku)^{-1} = 1 + \frac{ku}{1 - ku} \leq 1 + \frac{ku}{1 - \alpha} = 1 + ku'.$$

Letting $\alpha = \frac{1}{11}$ in (B.10) we obtain (B.9). \square

The number $u'_M := 1.1u_M$, corresponding to $\alpha = 1/11$, is called the **adjusted rounding unit**. In the literature many values of α can be found. [22] uses $\alpha = 1/10$ giving $u'_M = 1.12u_M$, while in [11] the value $\alpha = 0.01$ can be found. In the classical work [30] one finds $1/(1 - \alpha) = 1.06$.

Let us return to the backward error (B.8) in a sum of n numbers. Since $\delta_1 = 0$ we see that

$$|\eta_1| \leq (n - 1)u'_M, \quad |\eta_i| \leq (n - i + 1)u'_M, \text{ for } i = 2, \dots, n.$$

or more simply

$$|\eta_i| \leq (n - 1)u'_M, \text{ for } i = 1, \dots, n. \quad (\text{B.11})$$

This shows that the algorithm for computing a sum is backward stable.

The bounds from a backward rounding-error analysis can be used together with a condition number to bound the actual error in the computed result. To see this for the sum, we subtract the exact sum $s = \alpha_1 + \dots + \alpha_n$ from the computed sum $\tilde{s} = \alpha_1(1 + \eta_1) + \dots + \alpha_n(1 + \eta_n)$, to get

$$|\tilde{s} - s| = |\alpha_1\eta_1 + \dots + \alpha_n\eta_n| \leq (|\alpha_1| + \dots + |\alpha_n|)(n - 1)u'_M.$$

Thus the relative error in the computed sum of n numbers is bounded as follows

$$\left| \frac{\tilde{s} - s}{s} \right| \leq \kappa(n - 1)u'_M, \text{ where } \kappa := \frac{|\alpha_1| + \dots + |\alpha_n|}{\alpha_1 + \dots + \alpha_n}. \quad (\text{B.12})$$

This bound shows that the backward error can be magnified by at most κ . The number κ is called the **condition number** for the sum.

The condition number measures how much a relative error in each of the components in a sum can be magnified in the final sum. The backward error shows how large these relative perturbations can be in the actual algorithm we used to compute the sum. Using backward error analysis and condition number

separates the process of estimating the error in the final result into two distinct jobs.

A problem where small relative changes in the data leads to large relative changes in the exact result is called **ill conditioned**. We see that computing a sum can be ill-conditioned if the exact value of the sum is close to zero and some of the individual terms have large absolute values with opposite signs.

B.4.2 Computing an Inner Product

Computing an inner product $p := \alpha_1\gamma_1 + \cdots + \alpha_n\gamma_n$ is also backward stable using the standard algorithm

```

 $p_1 := \text{fl}(\alpha_1\gamma_1)$ 
for  $k = 2 : n$ 
     $p_k := \text{fl}(p_{k-1} + \text{fl}(\alpha_k\gamma_k))$ 
end
 $\tilde{p} := p_n$ 
```

For a backward error analysis of this algorithm we only need to modify (B.7) slightly. All we have to do is to add terms $\text{fl}(\alpha_k\gamma_k) = \alpha_k\gamma_k(1 + \pi_k)$ to the terms of the sum. The result is

$$\tilde{p} = \sum_{k=1}^n \alpha_k\gamma_k(1 + \eta_k), \quad \eta_k = (1 + \pi_k)(1 + \delta_k)\cdots(1 + \delta_n), \quad k = 1, \dots, n,$$

where $\delta_1 = 0$. Thus for the inner product of n terms we obtain

$$\left| \frac{\tilde{p} - p}{p} \right| \leq \kappa n u_M, \quad \kappa := \frac{|\alpha_1\gamma_1| + \cdots + |\alpha_n\gamma_n|}{|\alpha_1\gamma_1 + \cdots + \alpha_n\gamma_n|}. \quad (\text{B.13})$$

The computation can be ill conditioned if the exact value is close to zero and some of the components are large in absolute value.

B.4.3 Computing a Matrix Product

Using matrix norms we can bound the backward error in matrix algorithms. Suppose we want to compute the matrix product $\mathbf{C} = \mathbf{A} * \mathbf{B}$. Let n be the number of columns of \mathbf{A} and the number of rows of \mathbf{B} . Each element in \mathbf{C} is the inner product of a row of \mathbf{A} and a column of \mathbf{B} . Thus if $\tilde{\mathbf{C}}$ is the computed product then from (B.13)

$$\left| \frac{\tilde{c}_{ij} - c_{ij}}{c_{ij}} \right| \leq \kappa_{ij} n u'_M, \quad \kappa_{ij} := \frac{|a_1 b_1| + \cdots + |a_n b_n|}{|a_1 b_1 + \cdots + a_n b_n|}, \quad \text{all } i, j. \quad (\text{B.14})$$

We write this as $|\tilde{c}_{ij} - c_{ij}| \leq \kappa_{ij}|c_{ij}|nu'_M$. Using the infinity matrix norm we find

$$\sum_j |\tilde{c}_{ij} - c_{ij}| \leq nu'_M \sum_j \kappa_{ij}|c_{ij}| \leq \kappa n u'_M \sum_j |c_{ij}| \leq \kappa n u'_M \|\mathbf{C}\|_\infty, \text{ all } i,$$

where $\kappa := \max_{ij} \kappa_{ij}$. Maximizing over i we obtain

$$\frac{\|\tilde{\mathbf{C}} - \mathbf{C}\|_\infty}{\|\mathbf{C}\|_\infty} \leq \kappa n u'_M. \quad (\text{B.15})$$

The calculation of a matrix product can be ill conditioned if one or more of the product elements are small and the corresponding inner products have large terms of opposite signs.

Appendix C

Differentiation of Vector Functions

For any sufficiently differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R}$ we recall that the partial derivative with respect to the i th variable of f is defined by

$$D_i f(\mathbf{x}) := \frac{\partial f(\mathbf{x})}{\partial x_i} := \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}, \quad \mathbf{x} \in \mathbb{R}^n,$$

where \mathbf{e}_i is the i th unit vector in \mathbb{R}^n . For each $\mathbf{x} \in \mathbb{R}^n$ we define the **gradient** $\nabla f(\mathbf{x}) \in \mathbb{R}^n$, and the **hessian** $\nabla \nabla^T f(\mathbf{x}) \in \mathbb{R}^{n,n}$ of f by

$$\nabla f := \begin{bmatrix} D_1 f \\ \vdots \\ D_n f \end{bmatrix}, \quad \mathbf{H} f := \nabla \nabla^T f := \begin{bmatrix} D_1 D_1 f & \cdots & D_1 D_n f \\ \vdots & & \vdots \\ D_n D_1 & \cdots & D_n D_n f \end{bmatrix}, \quad (\text{C.1})$$

where $\nabla^T f := (\nabla f)^T$ is the row vector gradient. The operators $\nabla \nabla^T$ and $\nabla^T \nabla$ are quite different. Indeed, $\nabla^T \nabla f = D_1^2 f + \cdots + D_n^2 f =: \nabla^2$ the **Laplacian** of f , while $\nabla \nabla^T$ can be thought of as an outer product resulting in a matrix.

Lemma C.1 (Product rules)

For $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ we have the product rules

1. $\nabla(fg) = f\nabla g + g\nabla f, \quad \nabla^T(fg) = f\nabla^T g + g\nabla^T f,$
2. $\nabla \nabla^T(fg) = \nabla f \nabla^T g + \nabla g \nabla^T f + f \nabla \nabla^T g + g \nabla \nabla^T f.$
3. $\nabla^2(fg) = 2\nabla^T f \nabla g + f \nabla^2 g + g \nabla^2 f.$

We define the **Jacobian** of a vector function $\mathbf{f} = [f_1, \dots, f_m]^T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ as the m, n matrix

$$\nabla^T \mathbf{f} := \begin{bmatrix} D_1 f_1 & \cdots & D_n f_1 \\ \vdots & & \vdots \\ D_1 f_m & \cdots & D_n f_m \end{bmatrix}.$$

As an example, if $f(\mathbf{x}) = f(x, y) = x^2 - xy + y^2$ and $\mathbf{g}(x, y) := [f(x, y), x - y]^T$ then

$$\begin{aligned} \nabla f(x, y) &= \begin{bmatrix} 2x - y \\ -x + 2y \end{bmatrix}, & \nabla^T \mathbf{g}(x, y) &= \begin{bmatrix} 2x - y & -x + 2y \\ 1 & -1 \end{bmatrix}, \\ \mathbf{H}f(x, y) &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}. \end{aligned}$$

The second order Taylor expansion in n variables can be expressed in terms of the gradient and the hessian.

Lemma C.2 (Second order Taylor expansion)

Suppose $f \in C^2(\Omega)$, where $\Omega \subset \mathbb{R}^n$ contains two points $\mathbf{x}, \mathbf{x} + \mathbf{h} \in \Omega$, such that the line segment $L := \{\mathbf{x} + t\mathbf{h} : t \in (0, 1)\} \subset \Omega$. Then

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \mathbf{h}^T \nabla f(\mathbf{x}) + \frac{1}{2} \mathbf{h}^T \nabla \nabla^T f(\mathbf{c}) \mathbf{h}, \text{ for some } \mathbf{c} \in L. \quad (\text{C.2})$$

Proof. Let $g : [0, 1] \rightarrow \mathbb{R}$ be defined by $g(t) := f(\mathbf{x} + t\mathbf{h})$. Then $g \in C^2[0, 1]$ and by the chain rule

$$\begin{aligned} g(0) &= f(\mathbf{x}), & g(1) &= f(\mathbf{x} + \mathbf{h}), \\ g'(t) &= \sum_{i=1}^n h_i \frac{\partial f(\mathbf{x} + t\mathbf{h})}{\partial x_i} = \mathbf{h}^T \nabla f(\mathbf{x} + t\mathbf{h}), \\ g''(t) &= \sum_{i=1}^n \sum_{j=1}^n h_i h_j \frac{\partial^2 f(\mathbf{x} + t\mathbf{h})}{\partial x_i \partial x_j} = \mathbf{h}^T \nabla \nabla^T f(\mathbf{x} + t\mathbf{h}) \mathbf{h}. \end{aligned}$$

Inserting these expressions in the second order Taylor expansion

$$g(1) = g(0) + g'(0) + \frac{1}{2} g''(u), \text{ for some } u \in (0, 1),$$

we obtain (C.2) with $\mathbf{c} = \mathbf{x} + u\mathbf{h}$. \square

The gradient and hessian of some functions involving matrices can be found from the following lemma.

Lemma C.3 (Functions involving matrices)

For any $m, n \in \mathbb{N}$, $\mathbf{B} \in \mathbb{R}^{n,n}$, $\mathbf{C} \in \mathbb{R}^{m,n}$, and $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$ we have

1. $\nabla(\mathbf{y}^T \mathbf{C}) = \nabla^T(\mathbf{C}\mathbf{x}) = \mathbf{C}$,
2. $\nabla(\mathbf{x}^T \mathbf{B}\mathbf{x}) = (\mathbf{B} + \mathbf{B}^T)\mathbf{x}$, $\nabla^T(\mathbf{x}^T \mathbf{B}\mathbf{x}) = \mathbf{x}^T(\mathbf{B} + \mathbf{B}^T)$,
3. $\nabla\nabla^T(\mathbf{x}^T \mathbf{B}\mathbf{x}) = \mathbf{B} + \mathbf{B}^T$.

Proof.

1. We find $D_i(\mathbf{y}^T \mathbf{C}) = \lim_{h \rightarrow 0} \frac{1}{h} ((\mathbf{y} + h\mathbf{e}_i)^T \mathbf{C} - \mathbf{y}^T \mathbf{C}) = \mathbf{e}_i^T \mathbf{C}$ and $D_i(\mathbf{C}\mathbf{x}) = \lim_{h \rightarrow 0} \frac{1}{h} (\mathbf{C}(\mathbf{x} + h\mathbf{e}_i) - \mathbf{C}\mathbf{x}) = \mathbf{C}\mathbf{e}_i$ and 1. follows.

2. Here we find

$$\begin{aligned} D_i(\mathbf{x}^T \mathbf{B}\mathbf{x}) &= \lim_{h \rightarrow 0} \frac{1}{h} ((\mathbf{x} + h\mathbf{e}_i)^T \mathbf{B}(\mathbf{x} + h\mathbf{e}_i) - \mathbf{x}^T \mathbf{B}\mathbf{x}) \\ &= \lim_{h \rightarrow 0} (\mathbf{e}_i^T \mathbf{B}\mathbf{x} + \mathbf{x}^T \mathbf{B}\mathbf{e}_i + h\mathbf{e}_i^T \mathbf{e}_i) = \mathbf{e}_i^T (\mathbf{B} + \mathbf{B}^T) \mathbf{x}, \end{aligned}$$

and the first part of 2. follows. Taking transpose we obtain the second part.

3. Combining 1. and 2. we obtain 3.

□

Bibliography

- [1] Beckenbach, E. F, and R. Bellman, *Inequalities*, Springer Verlag, Berlin, Fourth Printing, 1983.
- [2] Björck, Åke, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1995.
- [3] E. Cohen, R. F. Riesenfeld, G. Elber, *Geometric Modeling with Splines: An Introduction*, A.K. Peters, Ltd., 2001,
- [4] Golub, G. H., and C. F. Van Loan, *Matrix Computations*, John Hopkins University Press, Baltimore, MD, third edition, 1996.
- [5] Grcar, Joseph F., Mathematicians of Gaussian elimination, Notices of the AMS, **58** (2011), 782–792.
- [6] Greenbaum, Anne, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [7] Hackbusch, Wolfgang, *Iterative Solution of Large Sparse Systems of Equations*, Springer-Verlag, Berlin, 1994.
- [8] Hall, C. A. and W. W. Meyer, Optimal error bounds for cubic spline interpolation. J. Approx. Theory, **16** (1976), 105122.
- [9] Hestenes, Magnus, *Conjugate Direction Methods in Optimization*, Springer-Verlag, Berlin, 1980.
- [10] Hestenes, M. and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, Journal of Research of the National Bureau of Standards **29**(1952), 409–439.
- [11] Higham, Nicloas J., *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.

- [12] Horn, Roger A. and Charles R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [13] Horn, Roger A. and Charles R. Johnson, *Topics in Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1991.
- [14] Kato, *Perturbation Theory for Linear Operators*, Pringer.
- [15] Lancaster, P., and Rodman, L., Canonical forms for hermitian matrix pairs under strict equivalence and congruence”, SIAM Review, vol. 47, 2005, 407-443.
- [16] Lawson, C.L. and R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, Englewood Cliffs, N.J, 1974.
- [17] Lax, Peter D., *Linear Algebra*, John Wiley & Sons, New York, 1997.
- [18] Lay, D.C: Linear algebra and its applications, 2012. Addison Wesley / Pearson. Fourth edition.
- [19] Leon, Steven J., *Linear Algebra with Applications*, Prentice Hall, NJ, Seventh Edition, 2006.
- [20] Meyer, Carl D., *Matrix Analysis and Applied Linear Algebra* , Siam Philadelphia, 2000.
- [21] Steel, J. Michael, *The Cauchy-Schwarz Master Class*, Cambridge University Press, Cambridge, UK, 2004.
- [22] Stewart, G. G., *Matrix Algorithms Volume I: Basic Decompositions*, Siam Philadelphia, 1998.
- [23] Stewart, G. G., *Matrix Algorithms Volume II: Eigensystems*, Siam Philadelphia, 2001.
- [24] Stewart, G. G. and Ji-guang Sun, *Matrix Perturbation Theory*, Academic Press, San Diego, 1990.
- [25] Stewart, G. G., *Introduction to Matrix Computations*, Academic press, New York, 1973.
- [26] Trefethen, Lloyd N., and David Bau III, *Numerical Linear Algebra*, Siam Philadelphia, 1997.
- [27] Tveito, A., and R. Winther, *Partial Differential Equations*, Springer, Berlin.
- [28] Van Loan, Charles, *Computational Frameworks for the Fast Fourier Transform*, Siam Philadelphia, 1992.

- [29] Varga, R. S., *Matrix Iterative Analysis/ 2nd Edn.*, Springer Verlag, New York, 2000.
- [30] Wilkinson, J. H., *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.
- [31] Young, D. M., *Iterative Solution of Large Linear Systems*, Academic Press, New York, 1971.
- [32] Zhang, F., *Matrix Theory*, Springer, New York, 1999.

List of Exercises

0.25	Linear combinations of convergent sequences	17
0.26	Coefficient norm	17
0.32	The $\mathbf{A}^T \mathbf{A}$ inner product	20
0.33	Complex inner product as sums of norms	20
0.34	Angle between vectors in complex case	20
0.53	The inverse of a general 2×2 matrix	28
0.54	The inverse of a 2×2 matrix	28
0.55	Sherman-Morrison formula	28
0.56	Cramer's rule; special case	29
0.57	Adjoint matrix; special case	29
0.59	Determinant equation for a plane	30
0.60	Signed area of a triangle	31
0.61	Vandermonde matrix	31
0.62	Cauchy determinant (1842)	31
0.63	Inverse of the Hilbert matrix	32
1.1	Matrix element as a quadratic form	42
1.2	Outer product expansion of a matrix	43
1.3	The product $\mathbf{A}^T \mathbf{A}$	43
1.4	Outer product expansion	43
1.5	System with many right hand sides; compact form	43
1.6	Block multiplication example	43
1.7	Another block multiplication example	43
1.14	Column oriented backsolve	47
1.17	Computing the inverse of a triangular matrix	48
1.21	Gaussian elimination example	53
1.22	Finite sums of integers	53
1.23	Operations	53
1.24	Multiplying triangular matrices	53
1.25	Matrix formulation of Gaussian elimination	53
1.31	Using PLU of \mathbf{A} to solve $\mathbf{A}^T \mathbf{x} = \mathbf{b}$	58
1.32	Using PLU to compute the determinant	58

1.33	Using PLU to compute the inverse	58
2.6	LU factorization of 2. derivative matrix	66
2.7	Inverse of 2. derivative matrix	66
2.8	Central difference approximation of 2. derivative	66
2.9	Two point boundary value problem	67
2.10	Two point boundary value problem; computation	67
2.16	Arctan example	75
2.18	Splineevaluation	76
2.20	Bounding the moments	77
2.21	Moment equations for 1. derivative boundary conditions	77
2.22	Proof of minimal 2. derivative property	77
3.9	Row interchange	86
3.10	LU of singular matrix	86
3.11	LU and determinant	86
3.12	Diagonal elements in U	86
3.16	Making a block LU into an LU	88
4.2	2×2 Poisson matrix	106
4.5	Properties of Kronecker products	109
4.14	2. derivative matrix is positive definite	114
4.15	1D test matrix is positive definite?	114
4.16	Eigenvalues 2×2 for 2D test matrix	114
4.17	Nine point scheme for Poisson problem	115
4.18	Matrix equation for nine point scheme	115
4.19	Biharmonic equation	115
5.5	Fourier matrix	126
5.6	Sine transform as Fourier transform	126
5.7	Explicit solution of the discrete Poisson equation	126
5.8	Improved version of Algorithm 5.1	126
5.9	Fast solution of 9 point scheme	127
5.10	Algorithm for fast solution of 9 point scheme	127
5.11	Fast solution of biharmonic equation	127
5.12	Algorithm for fast solution of biharmonic equation	128
5.13	Check algorithm for fast solution of biharmonic equation	128
5.14	Fast solution of biharmonic equation using 9 point rule	128
6.5	Unitary matrix	135
6.23	Find eigenpair example	140
6.24	Idempotent matrix	141
6.25	Idempotent matrix	141
6.26	Eigenvalues of a unitary matrix	141
6.27	Nonsingular approximation of a singular matrix	141
6.28	Companion matrix	141
6.32	Schur decomposition example	143

6.35	Skew-Hermitian matrix	144
6.36	Eigenvalues of a skew-Hermitian matrix	145
6.47	Eigenvalue perturbation for Hermitian matrices	149
6.49	Hoffman-Wielandt	149
6.54	Biorthogonal expansion	151
6.55	Generalized Rayleigh quotient	151
6.58	Jordan example	153
6.59	Big Jordan example	153
6.61	Jordan block example	154
6.62	Powers of a Jordan block	154
6.64	Minimal polynomial example	155
6.65	Similar matrix polynomials	155
6.66	Minimal polynomial of a diagonalizable matrix	155
7.11	SVD examples	165
7.12	More SVD examples	165
7.13	Singular values of a normal matrix	165
7.18	Orthonormal bases example	168
7.19	Some spanning sets	168
7.20	Singular values and eigenpair of composite matrix	168
7.25	Rank example	172
7.26	Another rank example	172
8.4	Consistency of sum norm?	178
8.5	Consistency of max norm?	178
8.6	Consistency of modified max norm?	178
8.8	The sum norm is subordinate to?	179
8.9	The max norm is subordinate to?	179
8.15	Spectral norm of the inverse	183
8.16	p -norm example	183
8.20	Univariance of spectral norm	184
8.21	$\ AU\ _2$ rectangular A	184
8.22	p -norm of diagonal matrix	184
8.23	spectral norm of a column vector	184
8.24	Norm of absolute value matrix	184
8.25	Spectral norm	185
8.26	Absolute norms	185
8.27	Is the spectral norm an absolute norm?	185
8.34	Sharpness of perturbation bounds	190
8.35	Condition number of 2. derivative matrix	190
8.43	p norm for $p = 1$ and $p = \infty$	194
8.44	The p - norm unit sphere	194
8.45	Sharpness of p -norm inequaltiy	194
8.46	p -norm inequaltiies for arbitrary p	194

9.9	Slow spectral radius convergence	208
9.11	A special norm	210
9.12	When is $\mathbf{A} + \mathbf{E}$ nonsingular?	210
9.17	Divergence example for J and GS	211
9.18	J and GS on spline matrix	211
9.19	Strictly diagonally dominance; The J method	211
9.20	Strictly diagonally dominance; The GS method	211
9.22	Estimate in Lemma 9.21 can be exact	212
9.25	The GS method converges, but not the J method	214
9.28	Convergence example for fix point iteration	217
10.1	Paraboloid	224
10.4	Steepest descent iteration	226
10.6	Maximum of a convex function	228
10.7	The \mathbf{A} -inner product	228
10.9	A test for the error bound	229
10.10	Orthogonality in steepest descent	229
10.12	Conjugate gradient iteration,II	231
10.13	Conjugate gradient iteration,III	232
10.16	The cg step length is optimal	233
10.17	Starting value in cg	233
10.23	Krylov space and cg iterations	237
10.24	Program code for testing steepest descent	238
10.25	Compare Richardson and steepest descent	238
10.26	Using cg to solve normal equations	238
10.31	An explicit formula for the Chebyshev polynomial	242
11.2	Reflector	256
11.5	What does algorithm housegen do when $\mathbf{x} = \mathbf{e}_1$?	258
11.6	Examples of Householder transformations	258
11.7	2×2 Householder transformation	258
11.15	QR decomposition	264
11.16	Householder triangulation	264
11.19	QR using Gram-Schmidt, II	265
11.21	Plane rotation	266
11.22	Solving upper Hessenberg system using roations	267
12.10	Straight line fit (linear regression)	274
12.11	Straight line fit using shifted power form	275
12.12	Fitting a circle to points	275
12.16	The generalized inverse	277
12.17	Uniqueness of generalized inverse	277
12.18	Verify that a matrix is a generalized inverse	278
12.19	Linearly independent columns and generalized inverse	278
12.20	The generalized inverse of a vector	278

12.21	The generalized inverse of an outer product	278
12.22	The generalized inverse of a diagonal matrix	278
12.23	Properties of the generalized inverse	278
12.24	The generalized inverse of a product	279
12.25	The generalized inverse of the conjugate transpose	279
12.26	Linearly independent columns	279
12.27	Analaysis of the general linear system	279
12.28	Fredholm's alternative	279
12.36	Condition number	285
12.37	Equality in perturbation bound	285
12.39	Problem using normal equations	286
13.7	Continuity of eigenvalues	296
13.9	Nonsingularity using Gerschgorin	296
13.10	Gerschgorin, strictly diagonally dominant matrix	296
13.12	Number of arithmetic operations	298
13.14	Number of arithmetic operations	298
13.15	Tridiagonalizea symmetric matrix	299
13.19	Counitng eigenvalues	302
13.20	Overflow in LDLT factorization	302
13.21	Simultaneous diagonalization	302
13.22	Program code for one eigenvalue	303
13.23	Determinant of upper Hessenberg matrix	303
13.25	∞ -norm of a diagonal matrix	304
14.3	Orthogonal vectors	309
14.13	QR convergence detail	318

List of everything

0.1	<i>Definition</i> Real vector space	4
0.2	<i>Definition</i> Linear combination	6
0.3	<i>Example</i> Linear combinations	6
0.4	<i>Definition</i> Linear independence	7
0.5	<i>Lemma</i> Linear independence and span	7
0.6	<i>Definition</i> basis	7
0.7	<i>Theorem</i> Basis subset of a spanning set	7
0.8	<i>Corollary</i> Existence of a basis	8
0.9	<i>Theorem</i> Dimension of a vector space	8
0.10	<i>Theorem</i> Enlarging vectors to a basis	8
0.11	<i>Definition</i> Subspace	8
0.12	<i>Example</i> Examples of subspaces	9
0.13	<i>Theorem</i> Dimension formula for sums of subspaces	9
0.14	<i>Theorem</i> Direct sum decomposition	10
0.15	<i>Lemma</i> Change of basis matrix	10
0.16	<i>Definition</i> Column space and null space	11
0.17	<i>Definition</i> Vector norm	12
0.18	<i>Definition</i> Vector p-norms	12
0.19	<i>Definition</i> Equivalent norms	13
0.20	<i>Theorem</i> Basic properties of vector norms	14
0.21	<i>Definition</i> Convergence of vectors	14
0.22	<i>Theorem</i> Norm convergence	15
0.23	<i>Definition</i> Cauchy sequence	15
0.24	<i>Theorem</i> Complete vector space	16
0.25	<i>Exercise</i> Linear combinations of convergent sequences	17
0.26	<i>Exercise</i> Coefficient norm	17
0.27	<i>Definition</i> Real inner product	17
0.28	<i>Definition</i> Complex inner product	18
0.29	<i>Theorem</i> Cauchy-Schwarz inequality	18
0.30	<i>Theorem</i> Inner product norm	19
0.31	<i>Theorem</i> Parallelogram Identity	19

0.32	<i>Exercise</i> The $\mathbf{A}^T \mathbf{A}$ inner product	20
0.33	<i>Exercise</i> Complex inner product as sums of norms	20
0.34	<i>Exercise</i> Angle between vectors in complex case	20
0.35	<i>Definition</i> Orthogonality	20
0.36	<i>Theorem</i> Pythagoras	20
0.37	<i>Definition</i> Orthogonal- and Orthonormal Bases	21
0.38	<i>Theorem</i> Gram-Schmidt	21
0.39	<i>Theorem</i> Orthogonal Extension of basis	22
0.40	<i>Corollary</i> Extending orthogonal vectors to a basis	22
0.41	<i>Theorem</i> Orthogonal Projection	22
0.42	<i>Definition</i> Orthogonal sum	23
0.43	<i>Theorem</i> Column space decomposition	24
0.44	<i>Lemma</i> Underdetermined system	25
0.45	<i>Definition</i> Real nonsingular matrix	25
0.46	<i>Theorem</i> Linear systems; existence and uniqueness	25
0.47	<i>Lemma</i> Complex underdetermined system	26
0.48	<i>Definition</i> Complex nonsingular matrix	26
0.49	<i>Theorem</i> Complex linear system; existence and uniqueness	26
0.50	<i>Theorem</i> Product of nonsingular matrices	26
0.51	<i>Theorem</i> When is a square matrix invertible?	27
0.52	<i>Corollary</i> Basic properties of the inverse matrix	27
0.53	<i>Exercise</i> The inverse of a general 2×2 matrix	28
0.54	<i>Exercise</i> The inverse of a 2×2 matrix	28
0.55	<i>Exercise</i> Sherman-Morrison formula	28
0.56	<i>Exercise</i> Cramer's rule; special case	29
0.57	<i>Exercise</i> Adjoint matrix; special case	29
0.58	<i>Example</i> Determinant equation for a straight line	30
0.59	<i>Exercise</i> Determinant equation for a plane	30
0.60	<i>Exercise</i> Signed area of a triangle	31
0.61	<i>Exercise</i> Vandermonde matrix	31
0.62	<i>Exercise</i> Cauchy determinant (1842)	31
0.63	<i>Exercise</i> Inverse of the Hilbert matrix	32
0.64	<i>Lemma</i> Characteristic equation	33
0.65	<i>Definition</i> Characteristic polynomial of a matrix	33
0.66	<i>Theorem</i> Derived eigenpairs	33
0.67	<i>Theorem</i> Eigenvalues of a triangular matrix	34
0.68	<i>Theorem</i> Sums and products of eigenvalues; trace	34
0.69	<i>Theorem</i> Zero eigenvalue	35
1.1	<i>Exercise</i> Matrix element as a quadratic form	42
1.2	<i>Exercise</i> Outer product expansion of a matrix	43
1.3	<i>Exercise</i> The product $\mathbf{A}^T \mathbf{A}$	43
1.4	<i>Exercise</i> Outer product expansion	43

1.5	<i>Exercise</i> System with many right hand sides; compact form	43
1.6	<i>Exercise</i> Block multiplication example	43
1.7	<i>Exercise</i> Another block multiplication example	43
1.8	<i>Lemma</i> Inverse of a block triangular matrix	43
1.9	<i>Lemma</i> Inverse of a triangular matrix	44
1.10	<i>Lemma</i> Product of triangular matrices	45
1.11	<i>Lemma</i> Unit triangular matrices	45
1.12	<i>Algorithm</i> forwardsolve	46
1.13	<i>Algorithm</i> backsolve	47
1.14	<i>Exercise</i> Column oriented backsolve	47
1.15	<i>Algorithm</i> Forward Solve (column oriented)	48
1.16	<i>Algorithm</i> Backsolve (column oriented)	48
1.17	<i>Exercise</i> Computing the inverse of a triangular matrix	48
1.18	<i>Theorem</i> When is naive Gaussian elimination possible?	50
1.19	<i>Theorem</i> Gauss=LU	50
1.20	<i>Definition</i> $G_n := \frac{2}{3}n^3$	51
1.21	<i>Exercise</i> Gaussian elimination example	53
1.22	<i>Exercise</i> Finite sums of integers	53
1.23	<i>Exercise</i> Operations	53
1.24	<i>Exercise</i> Multiplying triangular matrices	53
1.25	<i>Exercise</i> Matrix formulation of Gaussian elimination	53
1.26	<i>Definition</i> Elementary lower triangular matrix	54
1.27	<i>Definition</i>	55
1.28	<i>Definition</i> Interchange matrix	55
1.29	<i>Theorem</i> PLU theorem	56
1.30	<i>Algorithm</i> PLU factorization	58
1.31	<i>Exercise</i> Using PLU of \mathbf{A} to solve $\mathbf{A}^T \mathbf{x} = \mathbf{b}$	58
1.32	<i>Exercise</i> Using PLU to compute the determinant	58
1.33	<i>Exercise</i> Using PLU to compute the inverse	58
1.34	<i>Example</i> Row pivoting	58
2.1	<i>Algorithm</i> trifactor	63
2.2	<i>Algorithm</i> trisolve	64
2.3	<i>Definition</i> Diagonal dominance	64
2.4	<i>Theorem</i> Strict diagonal dominance	64
2.5	<i>Theorem</i> Weak diagonal dominance	65
2.6	<i>Exercise</i> LU factorization of 2. derivative matrix	66
2.7	<i>Exercise</i> Inverse of 2. derivative matrix	66
2.8	<i>Exercise</i> Central difference approximation of 2. derivative	66
2.9	<i>Exercise</i> Two point boundary value problem	67
2.10	<i>Exercise</i> Two point boundary value problem; computation	67
2.11	<i>Example</i> A cubic spline interpolant	70
2.12	<i>Theorem</i> Cubic spline; minimal 2. derivative	70

2.13	<i>Theorem</i> Cubic spline with not-a-knot boundary conditions	72
2.14	<i>Example</i> Not-a-knot	74
2.15	<i>Algorithm</i> findsubintervals	75
2.16	<i>Exercise</i> Arctan example	75
2.17	<i>Algorithm</i> splineint	76
2.18	<i>Exercise</i> Splineevaluation	76
2.19	<i>Algorithm</i> splineeval	77
2.20	<i>Exercise</i> Bounding the moments	77
2.21	<i>Exercise</i> Moment equations for 1. derivative bounary conditions	77
2.22	<i>Exercise</i> Proof of minimal 2. derivative property	77
3.1	<i>Example</i> LU of 2×2 matrix	82
3.2	<i>Example</i> LU of 3×3 matrices	83
3.3	<i>Definition</i> Principal submatrix	83
3.4	<i>Example</i> Principal submatrices	84
3.5	<i>Lemma</i> LU of leading principal sub matrices	84
3.6	<i>Theorem</i> LU Theorem	84
3.7	<i>Remark</i> LU of upper triangular matrix	85
3.8	<i>Remark</i> PLU factorization	86
3.9	<i>Exercise</i> Row interchange	86
3.10	<i>Exercise</i> LU of singular matrix	86
3.11	<i>Exercise</i> LU and determinant	86
3.12	<i>Exercise</i> Diagonal elements in U	86
3.13	<i>Theorem</i> Block LU theorem	87
3.14	<i>Remark</i> Comparing LU and block LU	87
3.15	<i>Remark</i> A block LU is not an LU	88
3.16	<i>Exercise</i> Making a block LU into an LU	88
3.17	<i>Definition</i> Symmetric LU	88
3.18	<i>Example</i> 2×2 symmetric LU	88
3.19	<i>Lemma</i> Symmetric LU of leading principal sub matrices	89
3.20	<i>Theorem</i>	89
3.21	<i>Example</i> 2×2 positive definite	91
3.22	<i>Lemma</i> \mathbf{T} is symmetric positive definite	91
3.23	<i>Example</i> Gradient and Hessian	91
3.24	<i>Theorem</i> A general criterium	92
3.25	<i>Theorem</i> Principal submatrices	92
3.26	<i>Corollary</i> positive (semi)definite criteria	93
3.27	<i>Corollary</i> $\mathbf{A}^T \mathbf{A}$ is symmetric positive semidefinite	93
3.28	<i>Lemma</i> Eigenvalues of a Hermitian matrix	93
3.29	<i>Lemma</i> Symmetry and positive eigenvalues	93
3.30	<i>Lemma</i> Symmetric posotive definite and symmetric LU	94
3.31	<i>Theorem</i> Symmetric positive definite characterization	94
3.32	<i>Definition</i> Cholesky	95

3.33	<i>Theorem</i> Cholesky	95
3.34	<i>Example</i> 2×2	95
3.35	<i>Lemma</i> Banded Cholesky factor	96
3.36	<i>Algorithm</i> bandcholesky	97
3.37	<i>Lemma</i> Criteria symmetric semidefinite	97
3.38	<i>Definition</i> Semi-Cholesky factorization	98
3.39	<i>Theorem</i> Characterization, semi-Cholesky factorization	98
3.40	<i>Theorem</i> Positive symmetric semidefinite characterization	99
3.41	<i>Theorem</i> Bandwidth semi Cholesky factor	100
3.42	<i>Algorithm</i> bandsemicholesky	101
4.1	<i>Definition</i> vec operation	105
4.2	<i>Exercise</i> 2×2 Poisson matrix	106
4.3	<i>Definition</i> Kronecker Product	108
4.4	<i>Definition</i> Kronecker sum	109
4.5	<i>Exercise</i> Properties of Kronecker products	109
4.6	<i>Lemma</i> Mixed Product Rule	109
4.7	<i>Lemma</i> Eigenvalues of Kronecker products	110
4.8	<i>Lemma</i> Eigenvalues of Kronecker sums	110
4.9	<i>Lemma</i> Kronecker product;inverse and positive definite	111
4.10	<i>Lemma</i> Conversion Kronecker product to matrix equation	111
4.11	<i>Lemma</i> Eigenpairs of 2. derivative matrix	112
4.12	<i>Lemma</i> Eigenpairs of a Hermitian matrix	113
4.13	<i>Theorem</i> Eigenpairs of 2D test matrix	113
4.14	<i>Exercise</i> 2. derivative matrix is positive definite	114
4.15	<i>Exercise</i> 1D test matrix is positive definite?	114
4.16	<i>Exercise</i> Eigenvalues 2×2 for 2D test matrix	114
4.17	<i>Exercise</i> Nine point scheme for Poisson problem	115
4.18	<i>Exercise</i> Matrix equation for nine point scheme	115
4.19	<i>Exercise</i> Biharmonic equation	115
5.1	<i>Algorithm</i> Fast Poisson Solver	120
5.2	<i>Lemma</i> Sine transform as Fourier transform	122
5.3	<i>Theorem</i> Fast Fourier Transform	124
5.4	<i>Algorithm</i> Recursive FFT	125
5.5	<i>Exercise</i> Fourier matrix	126
5.6	<i>Exercise</i> Sine transform as Fourier transform	126
5.7	<i>Exercise</i> Explicit solution of the discrete Poisson equation	126
5.8	<i>Exercise</i> Improved version of Algorithm 5.1	126
5.9	<i>Exercise</i> Fast solution of 9 point scheme	127
5.10	<i>Exercise</i> Algorithm for fast solution of 9 point scheme	127
5.11	<i>Exercise</i> Fast solution of biharmonic equation	127
5.12	<i>Exercise</i> Algorithm for fast solution of biharmonic equation	128
5.13	<i>Exercise</i> Check algorithm for fast solution of biharmonic equation	128

5.14	<i>Exercise</i> Fast solution of biharmonic equation using 9 point rule	128
6.1	<i>Definition</i> Orthonormal matrix	133
6.2	<i>Theorem</i> Orthonormal matrix	133
6.3	<i>Definition</i> Unitary matrix	134
6.4	<i>Theorem</i> Unitary matrix	134
6.5	<i>Exercise</i> Unitary matrix	135
6.6	<i>Definition</i> Similar matrices	135
6.7	<i>Theorem</i> Eigenvalues of similar matrices	135
6.8	<i>Corollary</i> Spectra of \mathbf{AB} and \mathbf{BA}	136
6.9	<i>Theorem</i> Eigenvectors of nondefective matrices	136
6.10	<i>Definition</i> Defective and nondefective matrix	136
6.11	<i>Corollary</i> Diagonalizable matrix	136
6.12	<i>Corollary</i> Eigenvectors of \mathbf{A}^*	137
6.13	<i>Theorem</i> Distinct eigenvalues	137
6.14	<i>Corollary</i> Nondefective matrix	137
6.15	<i>Example</i> Two upper triangular matrices	138
6.16	<i>Definition</i> Geometric multiplicity	139
6.17	<i>Example</i> Geometric multiplicity	139
6.18	<i>Theorem</i> $g \leq a$	139
6.19	<i>Definition</i> Defective eigenvalue	139
6.20	<i>Theorem</i> The number of linearly independent eigenvectors	139
6.21	<i>Corollary</i> Linearly independent eigenvectors characterization	140
6.22	<i>Theorem</i> Geometric multiplicity of similar matrices	140
6.23	<i>Exercise</i> Find eigenpair example	140
6.24	<i>Exercise</i> Idempotent matrix	141
6.25	<i>Exercise</i> Idempotent matrix	141
6.26	<i>Exercise</i> Eigenvalues of a unitary matrix	141
6.27	<i>Exercise</i> Nonsingular approximation of a singular matrix	141
6.28	<i>Exercise</i> Companion matrix	141
6.29	<i>Theorem</i> Schur decomposition	142
6.30	<i>Example</i> Deflation example	142
6.31	<i>Theorem</i> Schur form, real eigenvalues	143
6.32	<i>Exercise</i> Schur decomposition example	143
6.33	<i>Definition</i> Quasi-triangular matrix	143
6.34	<i>Definition</i> Normal Matrix	144
6.35	<i>Exercise</i> Skew-Hermitian matrix	144
6.36	<i>Exercise</i> Eigenvalues of a skew-Hermitian matrix	145
6.37	<i>Theorem</i> Orthonormal eigenpairs characterization	145
6.38	<i>Theorem</i> Spectral theorem, complex form	146
6.39	<i>Theorem</i> Spectral Theorem (real form)	146
6.40	<i>Example</i>	146
6.41	<i>Definition</i> Rayleigh quotient	146

6.42	<i>Lemma</i> Convex combination of the eigenvalues	146
6.43	<i>Theorem</i> Minmax	147
6.44	<i>Theorem</i> Maxmin	148
6.45	<i>Corollary</i> The Courant-Fischer Theorem	148
6.46	<i>Theorem</i> Eigenvalue perturbation for Hermitian matrices	148
6.47	<i>Exercise</i> Eigenvalue perturbation for Hermitian matrices	149
6.48	<i>Theorem</i> Hoffman-Wielandt Theorem	149
6.49	<i>Exercise</i> Hoffman-Wielandt	149
6.50	<i>Definition</i> Left eigenpair	149
6.51	<i>Theorem</i> Biorthogonality	149
6.52	<i>Theorem</i> Simple eigenvalue	150
6.53	<i>Theorem</i> Biorthogonal eigenvector expansion	150
6.54	<i>Exercise</i> Biorthogonal expansion	151
6.55	<i>Exercise</i> Generalized Rayleigh quotient	151
6.56	<i>Definition</i> Jordan block	151
6.57	<i>Theorem</i> The Jordan form of a matrix	152
6.58	<i>Exercise</i> Jordan example	153
6.59	<i>Exercise</i> Big Jordan example	153
6.60	<i>Lemma</i> Properties of the Jordan form	153
6.61	<i>Exercise</i> Jordan block example	154
6.62	<i>Exercise</i> Powers of a Jordan block	154
6.63	<i>Definition</i> Minimal polynomial of a matrix	154
6.64	<i>Exercise</i> Minimal polynomial example	155
6.65	<i>Exercise</i> Similar matrix polynomials	155
6.66	<i>Exercise</i> Minimal polynomial of a diagonalizable matrix	155
6.67	<i>Theorem</i> Proof of real Schur form	156
7.1	<i>Lemma</i> Eigenpairs of $\mathbf{A}^* \mathbf{A}$	160
7.2	<i>Theorem</i> Orthogonal bases for column- and null space of \mathbf{A}	160
7.3	<i>Definition</i> Singular values	161
7.4	<i>Theorem</i> rank = # positive singular values	161
7.5	<i>Theorem</i> Existence of singular value decomposition	161
7.6	<i>Corollary</i> Singular value factorization	162
7.7	<i>Example</i> Nonsingular matrix	163
7.8	<i>Example</i> Full row rank	163
7.9	<i>Example</i> Full column rank	164
7.10	<i>Example</i> $r < n < m$	164
7.11	<i>Exercise</i> SVD examples	165
7.12	<i>Exercise</i> More SVD examples	165
7.13	<i>Exercise</i> Singular values of a normal matrix	165
7.14	<i>Lemma</i> SVD of $\mathbf{A}^* \mathbf{A}$ and $\mathbf{A} \mathbf{A}^*$	166
7.15	<i>Theorem</i> Singular vectors and orthonormal bases	166
7.16	<i>Corollary</i> Counting dimensions of fundamental subspaces	167

7.17	<i>Theorem</i> Rank and nullity relations	167
7.18	<i>Exercise</i> Orthonormal bases example	168
7.19	<i>Exercise</i> Some spanning sets	168
7.20	<i>Exercise</i> Singular values and eigenpair of composite matrix	168
7.21	<i>Example</i> Ellipse	169
7.22	<i>Lemma</i> Frobenius norm properties	170
7.23	<i>Theorem</i> Frobenius norm and singular values	171
7.24	<i>Theorem</i> Best low rank approximation	171
7.25	<i>Exercise</i> Rank example	172
7.26	<i>Exercise</i> Another rank example	172
7.27	<i>Theorem</i> The Courant-Fischer Theorem for Singular Values	172
7.28	<i>Theorem</i> Hoffman-Wielandt Theorem for singular values	173
8.1	<i>Definition</i> Matrix Norms	177
8.2	<i>Theorem</i> Matrix norm equivalence	177
8.3	<i>Definition</i> Consistent Matrix Norms	178
8.4	<i>Exercise</i> Consistency of sum norm?	178
8.5	<i>Exercise</i> Consistency of max norm?	178
8.6	<i>Exercise</i> Consistency of modified max norm?	178
8.7	<i>Definition</i> Subordinate Matrix Norms	179
8.8	<i>Exercise</i> The sum norm is subordinate to?	179
8.9	<i>Exercise</i> The max norm is subordinate to?	179
8.10	<i>Definition</i> Operator Norm	179
8.11	<i>Lemma</i> The operator norm is a matrix norm	180
8.12	<i>Theorem</i> onetwoinfnorms	181
8.13	<i>Example</i> Compare onetwoinfnorms	182
8.14	<i>Theorem</i> Spectral norm	182
8.15	<i>Exercise</i> Spectral norm of the inverse	183
8.16	<i>Exercise</i> p -norm example	183
8.17	<i>Theorem</i> Spectral norm bound	183
8.18	<i>Definition</i> Unitary invariant norm	183
8.19	<i>Theorem</i> Unitary invariant norms	184
8.20	<i>Exercise</i> Univariance of spectral norm	184
8.21	<i>Exercise</i> $\ AU\ _2$ rectangular A	184
8.22	<i>Exercise</i> p -norm of diagonal matrix	184
8.23	<i>Exercise</i> spectral norm of a column vector	184
8.24	<i>Exercise</i> Norm of absolute value matrix	184
8.25	<i>Exercise</i> Spectral norm	185
8.26	<i>Exercise</i> Absolute norms	185
8.27	<i>Exercise</i> Is the spectral norm an absolute norm?	185
8.28	<i>Theorem</i> Perturbation in the right-hand side	186
8.29	<i>Theorem</i> Spectral condition number	187
8.30	<i>Theorem</i> Nonsingularity of perturbation of identity	187

8.31	<i>Theorem</i> Nonsingularity of perturbation	188
8.32	<i>Theorem</i> Perturbation and residual	189
8.33	<i>Theorem</i> Perturbation of inverse matrix	189
8.34	<i>Exercise</i> Sharpness of perturbation bounds	190
8.35	<i>Exercise</i> Condition number of 2. derivative matrix	190
8.36	<i>Theorem</i> The p norms are norms	191
8.37	<i>Definition</i> Convex function	191
8.38	<i>Lemma</i> A sufficient condition for convexity	191
8.39	<i>Theorem</i> Jensen's Inequality	192
8.40	<i>Corollary</i> Weighted geometric/arithmetic mean inequality	192
8.41	<i>Corollary</i> Hölder's inequality	193
8.42	<i>Corollary</i> Minkowski's inequality	193
8.43	<i>Exercise</i> p norm for $p = 1$ and $p = \infty$	194
8.44	<i>Exercise</i> The p -norm unit sphere	194
8.45	<i>Exercise</i> Sharpness of p -norm inequaltiy	194
8.46	<i>Exercise</i> p -norm inequalties for arbitrary p	194
9.1	<i>Algorithm</i> Jacobi	203
9.2	<i>Algorithm</i> SOR	203
9.3	<i>Proposition</i> Splitting matrices for J, GS, and SOR	205
9.4	<i>Example</i> Splitting matrices	205
9.5	<i>Theorem</i> When is $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$?	206
9.6	<i>Theorem</i> Any consistent norm majorizes the spectral radius	207
9.7	<i>Theorem</i> The spectral radius can be approximated by a norm	207
9.8	<i>Theorem</i> Spectral radius convergence	208
9.9	<i>Exercise</i> Slow spectral radius convergence	208
9.10	<i>Theorem</i> Neumann Series	209
9.11	<i>Exercise</i> A special norm	210
9.12	<i>Exercise</i> When is $\mathbf{A} + \mathbf{E}$ nonsingular?	210
9.13	<i>Definition</i> Convergence of fixed-point iteration	210
9.14	<i>Lemma</i> Convergence of an iterative method	211
9.15	<i>Theorem</i> When does an iterative method converge?	211
9.16	<i>Corollary</i> Sufficient condition for convergence	211
9.17	<i>Exercise</i> Divergence example for J and GS	211
9.18	<i>Exercise</i> J and GS on spline matrix	211
9.19	<i>Exercise</i> Strictly diagonally dominance; The J method	211
9.20	<i>Exercise</i> Strictly diagonally dominance; The GS method	211
9.21	<i>Lemma</i> Number of iterations	212
9.22	<i>Exercise</i> Estimate in Lemma 9.21 can be exact	212
9.23	<i>Lemma</i> Be careful when stopping	212
9.24	<i>Proposition</i> Convergence of Richardson's method	213
9.25	<i>Exercise</i> The GS method converges, but not the J method	214
9.26	<i>Theorem</i> The spectral radius of SOR matrix	215

9.28	<i>Exercise</i> Convergence example for fix point iteration	217
9.29	<i>Lemma</i> SOR iteration matrix	217
9.30	<i>Theorem</i> Necessay condition for convergence of SOR	218
9.31	<i>Theorem</i> SOR on positive definite matrix	218
9.32	<i>Theorem</i> The optimal ω	219
10.1	<i>Exercise</i> Paraboloid	224
10.2	<i>Lemma</i> Quadratic function	224
10.3	<i>Example</i> Steepest descent iteration	226
10.4	<i>Exercise</i> Steepest descent iteration	226
10.5	<i>Theorem</i> Kantorovich inequality	227
10.6	<i>Exercise</i> Maximum of a convex function	228
10.7	<i>Exercise</i> The \mathbf{A} -inner product	228
10.8	<i>Theorem</i> Error bound for steepest descent	228
10.9	<i>Exercise</i> A test for the error bound	229
10.10	<i>Exercise</i> Orthogonality in steepest descent	229
10.11	<i>Example</i> Conjugate gradient iteration	231
10.12	<i>Exercise</i> Conjugate gradient iteration,II	231
10.13	<i>Exercise</i> Conjugate gradient iteration,III	232
10.14	<i>Lemma</i> Krylov space	232
10.15	<i>Theorem</i> Best approximation property	232
10.16	<i>Exercise</i> The cg step length is optimal	233
10.17	<i>Exercise</i> Starting value in cg	233
10.18	<i>Algorithm</i> Conjugate Gradient Iteration	234
10.19	<i>Algorithm</i> Testing Conjugate Gradient	235
10.22	<i>Theorem</i> Error bounds for cg	236
10.23	<i>Exercise</i> Krylov space and cg iterations	237
10.24	<i>Exercise</i> Program code for testing steepest descent	238
10.25	<i>Exercise</i> Compare Richardson and steepest descent	238
10.26	<i>Exercise</i> Using cg to solve normal equations	238
10.27	<i>Lemma</i> Krylov space and polyomials	238
10.28	<i>Theorem</i> cg and best polynomial approximation	239
10.29	<i>Lemma</i> Closed forms of Chebyshev polynomials	240
10.30	<i>Theorem</i> A minimal norm problem	241
10.31	<i>Exercise</i> An explicit formula for the Chebyshev polynomial	242
10.32	<i>Theorem</i> The error in cg is strictly decreasing	243
10.33	<i>Algorithm</i> Preconditioned conjugate gradient	246
10.34	<i>Theorem</i> Error bound preconditioned cg	246
10.1	<i>Theorem</i> Positive definite matrix	249
10.3	<i>Theorem</i> Eigevalues of preconditioned matrix	251
11.1	<i>Definition</i> Householder Transformation	255
11.2	<i>Exercise</i> Reflector	256
11.3	<i>Theorem</i> Zeros in vectors	257

11.4	<i>Algorithm</i> Generate a Householder transformation	257
11.5	<i>Exercise</i> What does algorithm housegen do when $\mathbf{x} = \mathbf{e}_1$?	258
11.6	<i>Exercise</i> Examples of Householder transformations	258
11.7	<i>Exercise</i> 2×2 Householder transformation	258
11.8	<i>Algorithm</i> Householder Triangulation	260
11.9	<i>Lemma</i> Updating a Householder transformation	261
11.10	<i>Definition</i> QR decomposition	262
11.11	<i>Example</i> QR decomposition and factorization	262
11.12	<i>Theorem</i> Existence of QR decomposition	263
11.13	<i>Example</i> QR decomposition and factorization	263
11.14	<i>Theorem</i> Hadamard's Inequality	263
11.15	<i>Exercise</i> QR decomposition	264
11.16	<i>Exercise</i> Householder triangulation	264
11.17	<i>Theorem</i> QR and Gram-Schmidt	264
11.18	<i>Example</i> QR using Gram-Schmidt	265
11.19	<i>Exercise</i> QR using Gram-Schmidt, II	265
11.20	<i>Definition</i> Givens rotation, plane rotation	266
11.21	<i>Exercise</i> Plane rotation	266
11.22	<i>Exercise</i> Solving upper Hessenberg system using rotations	267
11.23	<i>Algorithm</i> Upper Hessenberg linear system	268
12.1	<i>Definition</i> Least squares problem	269
12.2	<i>Theorem</i> Existence	269
12.3	<i>Theorem</i> Uniqueness	270
12.4	<i>Theorem</i> Characterization	270
12.5	<i>Example</i> Average	271
12.6	<i>Example</i> Input/output model	271
12.7	<i>Lemma</i> Curve fitting	273
12.8	<i>Example</i> Straight line fit	273
12.9	<i>Example</i> Ill conditioning and the Hilbert matrix	274
12.10	<i>Exercise</i> Straight line fit (linear regression)	274
12.11	<i>Exercise</i> Straight line fit using shifted power form	275
12.12	<i>Exercise</i> Fitting a circle to points	275
12.13	<i>Theorem</i> Orthogonal projections	276
12.14	<i>Example</i> Orthogonal projections	276
12.15	<i>Theorem</i> General least squares solution	277
12.16	<i>Exercise</i> The generalized inverse	277
12.17	<i>Exercise</i> Uniqueness of generalized inverse	277
12.18	<i>Exercise</i> Verify that a matrix is a generalized inverse	278
12.19	<i>Exercise</i> Linearly independent columns and generalized inverse	278
12.20	<i>Exercise</i> The generalized inverse of a vector	278
12.21	<i>Exercise</i> The generalized inverse of an outer product	278
12.22	<i>Exercise</i> The generalized inverse of a diagonal matrix	278

12.23	<i>Exercise</i> Properties of the generalized inverse	278
12.24	<i>Exercise</i> The generalized inverse of a product	279
12.25	<i>Exercise</i> The generalized inverse of the conjugate transpose	279
12.26	<i>Exercise</i> Linearly independent columns	279
12.27	<i>Exercise</i> Analaysis of the general linear system	279
12.28	<i>Exercise</i> Fredholm's alternative	279
12.29	<i>Example</i> Unique solution of normal equations	280
12.30	<i>Example</i> Nonunique solution of normal equations	280
12.31	<i>Example</i> Solution using QR factorization	282
12.32	<i>Example</i> Solution using singular value factorization	283
12.33	<i>Theorem</i> Minimal solution	283
12.34	<i>Theorem</i> Perturbing the Right Hand Side	284
12.35	<i>Example</i> Perturbing the Right Hand Side	284
12.36	<i>Exercise</i> Condition number	285
12.37	<i>Exercise</i> Equality in perturbation bound	285
12.38	<i>Theorem</i> Perturbing the Matrix	286
12.39	<i>Exercise</i> Problem using normal equations	286
12.40	<i>Theorem</i> Perturbation of singular values	287
12.41	<i>Theorem</i> Generalized inverse when perturbing the matrix	287
13.1	<i>Theorem</i> Elsner's Theorem	292
13.2	<i>Theorem</i> Linearly independent eigenvectors	293
13.3	<i>Theorem</i> Perturbations, normal matrix	293
13.4	<i>Theorem</i> Gershgorin's Circle Theorem	294
13.5	<i>Example</i> Gershgorin	294
13.6	<i>Corollary</i> Disjoint Gershgorin disks	295
13.7	<i>Exercise</i> Continuity of eigenvalues	296
13.8	<i>Example</i>	296
13.9	<i>Exercise</i> Nonsingularity using Gerschgorin	296
13.10	<i>Exercise</i> Gerschgorin, strictly diagonally dominant matrix	296
13.11	<i>Algorithm</i> Householder reduction to Hessenberg form	298
13.12	<i>Exercise</i> Number of arithmetic operations	298
13.13	<i>Algorithm</i> Assemble Householder transformations	298
13.14	<i>Exercise</i> Number of arithmetic operations	298
13.15	<i>Exercise</i> Tridiagonalizea symmetric matrix	299
13.16	<i>Lemma</i> Distinct eigevalues of a tridiagonal matrix	300
13.17	<i>Theorem</i> Sylvester's Inertia Theorem	301
13.18	<i>Corollary</i> Counting eigenvalues using the LDLT factorization	301
13.19	<i>Exercise</i> Counting eigenvalues	302
13.20	<i>Exercise</i> Overflow in LDLT factorization	302
13.21	<i>Exercise</i> Simultaneous diagonalization	302
13.22	<i>Exercise</i> Program code for one eigenvalue	303
13.23	<i>Exercise</i> Determinant of upper Hessenberg matrix	303

13.24	<i>Lemma</i> p -norm of a diagonal matrix	304
13.25	<i>Exercise</i> ∞ -norm of a diagonal matrix	304
13.26	<i>Theorem</i> Absolute errors	304
13.27	<i>Theorem</i> Relative errors	305
14.1	<i>Lemma</i> Convergence of the power method	308
14.2	<i>Theorem</i> The Rayleigh quotient minimizes the residual	308
14.3	<i>Exercise</i> Orthogonal vectors	309
14.4	<i>Algorithm</i> The Power Method	309
14.5	<i>Example</i> Power method	309
14.6	<i>Algorithm</i> Rayleigh quotient iteration	311
14.7	<i>Example</i> Rayleigh quotient iteration	312
14.9	<i>Example</i> QR iteration; real eigenvalues	313
14.10	<i>Example</i> QR iteration; complex eigenvalues	313
14.11	<i>Theorem</i> QR and power	314
14.12	<i>Theorem</i> Convergence of basis QR	317
14.13	<i>Exercise</i> QR convergence detail	318
A.1	<i>Example</i> Properties of permutations	323
A.2	<i>Theorem</i> Cramer's rule (1750)	327
A.3	<i>Definition</i> Cofactor and Adjoint	328
A.4	<i>Theorem</i> The inverse as an adjoint	328
A.5	<i>Corollary</i> The adjoint and the inverse	329
A.6	<i>Corollary</i> Cofactor expansion	329
A.7	<i>Theorem</i> Cauchy-Binet formula	330
B.1	<i>Definition</i> Absolute Error	333
B.2	<i>Definition</i> Relative Error	333
B.3	<i>Lemma</i> Relative errors	333
B.4	<i>Example</i> Floating numbers	335
B.5	<i>Theorem</i> Relative error in rounding	337
B.6	<i>Lemma</i> Bound on factors	339
C.1	<i>Lemma</i> Product rules	343
C.2	<i>Lemma</i> Second order Taylor expansion	344
C.3	<i>Lemma</i> Functions involving matrices	344

Index

- 1D test matrix, 107
- 2D test matrix, 107
- convex combinations, 227
- cubic spline
 - C^2 , 69
- eigenvector expansion, 136
- singular values (SVD), 159
- absolute convergence, 17
- absolute error, 186, 333
- adjoint matrix, 328
- adjusted rounding unit, 340
- algebraic multiplicity, 138
- algorithms
 - assemble Householder transformations, 298
 - backsolve, 47
 - backsolve column oriented, 48
 - bandcholesky, 97
 - cg, 234
 - fastpoisson, 120
 - findsubintervals, 75
 - forwardsolve, 46
 - forwardsolve column oriented, 48
 - housegen, 257
 - Householder reduction to Hessenberg form, 298
 - Householder triangulation, 260
 - Jacobi, 203
 - PLU factorization, 58
- preconditioned cg, 246
- Rayleigh quotient iteration, 311
- SOR, 203
- splineevaluation, 76
- splineint, 76
- testing Conjugate Gradient, 235
- the Power Method, 309
- trifactor, 63
- trisolve, 64
- upper Hessenberg linear system, 268
- averaging matrix, 108
- backward error, 338
- backward stable, 338
- banded matrix, 4
 - symmetric LU factorization, 97
- banded symmetric LU factorization, 97
- bandsemicholesky, 101
- basis coefficients, 15
- biharmonic equation, 115
 - fast solution method, 128
 - nine point rule, 128
- block LU theorem, 87
- Cauchy determinant, 31
- Cauchy sequence, 16
- Cauchy-Binet formula, 330
- Cauchy-Schwarz inequality, 18

- Cayley Hamilton Theorem, 155
central difference, 66
central difference approximation
 second derivative, 67
change of basis matrix, 10
characteristic equation, 33
characteristic polynomial, 33, 138
Chebyshev polynomial, 240
Cholesky factor, 95
Cholesky factorization, 95
coefficient norm, 15
cofactor, 328
column operations, 330
column space (span), 12
column space decomposition, 22,
 24
companion matrix, 141
complete pivoting, 60
computer arithmetic, 333
condition number, 340
 ill-conditioned, 186
congruent matrices, 301
conjugate gradient method, 223
 convergence, 238
 derivation, 230
 Krylov spaces, 232
 least squares problem, 238
 preconditioning, 244
 preconditioning algorithm, 246
 preconditioning convergence,
 246
convergence
 absolute, 209
convex combination, 146, 191
convex function, 191
Courant-Fischer theorem, 148
Cramer's rule, 29, 328
Crout factorization, 82
cubic B-spline, 74
cubic spline
 first derivative boundary con-
 ditions, 70
minimal 2. derivative, 70
not-a-knot, 72
defective eigenvalue, 139
defective matrix, 136
deflation, 142
determinant, 323
 additivity, 324
 area of a triangle, 31
 block triangular, 324
 Cauchy, 31
 Cauchy-Binet, 330
 cofactor, 328
 cofactor expansion, 29, 329
 homogeneity, 324
 permutation of columns, 324
 plane equation, 30
 product rule, 324
 singular matrix, 324
 straight line equation, 30
 transpose, 324
 triangular matrix, 324
 Vandermonde, 31
dirac delta, 3
direct sum decomposition, 10
discrete cosine transform, 121
discrete Fourier transform, 121,
 122
 Fourier matrix, 122
discrete sine transform, 121
double precision, 337
eigenpair, 32
 left eigenpair, 149
eigenvalue, 32
 algebraic multiplicity, 138
 characteristic equation, 33
 characteristic polynomial, 33
 Courant-Fischer theorem, 148
 defective, 139
 derived eigenpairs, 34
 geometric multiplicity, 139

-
- Hoffman-Wielandt theorem, 149
 - Kronecker sum, 110
 - location, 294
 - Rayleigh quotient, 146
 - Schur form, real, 156
 - spectral theorem, 146
 - spectrum, 33
 - triangular matrix, 34
 - eigenvector, 32
 - Kronecker sum, 110
 - left eigenvector, 149
 - elementary divisors, 154
 - elementary lower triangular matrix, 54
 - elementary reflector, 255
 - Elsner's theorem, 292
 - equivalent norms, 13
 - extension of basis, 22
 - fast Fourier transform, 121, 123
 - recursive FFT, 125
 - fill-inn, 118
 - finite difference method, 61
 - finite dimensional vector space, 6
 - fixed-point, 210
 - fixed-point iteration, 210
 - floating-point number
 - bias, 336
 - denormalized, 336
 - double precision, 337
 - exponent part, 335
 - guard digits, 338
 - half precision, 337
 - Inf, 336
 - mantissa, 335
 - NaN, 336
 - normalized, 335
 - overflow, 336
 - quadruple precision, 337
 - round away from zero, 337
 - round to zero, 337
 - rounding, 337
 - rounding unit, 337
 - single precision, 337
 - subnormal, 336
 - Fourier matrix, 122
 - Fredholm's alternative, 279
 - Frobenius norm, 170
 - Gaussian elimination
 - complete pivoting, 60
 - elementary lower triangular matrix, 54
 - Gauss=LU, 50
 - interchange matrix, 55
 - naive, 48
 - partial pivoting, 58
 - pivot, 55
 - pivot vector, 56
 - pivoting, 54
 - PLU factorization, 56
 - row pivoting, 59
 - generalized inverse, 277
 - geometric multiplicity, 139
 - Gershgorin's theorem, 294
 - Given's rotation, 266
 - gradient, 91, 343
 - gradient method, 226
 - Gram-Schmidt, 21
 - guard digits, 338
 - Hölder's inequality, 13, 193
 - Hadamard's inequality, 263
 - half precision, 337
 - hessian, 91, 343
 - Hilbert matrix, 32, 274
 - Hoffman-Wielandt theorem, 149
 - Householder transformation, 255
 - identity matrix, 3
 - ill-conditioned, 341
 - ill-conditioned problem, 186
 - inequality

- geometric/arithmetic mean, 193
- Hölder, 193
- Kantorovich, 227
- Minkowski, 194
- Inf, 336
- inner product, 17
 - complex, 18
 - inner product norm, 17, 18
 - Pythagoras' theorem, 20
 - standard inner product in \mathbb{C}^n , 18
 - standard inner product in \mathbb{R}^n , 18
- inner product space
 - orthogonal basis, 21
 - orthonormal basis, 21
- interchange matrix, 55
- inverse power method, 310
- inverse triangle inequality, 14
- iterative method
 - convergence, 211
 - Gauss-Seidel, 200
 - Jacobi, 200
 - SOR, 200
 - SOR, convergence, 217
 - SSOR, 201
- iterative methods, 199
- Jacobian, 344
- Jordan form, 152
 - elementary divisors, 154
 - Jordan block, 151
 - Jordan canonical form, 152
 - principal vectors, 153
- Kronecker product, 108
 - eigenvalues, 110
 - eigenvectors, 110
 - inverse, 111
 - left product, 108
 - mixed product rule, 109
- nonsingular, 111
- positive definite, 111
- right product, 108
- symmetry, 111
- transpose, 109
- Kronecker sum, 109
 - eigenvalues, 110
 - eigenvectors, 110
 - nonsingular, 111
 - positive definite, 111
 - symmetry, 111
- Krylov space, 232
- Laplacian, 343
- leading principal block submatrices, 87
- leading principal minor, 84
- leading principal submatrices, 84
- least squares
 - error analysis, 283
- least squares problem, 269
- least squares solution, 269
- left eigenpair, 149
- left eigenvector, 149
- left triangular, 81
- linear combination, 6
- linear interpolation polynomial, 68
- linear system
 - existence and uniqueness, 25, 26
 - homogeneous, 24
 - overdetermined, 24
 - residual vector, 189
 - square, 24
 - underdetermined, 24
- linearly dependent, 7
- linearly independent, 7
- LLT factorization, 95
- LU factorization, 81
 - LDLT, 88
 - symmetric, 88

- symmetric LU , 89
- LU theorem, 84
- mantissa, 335
- matrix
 - addition, 3
 - adjoint, 29, 328
 - adjoint formula for the inverse, 29
 - block lower triangular, 4
 - block matrix, 41
 - block triangular, 44
 - block upper triangular, 4
 - blocks, 41
 - cofactor, 29
 - column space (span), 12
 - companion matrix, 141
 - computing inverse, 48
 - defective, 136
 - deflation, 142
 - diagonal, 3
 - element-by-element operations, 3
 - Hadamard product, 3
 - Hilbert, 32
 - idempotent, 141
 - ill-conditioned, 187
 - inverse, 26
 - inverse Hilbert matrix, 32
 - invertible, 26
 - Kronecker product, 108
 - leading principal minor, 83
 - leading principal submatrices, 83, 84
 - left inverse, 26
 - left triangular, 3
 - lower Hessenberg, 4
 - lower triangular, 3
 - LU theorem, 84
 - multiplication, 3
 - negative (semi)definite, 90
 - Neumann series, 209
 - nilpotent, 141
 - nonsingular, 25, 26
 - norm, 177
 - normal, 144
 - null space (\ker), 12
 - outer product, 43
 - outer product expansion, 43
 - permutation, 55
 - positive definite, 90
 - positive semidefinite, 90
 - principal minor, 83
 - principal submatrix, 83
 - product of triangular matrices, 45
 - quasi-triangular, 144
 - right inverse, 26
 - right triangular, 3
 - row space, 12
 - scalar multiplication, 3
 - Schur product, 3
 - second derivative, 62
 - similar matrices, 135
 - similarity transformation, 135
 - singular, 25, 26
 - spectral radius, 206
 - strictly diagonally dominant, 64
 - symmetric positive semidefinite, 90
 - test matrix, 2D, 107
 - test matrix, 1D , 107
 - trace, 35
 - triangular, 44
 - tridiagonal, 4
 - unit triangular, 45
 - upper Hessenberg, 3
 - upper trapezoidal, 259
 - upper triangular, 3
 - vec operation, 105
 - weakly diagonally dominant, 64
 - well-conditioned, 187

- matrix norm
 consistent norm, 178
 Frobenius norm, 170, 178
 max norm, 178
 operator norm, 179
 spectral norm, 181
 subordinate norm, 179
 sum norm, 178
 two-norm, 181
minimal polynomial, 154
Minkowski's inequality, 13, 194
mixed product rule, 109
moments, 72
- NaN, 336
natural ordering, 105
negative (semi)definite, 90
Neumann series, 209
nilpotent matrix, 141
nonsingular matrix, 25
nontrivial subspaces, 9
norm, 12
 l_1 -norm, 12
 l_2 -norm, 13
 l_∞ -norm, 13
 absolute norm, 185
 continuity, 14
 Euclidian norm, 13
 infinity-norm, 13
 max norm, 13
 monotone norm, 185
 one-norm, 12
 triangle inequality, 12
 two-norm, 13
norm convergence, 15
normal equations, 270
normal matrix, 144
null space (\ker), 12
- operation count, 51
optimal relaxation parameter, 216
optimal step length, 225
- orthogonal matrix, see orthonormal matrix, 133
orthogonal projection, 22
orthogonal sum, 23
orthonormal matrix, 133
outer product, 43
overflow, 336
- p-norms, 12
paraboloid, 224
parallelogram identity, 19
partial pivoting, 58
permutation, 321
 identity, 322
 inversion, 322
 sign, 322
 symmetric group, 323
permutation matrix, 55
perpendicular vectors, 20
pivot row, 55
pivot vector, 55, 56
plane rotation, 266
PLU factorization, 39, 57, 81, 86
Poisson matrix, 106
Poisson problem, 103
 five point stencil, 105
 nine point scheme, 115
 Poisson matrix, 106
 variable coefficients, 248
Poisson problem (1D), 61
positive definite, 90
positive semidefinite, 90
power method, 307
 inverse, 310
 Rayleigh quotient iteration, 311
 shifted, 310
preconditioned conjugate gradient method, 223
preconditioning, 244
principal minor, 84, 94
principal submatrix, 84

-
- principal vectors, 153
 - pseudo inverse, 277
 -
 - QR algorithm
 - implicit shift, 316
 - Rayleigh quotient shift, 316
 - shifted, 316
 - Wilkinson shift, 316
 - QR decomposition, 262
 - QR factorization, 262
 - quadratic form, 90
 - quadruple precision, 337
 -
 - rate of convergence, 212
 - Rayleigh quotient, 146
 - generalized, 151
 - Rayleigh quotient iteration, 311
 - relative error, 186, 333
 - residual vector, 189
 - right triangular, 81
 - rotation in the i, j -plane, 267
 - rounding unit, 337
 - rounding-error analysis
 - adjusted rounding unit, 340
 - backward error, 338
 - backward stable, 338
 - condition number, 340
 - ill-conditioned, 341
 - row operations, 330
 - row space, 12
 -
 - scalar product, 17
 - scaled partial pivoting, 59
 - Schur form, real, 156
 - second derivative matrix, 62
 - semi-Cholesky factorization, 98
 - Sherman-Morrison formula, 28
 - shifted power method, 310
 - similar matrices, 135
 - similarity transformation, 135
 - single precision, 337
 - singular value factorization(SVF), 162
 -
 - singular values, 161
 - Courant-Fischer theorem, 172
 - error analysis, 287
 - Hoffman-Wielandt theorem, 173
 - singular vectors, 159
 - span, 6
 - spectral radius, 206
 - spectral theorem, 146
 - spectrum, 33
 - splitting matrices for J, GS, and SOR, 205
 - steepest descent, 226
 - stencil, 105
 - sums of integers, 53
 - Sylvester's inertia theorem, 301
 - symmetric positive semidefinite, 90
 -
 - trace, 35
 - triangle inequality, 12
 - triangular matrix
 - left triangular, 81
 - right triangular, 81
 - unit triangular, 82
 - trivial subspace, 9
 - two point boundary value problem, 61
 -
 - unit triangular, 82
 - unit vectors, 3
 - upper trapezoidal matrix, 259
 -
 - vector
 - angle, 20
 - linearly dependent, 7
 - linearly independent, 7
 - nontrivial subspaces, 9
 - orthogonal, 20
 - orthonormal, 20
 - vector norm, 12
 - vector space
 - absolute convergence, 17

basis, 7
basis coefficients, 15
bounded sequence, 16
Cauchy sequence, 16
change of basis matrix, 10
coefficient norm, 15
complementary, 9
complete, 16
complex inner product space,
 18
convergent series, 17
dimension, 8
dimension formula for sums
 of subspaces, 9
direct sum, 9
direct sum decomposition, 10
enlarging vectors to a basis,
 8
examples of subspaces, 9
existence of basis, 8
intersection, 9
norm convergence, 15
normed, 12
orthogonal vectors, 20
real, 5
real inner product space, 17
span, 7
subsequence, 16
subspace, 8
sum, 9
union, 9
vectorization, 105