

# Lecture Notes for Inf-Mat 4350, 2009

Tom Lyche

August 11, 2009



# Contents

<b>Preface</b>	<b>vii</b>
<b>I A Review of Linear Algebra</b>	<b>1</b>
<b>1 Introduction</b>	<b>3</b>
1.1 Notation . . . . .	3
<b>2 Vectors</b>	<b>7</b>
2.1 Vector Spaces . . . . .	7
2.2 Linear Independence and Bases . . . . .	10
2.3 Operations on Subspaces . . . . .	12
2.3.1 Sums and intersections of subspaces . . . . .	12
2.3.2 The quotient space . . . . .	14
2.4 Norms . . . . .	14
2.5 Convergence of Vectors . . . . .	16
2.5.1 Convergence of Series of Vectors . . . . .	18
2.6 Inner Products . . . . .	19
2.7 Orthogonality . . . . .	21
2.8 Projections and Orthogonal Complements . . . . .	24
<b>3 Matrices</b>	<b>25</b>
3.1 Arithmetic Operations and Block Multiplication . . . . .	25
3.1.1 Block Multiplication . . . . .	26
3.2 The Transpose Matrix . . . . .	26
3.3 Linear Systems . . . . .	27
3.4 The Inverse matrix . . . . .	29
3.5 Rank, Nullity, and the Fundamental Subspaces . . . . .	31
3.6 Linear Transformations and Matrices . . . . .	33
3.7 Orthonormal and Unitary Matrices . . . . .	34
<b>4 Determinants</b>	<b>37</b>
4.1 Permutations . . . . .	37
4.2 Basic Properties of Determinants . . . . .	39

4.3	The Adjoint Matrix and Cofactor Expansion . . . . .	43
4.4	Computing Determinants . . . . .	45
4.5	Some Useful Determinant Formulas . . . . .	48
<b>5</b>	<b>Eigenvalues and Eigenvectors</b>	<b>49</b>
5.1	The Characteristic Polynomial . . . . .	49
5.1.1	The characteristic equation . . . . .	49
5.2	Similarity Transformations . . . . .	53
5.3	Linear Independence of Eigenvectors . . . . .	54
5.4	Left Eigenvectors . . . . .	57
<b>II</b>	<b>Some Linear Systems with a Special Structure</b>	<b>59</b>
<b>6</b>	<b>Examples of Linear Systems</b>	<b>61</b>
6.1	Cubic Spline Interpolation . . . . .	61
6.1.1	Cubic $C^2$ Splines . . . . .	61
6.1.2	Determining the Interpolant . . . . .	64
6.1.3	Strictly Diagonally Dominant Matrices . . . . .	67
6.2	The Second Derivative Matrix . . . . .	69
6.3	LU Factorization of a Tridiagonal System . . . . .	69
6.3.1	Diagonal Dominance . . . . .	71
6.3.2	Periodic Boundary Conditions . . . . .	72
6.4	Block Multiplication and Triangular Matrices . . . . .	76
6.4.1	Block Multiplication . . . . .	76
6.4.2	Triangular matrices . . . . .	79
<b>7</b>	<b>LU Factorizations</b>	<b>81</b>
7.1	The LU Factorization . . . . .	81
7.2	Block LU Factorization . . . . .	84
7.3	The Symmetric LU Factorization . . . . .	85
7.4	Positive Definite- and Positive Semidefinite Matrices . . . . .	86
7.4.1	Definition and Examples . . . . .	86
7.4.2	Some Criteria for the Nonsymmetric Case . . . . .	87
7.5	The Symmetric Case and Cholesky Factorization . . . . .	89
7.6	An Algorithm for SemiCholesky Factorization of a Banded Matrix . . . . .	92
7.7	The PLU Factorization . . . . .	95
<b>8</b>	<b>The Kronecker Product</b>	<b>97</b>
8.1	Test Matrices . . . . .	97
8.1.1	The 2D Poisson Problem . . . . .	97
8.1.2	The test Matrices . . . . .	100
8.2	The Kronecker Product . . . . .	101
8.3	Properties of the 1D and 2D Test Matrices . . . . .	104
<b>9</b>	<b>Fast Direct Solution of a Large Linear System</b>	<b>109</b>

9.1	Algorithms for a Banded Positive Definite System . . . . .	109
9.1.1	Cholesky Factorization . . . . .	109
9.1.2	Block LU Factorization of a Block Tridiagonal Matrix	110
9.1.3	Other Methods . . . . .	110
9.2	A Fast Poisson Solver based on Diagonalization . . . . .	111
9.3	A Fast Poisson Solver based on the Discrete Sine and Fourier Transforms . . . . .	112
9.3.1	The Discrete Sine Transform (DST) . . . . .	112
9.3.2	The Discrete Fourier Transform (DFT) . . . . .	113
9.3.3	The Fast Fourier Transform (FFT) . . . . .	114
9.3.4	A Poisson Solver based on the FFT . . . . .	116
9.4	Problems . . . . .	117
<b>III Some Matrix Theory</b>		<b>119</b>
<b>10</b>	<b>Orthonormal Eigenpairs and the Schur Form</b>	<b>121</b>
10.1	The Schur Form . . . . .	121
10.2	Hermitian and Normal Matrices . . . . .	123
10.2.1	The Spectral Theorem . . . . .	124
10.3	The Rayleigh Quotient and Minmax Theorems . . . . .	125
10.3.1	The Rayleigh Quotient . . . . .	125
10.3.2	Minmax and Maxmin Theorems . . . . .	126
10.3.3	The Hoffman-Wielandt Theorem . . . . .	128
10.4	Proof of the Real Schur Form . . . . .	128
<b>11</b>	<b>The Singular Value Decomposition</b>	<b>131</b>
11.1	Singular Values and Singular Vectors . . . . .	131
11.1.1	SVD and SVF . . . . .	131
11.1.2	Examples . . . . .	134
11.1.3	Singular Values of Normal and Positive Semidefinite Matrices . . . . .	137
11.1.4	A Geometric Interpretation . . . . .	137
11.2	Singular Vectors . . . . .	138
11.2.1	The SVD of $\mathbf{A}^H \mathbf{A}$ and $\mathbf{A} \mathbf{A}^H$ . . . . .	139
11.3	The Pseudo-Inverse and Orthogonal Projections . . . . .	140
11.3.1	The Pseudo-Inverse . . . . .	140
11.3.2	Orthogonal Projections . . . . .	142
11.4	The Minmax Theorem for Singular Values and the Hoffman- Wielandt Theorem . . . . .	144
<b>12</b>	<b>Matrix Norms</b>	<b>147</b>
12.1	Matrix Norms . . . . .	147
12.1.1	The Frobenius Norm . . . . .	147
12.1.2	Consistent and Subordinate Matrix Norms . . . . .	149
12.1.3	Operator Norms . . . . .	150

12.1.4	The $p$ -Norms . . . . .	151
12.1.5	Unitary Invariant Matrix Norms . . . . .	153
12.1.6	Absolute and Monotone Norms . . . . .	155
12.2	The Condition Number with Respect to Inversion . . . . .	155
12.3	Determining the Rank of a Matrix . . . . .	159
12.4	Convergence and Spectral Radius . . . . .	160
12.4.1	Convergence in $\mathbb{R}^{m,n}$ and $\mathbb{C}^{m,n}$ . . . . .	160
12.4.2	The Spectral Radius . . . . .	161
12.4.3	Neumann Series . . . . .	163
<b>IV</b>	<b>Iterative Methods for Large Linear Systems</b>	<b>165</b>
<b>13</b>	<b>The Classical Iterative Methods</b>	<b>167</b>
13.1	Classical Iterative Methods; Component Form . . . . .	167
13.2	The Discrete Poisson System . . . . .	169
13.3	Matrix Formulations of the Classical Methods . . . . .	171
13.3.1	The Splitting Matrices for the Classical Methods . . . . .	171
13.4	Convergence of Fixed-point Iteration . . . . .	173
13.4.1	Stopping the Iteration . . . . .	175
13.4.2	Richardson's Method (R method) . . . . .	175
13.5	Convergence of the Classical Methods for the Discrete Poisson Matrix . . . . .	176
13.5.1	Number of Iterations . . . . .	178
13.6	Convergence Analysis for SOR . . . . .	179
<b>14</b>	<b>The Conjugate Gradient Method</b>	<b>185</b>
14.1	The Conjugate Gradient Algorithm . . . . .	186
14.2	Numerical Example . . . . .	187
14.3	Derivation and Basic Properties . . . . .	189
14.4	Convergence . . . . .	193
<b>15</b>	<b>Minimization and Preconditioning</b>	<b>199</b>
15.1	Minimization . . . . .	199
15.2	Preconditioning . . . . .	201
15.3	Preconditioning Example . . . . .	204
15.3.1	A Banded Matrix . . . . .	204
15.3.2	Preconditioning . . . . .	207
<b>V</b>	<b>Orthonormal Transformations and Least Squares</b>	<b>209</b>
<b>16</b>	<b>Orthonormal Transformations</b>	<b>211</b>
16.1	The QR Decomposition and QR Factorization. . . . .	211
16.1.1	QR and Gram-Schmidt . . . . .	213
16.2	The Householder Transformation . . . . .	214

16.3	Householder Triangulation . . . . .	217
16.4	Givens Rotations . . . . .	218
<b>17</b>	<b>Least Squares</b>	<b>221</b>
17.1	Examples . . . . .	222
17.2	Numerical Solution using the Normal Equations . . . . .	225
17.3	Numerical Solution using the QR Factorization . . . . .	226
17.3.1	QR and Linear Systems . . . . .	228
17.4	Numerical Solution using the Singular Value Factorization . . . . .	228
17.5	Perturbation Theory for Least Squares . . . . .	229
17.5.1	Perturbing the Right Hand Side . . . . .	229
17.5.2	Perturbing the Matrix . . . . .	231
17.6	Perturbation Theory for Singular Values . . . . .	232
<b>VI</b>	<b>Eigenvalues and Eigenvectors</b>	<b>235</b>
<b>18</b>	<b>Numerical Eigenvalue Problems</b>	<b>237</b>
18.1	Perturbation of Eigenvalues . . . . .	237
18.1.1	Gerschgorin's Theorem . . . . .	239
18.2	Unitary Similarity Transformation of a Matrix into Upper Hessenberg or Tridiagonal Form . . . . .	242
18.3	Computing a Selected Eigenvalue of a Symmetric Matrix . . . . .	244
18.3.1	The Inertia Theorem . . . . .	246
18.3.2	Approximating $\lambda_m$ . . . . .	248
18.4	Perturbation Proofs . . . . .	249
<b>19</b>	<b>The Power and QR Methods</b>	<b>251</b>
19.1	The Power Method . . . . .	251
19.1.1	The Inverse Power Method . . . . .	254
19.2	The QR Algorithm . . . . .	255
19.2.1	The Relation to the Power Method . . . . .	257
19.2.2	A convergence theorem . . . . .	258
19.2.3	The Shifted QR Algorithms . . . . .	259
<b>VII</b>	<b>Appendix</b>	<b>261</b>
<b>A</b>	<b>Gaussian Elimination</b>	<b>263</b>
A.1	Gaussian Elimination and LU factorization . . . . .	264
A.1.1	Algoritms . . . . .	266
A.1.2	Operation Count . . . . .	268
A.2	Pivoting . . . . .	269
A.2.1	Permutation Matrices . . . . .	270
A.2.2	Gaussian Elimination Works Mathematically . . . . .	270
A.2.3	Pivot Strategies . . . . .	271

---

A.3	The PLU Factorization . . . . .	272
A.4	An Algorithm for Finding the PLU Factorization . . . . .	273
<b>B</b>	<b>Computer Arithmetic</b>	<b>277</b>
B.1	Absolute and Relative Errors . . . . .	277
B.2	Floating Point Numbers . . . . .	278
B.3	Rounding and Arithmetic Operations . . . . .	280
B.3.1	Rounding . . . . .	281
B.3.2	Arithmetic Operations . . . . .	281
B.4	Backward Rounding-Error Analysis . . . . .	281
B.4.1	Computing a Sum . . . . .	282
B.4.2	Computing an Inner Product . . . . .	284
B.4.3	Computing a Matrix Product . . . . .	285
<b>C</b>	<b>Differentiation of Vector Functions</b>	<b>287</b>
<b>D</b>	<b>Some Inequalities</b>	<b>291</b>
D.1	Convexity . . . . .	291
D.2	Inequalities . . . . .	292
<b>E</b>	<b>The Jordan Form</b>	<b>295</b>
E.1	The Jordan Form . . . . .	295
E.1.1	The Minimal Polynomial . . . . .	298
	<b>Bibliography</b>	<b>301</b>
	<b>Index</b>	<b>303</b>



# Preface

These lecture notes are the result of an ongoing project to write a text for a course in matrix analysis and numerical linear algebra given at the advanced undergraduate and beginning graduate level at the University of Oslo. In the first 5 chapters we give a quick review of basis linear algebra. I usually start with Chapter 6 and then each of the remaining Chapters 7-19 should correspond approximately to three hours of weekly lectures.

Oslo, 12 August, 2009

Tom Lyche



# List of Figures

2.1	The orthogonal projection of $\mathbf{x}$ into $\mathcal{S}$ . . . . .	23
4.12	The triangle $T$ defined by the three points $P_1$ , $P_2$ and $P_3$ . . . . .	47
6.1	A physical spline with ducks. . . . .	62
6.2	A two piece cubic spline interpolant to $f(x) = x^4$ . . . . .	63
6.3	Cubic spline interpolation to the data in Example 6.6. The break-points $(x_i, y_i)$ , $i = 2, 3, 4$ are marked with dots on the curve. . . .	66
6.4	A 4 piece periodic cubic spline interpolant to the unit circle. . . .	74
8.1	Numbering of grid points . . . . .	99
8.2	The 5-point stencil . . . . .	99
8.3	Band structure of the 2D test matrix, $n = 9$ , $n = 25$ , $n = 100$ . . .	100
9.1	Fill-inn in the Cholesky factor of the Poisson matrix ( $n = 100$ ). . .	110
11.1	The ellipse $y_1^2/9 + y_2^2 = 1$ (left) and the rotated ellipse $\mathbf{A}\mathbf{S}$ (right). .	138
13.1	$\rho(\mathbf{G}_\omega)$ with $\omega \in [0, 2]$ for $n = 100$ , (lower curve) and $n = 2500$ (upper curve). . . . .	177
14.10	Orthogonality in the conjugate gradient algorithm. . . . .	190
16.1	The Householder transformation . . . . .	215
16.2	A plane rotation. . . . .	219
17.1	A least squares fit to data. . . . .	224
17.2	Graphical interpretation of the bounds in Theorem 17.13. . . . .	230
18.1	The Gerschgorin disk $R_i$ . . . . .	240
A.1	Gaussian elimination . . . . .	264
B.1	Distribution of some positive floating-point numbers . . . . .	279
D.1	A convex function. . . . .	291



# List of Tables

13.1	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 $10^{-8}$ . . . . .	170
13.2	Spectral radii for $\mathbf{G}_J$ , $\mathbf{G}_1$ , $\mathbf{G}_{\omega^*}$ and the smallest integer $k_n$ such that $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$ . . . . .	178
14.6	The number of iterations $K$ for the averaging problem on a $\sqrt{n} \times \sqrt{n}$ grid for various $n$ . . . . .	188
14.7	The number of iterations $K$ for the Poisson problem on a $\sqrt{n} \times \sqrt{n}$ grid for various $n$ . . . . .	188
15.2	The number of iterations $K$ (no preconditioning) and $K_{pre}$ (with preconditioning) for the problem (15.14) using the discrete Poisson problem as a preconditioner. . . . .	207
19.7	Quadratic convergence of Rayleigh quotient iteration. . . . .	256



# List of Algorithms

6.2	findsubintervals . . . . .	64
6.3	cubppeval . . . . .	64
6.12	trifactor . . . . .	70
6.13	trisolve . . . . .	71
7.37	bandcholesky . . . . .	93
7.38	bandforwardsolve . . . . .	94
7.39	bandbacksolve . . . . .	94
9.1	Fast Poisson Solver . . . . .	112
9.4	Recursive FFT . . . . .	116
13.1	Jacobi . . . . .	169
13.2	SOR . . . . .	170
14.4	Conjugate Gradient Iteration . . . . .	187
14.5	Testing Conjugate Gradient . . . . .	188
15.3	Preconditioned Conjugate Gradient Algorithm . . . . .	203
16.12	Generate a Householder transformation . . . . .	216
16.20	Upper Hessenberg linear system . . . . .	220
17.9	Solving least squares by Householder $QR$ . . . . .	227
18.11	Householder reduction to Hessenberg form . . . . .	243
18.13	Assemble Householder transformations . . . . .	244
19.3	The Power Method . . . . .	253
19.5	Rayleigh quotient iteration . . . . .	255
A.5	lufactor . . . . .	267
A.6	forwardsolve . . . . .	267
A.7	backsolve . . . . .	268
A.12	PLU factorization . . . . .	274
A.14	Forward Substitution (column oriented) . . . . .	275
A.15	Backward Substitution (column oriented) . . . . .	275





# List of Exercises

2.9	.....	9
2.10	.....	9
2.11	.....	9
2.12	.....	9
2.21	.....	12
2.22	.....	12
2.23	.....	12
2.24	.....	12
2.26	.....	12
2.30	.....	14
2.35	.....	16
2.36	.....	16
2.37	.....	16
2.38	.....	16
2.43	.....	18
2.44	.....	18
2.50	.....	21
2.51	.....	21
2.52	.....	21
2.53	.....	21
2.54	.....	21
2.63	.....	24
3.4	.....	27
3.11	.....	30
3.12	.....	30
3.13	.....	30
3.18	.....	32
3.19	.....	32
3.20	.....	32
3.29	.....	35
4.1	.....	39
4.3	.....	43
4.5	.....	43
4.10	.....	46
4.11	.....	46

---

4.13	.....	46
4.14	.....	47
4.15	.....	48
5.8	.....	52
5.9	.....	52
5.10	.....	52
5.11	.....	52
5.12	.....	52
5.13	.....	52
6.5	.....	66
6.7	.....	67
6.9	.....	67
6.10	.....	67
6.11	.....	68
6.17	.....	72
6.18	.....	72
6.20	.....	73
6.21	.....	75
6.22	.....	75
6.23	.....	75
6.24	.....	75
6.25	.....	78
6.26	.....	78
6.27	.....	78
6.28	.....	79
6.29	.....	79
6.30	.....	79
7.11	.....	83
7.12	.....	84
7.13	.....	84
7.14	.....	84
7.18	.....	85
7.31	.....	89
8.2	.....	99
8.5	.....	102
8.14	.....	106
8.15	.....	106
8.16	.....	107
8.17	.....	107
8.18	.....	107
8.19	.....	107
9.1	.....	117
9.2	.....	117
9.3	.....	117
9.4	.....	117
9.5	.....	117

---

9.6	.....	118
9.7	.....	118
9.8	.....	118
9.9	.....	118
9.10	.....	118
10.3	.....	122
10.9	.....	125
10.10	.....	125
10.12	.....	125
10.15	.....	127
10.17	.....	128
10.19	.....	128
11.9	.....	136
11.10	.....	136
11.16	.....	139
11.19	.....	140
11.20	.....	141
11.21	.....	141
11.22	.....	141
11.23	.....	141
11.24	.....	141
11.25	.....	141
11.26	.....	141
11.27	.....	142
11.28	.....	142
11.29	.....	142
11.32	.....	143
11.33	.....	143
11.34	.....	144
12.6	.....	149
12.7	.....	149
12.8	.....	149
12.10	.....	149
12.11	.....	150
12.17	.....	153
12.18	.....	153
12.22	.....	154
12.23	.....	154
12.24	.....	154
12.25	.....	154
12.26	.....	154
12.27	.....	154
12.28	.....	155
12.29	.....	155
12.32	.....	156
12.36	.....	158

---

12.38	.....	159
12.39	.....	160
12.46	.....	162
12.48	.....	164
12.49	.....	164
13.9	.....	174
13.10	.....	174
13.11	.....	174
13.12	.....	174
13.14	.....	174
13.17	.....	176
13.20	.....	179
14.2	.....	186
14.3	.....	186
14.8	.....	189
14.12	.....	191
14.13	.....	192
14.14	.....	192
14.20	.....	195
15.1	.....	201
15.2	.....	201
16.5	.....	213
16.6	.....	214
16.9	.....	215
16.13	.....	216
16.14	.....	216
16.15	.....	217
16.16	.....	217
16.18	.....	219
16.19	.....	219
17.5	.....	224
17.6	.....	224
17.7	.....	225
17.15	.....	231
17.16	.....	231
17.18	.....	232
18.7	.....	241
18.9	.....	241
18.10	.....	242
18.12	.....	243
18.14	.....	244
18.15	.....	244
18.19	.....	247
18.20	.....	247
18.21	.....	247
18.22	.....	248

---

18.23	.....	248
18.25	.....	249
19.12	.....	259
A.2	.....	265
A.13	.....	274
B.4	.....	278
B.6	.....	280
B.8	.....	281
C.1	.....	287
E.3	.....	296
E.4	.....	296
E.6	.....	297
E.7	.....	297
E.9	.....	298
E.10	.....	299
E.11	.....	299



**Part I**

**A Review of Linear Algebra**





## Chapter 1

# Introduction

### 1.1 Notation

The following sets will be used throughout these notes.

1. The set 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.  $\mathbb{R}^{m,n}$  is the set of  $m \times n$  matrices with real elements represented as

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

The element in the  $i$ th row and  $j$ th column of a matrix  $\mathbf{A}$  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 dots 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  $\mathbf{x}, \mathbf{y}, \mathbf{z}, \dots$ .

5. 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  $z = x + iy = re^{i\phi} = r(\cos \phi + i \sin \phi)$  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  $z$ .
6. For matrices and vectors with complex elements we use the notation  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{x} \in \mathbb{C}^n$ . We identify complex row vectors using either the transpose  $T$  or the Hermitian (conjugate) transpose operation  $\mathbf{x}^H := \bar{\mathbf{x}}^T = [\bar{x}_1, \dots, \bar{x}_n]$ .
7. 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

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

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

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

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

$$\mathbf{B} = \text{tridiag}(a_i, d_i, c_i) = \text{tridiag}(\mathbf{a}, \mathbf{d}, \mathbf{c}) := \begin{bmatrix} d_1 & c_1 & & & \\ a_2 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-1} & d_{n-1} & c_{n-1} \\ & & & a_n & d_n \end{bmatrix}.$$

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

9. Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  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  $\mathbf{A}(\mathbf{i}, \mathbf{j}) \in \mathbb{C}^{r,c}$  is the submatrix of  $\mathbf{A}$  consisting of rows  $\mathbf{i} := [i_1, \dots, i_r]$  and columns  $\mathbf{j} := [j_1, \dots, j_c]$

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

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

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



## Chapter 2

# Vectors

This chapter contains a review of vector space concepts that will be useful in this text. we start by introducing a vector space. To define a vector space we need a field  $\mathbb{F}$ , a set of vectors  $\mathcal{V}$ , a way to combine vectors called **vector addition**, and a way to combine elements of  $\mathbb{F}$  and  $\mathcal{V}$  called **scalar multiplication**. In the first part of this section  $\mathbb{F}$  will be an arbitrary field, but later the field will be the set of real or complex numbers with the usual arithmetic operations.

### 2.1 Vector Spaces

**Definition 2.1** A **field** is a set  $\mathbb{F}$  together with two operations  $+, \cdot : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$  such that for all  $a, b, c \in \mathbb{F}$  the following arithmetic rules hold

(A0) there exists an element  $0 \in \mathbb{F}$  such that  $a + 0 = a$ .

(Am) there exists an element  $(-a) \in \mathbb{F}$  such that  $a + (-a) = 0$ . We define subtraction as  $a - b := a + (-b)$ .

(Aa)  $a + (b + c) = (a + b) + c$ .

(Ac)  $a + b = b + a$ .

(M1) there exists an element  $1 \in \mathbb{F}$  such that  $a \cdot 1 = a$ .

(Mi) if  $a \neq 0$  then there exists an element  $a^{-1} \in \mathbb{F}$  such that  $a \cdot a^{-1} = 1$ .

(Ma)  $a \cdot (b \cdot c) = (a \cdot b) \cdot c$ .

(Mc)  $a \cdot b = b \cdot a$ .

(D)  $a \cdot (b + c) = a \cdot b + a \cdot c$ .

The requirements (A0), (Am), (Aa) are the axioms for a group. They state that  $(\mathbb{F}, +)$  is a **group**, and since in addition (Ac) holds then  $(\mathbb{F}, +)$  is by definition an

**abelian group.** The axioms (M1), (Mi), (Ma), (Mc) state that  $(\mathbb{F} \setminus \{0\}, \cdot)$  is an abelian group. Often we drop the dot and write  $ab$  for the product  $a \cdot b$ . Examples of fields are  $\mathbb{R}$  or  $\mathbb{C}$  with ordinary addition and multiplication.

**Definition 2.2** A **vector space** over a field  $\mathbb{F}$  is a set  $\mathcal{V}$  together with two operations vector addition,  $+: \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$  and scalar multiplication,  $\cdot: \mathbb{F} \times \mathcal{V} \rightarrow \mathcal{V}$  such that for all  $a, b \in \mathbb{F}$  and  $\mathbf{v}, \mathbf{w} \in \mathcal{V}$  the following hold

(V)  $(\mathcal{V}, +)$  is an abelian group.

(Va)  $(a \cdot b) \cdot \mathbf{v} = a \cdot (b \cdot \mathbf{v})$ .

(Vd1)  $(a + b) \cdot \mathbf{v} = a \cdot \mathbf{v} + b \cdot \mathbf{v}$ .

(Vd2)  $a \cdot (\mathbf{v} + \mathbf{w}) = a \cdot \mathbf{v} + a \cdot \mathbf{w}$ .

(M1)  $1 \cdot \mathbf{v} = \mathbf{v}$ .

We denote a vector space by  $(\mathcal{V}, \mathbb{F})$  or by  $\mathcal{V}$  if the underlying field is clear from the context.

**Definition 2.3** Let  $(\mathcal{V}, \mathbb{F})$  be a vector space and  $\mathcal{S}$  a nonempty subset of  $\mathcal{V}$ . Then  $(\mathcal{S}, \mathbb{F})$  is a **subspace** of  $(\mathcal{V}, \mathbb{F})$  if  $(\mathcal{S}, \mathbb{F})$  is itself a vector space.

It follows that  $(\mathcal{S}, \mathbb{F})$  is a subspace of  $(\mathcal{V}, \mathbb{F})$  if  $\mathcal{S}$  is closed under vector addition and scalar multiplication, i.e.  $as_1 + bs_2 \in \mathcal{S}$  for all  $a, b \in \mathbb{F}$  and all  $s_1, s_2 \in \mathcal{S}$ . For any vector space  $(\mathcal{V}, \mathbb{F})$  the two sets  $\{\mathbf{0}\}$ , consisting only of the zero element in  $\mathcal{V}$ , and  $\mathcal{V}$  itself are subspaces. They are called the **trivial subspaces**.

Here are some examples of vector spaces.

**Example 2.4 (The Vector Spaces  $\mathbb{R}^n$  and  $\mathbb{C}^n$ )** In the following chapters we will deal almost exclusively with the vector spaces  $\mathbb{R}^n = (\mathbb{R}^n, \mathbb{R})$ ,  $\mathbb{C}^n = (\mathbb{C}^n, \mathbb{C})$  and their subspaces. Addition and scalar multiplication are defined by

$$\mathbf{v} + \mathbf{w} = \begin{bmatrix} v_1 + w_1 \\ \vdots \\ v_n + w_n \end{bmatrix}, \quad a\mathbf{v} = \begin{bmatrix} av_1 \\ \vdots \\ av_n \end{bmatrix}.$$

**Example 2.5 (Subspaces of  $\mathbb{R}^2$  and  $\mathbb{R}^3$ )** For a given vector  $\mathbf{x} \in \mathbb{R}^n$  let  $\mathcal{S} = \{t\mathbf{x} : t \in \mathbb{R}\}$ . Then  $\mathcal{S}$  is a subspace of  $\mathbb{R}^n$ , in fact it represents a straight line passing through the origin. For  $n = 2$  it can be shown that all nontrivial subspaces of  $\mathbb{R}^2$  are of this form. For  $n = 3$  the nontrivial subspaces are all lines and all planes containing  $\{\mathbf{0}\}$ .

**Example 2.6 (The Vector Space  $C(I)$ )** Let  $\mathbb{F} = \mathbb{R}$  and let  $C(I)$  be the set of all real valued functions  $f : I \rightarrow \mathbb{R}$  which are defined and continuous on an interval  $I \subset \mathbb{R}$ . Here the vectors are functions in  $C(I)$ . Vector addition and scalar multiplication are defined for all  $f, g \in C(I)$  and all  $a \in \mathbb{R}$  by

$$(f + g)(x) := f(x) + g(x), \quad (af)(x) := af(x), \quad \text{for all } x \in I.$$

$C(I) = (C(I), \mathbb{R})$  is a vector space since

- the sum of two continuous functions is continuous,
- a constant times a continuous function is continuous
- vector addition and scalar multiplication are defined point-wise, so the axioms for a vector space follows from properties of real numbers.

**Example 2.7 (The Vector Space  $\Pi_n$ )** Let  $\Pi_n(I)$  be the set of all polynomials of degree at most  $n$  defined on a subset  $I \subset \mathbb{R}$  or  $I \subset \mathbb{C}$ . We write simply  $\Pi_n$  if  $I = \mathbb{R}$  or  $I = \mathbb{C}$ . With pointwise addition and scalar multiplication defined as in Example 2.6 the set  $(\Pi_n(I), \mathbb{R})$  is a subspace of  $(C(I), \mathbb{R})$ .

**Definition 2.8 (Linear Combinations)** The sum  $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \cdots + c_n\mathbf{v}_n$  with  $c_i \in \mathbb{F}$  and  $\mathbf{v}_i \in \mathcal{V}$  for  $i = 1, \dots, n$  is called a **linear combination** of  $\mathbf{v}_1, \dots, \mathbf{v}_n$ . We say that the linear combination is nontrivial if at least one of the  $c_i$ 's is nonzero. The set

$$\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\} := \{c_1\mathbf{v}_1 + \cdots + c_n\mathbf{v}_n : c_i \in \mathbb{F}, i = 1, \dots, n\}$$

spanned by  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathcal{V}$  is a subspace of  $(\mathcal{V}, \mathbb{F})$ . A vector space  $\mathcal{V}$  is called **finite dimensional** if it has a finite spanning set; i.e. there exist  $n \in \mathbb{N}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  in  $\mathcal{V}$  such that  $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ .

**Exercise 2.9** Show that the  $\mathbf{0}$  of vector addition is unique and that  $\{\mathbf{0}\}$  is a subspace.

**Exercise 2.10** Show that  $0 \cdot \mathbf{x} = \mathbf{0}$  for any  $\mathbf{x} \in \mathcal{V}$ .

**Exercise 2.11** Show that  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a subspace.

**Exercise 2.12** Show that  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is the smallest subspace containing the vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$ .

## 2.2 Linear Independence and Bases

**Definition 2.13** Let  $\mathcal{X} := \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  be a set of vectors in a vector space  $(\mathcal{V}, \mathbb{F})$ . We say that  $\mathcal{X}$  is **linearly dependent** if we can find a nontrivial linear combination which is equal to zero. We say that  $\mathcal{X}$  is **linearly independent** if it is not linearly dependent. In other words

$$c_1\mathbf{v}_1 + \dots + c_n\mathbf{v}_n = \mathbf{0} \text{ for some } c_1, \dots, c_n \in \mathbb{F} \implies c_1 = \dots = c_n = 0.$$

The elements in a set of linearly independent vectors must all be nonzero and we have

**Lemma 2.14** Suppose  $\mathbf{v}_1, \dots, \mathbf{v}_n$  span a vector space  $V$  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  $\mathbf{w}_1 = c_1\mathbf{v}_1 + \dots + c_n\mathbf{v}_n$ . 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  $i_2$ . Moreover  $i_2 \neq i_1$  for otherwise  $\mathbf{w}_2 = d_{i_1}\mathbf{w}_1$  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 all the  $\mathbf{v}$ 's 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 2.15** A finite set of vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  in a vector space  $(\mathcal{V}, \mathbb{F})$  is a **basis** for  $(\mathcal{V}, \mathbb{F})$  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 2.16** Suppose  $(\mathcal{V}, \mathbb{F})$  is a vector space and that  $S := \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a spanning set for  $\mathcal{V}$ . Then we can find a subset  $\{\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}\}$  of  $S$  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 2.17** 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 2.16  $\mathcal{V}$  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 2.18** *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 2.14 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  $\mathbb{C}^n$ . The dimension of the trivial subspace  $\{\mathbf{0}\}$  is defined to be zero.

**Theorem 2.19** *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.  $\square$

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

**Lemma 2.20** *Suppose  $\mathcal{S}$  is a subspace of a finite dimensional vector space  $(\mathcal{V}, \mathbb{F})$  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*

$$\mathbf{s}_j = \sum_{i=1}^m a_{ij} \mathbf{v}_i \text{ for } j = 1, \dots, n. \quad (2.1)$$

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  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{m,n}$ . The matrix  $\mathbf{A}$  has linearly independent columns.

**Proof.** (2.1) holds 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

$$\mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j \stackrel{(2.1)}{=} \sum_{j=1}^n c_j \left( \sum_{i=1}^m a_{ij} \mathbf{v}_i \right) = \sum_{i=1}^m \left( \sum_{j=1}^n a_{ij} c_j \right) \mathbf{v}_i = \sum_{i=1}^m b_i \mathbf{v}_i.$$

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 we have  $\mathbf{c} = \mathbf{0}$ .  $\square$

The matrix  $\mathbf{A}$  in Lemma 2.20 is called a **change of basis matrix**.

**Exercise 2.21** Show that the elements in a linearly independent set must be nonzero.

**Exercise 2.22** Show that the set of unit vectors  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  form a basis both for  $\mathbb{R}^n$  and for  $\mathbb{C}^n$ . Why does this show that the dimension of  $\mathbb{R}^n$  and  $\mathbb{C}^n$  is  $n$ ?

## 2.3 Operations on Subspaces

Let  $\mathcal{R}$  and  $\mathcal{S}$  be two subsets of a vector space  $(\mathcal{V}, \mathbb{F})$  and let  $a$  be a scalar. The **sum**, **multiplication by scalar**, **union**, and **intersection** of  $\mathcal{R}$  and  $\mathcal{S}$  are defined by

$$\mathcal{R} + \mathcal{S} := \{\mathbf{r} + \mathbf{s} : \mathbf{r} \in \mathcal{R} \text{ and } \mathbf{s} \in \mathcal{S}\}, \quad (2.2)$$

$$a\mathcal{S} := \{a\mathbf{s} : \mathbf{s} \in \mathcal{S}\}, \quad (2.3)$$

$$\mathcal{R} \cup \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ or } \mathbf{x} \in \mathcal{S}\}. \quad (2.4)$$

$$\mathcal{R} \cap \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ and } \mathbf{x} \in \mathcal{S}\}. \quad (2.5)$$

**Exercise 2.23** Let  $\mathcal{R} = \{(x, y) : x^2 + y^2 \leq 1\}$  be the unit disc in  $\mathbb{R}^2$  and set  $\mathcal{S} = \{(x, y) : (x - \frac{1}{2})^2 + y^2 \leq 1\}$ . Find  $\mathcal{R} + \mathcal{S}$ ,  $2\mathcal{S}$ ,  $\mathcal{R} \cup \mathcal{S}$ , and  $\mathcal{R} \cap \mathcal{S}$ .

### 2.3.1 Sums and intersections of subspaces

In many cases  $\mathcal{R}$  and  $\mathcal{S}$  will be subspaces. Then  $a\mathcal{S} = \mathcal{S}$  and both the sum and intersection of two subspaces is a subspace of  $(\mathcal{V}, \mathbb{F})$ . Note however that the union  $\mathcal{R} \cup \mathcal{S}$  of two subspaces is not necessarily a subspace.

**Exercise 2.24** Let  $\mathcal{R}$  and  $\mathcal{S}$  be two subspaces of a vector space  $(\mathcal{V}, \mathbb{F})$ . Show that  $a\mathcal{S} = \mathcal{S}$  and that both  $\mathcal{R} + \mathcal{S}$  and  $\mathcal{R} \cap \mathcal{S}$  are subspaces of  $(\mathcal{V}, \mathbb{F})$ .

**Example 2.25** For given vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $\mathbf{x}$  and  $\mathbf{y}$  linearly independent let  $\mathcal{R} = \text{span}\{\mathbf{x}\}$  and  $\mathcal{S} = \text{span}\{\mathbf{y}\}$ . Then  $\mathcal{R}$  and  $\mathcal{S}$  are subspaces of  $\mathbb{R}^n$ . For  $n = 2$  we have  $\mathcal{R} + \mathcal{S} = \mathbb{R}^2$ , while for  $n = 3$  the sum represents a plane passing through the origin. We also see that  $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$  and that  $\mathcal{R} \cup \mathcal{S}$  is not a subspace.

**Exercise 2.26** Show the statements made in Example 2.25.

**Theorem 2.27** Let  $\mathcal{R}$  and  $\mathcal{S}$  be two subspaces of a vector space  $(\mathcal{V}, \mathbb{F})$ . Then

$$\dim(\mathcal{R} + \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S}) - \dim(\mathcal{R} \cap \mathcal{S}). \quad (2.6)$$

**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 2.19 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  $\mathbf{r} := \sum_{j=1}^p \beta_j \mathbf{u}_j$  and at the same time as a linear combination of  $\mathbf{r}_1, \dots, \mathbf{r}_q$ . 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 now 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  $\mathcal{R} + \mathcal{S}$ . The result now follows from a simple calculation

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

□

From this theorem it follows that  $\dim(\mathcal{R} + \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S})$  provided  $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$ .

**Definition 2.28 (Direct Sum)** Let  $\mathcal{R}$  and  $\mathcal{S}$  be two subspaces of a vector space  $(\mathcal{V}, \mathbb{F})$ . If  $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$  then the subspace  $\mathcal{R} + \mathcal{S}$  is called a **direct sum** and denoted  $\mathcal{R} \oplus \mathcal{S}$ . The subspaces  $\mathcal{R}$  and  $\mathcal{S}$  are called **complementary** in the subspace  $\mathcal{R} \oplus \mathcal{S}$ .

**Theorem 2.29** Let  $\mathcal{R}$  and  $\mathcal{S}$  be two subspaces of a vector space  $(\mathcal{V}, \mathbb{F})$  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. □

### 2.3.2 The quotient space

For the sum of two sets we write  $\mathbf{x} + \mathcal{S} := \{\mathbf{x} + \mathbf{s} : \mathbf{s} \in \mathcal{S}\}$  when one of the sets is a singleton set  $\{\mathbf{x}\}$ . Suppose  $\mathcal{S}$  is a subspace of a vector space  $(\mathcal{X}, \mathbb{F})$ . Since  $a\mathcal{S} = \mathcal{S}$  we have

$$a(\mathbf{x} + \mathcal{S}) + b(\mathbf{y} + \mathcal{S}) = (a\mathbf{x} + b\mathbf{y}) + \mathcal{S}, \text{ for all } a, b \in \mathbb{F} \text{ and all } \mathbf{x}, \mathbf{y} \in \mathcal{S}.$$

The set

$$\mathcal{X}/\mathcal{S} := \{\mathbf{x} + \mathcal{S} : \mathbf{x} \in \mathcal{X}\} \quad (2.7)$$

is a vector space if we define

$$a(\mathbf{x} + \mathcal{S}) + b(\mathbf{y} + \mathcal{S}) := (a\mathbf{x} + b\mathbf{y}) + \mathcal{S}, \text{ for all } a, b \in \mathbb{F} \text{ and all } \mathbf{x}, \mathbf{y} \in \mathcal{S}.$$

The space  $\mathcal{X}/\mathcal{S}$  is called the **quotient space** of  $\mathcal{X}$  by  $\mathcal{S}$ . The zero element in  $\mathcal{X}/\mathcal{S}$  is  $\mathcal{S}$  itself. Moreover, if  $\mathbf{x} + \mathcal{S} = \mathbf{y} + \mathcal{S}$  then  $\mathbf{x} - \mathbf{y} \in \mathcal{S}$ .

**Exercise 2.30** Show that  $\mathcal{X}/\mathcal{S}$  is a vector space.

**Theorem 2.31** Suppose  $\mathcal{S}$  is a subspace of a finite dimensional vector space  $(\mathcal{X}, \mathbb{F})$ . Then

$$\dim(\mathcal{S}) + \dim(\mathcal{X}/\mathcal{S}) = \dim(\mathcal{X}). \quad (2.8)$$

**Proof.** Let  $n := \dim(\mathcal{X})$ ,  $k = \dim(\mathcal{S})$ , and let  $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$  be a basis for  $\mathcal{S}$ . By Theorem 2.19 we can extend it to a basis  $\{\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{t}_{k+1}, \dots, \mathbf{t}_n\}$  for  $\mathcal{X}$ . The result will follow if we can show that  $\{\mathbf{t}_{k+1} + \mathcal{S}, \dots, \mathbf{t}_n + \mathcal{S}\}$  is a basis for  $\mathcal{X}/\mathcal{S}$ . Recall that the zero element in  $\mathcal{X}/\mathcal{S}$  is  $\mathcal{S}$ . To show linear independence suppose  $\sum_{j=k+1}^n a_j(\mathbf{t}_j + \mathcal{S}) = \mathcal{S}$  for some  $a_{k+1}, \dots, a_n$  in  $\mathbb{F}$ . Since  $\sum_{j=k+1}^n a_j\mathcal{S} = \mathcal{S}$  and the zero element in  $\mathcal{X}/\mathcal{S}$  is unique we must have  $\sum_{j=k+1}^n a_j\mathbf{t}_j = \mathbf{0}$  which implies that  $a_{k+1} = \dots = a_n = 0$  by linear independence of the  $\mathbf{t}$ 's. It remains to show that  $\text{span}\{\mathbf{t}_{k+1} + \mathcal{S}, \dots, \mathbf{t}_n + \mathcal{S}\} = \mathcal{X}/\mathcal{S}$ . Suppose  $\mathbf{x} + \mathcal{S} \in \mathcal{X}/\mathcal{S}$ . For some  $a_1, \dots, a_n$  we have  $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$ , where  $\mathbf{x}_1 = \sum_{j=1}^k a_j\mathbf{s}_j$  and  $\mathbf{x}_2 = \sum_{j=k+1}^n a_j\mathbf{t}_j$ . Since  $\mathbf{x}_1 + \mathcal{S} = \mathcal{S}$  we have  $\mathbf{x} + \mathcal{S} = \mathbf{x}_2 + \mathcal{S} = \sum_{j=k+1}^n a_j\mathbf{t}_j + \mathcal{S} = \sum_{j=k+1}^n a_j(\mathbf{t}_j + \mathcal{S}) \in \mathcal{X}/\mathcal{S}$ .  $\square$

## 2.4 Norms

To measure the size of a vector in a vector space  $(\mathcal{V}, \mathbb{F})$  we use norms.

**Definition 2.32 (Norm)** A norm in a vector space  $(\mathcal{V}, \mathbb{F})$ , where  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ , is a function  $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$  that satisfies for all  $\mathbf{x}, \mathbf{y}$  in  $\mathcal{V}$  and all  $a$  in  $\mathbb{F}$

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  $(\mathcal{V}, \mathbb{F}, \|\cdot\|)$  is said to be a **normed vector space** and the inequality 3. is called the **triangle inequality**.

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

$$\|\mathbf{x} - \mathbf{y}\| \geq |\|\mathbf{x}\| - \|\mathbf{y}\||, \quad \mathbf{x}, \mathbf{y} \in \mathcal{V}. \quad (2.9)$$

Consider now some specific vector norms. For the vector spaces  $\mathbb{R}^n$  and  $\mathbb{C}^n$  we define for  $p \geq 1$  the  $p$ -norms by

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad (2.10)$$

$$\|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|. \quad (2.11)$$

The most important cases are:

1.  $\|\mathbf{x}\|_1 = \sum_{j=1}^n |x_j|$ , (**the one-norm** or  **$l_1$ -norm**)
2.  $\|\mathbf{x}\|_2 = \left( \sum_{j=1}^n |x_j|^2 \right)^{1/2}$ , **the two-norm**,  **$l_2$ -norm**, or **Euclidian norm**)
3.  $\|\mathbf{x}\|_\infty = \max_{1 \leq j \leq n} |x_j|$ , (**the infinity-norm**,  **$l_\infty$ -norm**, or **max norm**.)

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

$$\lim_{p \rightarrow \infty} \|\mathbf{x}\|_p = \|\mathbf{x}\|_\infty \text{ for all } \mathbf{x} \in \mathbb{C}^n. \quad (2.12)$$

This 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

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_p \leq n^{1/p} \|\mathbf{x}\|_\infty, \quad p \geq 1. \quad (2.13)$$

Since  $\lim_{p \rightarrow \infty} n^{1/p} = 1$  for any  $n \in \mathbb{N}$  we see that (2.12) follows.

It can be shown (cf. Appendix D) that the  $p$ -norm are norms in  $\mathbb{R}^n$  and in  $\mathbb{C}^n$  for any  $p$  with  $1 \leq p \leq \infty$ . The triangle inequality  $\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$  is called **Minkowski's inequality**. To prove it one first establishes **Hölder's inequality**

$$\sum_{j=1}^n |x_j y_j| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1. \quad (2.14)$$

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

(2.13) shows that the infinity norm and any other  $p$ -norm can be bounded in terms of each other. We define

**Definition 2.33** Two norms  $\|\cdot\|$  and  $\|\cdot\|'$  in a finite dimensional vector space  $(\mathcal{V}, \mathbb{F})$  are equivalent if there are positive constants  $m$  and  $M$  (depending only on the dimension of  $\mathcal{V}$ ) such that for all vectors  $\mathbf{x} \in \mathcal{V}$  we have

$$m\|\mathbf{x}\| \leq \|\mathbf{x}\|' \leq M\|\mathbf{x}\|. \quad (2.15)$$

The following result is proved in Appendix D.

**Theorem 2.34** All vector norms in a finite dimensional vector space are equivalent.

The inverse triangle inequality (2.9) shows that a norm is a continuous function  $\mathcal{V} \rightarrow \mathbb{R}$ .

**Exercise 2.35** Show that  $\|\cdot\|_p$  is a vector norm in  $\mathbb{R}^n$  for  $p = 1, p = \infty$ .

**Exercise 2.36** 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 2.37** 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  $\|\mathbf{x}_u\|_p = n^{1/p}\|\mathbf{x}_u\|_\infty$ . Thus the inequalities in (2.12) are sharp.

**Exercise 2.38** 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 D.2 with  $f(z) = z^{p/q}$  and  $z_i = |x_i|^q$ . For the left inequality consider first  $y_i = x_i/\|\mathbf{x}\|_\infty$ ,  $i = 1, 2, \dots, n$ .

## 2.5 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.

Consider now a more general vector space.

**Definition 2.39** Let  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  be a basis for a finite dimensional vector space  $(\mathcal{V}, \mathbb{F})$ , where  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ , 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 2.20  $\mathbf{b}_k = \mathbf{A} \mathbf{c}_k$  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

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

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  $\alpha, \beta$  such that any  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{v}_j$  satisfy

$$\|\mathbf{x}\| \leq \alpha \max_{1 \leq j \leq n} |c_j| \text{ and } |c_j| \leq \beta \|\mathbf{x}\| \text{ for } j = 1, \dots, n. \quad (2.16)$$

Suppose now  $(\mathcal{V}, \mathbb{F}, \|\cdot\|)$  is a normed vector space with  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . The notion of limit can then be stated in terms of convergence in norm.

**Theorem 2.40** *In a normed vector space 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 (2.16) 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 2.41** *Let  $(\mathcal{V}, \mathbb{F}, \|\cdot\|)$  be a normed vector space and let  $\{\mathbf{x}_k\}$  in  $\mathcal{V}$  be an infinite sequence.*

1.  $\{\mathbf{x}_k\}$  is a **Cauchy sequence** if  $\lim_{k,l \rightarrow \infty} (\mathbf{x}_k - \mathbf{x}_l) = \mathbf{0}$  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  $0 \leq n_0 < n_1 < n_2 \cdots$ .

**Theorem 2.42** 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 (2.16) 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  $\{\mathbf{x}_k\}$ .

□

### 2.5.1 Convergence of Series of Vectors

Consider now an infinite series  $\sum_{m=0}^{\infty} \mathbf{y}_m$  of vectors in a vector space  $(\mathcal{V}, \mathbb{F})$  with  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . 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 Theorem 2.34. In an absolute convergent series we may change the order of the terms without changing the value of the sum.

**Exercise 2.43** 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 2.44** Show that  $\|\cdot\|_c$  is a norm.



## 2.6 Inner Products

An **inner product** or **scalar product** in a vector space  $(\mathcal{V}, \mathbb{F})$ , where  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ , is a function  $\langle \cdot, \cdot \rangle$  mapping pairs of vectors into a scalar. We consider first the case where  $\mathbb{F} = \mathbb{R}$ .

**Definition 2.45** *An inner product in a vector space  $(\mathcal{V}, \mathbb{R})$  is a function  $\mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  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 triple  $(\mathcal{V}, \mathbb{R}, \langle \cdot, \cdot \rangle)$  is called a **real inner product space**

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  $\mathbb{R}^n$ .

When the field of scalars is  $\mathbb{C}$  the inner product is complex valued and properties 2. and 3. are altered as follows:

**Definition 2.46** *An inner product in a vector space  $(\mathcal{V}, \mathbb{C})$  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 triple  $(\mathcal{V}, \mathbb{C}, \langle \cdot, \cdot \rangle)$  is called a **complex inner product space**

Note the complex conjugate in 2. and that (Cf. Exercise 2.52)

$$\langle \mathbf{x}, a\mathbf{y} + b\mathbf{z} \rangle = \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle + \bar{b}\langle \mathbf{x}, \mathbf{z} \rangle. \quad (2.17)$$

The **standard inner product** in  $\mathbb{C}^n$  is given by  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^H \mathbf{y} = \sum_{j=1}^n \bar{x}_j y_j$ . It is clearly an inner product in  $\mathbb{C}^n$ .

Suppose now  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$  is an inner product space with  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . We define the **inner product norm** by

$$\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}, \quad \mathbf{x} \in \mathcal{V}.$$

For any vectors  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$  and scalar  $a \in \mathbb{F}$  we have (Cf. Exercises 2.51 and 2.52) by linearity and symmetry the expansion

$$\|\mathbf{x} + a\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2a\langle \mathbf{x}, \mathbf{y} \rangle + a^2\|\mathbf{y}\|^2 \quad (\text{real case}), \quad (2.18)$$

$$= \|\mathbf{x}\|^2 + 2\operatorname{Re}\langle \mathbf{x}, a\mathbf{y} \rangle + |a|^2\|\mathbf{y}\|^2 \quad (\text{complex case}), \quad (2.19)$$

where  $\operatorname{Re} z$  and  $\operatorname{Im} z$  denotes the real- and imaginary part of the complex number  $z$ .

In the complex case we can write the inner product of two vectors as a sum of inner product norms. For any  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$  it follows from (2.19) that

$$4\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2 + i\|\mathbf{x} - i\mathbf{y}\|^2 - i\|\mathbf{x} + i\mathbf{y}\|^2, \quad (2.20)$$

where  $i = \sqrt{-1}$  and we used that  $\text{Im}(z) = \text{Re}(-iz)$  for any  $z \in \mathbb{C}$ .

To show that the inner product norm is a norm in  $(\mathcal{V}, \mathbb{R})$  we need the triangle inequality. To show it we start with a famous inequality.

**Theorem 2.47 (Cauchy-Schwarz inequality)** *For any  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space*

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

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

**Proof.** The inequality is trivial if  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$  so assume  $\langle \mathbf{x}, \mathbf{y} \rangle \neq 0$ . Suppose first  $\langle \mathbf{x}, \mathbf{y} \rangle \in \mathbb{R}$ . We define the scalar  $a := -\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2}$ , and use (2.18) to obtain  $0 \leq \|\mathbf{x} + a\mathbf{y}\|^2 = \|\mathbf{x}\|^2 - (\langle \mathbf{x}, \mathbf{y} \rangle)^2 / \|\mathbf{y}\|^2$ . Thus the inequality follows in the real case. Suppose next  $\langle \mathbf{x}, \mathbf{y} \rangle$  is complex valued, say  $\langle \mathbf{x}, \mathbf{y} \rangle = re^{i\phi}$ . We define  $b := e^{-i\phi}$  and observe that  $b\langle \mathbf{x}, \mathbf{y} \rangle = r$  is real valued and  $|b| = 1$ . Using the real case of the Cauchy-Schwarz inequality we find

$$|\langle \mathbf{x}, \mathbf{y} \rangle| = |b\langle \mathbf{x}, \mathbf{y} \rangle| = |\langle b\mathbf{x}, \mathbf{y} \rangle| \leq \|b\mathbf{x}\| \|\mathbf{y}\| = \|\mathbf{x}\| \|\mathbf{y}\|$$

which proves the inequality also in the complex case. We have equality if and only if  $\mathbf{x} + a\mathbf{y} = \mathbf{0}$  which means that  $\mathbf{x}$  and  $\mathbf{y}$  are linearly dependent.  $\square$

**Theorem 2.48 (Triangle Inequality)** *For any  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space*

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$$

**Proof.** From the Cauchy-Schwarz inequality it follows that  $\text{Re}\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \|\mathbf{y}\|$ . Using this on the inner product term in (2.19) with  $a = 1$  we get

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

Taking square roots completes the proof.  $\square$

**Theorem 2.49 (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 the inner product expansion (2.19) and add the two equations.  $\square$

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  $\theta$  in  $[0, \pi]$  such that

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}. \quad (2.21)$$

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

**Exercise 2.50** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$  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 2.51** Show (2.18)

**Exercise 2.52** Show (2.17) and (2.19).

**Exercise 2.53** Show (2.20)

**Exercise 2.54** Show that in the complex case there is a unique angle  $\theta$  in  $[0, \pi/2]$  such that

$$\cos \theta = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}. \quad (2.22)$$

## 2.7 Orthogonality

As in the previous section we assume that  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$  is an inner product space with  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . Also  $\|\cdot\|$  denotes the inner product norm.

**Definition 2.55 (Orthogonality)** Two vectors  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space are called **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  $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$ .

For orthogonal vectors it follows from (2.19) that the Pythagorean theorem holds

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

**Definition 2.56 (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 called 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 an inner product space can be turned into an orthogonal- or orthonormal basis for the subspace by the following construction.

**Theorem 2.57 (Gram-Schmidt)** Let  $\{s_1, \dots, s_k\}$  be a basis for a real or complex inner product space  $(\mathcal{S}, \mathbb{F}, \langle \cdot, \cdot \rangle)$ . Define

$$v_1 := s_1, \quad v_j := s_j - \sum_{i=1}^{j-1} \frac{\langle s_j, v_i \rangle}{\langle v_i, v_i \rangle} v_i, \quad j = 2, \dots, k. \quad (2.23)$$

Then  $\{v_1, \dots, v_k\}$  is an orthogonal basis for  $\mathcal{S}$  and the normalized vectors

$$\{u_1, \dots, u_k\} := \left\{ \frac{v_1}{\|v_1\|}, \dots, \frac{v_k}{\|v_k\|} \right\}$$

is an orthonormal basis for  $\mathcal{S}$ .

**Proof.** To show that  $\{v_1, \dots, v_k\}$  is an orthogonal basis for  $\mathcal{S}$  we use induction on  $k$ . Let  $S_j := \text{span}\{s_1, \dots, s_j\}$  for  $j = 1, \dots, k$ . Clearly  $v_1 = s_1$  is an orthogonal basis for  $S_1$ . Suppose for some  $j \geq 2$  that  $v_1, \dots, v_{j-1}$  is an orthogonal basis for  $S_{j-1}$  and let  $v_j$  be given by (2.23) as a linear combination of  $s_j$  and  $v_1, \dots, v_{j-1}$ . Replacing each of these  $v_i$  by a linear combination of  $s_1, \dots, s_{j-1}$  we obtain  $v_j = \sum_{i=1}^j a_i s_i$  for some  $a_0, \dots, a_j$  with  $a_j = 1$ . Since  $s_1, \dots, s_j$  are linearly independent and  $a_j \neq 0$  we deduce that  $v_j \neq 0$ . By the induction hypothesis

$$\langle v_j, v_l \rangle = \langle s_j, v_l \rangle - \sum_{i=1}^{j-1} \frac{\langle s_j, v_i \rangle}{\langle v_i, v_i \rangle} \langle v_i, v_l \rangle = \langle s_j, v_l \rangle - \frac{\langle s_j, v_l \rangle}{\langle v_l, v_l \rangle} \langle v_l, v_l \rangle = 0$$

for  $l = 1, \dots, j-1$ . Thus  $v_1, \dots, v_j$  is an orthogonal basis for  $S_j$ .

If  $\{v_1, \dots, v_k\}$  is an orthogonal basis for  $\mathcal{S}$  then clearly  $\{u_1, \dots, u_k\}$  is an orthonormal basis for  $\mathcal{S}$ .  $\square$

**Theorem 2.58 (Orthogonal Projection)** Let  $\mathcal{S}$  be a subspace of a finite dimensional real or complex inner product space  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$ . To each  $x \in \mathcal{V}$  there is a unique vector  $p \in \mathcal{S}$  such that

$$\langle x - p, s \rangle = 0, \quad \text{for all } s \in \mathcal{S}. \quad (2.24)$$

If  $(v_1, \dots, v_k)$  is an orthogonal basis for  $\mathcal{S}$  then

$$p = \sum_{i=1}^k \frac{\langle x, v_i \rangle}{\langle v_i, v_i \rangle} v_i. \quad (2.25)$$

**Proof.** Define  $p$  by (2.25). Then

$$\langle p, v_j \rangle = \sum_{i=1}^k \frac{\langle x, v_i \rangle}{\langle v_i, v_i \rangle} \langle v_i, v_j \rangle = \frac{\langle x, v_j \rangle}{\langle v_j, v_j \rangle} \langle v_j, v_j \rangle = \langle x, v_j \rangle$$

so that by linearity  $\langle x - p, v_j \rangle = 0$  for  $j = 1, \dots, k$ . But then  $\langle x - p, s \rangle = 0$  for all  $s \in \mathcal{S}$ . This shows existence of a  $p$  satisfying (2.24). 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  $\mathbf{p}_1 = \mathbf{p}_2$ .  $\square$

**Theorem 2.59 (Best Approximation)** *Let  $\mathcal{S}$  be a subspace of a finite dimensional real or complex inner product space  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$ . Let  $\mathbf{x} \in \mathcal{V}$ , and  $\mathbf{p} \in \mathcal{S}$ . The following statements are equivalent*

1.  $\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0$ , for all  $\mathbf{s} \in \mathcal{S}$ .
2.  $\|\mathbf{x} - \mathbf{s}\| > \|\mathbf{x} - \mathbf{p}\|$  for all  $\mathbf{s} \in \mathcal{S}$  with  $\mathbf{s} \neq \mathbf{p}$ .

**Proof.** Suppose 1. holds and that  $\mathbf{s} \neq \mathbf{p}$ . Using Pythagoras for inner products we have

$$\|\mathbf{x} - \mathbf{s}\|^2 = \|(\mathbf{x} - \mathbf{p}) + (\mathbf{p} - \mathbf{s})\|^2 = \|\mathbf{x} - \mathbf{p}\|^2 + \|\mathbf{p} - \mathbf{s}\|^2 > \|\mathbf{x} - \mathbf{p}\|^2.$$

Conversely, suppose 2. holds. Pick any nonzero  $\mathbf{s} \in \mathcal{S}$  and define the scalar  $a := -\operatorname{Re}\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle / \|\mathbf{s}\|^2$ . Using (2.19) and the minimality of  $\mathbf{p}$  we obtain

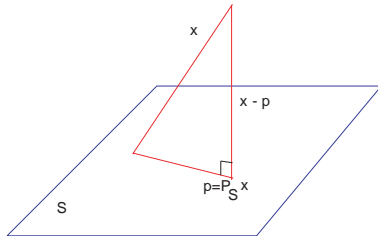
$$\begin{aligned} \|\mathbf{x} - \mathbf{p}\|^2 &\leq \|\mathbf{x} - \mathbf{p} + a\mathbf{s}\|^2 = \|\mathbf{x} - \mathbf{p}\|^2 + 2a\operatorname{Re}\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle + a^2\|\mathbf{s}\|^2 \\ &= \|\mathbf{x} - \mathbf{p}\|^2 - (\operatorname{Re}\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle) / \|\mathbf{s}\|^2. \end{aligned}$$

This can only be true if  $\operatorname{Re}\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0$  for all  $\mathbf{s} \in \mathcal{S}$ . Since  $\mathbf{s} \in \mathcal{S}$  implies that  $i\mathbf{s} \in \mathcal{S}$ , where  $i = \sqrt{-1}$ , we have

$$0 = \operatorname{Re}\langle \mathbf{x} - \mathbf{p}, i\mathbf{s} \rangle = \operatorname{Re}(-i\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle) = \operatorname{Im}\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle$$

and hence  $\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0$  for all  $\mathbf{s} \in \mathcal{S}$ .  $\square$

The vector  $\mathbf{p}$  is called the **orthogonal projection** of  $\mathbf{x}$  into  $\mathcal{S}$  with respect to  $\langle \cdot, \cdot \rangle$ , and denoted by  $\mathbf{p} = P_{\mathcal{S}}\mathbf{x}$ .



**Figure 2.1.** The orthogonal projection of  $\mathbf{x}$  into  $\mathcal{S}$ .

In terms of an orthogonal basis  $(\mathbf{v}_1, \dots, \mathbf{v}_k)$  for  $\mathcal{S}$  we have the representation

$$\mathbf{s} = \sum_{i=1}^k \frac{\langle \mathbf{s}, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \mathbf{v}_i, \text{ all } \mathbf{s} \in \mathcal{S}. \quad (2.26)$$

## 2.8 Projections and Orthogonal Complements

**Theorem 2.60** *Let  $\mathcal{S}$  be a subspace in a real or complex inner product space  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$  and let  $P_{\mathcal{S}} : \mathbb{R}^n \rightarrow \mathcal{S}$  be the operator mapping a vector  $\mathbf{x} \in \mathcal{V}$  into the orthogonal projection  $\mathbf{p}$  in  $\mathcal{S}$ . Then  $P_{\mathcal{S}}$  is a **linear projection operator**, i.e.*

1.  $P_{\mathcal{S}}(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha P_{\mathcal{S}}\mathbf{x} + \beta P_{\mathcal{S}}\mathbf{y}$  for all  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$  and all  $\alpha, \beta \in \mathbb{F}$ .
2.  $P_{\mathcal{S}}^2 = P_{\mathcal{S}}$ , i.e.  $P_{\mathcal{S}}(P_{\mathcal{S}}\mathbf{x}) = P_{\mathcal{S}}\mathbf{x}$  for all  $\mathbf{x} \in \mathcal{V}$ .

**Proof.**

1. Let  $\mathbf{p} := P_{\mathcal{S}}\mathbf{x}$  and  $\mathbf{q} := P_{\mathcal{S}}\mathbf{y}$ . Then  $\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle = 0$  and  $\langle \mathbf{y} - \mathbf{q}, \mathbf{s} \rangle = 0$  for all  $\mathbf{s} \in \mathcal{S}$ , and by linearity of the inner product

$$\langle \alpha\mathbf{x} + \beta\mathbf{y} - (\alpha\mathbf{p} + \beta\mathbf{q}), \mathbf{s} \rangle = \alpha\langle \mathbf{x} - \mathbf{p}, \mathbf{s} \rangle + \beta\langle \mathbf{y} - \mathbf{q}, \mathbf{s} \rangle = 0.$$

But then  $\alpha\mathbf{p} + \beta\mathbf{q} = \alpha P_{\mathcal{S}}\mathbf{x} + \beta P_{\mathcal{S}}\mathbf{y}$  is the orthogonal projection of  $\alpha\mathbf{x} + \beta\mathbf{y}$  into  $\mathcal{S}$  and 1. follows.

2. Since  $\mathbf{p} = P_{\mathcal{S}}\mathbf{x} \in \mathcal{S}$  the uniqueness implies that  $P_{\mathcal{S}}\mathbf{p} = \mathbf{p}$  which gives 2.

□

**Definition 2.61 (Orthogonal Complement)** *Let  $\mathcal{S}$  be a subspace in a real or complex inner product space  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$ . The Orthogonal Complement of  $\mathcal{S}$ , which is denoted by  $\mathcal{S}^{\perp}$ , consists of all vectors in  $\mathcal{V}$  that are orthogonal to every  $\mathbf{s} \in \mathcal{S}$ . In other words*

$$\mathbf{x} \in \mathcal{S}^{\perp} \iff \langle \mathbf{x}, \mathbf{s} \rangle = 0, \text{ for all } \mathbf{s} \in \mathcal{S}.$$

Clearly  $\mathcal{S}^{\perp}$  is a subspace of  $\mathcal{V}$ .

**Theorem 2.62 (Orthogonal Decomposition)** *For each subspace  $\mathcal{S}$  of a real or complex inner product space we have the direct sum decomposition  $\mathcal{V} = \mathcal{S} \oplus \mathcal{S}^{\perp}$ . If  $(\mathbf{s}_1, \dots, \mathbf{s}_k)$  is a basis for  $\mathcal{S}$  and  $(\mathbf{t}_1, \dots, \mathbf{t}_n)$  is a basis for  $\mathcal{S}^{\perp}$  then  $(\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{t}_1, \dots, \mathbf{t}_n)$  is a basis for  $\mathcal{S} \oplus \mathcal{S}^{\perp}$ . In particular, any orthonormal basis for  $\mathcal{S}$  can be extended to an orthonormal basis for  $\mathcal{V}$ .*

**Proof.** If  $\mathbf{x} \in \mathcal{S} \cap \mathcal{S}^{\perp}$  then  $\langle \mathbf{x}, \mathbf{x} \rangle = 0$  so  $\mathbf{x} = \mathbf{0}$ . This means that  $\mathcal{S} \cap \mathcal{S}^{\perp} = \{\mathbf{0}\}$  and  $\mathcal{S} \oplus \mathcal{S}^{\perp}$  is a direct sum. Every  $\mathbf{x} \in \mathbb{R}^n$  can be decomposed as  $\mathbf{x} = P_{\mathcal{S}}\mathbf{x} + (\mathbf{x} - P_{\mathcal{S}}\mathbf{x})$  where  $P_{\mathcal{S}}\mathbf{x} \in \mathcal{S}$  and  $\mathbf{x} - P_{\mathcal{S}}\mathbf{x} \in \mathcal{S}^{\perp}$ . Since we are dealing with a direct sum it follows from Theorem 2.29 that any basis  $(\mathbf{s}_1, \dots, \mathbf{s}_k)$  for  $\mathcal{S}$  and any basis  $(\mathbf{t}_1, \dots, \mathbf{t}_n)$  for  $\mathcal{S}^{\perp}$  can be combined into a basis for  $\mathcal{V}$ . If  $(\mathbf{s}_1, \dots, \mathbf{s}_k)$  is an orthonormal basis for  $\mathcal{S}$  then we apply the Gram-Schmidt process to  $(\mathbf{t}_1, \dots, \mathbf{t}_n)$  to obtain a combined orthonormal basis for  $\mathcal{V}$ . □

**Exercise 2.63** *Show that  $(\mathcal{S}^{\perp})^{\perp} = \mathcal{S}$  for any subspace  $\mathcal{S}$  of a real or complex inner product space.*

## Chapter 3

# Matrices

In this chapter we review some topics related to matrices. In Section 3.1 we briefly study block-multiplication, a basic tool in matrix analysis. We then review the transpose matrix, linear systems, and inverse matrices. We end the chapter with some basic facts about orthonormal-, and unitary matrices.

Some matrices with many zeros have names indicating their "shape". Suppose  $A \in \mathbb{R}^{n,n}$  or  $A \in \mathbb{C}^{n,n}$ . Then  $A$  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  $|i - j| > 1$ .
- **lower banded** with bandwidth  $p$  if  $a_{ij} = 0$  for  $i > j + p$ .
- **upper banded** with bandwidth  $q$  if  $a_{ij} = 0$  for  $i < j - q$ .
- **banded** with bandwidth  $p + q + 1$  if  $A$  is both lower banded with bandwidth  $p$  and upper banded with bandwidth  $q$ .
- **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  $A^T$  is block upper triangular.

## 3.1 Arithmetic Operations and Block Multiplication

The arithmetic operations on rectangular matrices are

- **matrix addition**  $C = A + B$  if  $c_{ij} = a_{ij} + b_{ij}$  for all  $i, j$  and  $A, B, C$  are matrices of the same dimension.
- **multiplication by a scalar**  $C = \alpha A$ , where  $c_{ij} = \alpha a_{ij}$  for all  $i, j$ .

- **multiplication by another matrix**  $C = AB$ ,  $C = A \cdot B$  or  $C = A * B$ , where  $A \in \mathbb{C}^{m,p}$ ,  $B \in \mathbb{C}^{p,n}$ ,  $C \in \mathbb{C}^{m,n}$ , 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**(add a dot)  $C = A \cdot B$  and  $D = A ./ B$ , and  $E = A \wedge r$  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  $C = A \cdot B$  is known as the **Schur product** and also the **Hadamard product**.

**Example 3.1 (The Vector Space of  $m \times n$  matrices)** On the set  $\mathbb{C}^{m,n}$  of  $m \times n$  matrices we define vector addition as matrix addition and scalar multiplication as a scalar times a matrix. Then  $\mathbb{C}^{m,n} = (\mathbb{C}^{m,n}, \mathbb{C})$  is a vector space. Of course  $\mathbb{R}^{m,n} = (\mathbb{R}^{m,n}, \mathbb{R})$  is also a vector space.

### 3.1.1 Block Multiplication

A rectangular matrix  $A$  can be partitioned into submatrices by drawing horizontal lines between selected rows and vertical lines between selected columns. Such a matrix is often referred to as a **block matrix**. Under certain mild restriction we can multiply two block matrices by applying the matrix multiplication process to the blocks, treating them as elements. In particular, suppose  $C = AB$ . If  $A, B, C$  are block matrices

$$A = \begin{bmatrix} A_{11} & \cdots & A_{1s} \\ \vdots & & \vdots \\ A_{p1} & \cdots & A_{ps} \end{bmatrix}, \quad B = \begin{bmatrix} B_{11} & \cdots & B_{1q} \\ \vdots & & \vdots \\ B_{s1} & \cdots & B_{sq} \end{bmatrix}, \quad C = \begin{bmatrix} C_{11} & \cdots & C_{1q} \\ \vdots & & \vdots \\ C_{p1} & \cdots & C_{pq} \end{bmatrix}$$

then

$$C_{ij} = \sum_{k=1}^s A_{ik}B_{kj}, \quad i = 1, \dots, p, \quad j = 1, \dots, q$$

provided

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

## 3.2 The Transpose Matrix

The **transpose** of  $A \in \mathbb{C}^{m,n}$  is a matrix  $B \in \mathbb{C}^{n,m}$ , where  $b_{ij} = a_{ji}$  for all  $i, j$ . Thus the rows of  $A$  are the columns of  $B$  and vice versa. The transpose of  $A$  is denoted  $A^T$ . Three important properties of the transpose are

1.  $(A + B)^T = A^T + B^T$ .
2.  $(AC)^T = C^T A^T$ .



$$3. (\mathbf{A}^T)^T = \mathbf{A}.$$

Here  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m,n}$  and  $\mathbf{C} \in \mathbb{C}^{n,k}$ , where  $k, m, n$  are any positive integers.

Consider now the real case  $\mathbf{A} \in \mathbb{R}^{m,n}$ . A useful characterization is the following:

**Theorem 3.2** *Let  $\langle x, y \rangle := \mathbf{x}^T \mathbf{y} = \sum_{i=1}^m x_i y_i$  be the usual inner product on  $\mathbb{R}^m$ . For any  $\mathbf{A} \in \mathbb{R}^{m,n}$  we have  $\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \langle \mathbf{A}^T \mathbf{x}, \mathbf{y} \rangle$ , all  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{y} \in \mathbb{R}^n$ . If  $\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \langle \mathbf{B}\mathbf{x}, \mathbf{y} \rangle$  holds for some  $\mathbf{B} \in \mathbb{R}^{n,m}$  and all  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{y} \in \mathbb{R}^n$  then  $\mathbf{B} = \mathbf{A}^T$ .*

**Proof.** For any  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^n$

$$\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \sum_{i=1}^m x_i \left( \sum_{j=1}^n a_{ij} y_j \right) = \sum_{j=1}^n \left( \sum_{i=1}^m x_i a_{ij} \right) y_j = \sum_{j=1}^n (\mathbf{A}^T \mathbf{x})_j y_j = \langle \mathbf{A}^T \mathbf{x}, \mathbf{y} \rangle.$$

If we choose  $\mathbf{x} = \mathbf{e}_i$  and  $\mathbf{y} = \mathbf{e}_j$  then  $a_{ij} = \langle \mathbf{e}_i, \mathbf{A}\mathbf{e}_j \rangle = \langle \mathbf{B}\mathbf{e}_i, \mathbf{e}_j \rangle = b_{ji}$  for  $i = 1, \dots, m$  and  $j = 1, \dots, n$  so  $\mathbf{B} = \mathbf{A}^T$ .  $\square$

The **Hermitian transpose** or **conjugate transpose** of  $\mathbf{A} \in \mathbb{C}^{m,n}$  is the matrix  $\mathbf{B} \in \mathbb{C}^{n,m}$  given by  $\mathbf{B} = (\overline{\mathbf{A}})^T$ . Here  $\bar{z} = x - iy$  denotes the complex conjugate of  $z = x + iy$ , where  $i = \sqrt{-1}$  is the imaginary unit and  $x, y \in \mathbb{R}$ . Moreover  $\overline{\mathbf{A}}$  is obtained from  $\mathbf{A}$  by taking the complex conjugate of all its elements. The Hermitian transpose of  $\mathbf{A}$  is denoted  $\mathbf{A}^H$ . The Hermitian transpose enjoys the same properties as the transpose:

1.  $(\mathbf{A} + \mathbf{B})^H = \mathbf{A}^H + \mathbf{B}^H$ .
2.  $(\mathbf{A}\mathbf{C})^H = \mathbf{C}^H \mathbf{A}^H$ .
3.  $(\mathbf{A}^H)^H = \mathbf{A}$ .

Again  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m,n}$  and  $\mathbf{C} \in \mathbb{C}^{n,k}$ , where  $k, m, n$  are any positive integers.

We obtain the same characterization in the complex case.

**Theorem 3.3** *Let  $\langle x, y \rangle := \mathbf{x}^H \mathbf{y} = \sum_{i=1}^m \bar{x}_i y_i$  be the usual inner product in  $\mathbb{C}^m$ . For any  $\mathbf{A} \in \mathbb{C}^{m,n}$  we have  $\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \langle \mathbf{A}^H \mathbf{x}, \mathbf{y} \rangle$ , all  $\mathbf{x} \in \mathbb{C}^m$ ,  $\mathbf{y} \in \mathbb{C}^n$ . If  $\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle = \langle \mathbf{B}\mathbf{x}, \mathbf{y} \rangle$  holds for some  $\mathbf{B} \in \mathbb{C}^{n,m}$  and all  $\mathbf{x} \in \mathbb{C}^m$ ,  $\mathbf{y} \in \mathbb{C}^n$  then  $\mathbf{B} = \mathbf{A}^H$ .*

**Exercise 3.4** Use Theorem 3.3 to show that  $(\mathbf{A}\mathbf{C})^H = \mathbf{C}^H \mathbf{A}^H$  and  $(\mathbf{A}^H)^H = \mathbf{A}$ .

### 3.3 Linear Systems

Consider a linear system

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

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  $\mathbf{b} = [b_1, \dots, b_m]^T$ . It can also be written as a matrix equation

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

The system is **homogenous** 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 a homogenous underdetermined system.

**Lemma 3.5** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}(\mathbb{C}^{m,n})$  with  $m < n$ . Then there is a nonzero  $\mathbf{x} \in \mathbb{R}^n(\mathbb{C}^n)$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .

**Proof.** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}(\mathbb{C}^{m,n})$  with  $m < n$ . The  $n$  columns of  $\mathbf{A}$  span a subspace of  $\mathbb{R}^m(\mathbb{C}^m)$ . Since  $\mathbb{R}^m(\mathbb{C}^m)$  has dimension  $m$  the dimension of this subspace is at most  $m$ . By Lemma 2.14 the columns of  $\mathbf{A}$  must be linearly dependent. It follows that there is a nonzero  $\mathbf{x} \in \mathbb{R}^n(\mathbb{C}^n)$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .  $\square$

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

**Definition 3.6** A square matrix  $\mathbf{A}$  is said to be **nonsingular** if the only solution of the homogenous system  $\mathbf{A}\mathbf{x} = \mathbf{0}$  is  $\mathbf{x} = \mathbf{0}$ . The matrix is **singular** if it is not nonsingular.

**Theorem 3.7** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$ . The linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a unique solution  $\mathbf{x} \in \mathbb{R}^n(\mathbb{C}^n)$  for any  $\mathbf{b} \in \mathbb{R}^n(\mathbb{C}^n)$  if and only if the matrix  $\mathbf{A}$  is nonsingular.

**Proof.** Suppose  $\mathbf{A}$  is nonsingular. We define  $\mathbf{B} = [\mathbf{A} \ \mathbf{b}] \in \mathbb{R}^{n,n+1}(\mathbb{C}^{n,n+1})$  by adding a column to  $\mathbf{A}$ . By Lemma 3.5 there is a nonzero  $\mathbf{z} \in \mathbb{R}^{n+1}(\mathbb{C}^{n+1})$  such that  $\mathbf{B}\mathbf{z} = \mathbf{0}$ . If we write  $\mathbf{z} = \begin{bmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{bmatrix}$  where  $\tilde{\mathbf{z}} = [z_1, \dots, z_n]^T \in \mathbb{R}^n(\mathbb{C}^n)$  and  $z_{n+1} \in \mathbb{R}(\mathbb{C})$ , then

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

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

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

so  $\mathbf{x}$  is a solution.

Suppose  $\mathbf{Ax} = \mathbf{b}$  and  $\mathbf{Ay} = \mathbf{b}$  for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n(\mathbb{C}^n)$ . 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  $\mathbf{Ax} = \mathbf{b}$  has a unique solution for any  $\mathbf{b} \in \mathbb{R}^n(\mathbb{C}^n)$  then  $\mathbf{Ax} = \mathbf{0}$  has a unique solution which must be  $\mathbf{x} = \mathbf{0}$ . Thus  $\mathbf{A}$  is nonsingular.  $\square$

### 3.4 The Inverse matrix

Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$  is a square matrix. A matrix  $\mathbf{B} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$  is called a **right inverse** of  $\mathbf{A}$  if  $\mathbf{AB} = \mathbf{I}$ . A matrix  $\mathbf{C} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$  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  $\mathbf{A}^{-1}$ . 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.

**Lemma 3.8** *If  $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$  with  $\mathbf{AB} = \mathbf{C}$  then  $\mathbf{C}$  is nonsingular if and only if both  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular.*

**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  $\mathbf{x} \in \mathbb{R}^n(\mathbb{C}^n)$  be a nonzero vector so that  $\mathbf{Bx} = \mathbf{0}$ . But then  $\mathbf{Cx} = (\mathbf{AB})\mathbf{x} = \mathbf{A}(\mathbf{Bx}) = \mathbf{A0} = \mathbf{0}$  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 3.7 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 3.9** *A square matrix is invertible if and only if it is nonsingular.*

**Proof.** Suppose first  $\mathbf{A}$  is a nonsingular matrix. By Theorem 3.7 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 3.8  $\mathbf{B}$  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  $\mathbf{B}$  and since  $\mathbf{AB} = \mathbf{I}$  and  $\mathbf{I}$  is nonsingular we again use Lemma 3.8 to conclude that  $\mathbf{A}$  is nonsingular.  $\square$

The theorem shows that we can use the terms "nonsingular" and "invertible" interchangeably. If  $\mathbf{B}$  is a right inverse or a left inverse of  $\mathbf{A}$  then it follows from Lemma 3.8 that  $\mathbf{A}$  is nonsingular. 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 of 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 3.10** *Suppose  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,n}(\mathbb{C}^{n,n})$  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 3.11** *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 3.12** *Find the inverse of*

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

**Exercise 3.13** *Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$ , and  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{n,m}$  for some  $n, m \in \mathbb{N}$ . If  $(\mathbf{I} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{B})^{-1}$  exists then*

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

### 3.5 Rank, Nullity, and the Fundamental Subspaces

Recall that the column space (or span) and the null space (kernel) of a matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  are defined by

$$\begin{aligned}\text{span } \mathbf{A} &:= \{\mathbf{y} \in \mathbb{C}^m : \mathbf{y} = \mathbf{A}\mathbf{x}, \mathbf{x} \in \mathbb{R}^n\} \\ \ker \mathbf{A} &:= \{\mathbf{x} \in \mathbb{C}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}.\end{aligned}$$

These sets are subspaces of  $\mathbb{C}^m$  and  $\mathbb{C}^n$ , respectively. The four subspaces  $\text{span } \mathbf{A}$ ,  $\ker \mathbf{A}$ ,  $\text{span } \mathbf{A}^H$  and  $\ker \mathbf{A}^H$  are known as the four **fundamental subspaces** of a matrix. The dimension of the column space of  $\mathbf{A}$  is called the **rank** of  $\mathbf{A}$  and denoted  $\text{rank } \mathbf{A}$ . The dimension  $\dim \ker \mathbf{A}$  of the null space is called the **nullity** of  $\mathbf{A}$  and denoted  $\text{null } \mathbf{A}$ .

Recall that the orthogonal complement  $\mathcal{S}^\perp$  of a subspace  $\mathcal{S}$  of  $\mathbb{C}^n$  is  $\{\mathbf{t} \in \mathbb{C}^n : \langle \mathbf{s}, \mathbf{t} \rangle = 0 \text{ for all } \mathbf{s} \in \mathcal{S}\}$ . For  $\mathcal{S} = \text{span } \mathbf{A}$  we have

**Theorem 3.14** *The orthogonal complement of the column space of a matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  is the null space of  $\mathbf{A}^H$ . We have the orthogonal decomposition*

$$\mathbb{C}^m = \text{span } \mathbf{A} \oplus \ker \mathbf{A}^H. \quad (3.1)$$

**Proof.** We first show that

$$\text{span}(\mathbf{A})^\perp = \ker(\mathbf{A}^H) := \{\mathbf{y} \in \mathbb{R}^m : \mathbf{A}^H \mathbf{y} = \mathbf{0}\}.$$

Suppose  $\mathbf{c} \in \text{span}(\mathbf{A})$ . Then  $\mathbf{c} = \mathbf{A}\mathbf{x}$  for some  $\mathbf{x} \in \mathbb{R}^n$ . If  $\mathbf{y} \in \ker(\mathbf{A}^H)$  then  $\langle \mathbf{y}, \mathbf{c} \rangle = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle = \langle \mathbf{A}^H \mathbf{y}, \mathbf{x} \rangle = 0$ . Thus  $\ker(\mathbf{A}^H) \subset \text{span}(\mathbf{A})^\perp$ . To show that  $\text{span}(\mathbf{A})^\perp \subset \ker(\mathbf{A}^H)$  we pick any  $\mathbf{y} \in \text{span}(\mathbf{A})^\perp$ . Then  $\langle \mathbf{A}^H \mathbf{y}, \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle = 0$  for all  $\mathbf{x} \in \mathbb{R}^n$  which means that  $\mathbf{y} \in \ker(\mathbf{A}^H)$ . The orthogonal decomposition (3.1) now follows from Theorem 2.62.  $\square$

The following formula for the rank of a product of two matrices will also be useful.

**Lemma 3.15** *If  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{B} \in \mathbb{C}^{n,p}$  for some  $m, n, p \in \mathbb{N}$  then*

$$\text{rank}(\mathbf{AB}) = \text{rank } \mathbf{B} - \dim(\ker \mathbf{A} \cap \text{span } \mathbf{B}).$$

**Proof.** Pick a basis  $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$  for  $\ker \mathbf{A} \cap \text{span } \mathbf{B}$  and extend it to a basis  $\{\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_l\}$  for  $\text{span } \mathbf{B}$ . The result will follow if we can show that  $Y := \{\mathbf{A}\mathbf{x}_{k+1}, \dots, \mathbf{A}\mathbf{x}_l\}$  is a basis for  $\text{span}(\mathbf{AB})$ .

- (i)  $Y$  is linearly independent. For if  $\sum c_j \mathbf{A}\mathbf{x}_j := \sum_{j=k+1}^l c_j \mathbf{A}\mathbf{x}_j = \mathbf{0}$  then  $\mathbf{A}(\sum c_j \mathbf{x}_j) = \mathbf{0}$ , and hence  $\sum c_j \mathbf{x}_j \in \ker \mathbf{A} \cap \text{span } \mathbf{B}$ . But then  $\sum_{j=k+1}^l c_j \mathbf{x}_j = \sum_{j=1}^k c_j \mathbf{s}_j$  for some  $c_1, \dots, c_k$ , and by linear independence we have  $c_1 = \dots = c_l = 0$ .
- (ii)  $\text{span } Y \subset \text{span}(\mathbf{AB})$ . Suppose  $\mathbf{y} = \sum c_j \mathbf{A}\mathbf{x}_j \in \text{span } Y$ . Since  $\mathbf{x}_j \in \text{span}(\mathbf{B})$  we have  $\mathbf{x}_j = \mathbf{B}\mathbf{z}_j$ , for some  $\mathbf{z}_j$ ,  $j = k+1, \dots, l$ . But then  $\mathbf{y} = \sum_j c_j \mathbf{AB}\mathbf{z}_j \in \text{span}(\mathbf{AB})$ .

- (iii)  $\text{span}(\mathbf{AB}) \subset \text{span } Y$ . If  $\mathbf{y} \in \text{span}(\mathbf{AB})$  then  $\mathbf{y} = \mathbf{Ax}$  for some  $\mathbf{x} \in \text{span } \mathbf{B}$ . Since  $\{\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_l\}$  is a basis for  $\text{span } \mathbf{B}$  we have  $\mathbf{x} = \sum_{j=1}^k c_j \mathbf{s}_j + \sum_{j=k+1}^l c_j \mathbf{x}_j$  for some  $c_1, \dots, c_l$  and  $\mathbf{s}_j \in \ker(\mathbf{A})$ . But then

$$\mathbf{y} = \mathbf{Ax} = \sum_{j=1}^k c_j \mathbf{As}_j + \sum_{j=k+1}^l c_j \mathbf{Ax}_j = \sum_{j=k+1}^l c_j \mathbf{Ax}_j \in \text{span}(Y).$$

□

Consider now the four fundamental subspaces.

**Theorem 3.16** *For any matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  we have*

1.  $\text{rank } \mathbf{A} + \text{null } \mathbf{A} = n$ ,
2.  $\text{rank } \mathbf{A} + \text{null } \mathbf{A}^H = m$ ,
3.  $\text{rank } \mathbf{A} = \text{rank } \mathbf{A}^H$ .

**Proof.**

1. Taking  $\mathbf{B}$  to be the identity matrix in Lemma 3.15 we obtain  $\text{rank}(\mathbf{A}) = \text{rank } \mathbf{I} - \dim(\ker \mathbf{A} \cap \text{span } \mathbf{I}) = n - \dim(\ker \mathbf{A} \cap \mathbb{C}^n) = n - \dim \ker \mathbf{A}$ .
2. This follows from Theorems 2.62 and 3.14.
3. If we apply 2. to  $\mathbf{A}^H$  we obtain  $\dim \text{span } \mathbf{A}^H + \dim \ker \mathbf{A} = n$ . But then  $\text{rank } \mathbf{A} = \dim \text{span } \mathbf{A} \stackrel{1.}{=} n - \dim \ker \mathbf{A} = n - (n - \dim \text{span } \mathbf{A}^H) = \dim \text{span } \mathbf{A}^H = \text{rank}(\mathbf{A}^H)$ . □

To derive some further results about rank and nullity we start with a definition:

**Definition 3.17 (Equivalent matrices)** *Suppose  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m,n}$ . We say that  $\mathbf{A}$  is equivalent to  $\mathbf{B}$ , denoted  $\mathbf{A} \sim \mathbf{B}$ , if  $\mathbf{B} = \mathbf{XAY}$  for some nonsingular matrices  $\mathbf{X} \in \mathbb{C}^{m,m}$  and  $\mathbf{Y} \in \mathbb{C}^{n,n}$ .*

**Exercise 3.18** *Show that  $\sim$  is an equivalence relation, i. e.,*

- (i)  $\mathbf{A} \sim \mathbf{A}$ ,
- (ii) if  $\mathbf{A} \sim \mathbf{B}$  then  $\mathbf{B} \sim \mathbf{A}$ ,
- (iii) if  $\mathbf{A} \sim \mathbf{B}$  and  $\mathbf{B} \sim \mathbf{C}$  then  $\mathbf{A} \sim \mathbf{C}$ .

For any subspace  $\mathcal{S}$  of  $\mathbb{C}^n$  and  $\mathbf{B} \in \mathbb{C}^{m,n}$  we define  $\mathbf{BS} := \{\mathbf{Bs} : \mathbf{s} \in \mathcal{S}\}$ .

**Exercise 3.19** *Show that  $\mathbf{BS}$  is a subspace of  $\mathbb{C}^m$ .*

**Exercise 3.20** *Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  and that  $\mathbf{X} \in \mathbb{C}^{m,m}$  and  $\mathbf{Y} \in \mathbb{C}^{n,n}$  are nonsingular. Show that*

1.  $\text{span}(\mathbf{A}) = \text{span}(\mathbf{A}\mathbf{Y}) = \mathbf{X}^{-1} \text{span}(\mathbf{X}\mathbf{A})$ ,
2.  $\ker \mathbf{A} = \mathbf{Y} \ker(\mathbf{A}\mathbf{Y}) = \ker(\mathbf{X}\mathbf{A})$ ,
3.  $\text{rank}(\mathbf{X}\mathbf{A}\mathbf{Y}) = \text{rank}(\mathbf{A})$ ,
4.  $\text{null}(\mathbf{X}\mathbf{A}\mathbf{Y}) = \text{null } \mathbf{A}$ .

For the rank of a general product we have

**Theorem 3.21** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{B} \in \mathbb{C}^{n,p}$  for some  $m, n, p \in \mathbb{N}$ . Then

$$\text{rank}(\mathbf{AB}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}.$$

**Proof.** Since  $\text{span}(\mathbf{AB}) \subset \text{span}(\mathbf{A})$  we have  $\text{rank}(\mathbf{AB}) \leq \text{rank}(\mathbf{A})$ . Now  $\text{span}(\mathbf{B}^H \mathbf{A}^H) \subset \text{span}(\mathbf{B}^H)$ . Therefore  $\text{rank}(\mathbf{AB}) = \text{rank}((\mathbf{AB})^H) = \text{rank}(\mathbf{B}^H \mathbf{A}^H) \leq \text{rank}(\mathbf{B}^H) = \text{rank}(\mathbf{B})$ .  $\square$

We end this section with the following useful result.

**Theorem 3.22** If the matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  has rank  $r$  then there is at least one nonsingular  $r \times r$  submatrix in  $\mathbf{A}$ . Moreover there are no nonsingular submatrices of larger order.

**Proof.** We use Theorem 2.16 twice. There is a subset  $\{\mathbf{a}_{\cdot j_1}, \dots, \mathbf{a}_{\cdot j_r}\}$  of the columns of  $\mathbf{A}$  which forms a basis for  $\text{span}(\mathbf{A})$ . Consider the matrix  $\mathbf{B}^H \in \mathbb{C}^{m,r}$ , where  $\mathbf{B} = [\mathbf{a}_{\cdot j_1}, \dots, \mathbf{a}_{\cdot j_r}]$ . Since  $r = \text{rank}(\mathbf{B}) = \text{rank}(\mathbf{B}^H)$  there is a subset  $\{i_1, \dots, i_r\}$  of  $\{1, \dots, m\}$  such that columns  $i_1, \dots, i_r$  of  $\mathbf{B}^H$  form a basis for  $\text{span}(\mathbf{B}^H)$ . But then rows  $i_1, \dots, i_r$  of  $\mathbf{B}$  are linearly independent, defining a nonsingular  $r \times r$  submatrix in  $\mathbf{A}$ . Suppose  $M$  is a nonsingular submatrix in  $\mathbf{A}$  of order  $k$ . The columns in  $\mathbf{A}$  corresponding to the columns in  $M$  are linearly independent and hence  $k \leq r$ .  $\square$

## 3.6 Linear Transformations and Matrices

Let  $(\mathcal{X}, \mathbb{F})$  and  $(\mathcal{Y}, \mathbb{F})$  be vector spaces over the same field  $\mathbb{F}$ . A mapping  $T : \mathcal{X} \rightarrow \mathcal{Y}$  is called **linear** if for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$  and all  $a \in \mathbb{F}$  we have

1.  $T(\mathbf{x} + \mathbf{y}) = T\mathbf{x} + T\mathbf{y}$ , (additivity)
2.  $T(a\mathbf{x}) = aT\mathbf{x}$ . (homogeneity)

If  $\mathcal{Y}$  is the vector space of all functions  $f : \mathbb{R} \rightarrow \mathbb{R}$  and  $\mathcal{X}$  is the space of all differentiable functions  $f : \mathbb{R} \rightarrow \mathbb{R}$ , then the mapping  $T : \mathcal{X} \rightarrow \mathcal{Y}$  given by  $Tf := df/dx$  is a linear transformation from  $\mathcal{X}$  to  $\mathcal{Y}$ . The mapping  $T$  given by  $(Tf)(x) := \int_0^x f(t)dt$  is a linear transformation from the space  $\mathcal{X}$  of all continuous functions  $f : \mathbb{R} \rightarrow \mathbb{R}$  into  $\mathcal{X}$ .

Linear transformations are not the main emphasis of this text and we will only consider briefly the special case where  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}^m$ . The same results hold for the complex case  $\mathcal{X} = \mathbb{C}^n$  and  $\mathcal{Y} = \mathbb{C}^m$ . Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$ . The mapping

$T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  given by  $T\mathbf{x} = A\mathbf{x}$  is clearly additive and homogenous. Thus it is a linear mapping. It turns out that all linear mappings  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are of this form.

**Theorem 3.23** *Every linear map from  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  can be written in the form  $T = A\mathbf{x}$  for some  $A \in \mathbb{R}^{m,n}$ .*

**Proof.** Suppose  $\mathbf{x} \in \mathbb{R}^n$ . Then  $\mathbf{x} = \sum_{j=1}^n x_j \mathbf{e}_j$  and by linearity

$$T\mathbf{x} = T\left(\sum_{j=1}^n x_j \mathbf{e}_j\right) = \sum_{j=1}^n x_j T\mathbf{e}_j = \sum_{j=1}^n x_j \mathbf{a}_j = A\mathbf{x},$$

where  $A = [\mathbf{a}_1, \dots, \mathbf{a}_n] = [T\mathbf{e}_1, \dots, T\mathbf{e}_n] \in \mathbb{R}^{m,n}$ .  $\square$

Let

$$\begin{aligned} \text{span } T &:= \{\mathbf{y} \in \mathbb{R}^m : \mathbf{y} = T\mathbf{x} \text{ for some } \mathbf{x} \in \mathbb{R}^n\}, \\ \ker T &:= \{\mathbf{x} \in \mathbb{R}^n : T\mathbf{x} = \mathbf{0}\}, \end{aligned} \tag{3.2}$$

be the **span** and **kernel** of the linear transformation  $T$ . The sets  $\text{span } T$  and  $\ker T$  are subspaces of  $\mathbb{R}^m$  and  $\mathbb{R}^n$ , respectively.

**Theorem 3.24** *Suppose  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a linear transformation. Then For any matrix  $A \in \mathbb{C}^{m,n}$  we have*

$$\dim \text{span } T + \dim \ker T = n.$$

**Proof.** This follows from Theorem 3.16 since  $T\mathbf{x} = A\mathbf{x}$  for some matrix  $A$ .  $\square$

Much more can be said about linear transformations and matrices. We refer to any book on linear algebra.

### 3.7 Orthonormal and Unitary Matrices

**Definition 3.25** *A matrix  $Q \in \mathbb{R}^{n,n}$  is said to be **orthonormal** if  $Q^T Q = I$ .*

Since the columns of an orthonormal matrix are orthonormal, we have chosen the term "orthonormal matrix" although "orthogonal matrix" is more common in the classical literature.

**Theorem 3.26** *Suppose  $Q \in \mathbb{R}^{n,n}$ . The following is equivalent:*

1.  $Q$  is orthonormal,
2. the columns of  $Q$  form an orthonormal basis for  $\mathbb{R}^n$ .
3.  $Q^{-1} = Q^T$
4.  $QQ^T = I$
5. the columns of  $Q^T$  (rows of  $Q$ ) form an orthonormal basis for  $\mathbb{R}^n$ ,



6.  $\langle Q\mathbf{x}, Q\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , where  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y}$  is the usual inner product on  $\mathbb{R}^n$ .

We also have

- (i) The product  $Q_1 Q_2$  of two orthonormal matrices is orthonormal.  
(ii) If  $Q$  is orthonormal then  $\|Q\mathbf{x}\|_2 = \|\mathbf{x}\|_2$  for all  $\mathbf{x} \in \mathbb{R}^n$ .

**Proof.** Let  $\mathbf{q}_1, \dots, \mathbf{q}_n$  be the columns of  $Q$ .

$1 \Leftrightarrow 2$  This follows since  $(Q^T Q)_{ij} = \langle \mathbf{q}_i, \mathbf{q}_j \rangle$  for all  $i, j$ .

$1 \Leftrightarrow 3$  Since  $Q^T$  is a left inverse of  $Q$  it follows from the discussion after Theorem 3.9 that  $Q$  is invertible and  $Q^{-1} = Q^T$ .

$3 \Leftrightarrow 4$  Since  $Q^T = Q^{-1}$  the definition of the inverse matrix implies that  $Q^T$  is a right inverse of  $Q$  so that  $Q Q^T = I$ .

$4 \Leftrightarrow 5$  This follows since  $Q^T$  is orthonormal and  $(Q^T)^T = Q$ .

$1 \Leftrightarrow 6$  If  $Q$  is orthonormal then by Theorem 3.2 we have  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle Q^T Q \mathbf{x}, \mathbf{y} \rangle = \langle Q \mathbf{x}, Q \mathbf{y} \rangle$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Conversely, taking  $\mathbf{x} = \mathbf{e}_i$  and  $\mathbf{y} = \mathbf{e}_j$  we find  $(Q^T Q)_{ij} = \langle Q^T Q \mathbf{e}_i, \mathbf{e}_j \rangle = \langle Q \mathbf{e}_i, Q \mathbf{e}_j \rangle = \langle \mathbf{e}_i, \mathbf{e}_j \rangle = \delta_{ij}$  for all  $i, j = 1, \dots, n$ .

Suppose  $Q_1$  and  $Q_2$  are orthonormal. Then  $(Q_1 Q_2)^T Q_1 Q_2 = Q_2^T Q_1^T Q_1 Q_2 = I$  so the product  $Q_1 Q_2$  is orthonormal. Using 6. with  $\mathbf{y} = \mathbf{x}$  we obtain (ii).  $\square$

Consider now the complex case.

**Definition 3.27** A matrix  $U \in \mathbb{C}^{n,n}$  is said to be **unitary** if  $U^H U = I$ .

Note that a real unitary matrix is orthonormal.

**Theorem 3.28** Suppose  $U \in \mathbb{C}^{n,n}$ . The following is equivalent:

1.  $U$  is unitary,
2. the columns of  $U$  form an orthonormal basis for  $\mathbb{C}^n$ .
3.  $U^{-1} = U^H$
4.  $U U^H = I$
5. the columns of  $U^H$  (rows of  $U$ ) form an orthonormal basis for  $\mathbb{C}^n$ ,
6.  $\langle U\mathbf{x}, U\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ , where  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^H \mathbf{y}$  is the usual inner product on  $\mathbb{C}^n$ .
7.  $\|U\mathbf{x}\|_2 = \|\mathbf{x}\|_2$  for all  $\mathbf{x} \in \mathbb{C}^n$ .

The product  $U_1 U_2$  of two unitary matrices is unitary.

**Proof.** That 1-6 are equivalent is similar to the proof of the real case. Clearly 6 implies 7. That 7 implies 6 follows from the fact that we can write a complex inner product as a sum of norms. (Cf. (2.20).)  $\square$

**Exercise 3.29** *Show in Theorem 3.28 that 7. implies 6.*

## Chapter 4

# 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.

### 4.1 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  $\sigma, \tau$  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  $\sigma, \tau$  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  $\sigma^{-1}$ . 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  $\sigma$ . 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

$$\frac{g(\sigma\tau)}{g(\tau)} = \frac{\prod_{i=2}^n (\sigma(\tau(i)) - \sigma(\tau(1))) \cdots (\sigma(\tau(i)) - \sigma(\tau(i-1)))}{\prod_{i=2}^n (\tau(i) - \tau(1)) \cdots (\tau(i) - \tau(i-1))} = \prod_{i=2}^n \prod_{j=1}^{i-1} r_\sigma(\tau(i), \tau(j)).$$

$\tau$  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  $\sigma$  and  $\sigma^{-1}$  have the same sign.

**Exercise 4.1** Show 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  $\epsilon$ .
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  $n$ ).

## 4.2 Basic Properties of Determinants

For any  $A \in \mathbb{C}^{n,n}$  the determinant of  $A$  is defined the number

$$\det(A) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma(1),1} a_{\sigma(2),2} \cdots a_{\sigma(n),n}. \quad (4.1)$$

This sum ranges of all  $n!$  permutations of  $\{1, 2, \dots, n\}$ . We also denote the determinant by (Cayley, 1841)

$$\begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{vmatrix}.$$

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  $\epsilon$  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  $A, B \in \mathbb{C}^{n,n}$  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  $\sigma$  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 5.6, 5.7). But then the result follows.

7. To better understand the general proof, we do the  $2 \times 2$  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  $\sigma \in S_{n,k}$ ,  $\rho$  and  $\tau$  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  $\sigma$  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

$$\begin{aligned} \det(\mathbf{A}) &= \sum_{\hat{\rho} \in S_k} \sum_{\hat{\tau} \in S_{n-k}} \text{sign}(\hat{\rho})\text{sign}(\hat{\tau}) b_{\hat{\rho}(1),1} \cdots b_{\hat{\rho}(k),k} d_{\hat{\tau}(1),1} \cdots d_{\hat{\tau}(n-k),n-k} \\ &= \det(\mathbf{B}) \det(\mathbf{D}). \end{aligned}$$

□



### 4.3 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,

$$\begin{aligned} \mathbf{A} &= \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 3 \\ 6 \end{pmatrix}, \quad \mathbf{A}_1(\mathbf{b}) = \begin{pmatrix} 3 & 2 \\ 6 & 1 \end{pmatrix}, \quad \mathbf{A}_2(\mathbf{b}) = \begin{pmatrix} 1 & 3 \\ 2 & 6 \end{pmatrix}, \\ \mathbf{I} &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \mathbf{I}_1(\mathbf{x}) = \begin{pmatrix} x_1 & 0 \\ x_2 & 1 \end{pmatrix}, \quad \mathbf{I}_2(\mathbf{x}) = \begin{pmatrix} 1 & x_1 \\ 0 & x_2 \end{pmatrix}. \end{aligned}$$

**Theorem 4.2 (Cramers rule (1750))** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Then

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

**Proof.** Since  $1 = \det(\mathbf{I}) = \det(\mathbf{A}\mathbf{A}^{-1}) = \det(\mathbf{A})\det(\mathbf{A}^{-1})$  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.  $\square$

**Exercise 4.3** Solve the following system by Cramers rule:

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 6 \end{bmatrix}$$

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 4.4 (Cofactor and Adjoint)** For  $\mathbf{A} \in \mathbb{C}^{n,n}$  and  $1 \leq i, j \leq n$  the determinant  $\det(\mathbf{A}_{ij})$  is called the **cofactor** of  $a_{ij}$ . The matrix  $\text{adj}(\mathbf{A}) \in \mathbb{C}^{n,n}$  with elements  $(-1)^{i+j} \det(\mathbf{A}_{j,i})$  is called the **adjoint** of  $\mathbf{A}$ .

**Exercise 4.5** 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,

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

**Theorem 4.6 (The inverse as an adjoint)** If  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$  implies that  $\mathbf{A}\mathbf{x}_j = \mathbf{e}_j$  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 4.7** For any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have

$$\mathbf{A} \text{adj}(\mathbf{A}) = \text{adj}(\mathbf{A})\mathbf{A} = \det(\mathbf{A})\mathbf{I}. \quad (4.2)$$

**Proof.** If  $\mathbf{A}$  is nonsingular then (4.2) follows from Theorem 4.6. 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

$$(\mathbf{A} + \epsilon\mathbf{I}) \text{adj}(\mathbf{A} + \epsilon\mathbf{I}) = \text{adj}(\mathbf{A} + \epsilon\mathbf{I})(\mathbf{A} + \epsilon\mathbf{I}) = \det(\mathbf{A} + \epsilon\mathbf{I})\mathbf{I}. \quad (4.3)$$

Since the elements in  $\mathbf{A} + \epsilon\mathbf{I}$  and  $\text{adj}(\mathbf{A} + \epsilon\mathbf{I})$  depend continuously on  $\epsilon$  we can take limits in (4.3) to obtain (4.2). □

**Corollary 4.8 (Cofactor expansion)** *For any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have*

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } i = 1, \dots, n, \quad (4.4)$$

$$\det(\mathbf{A}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } j = 1, \dots, n. \quad (4.5)$$

**Proof.** By (4.2) we have  $\mathbf{A} \operatorname{adj}(\mathbf{A}) = \det(\mathbf{A}) \mathbf{I}$ . But then  $\det(\mathbf{A}) = \mathbf{e}_i^T \mathbf{A} \operatorname{adj}(\mathbf{A}) \mathbf{e}_i = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij})$  which is (4.4). 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 (4.5).  $\square$

## 4.4 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  $\alpha$ .
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  $\alpha$  times column  $k$  to column  $i$  for some  $k \neq i$ . Then using Properties 5. and 7. we find

$$\begin{aligned} \det(\mathbf{B}) &= \det \left( [\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{a}_i + \alpha \mathbf{a}_k, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n] \right) \\ &\stackrel{5.}{=} \det(\mathbf{A}) + \det \left( [\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \alpha \mathbf{a}_k, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n] \right) \stackrel{7.}{=} \det(\mathbf{A}) \end{aligned}$$

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 in Corollary 4.8.

**Example 4.9** 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 4.10** Show that the equation for the plane through the points  $(x_1, y_1, z_1)$ ,  $(x_2, y_2, z_2)$  and  $(x_3, y_3, z_3)$  is

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

**Exercise 4.11** 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  $T$  is

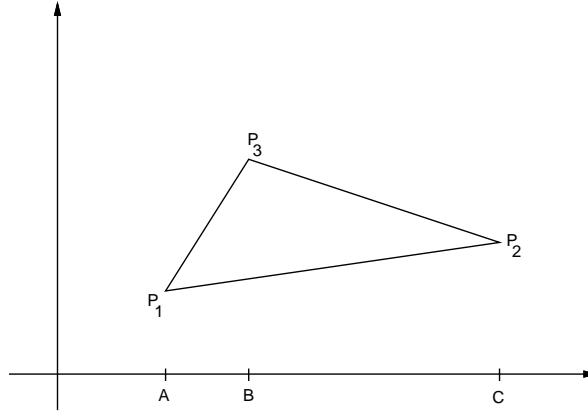
$$A(T) = \frac{1}{2} \begin{vmatrix} x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \\ 1 & 1 & 1 \end{vmatrix}$$

*Hint:*  $A(T) = A(ABP_3P_1) + A(P_3BCP_2) - A(P_1ACP_2)$ , c.f. Figure 4.12.

**Exercise 4.13** 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 Van der Monde determinant. *Hint:* Subtract  $x_n^k$  times column  $k$  from column  $k+1$  for  $k = n-1, n-2, \dots, 1$ .



**Figure 4.12.** The triangle  $T$  defined by the three points  $P_1$ ,  $P_2$  and  $P_3$ .

**Exercise 4.14** (Cauchy 1842). Let  $\alpha = [\alpha_1, \dots, \alpha_n]^T$ ,  $\beta = [\beta_1, \dots, \beta_n]^T$  be in  $\mathbb{R}^n$ .

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

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

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  $\alpha_i$  and  $\beta_j$ . 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  $\alpha$  and  $\beta$ . 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 4.15** Let  $\mathbf{H}_n = (h_{i,j})$  be the  $n \times n$  matrix with elements  $h_{i,j} = 1/(i+j-1)$ . Use Exercise 4.14 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

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

## 4.5 Some Useful Determinant Formulas

Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  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  $1 \leq j_1 < j_2 < \dots < j_r \leq n$ . 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 4.16 (Cauchy-Binet formula)** Let  $\mathbf{A} \in \mathbb{C}^{m,p}$ ,  $\mathbf{B} \in \mathbb{C}^{p,n}$  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

$$\det(\mathbf{C}(\mathbf{i}, \mathbf{j})) = \sum_{\mathbf{k}} \det(\mathbf{A}(\mathbf{i}, \mathbf{k})) \det(\mathbf{B}(\mathbf{k}, \mathbf{j})), \quad (4.6)$$

where we sum over all  $\mathbf{k} = \{k_1, \dots, k_r\}$  with  $1 \leq k_1 < k_2 < \dots < k_r \leq p$ .

## Chapter 5

# Eigenvalues and Eigenvectors

Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is a square matrix,  $\lambda \in \mathbb{C}$  and  $\mathbf{x} \in \mathbb{C}^n$ . We say that  $(\lambda, \mathbf{x})$  is an **eigenpair** for  $\mathbf{A}$  if  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbf{x}$  is nonzero. The scalar  $\lambda$  is called an **eigenvalue** and  $\mathbf{x}$  is said to be an **eigenvector**. If  $(\lambda, \mathbf{x})$  is an eigenpair then  $(\lambda, \alpha\mathbf{x})$  is an eigenpair for any  $\alpha \in \mathbb{C}$  with  $\alpha \neq 0$ . An eigenvector is a special vector that is mapped by  $\mathbf{A}$  into a vector parallel to itself. The length is increased if  $|\lambda| > 1$  and decreased if  $|\lambda| < 1$ . The set of distinct eigenvalues is called the **spectrum** of  $\mathbf{A}$  and is denoted by  $\sigma(\mathbf{A})$ .

## 5.1 The Characteristic Polynomial

### 5.1.1 The characteristic equation

**Lemma 5.1** For any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have  $\lambda \in \sigma(\mathbf{A}) \iff \det(\mathbf{A} - \lambda\mathbf{I}) = 0$ .

*Proof.* Suppose  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$ . The equation  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  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  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$ .  $\square$

We observe that  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$  if and only if  $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$ . The equation  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$  or equivalently  $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$  is called the **characteristic equation** of  $\mathbf{A}$ .

**Definition 5.2** 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  $\mathbf{A}$ .

To see that  $\pi_{\mathbf{A}}$  is in fact a polynomial let us take a closer look at this function. 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

$$\det(\mathbf{A} - \lambda \mathbf{I}) = (a_{11} - \lambda)(a_{22} - \lambda) \cdots (a_{nn} - \lambda) + r(\lambda), \quad (5.1)$$

where each term in  $r(\lambda)$  has at most  $n - 2$  factors containing  $\lambda$ . 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.

By the fundamental theorem of algebra an  $n \times n$  matrix has 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  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with  $\mathbf{A}$  real gives  $\mathbf{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$ .

The following result will be useful.

**Theorem 5.3** Suppose  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} \in \mathbb{C}^{n,n}$ . Then

1. If  $\mathbf{A}$  is nonsingular then  $(\mu^{-1}, \mathbf{x})$  is an eigenpair for  $\mathbf{A}^{-1}$ .
2.  $(\mu^k, \mathbf{x})$  is an eigenpair for  $\mathbf{A}^k$  for  $k \in \mathbb{N}$ .
3. If  $p$  given by  $p(t) = a_0 + a_1 t + a_2 t^2 + \cdots + a_k t^k$  is a polynomial, then  $(p(\mu), \mathbf{x})$  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.  $\mu$  is an eigenvalue for  $\mathbf{A}^T$ , in fact  $\pi_{\mathbf{A}^T} = \pi_{\mathbf{A}}$ .
5.  $\bar{\mu}$  is an eigenvalue for  $\mathbf{A}^H$ , in fact  $\pi_{\mathbf{A}^H}(\bar{\lambda}) = \overline{\pi_{\mathbf{A}}(\lambda)}$  for all  $\lambda \in \mathbb{C}$ .
6. If  $\mathbf{A} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}$  is block triangular then  $\pi_{\mathbf{A}} = \pi_{\mathbf{B}} \cdot \pi_{\mathbf{D}}$ .

**Proof.**

1.  $\mathbf{A}\mathbf{x} = \mu\mathbf{x} \implies \mathbf{A}^{-1}\mathbf{x} = \mu^{-1}\mathbf{x}$ .
2. We use induction on  $k$ . The case  $k = 1$  is trivial and if  $\mathbf{A}^{k-1}\mathbf{x} = \mu^{k-1}\mathbf{x}$  then  $\mathbf{A}^k\mathbf{x} = \mathbf{A}\mathbf{A}^{k-1}\mathbf{x} = \mu^{k-1}\mathbf{A}\mathbf{x} = \mu^k\mathbf{x}$ .
3.  $p(\mathbf{A})\mathbf{x} = \sum_{j=0}^k a_j \mathbf{A}^j \mathbf{x} \stackrel{2.}{=} \sum_{j=0}^k a_j \mu^j \mathbf{x} = p(\mu)\mathbf{x}$ .



4. Since  $\det(\mathbf{B}^T) = \det(\mathbf{B})$  for any matrix  $\mathbf{B}$  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  $\pi_{\mathbf{A}^H}(\bar{\lambda}) \stackrel{4}{=} \pi_{\overline{\mathbf{A}}}(\bar{\lambda}) = \det(\overline{\mathbf{A}} - \bar{\lambda} \mathbf{I}) = \overline{\det(\mathbf{A} - \lambda \mathbf{I})} = \overline{\pi_{\mathbf{A}}(\lambda)}$ . Thus  $\pi_{\mathbf{A}}(\lambda) = 0 \Leftrightarrow \pi_{\mathbf{A}^H}(\bar{\lambda}) = 0$  and the result follows.

6. By Property 8 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.

**Theorem 5.4** *The eigenvalues of a triangular matrix are given by its diagonal elements.*

**Proof.** If  $\mathbf{A} \in \mathbb{C}^{n,n}$  is triangular then  $\mathbf{A} - \lambda \mathbf{I}$  is also triangular with diagonal elements  $a_{ii} - \lambda$  for  $i = 1, \dots, n$ . But then the roots of  $\det(\mathbf{A} - \lambda \mathbf{I}) = \prod_{i=1}^n (a_{ii} - \lambda) = 0$  are  $\lambda_i = a_{ii}$  for  $i = 1, \dots, n$ . □

To find the eigenvectors of a triangular matrix requires more work. Indeed, the eigenvectors are nontrivial solutions of a homogenous triangular linear system with at least one zero on the diagonal.

**Example 5.5** *The  $3 \times 3$  matrix  $\mathbf{A} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$  has the eigenvalue  $\lambda = 1$ . The homogenous triangular linear system for an eigenvector  $\mathbf{x} = [x_1, x_2, x_3]^T$  is*

$$(\mathbf{A} - \mathbf{I})\mathbf{x} = \mathbf{0} \text{ or } \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

We find  $x_2 = x_3 = 0$  so any eigenvector must be a multiple of  $\mathbf{e}_1$ .

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

**Theorem 5.6** *For any  $\mathbf{A} \in \mathbb{C}^{n,n}$*

$$\text{trace}(\mathbf{A}) = \lambda_1 + \lambda_2 + \dots + \lambda_n, \quad \det(\mathbf{A}) = \lambda_1 \lambda_2 \dots \lambda_n, \quad (5.2)$$

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

$$\text{trace}(\mathbf{A}) := a_{11} + a_{22} + \dots + a_{nn} \quad (5.3)$$

and  $\det(\mathbf{A})$  is the determinant of  $\mathbf{A}$ .

**Proof.** We compare two different expansion of  $\pi_{\mathbf{A}}$ . On the one hand from (5.1) 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 (5.2).  $\square$

For a  $2 \times 2$  matrix the characteristic equation takes the convenient form

$$\lambda^2 - \text{trace}(\mathbf{A})\lambda + \det(\mathbf{A}) = 0. \quad (5.4)$$

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

In terms of eigenvalues we have an additional characterization of a singular matrix.

**Theorem 5.7** *The matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  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.  $\square$

**Exercise 5.8** Find eigenvalues and eigenvectors of  $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 2 \end{bmatrix}$ .

**Exercise 5.9** Let  $\lambda \in \sigma(A)$  where  $A^2 = A \in \mathbb{C}^{n,n}$ . Show that  $\lambda = 0$  or  $\lambda = 1$ . ( $A$  matrix is called idempotent if  $A^2 = A$ ).

**Exercise 5.10** Let  $\lambda \in \sigma(A)$  where  $A^k = 0$  for some  $k \in \mathbb{N}$ . Show that  $\lambda = 0$ . ( $A$  matrix  $A \in \mathbb{C}^{n,n}$  such that  $A^k = 0$  for some  $k \in \mathbb{N}$  is called nilpotent).

**Exercise 5.11** Let  $\lambda \in \sigma(A)$  where  $A^H A = I$ . Show that  $|\lambda| = 1$ .

**Exercise 5.12** Suppose  $A \in \mathbb{C}^{n,n}$  is singular. Then we can find  $\epsilon_0 > 0$  such that  $A + \epsilon I$  is nonsingular for all  $\epsilon \in (0, \epsilon_0)$ . Hint:  $\det(A) = \lambda_1 \lambda_2 \cdots \lambda_n$ , where  $\lambda_i$  are the eigenvalues of  $A$ .

**Exercise 5.13** For  $q_i \in \mathbb{C}$  let  $f(\lambda) = \lambda^n + q_{n-1} \lambda^{n-1} + \cdots + q_0$  be a polynomial of degree  $n$  in  $\lambda$ . We derive two matrices which have  $(-1)^n f$  as its characteristic polynomial.

a) Show that  $f = (-1)^n \pi_A$  where

$$A = \begin{bmatrix} -q_{n-1} & -q_{n-2} & \cdots & -q_1 & -q_0 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}.$$

$A$  is called the companion matrix of  $f$ .

b) Show that  $f = (-1)^n \pi_{A'}$  where

$$A' = \begin{bmatrix} 0 & 0 & \cdots & 0 & -q_0 \\ 1 & 0 & \cdots & 0 & -q_1 \\ 0 & 1 & \cdots & 0 & -q_2 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -q_{n-1} \end{bmatrix}.$$

Thus  $A'$  can also be regarded as a companion matrix for  $f$ .

## 5.2 Similarity Transformations

Row operations can be used 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 5.14** Two matrices  $A, B \in \mathbb{C}^{n,n}$  are said to be **similar** if there is a nonsingular matrix  $S \in \mathbb{C}^{n,n}$  such that  $B = S^{-1}AS$ . The transformation  $A \rightarrow B$  is called a **similarity transformation**.

A similarity transformation does not change the eigenvalues.

**Theorem 5.15** Similar matrices have the same characteristic polynomial and therefore the same eigenvalues.

**Proof.** Let  $B = S^{-1}AS$ . By properties of determinants

$$\begin{aligned} \pi_B(\lambda) &= \det(S^{-1}AS - \lambda I) = \det(S^{-1}(A - \lambda I)S) \\ &= \det(S^{-1}) \det(A - \lambda I) \det(S) = \det(S^{-1}S) \det(A - \lambda I) = \pi_A(\lambda). \end{aligned}$$

But then  $A$  and  $B$  have the same characteristic polynomial and hence the same eigenvalues.  $\square$

Consider next what a similarity transformation does to the eigenvectors.

**Theorem 5.16** 1.  $(\lambda, \mathbf{x})$  is an eigenpair for  $B = S^{-1}AS$  if and only if  $(\lambda, S\mathbf{x})$  is an eigenpair for  $A$ .

2. The columns of  $\mathbf{S}$  are eigenvectors of  $\mathbf{A}$  if and only if  $\mathbf{B}$  is diagonal.

**Proof.**

1.  $\mathbf{B}\mathbf{x} = \lambda\mathbf{x} \Leftrightarrow \mathbf{S}^{-1}\mathbf{A}\mathbf{S}\mathbf{x} = \lambda\mathbf{x} \Leftrightarrow \mathbf{A}(\mathbf{S}\mathbf{x}) = \lambda(\mathbf{S}\mathbf{x})$ , and  $\mathbf{S}\mathbf{x} \neq 0$  since  $\mathbf{S}$  is nonsingular.
2. Suppose  $\mathbf{A}$  has eigenvalues  $\lambda_1, \dots, \lambda_n$  and let  $\mathbf{s}_1, \dots, \mathbf{s}_n$  be the columns of  $\mathbf{S}$ . If  $\mathbf{B}$  is diagonal then  $(\lambda_i, \mathbf{e}_i)$  is an eigenpair for  $\mathbf{B}$  and  $(\lambda_i, \mathbf{S}\mathbf{e}_i) = (\lambda_i, \mathbf{s}_i)$  is an eigenpair for  $\mathbf{A}$  for  $i = 1, \dots, n$ . Conversely, if  $\mathbf{B} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}$  and the columns  $\mathbf{s}_1, \dots, \mathbf{s}_n$  of  $\mathbf{S}$  are eigenvectors of  $\mathbf{A}$  then  $\mathbf{A}\mathbf{s}_i = \lambda_i\mathbf{s}_i$  for  $i = 1, \dots, n$ . But then  $\mathbf{A}\mathbf{S} = \mathbf{S}\mathbf{C}$ , where  $\mathbf{C} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is diagonal. Thus  $\mathbf{C} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S} = \mathbf{B}$  is diagonal.

□

The following result is sometimes useful.

**Theorem 5.17** For any  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{B} \in \mathbb{C}^{n,m}$  the matrices  $\mathbf{AB}$  and  $\mathbf{BA}$  have the same spectrum. More precisely,

$$\lambda^n \pi_{\mathbf{AB}}(\lambda) = \lambda^m \pi_{\mathbf{BA}}(\lambda), \quad \lambda \in \mathbb{C}.$$

**Proof.** Define block matrices of order  $n + m$  by

$$\mathbf{E} = \begin{bmatrix} \mathbf{AB} & 0 \\ \mathbf{B} & 0 \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} 0 & 0 \\ \mathbf{B} & \mathbf{BA} \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} \mathbf{I}_m & \mathbf{A} \\ 0 & \mathbf{I}_n \end{bmatrix}.$$

By Property 6. of Theorem 5.3 we have  $\pi_{\mathbf{E}}(\lambda) = \lambda^n \pi_{\mathbf{AB}}(\lambda)$  and  $\pi_{\mathbf{F}}(\lambda) = \lambda^m \pi_{\mathbf{BA}}(\lambda)$ . But  $\mathbf{ES} = \mathbf{SF}$  so  $\mathbf{E}$  and  $\mathbf{F}$  are similar and have the same characteristic polynomial by the proof of Theorem 5.15. □

### 5.3 Linear Independence of Eigenvectors

**Definition 5.18** A square matrix  $\mathbf{A}$  is **diagonalizable** if it is similar to a diagonal matrix,  $\mathbf{S}^{-1}\mathbf{A}\mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$ .

Since  $\mathbf{S}$  is nonsingular its columns are eigenvectors of  $\mathbf{A}$ , and Theorem 5.16 implies the following result.

**Theorem 5.19** A matrix is diagonalizable if and only if its eigenvectors form a basis for  $\mathbb{R}^n$  or  $\mathbb{C}^n$ .

A matrix with distinct eigenvalues can be diagonalized.

**Theorem 5.20** Eigenvectors corresponding to distinct eigenvalues are linearly independent.

**Proof.** Suppose  $(\lambda_1, \mathbf{x}_1), \dots, (\lambda_k, \mathbf{x}_k)$  are eigenpairs for  $\mathbf{A} \in \mathbb{C}^{n,n}$  with  $\lambda_i \neq \lambda_j$  for  $i \neq j$ . Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_k$  are linearly dependent. Let  $m \leq k$  be the smallest positive integer so that  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are linearly dependent. Since  $\mathbf{x}_1 \neq \mathbf{0}$  we see that  $m \geq 2$ . For some nonzero  $(c_1, \dots, c_n)$  we have

$$\sum_{j=1}^m c_j \mathbf{x}_j = \mathbf{0}. \quad (5.5)$$

Applying  $\mathbf{A}$  to this equation we obtain by linearity  $\sum_{j=1}^m c_j \lambda_j \mathbf{x}_j = \mathbf{0}$ . From this relation we subtract  $\lambda_m$  times (5.5) 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 5.21** *If  $\mathbf{A} \in \mathbb{C}^{n,n}$  has distinct eigenvalues then the corresponding eigenvectors form a basis for  $\mathbb{C}^n$ .*

**Proof.** By the previous theorem the  $n$  eigenvectors are linearly independent. Since  $n$  is the dimension of  $\mathbb{C}^n$  the eigenvectors form a basis.  $\square$

For a matrix with multiple eigenvalues the situation is more complicated. We have seen that any eigenvector of the  $3 \times 3$  matrix  $\mathbf{A} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$  is a multiple of  $\mathbf{e}_1$ . Thus this matrix does not have a set of linearly independent eigenvectors. On the other hand the unit matrix has a basis of eigenvectors, namely the unit vectors.

In order to characterize the matrices with eigenvectors which form a basis we have to count carefully the multiplicity of the eigenvalues. We consider two kinds of multiplicities called algebraic and geometric multiplicities. The algebraic multiplicity of an eigenvalue  $\lambda$  is simply the multiplicity of  $\lambda$  as a root in the characteristic polynomial. More formally we state:

**Definition 5.22** *We say that an eigenvalue  $\lambda$  of  $\mathbf{A}$  has **algebraic multiplicity**  $a = a(\lambda) = a_{\mathbf{A}}(\lambda)$  if  $\pi_{\mathbf{A}}(z) = (z - \lambda)^a p(z)$ , where  $p(z) \neq 0$ . The eigenvalue  $\lambda$  is **simple** (**double**, **triple**) if  $a$  is equal to one (two, three). A complex number  $z$  which is not an eigenvalue is defined to have algebraic multiplicity  $a_{\mathbf{A}}(z) = 0$ .*

To define the second kind of multiplicity we consider for each  $\lambda \in \sigma(\mathbf{A})$  the nullspace

$$\ker(\mathbf{A} - \lambda \mathbf{I}) := \{\mathbf{x} \in \mathbb{R}^n : (\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = \mathbf{0}\} \quad (5.6)$$

of  $\mathbf{A} - \lambda \mathbf{I}$ . This set consists of all eigenvectors of  $\mathbf{A}$  corresponding to the eigenvalue  $\lambda$ . If  $\mathbf{x}, \mathbf{y} \in \ker(\mathbf{A} - \lambda \mathbf{I})$  and  $\alpha, \beta$  are scalars then  $\alpha \mathbf{x} + \beta \mathbf{y} \in \ker(\mathbf{A} - \lambda \mathbf{I})$ . So this nullspace is a subspace of  $\mathbb{C}^n$ . The dimension of the subspace must be at least one since  $\mathbf{A} - \lambda \mathbf{I}$  is singular.

**Definition 5.23** *The **geometric multiplicity**  $g = g(\lambda) = g_{\mathbf{A}}(\lambda)$  of an eigenvalue  $\lambda$  of  $\mathbf{A}$  is the dimension of the nullspace  $\ker(\mathbf{A} - \lambda \mathbf{I})$ .*

**Example 5.24** The  $n \times n$  identity matrix has the eigenvalue  $\lambda = 1$  with  $\pi_I(\lambda) = (1 - \lambda)^n$ . 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 the  $3 \times 3$  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 5.25** 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  $g := g_{\lambda}(\mathbf{A})$ , 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  $\mathbb{C}^n$ . Then the matrix  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n,n}$  is unitary and  $\mathbf{V}^{-1} = \mathbf{V}^H$ . 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  $\mathbf{A}\mathbf{V}_1 = \lambda\mathbf{V}_1$ ,  $\mathbf{V}_1^H\mathbf{V}_1 = \mathbf{I}_g$ ,  $\mathbf{V}_2^H\mathbf{V}_1 = 0$ , and

$$\mathbf{B} := \mathbf{V}^H \mathbf{A} \mathbf{V} = \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} \mathbf{A} [\mathbf{V}_1 \quad \mathbf{V}_2] = \begin{bmatrix} \mathbf{V}_1^H \mathbf{A} \mathbf{V}_1 & \mathbf{V}_1^H \mathbf{A} \mathbf{V}_2 \\ \mathbf{V}_2^H \mathbf{A} \mathbf{V}_1 & \mathbf{V}_2^H \mathbf{A} \mathbf{V}_2 \end{bmatrix} = \begin{bmatrix} \lambda \mathbf{I}_g & \mathbf{V}_1^H \mathbf{A} \mathbf{V}_2 \\ 0 & \mathbf{V}_2^H \mathbf{A} \mathbf{V}_2 \end{bmatrix}.$$

Since  $\mathbf{B}$  is block triangular Property 6 of Theorem 5.3 implies that  $\pi_{\mathbf{B}}(z) = (z - \lambda)^g \pi_{\mathbf{V}_2^H \mathbf{A} \mathbf{V}_2}(z)$ . 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 5.26** An eigenvalue where  $g_{\mathbf{A}}(\lambda) < a_{\mathbf{A}}(\lambda)$  is said to be **defective**. A matrix is **defective** if at least one of its eigenvalues is defective.

**Theorem 5.27** A matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  has  $n$  linearly independent eigenvectors if and only if the algebraic and geometric multiplicity of all eigenvalues are the same.

**Proof.** Suppose  $\mathbf{A}$  has distinct eigenvalues  $\mu_1, \dots, \mu_r$  with algebraic multiplicities  $a_1, \dots, a_r$  and geometric multiplicities  $g_1, \dots, g_r$ . Suppose  $\{\mathbf{v}_{j1}, \dots, \mathbf{v}_{jg_j}\}$  is a basis for  $\ker(\mathbf{A} - \mu_j\mathbf{I})$  for  $j = 1, \dots, r$ . We claim that the combined set  $\{\mathbf{v}_{jk}\}_{k=1, j=1}^{g_j, r}$  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

$$\sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} \mathbf{v}_{jk} = \mathbf{0} \text{ for some scalars } a_{jk}. \quad (5.7)$$

We multiply this equation by  $(\mathbf{A} - \mu_r\mathbf{I})$  and obtain by linearity

$$\mathbf{0} = \sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} (\mathbf{A} - \mu_r\mathbf{I}) \mathbf{v}_{jk} = \sum_{j=1}^r \sum_{k=1}^{g_j} a_{jk} (\mu_j - \mu_r) \mathbf{v}_{jk} = \sum_{j=1}^{r-1} \sum_{k=1}^{g_j} a_{jk} (\mu_j - \mu_r) \mathbf{v}_{jk}.$$

By the induction hypothesis all these  $a_{jk}$  vanish and in (5.7) we are left with  $\sum_{k=1}^{g_r} a_{rk} \mathbf{v}_{rk} = \mathbf{0}$ . Since these  $\mathbf{v}'$ s form a basis for  $\ker(\mathbf{A} - \mu_r \mathbf{I})$  we also have  $a_{rk} = 0$  for  $k = 1, \dots, g_r$ . (This also proves the induction hypothesis for  $r = 1$ .) 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  $\sum_j g_j$ . 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  $j = 1, \dots, r$ .  $\square$

## 5.4 Left Eigenvectors

**Definition 5.28** A nonzero vector  $\mathbf{y} \in \mathbb{C}^n$  corresponding to an eigenvalue  $\lambda$  of  $\mathbf{A}$  is called a **left eigenvector** of  $\mathbf{A}$  if  $\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H$ . We say that  $(\lambda, \mathbf{y})$  is a **left eigenpair** of  $\mathbf{A}$ .

Note that  $\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H$  if and only if  $\mathbf{A}^H \mathbf{y} = \bar{\lambda} \mathbf{y}$ . It follows from Theorem 5.3 that if  $\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H$  then  $\lambda$  must be an eigenvalue for  $\mathbf{A}$ , while a left eigenvector  $\mathbf{y}$  is an eigenvector for  $\mathbf{A}^H$ . If we need to make a distinction then an ordinary eigenvector, eigenpair is called a **right eigenvector** and **right eigenpair**, respectively.

Left- and right eigenvectors corresponding to distinct eigenvalues are orthogonal.

**Theorem 5.29** Suppose  $(\mu, \mathbf{y})$  and  $(\lambda, \mathbf{x})$  are left and right eigenpairs of  $\mathbf{A} \in \mathbb{C}^{n,n}$ . If  $\lambda \neq \mu$  then  $\mathbf{y}^H \mathbf{x} = 0$ .

**Proof.** Using the eigenpair relation in two ways we obtain  $\mathbf{y}^H \mathbf{A} \mathbf{x} = \lambda \mathbf{y}^H \mathbf{x} = \mu \mathbf{y}^H \mathbf{x}$  and we conclude that  $\mathbf{y}^H \mathbf{x} = 0$ .  $\square$

The case where  $\lambda = \mu$  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  $\mathbf{y}^H \mathbf{x} = 0$ . Two sufficient conditions guaranteeing that  $\mathbf{y}^H \mathbf{x} \neq 0$  are given in the following theorem.

**Theorem 5.30** Suppose  $\mathbf{y}$  and  $\mathbf{x}$  are left- and right eigenvectors corresponding to the same eigenvalue  $\lambda$  of  $\mathbf{A} \in \mathbb{C}^{n,n}$ . Then  $\mathbf{y}^H \mathbf{x} \neq 0$  in the following two cases:

1.  $\mathbf{A}$  can be diagonalized.
2. The algebraic multiplicity of  $\lambda$  is equal to one.

**Proof.**

1. Suppose  $\mathbf{Y} \mathbf{A} \mathbf{X} = \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ , where  $\mathbf{Y} = \mathbf{X}^{-1}$ . Partition  $\mathbf{Y}$  by rows and  $\mathbf{X}$  by columns as

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1^H \\ \vdots \\ \mathbf{y}_n^H \end{bmatrix}, \quad \mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n].$$

Since  $\mathbf{Y}\mathbf{A} = \mathbf{D}\mathbf{Y}$  and  $\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{D}$ , we see that  $\mathbf{y}_i$  is a left eigenvector and  $\mathbf{x}_i$  is a right eigenvector corresponding to  $\lambda_i$  for  $i = 1, \dots, n$ . But since  $\mathbf{Y}\mathbf{X} = \mathbf{I}$  we have  $\mathbf{y}_i^H \mathbf{x}_i = 1$  for all  $i$ .

2. Assume that  $\|\mathbf{x}\|_2 = 1$ . We have (cf. (10.1))

$$\mathbf{V}^H \mathbf{A} \mathbf{V} = \left[ \begin{array}{c|c} \lambda & \mathbf{z}^H \\ \hline \mathbf{0} & \mathbf{M} \end{array} \right],$$

where  $\mathbf{V}$  is unitary and  $\mathbf{V}\mathbf{e}_1 = \mathbf{x}$ . Let  $\mathbf{u} := \mathbf{V}^H \mathbf{y}$ . Then

$$(\mathbf{V}^H \mathbf{A}^H \mathbf{V})\mathbf{u} = \mathbf{V}^H \mathbf{A}^H \mathbf{y} = \bar{\lambda} \mathbf{V}^H \mathbf{y} = \bar{\lambda} \mathbf{u},$$

so  $(\bar{\lambda}, \mathbf{u})$  is an eigenpair of  $\mathbf{V}^H \mathbf{A}^H \mathbf{V}$ . But then  $\mathbf{y}^H \mathbf{x} = \mathbf{u}^H \mathbf{V}^H \mathbf{V} \mathbf{e}_1$ . Suppose that  $\mathbf{u}^H \mathbf{e}_1 = 0$ , i. e.,  $\mathbf{u} = \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix}$  for some nonzero  $\mathbf{v} \in \mathbb{C}^{n-1}$ . Then

$$\mathbf{V}^H \mathbf{A}^H \mathbf{V} \mathbf{u} = \left[ \begin{array}{c|c} \bar{\lambda} & \mathbf{0}^H \\ \hline \mathbf{z} & \mathbf{M}^H \end{array} \right] \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{M}^H \mathbf{v} \end{bmatrix} = \bar{\lambda} \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix}$$

and by Theorem 5.3 it follows that  $\lambda$  is an eigenvalue of  $\mathbf{M}$ . But this is impossible since  $\lambda$  has algebraic multiplicity one and the eigenvalues of  $\mathbf{A}$  are the union of  $\lambda$  and the eigenvalues of  $\mathbf{M}$ .

□

**Corollary 5.31** *If  $\mathbf{A} \in \mathbb{C}^{n,n}$  has linearly independent right eigenvectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$  then  $\mathbf{A}$ , also has linearly independent left eigenvectors  $\mathbf{y}_1, \dots, \mathbf{y}_n$ . For any  $\mathbf{v} \in \mathbb{C}^n$  we have*

$$\mathbf{v} = \sum_{j=1}^n (\mathbf{y}_j^H \mathbf{v}) \mathbf{x}_j = \sum_{k=1}^n (\mathbf{x}_k^H \mathbf{v}) \mathbf{y}_k. \quad (5.8)$$

**Proof.** From the proof of the previous theorem we have  $\mathbf{y}_k^H \mathbf{x}_j = \delta_{kj}$  for all  $j, k$ . So if  $\mathbf{v} = \sum_{j=1}^n c_j \mathbf{x}_j$ , then  $\mathbf{y}_k^H \mathbf{v} = \sum_{j=1}^n c_j \mathbf{y}_k^H \mathbf{x}_j = c_k$  for  $k = 1, \dots, n$ . The proof of the second formula is similar. □



## **Part II**

# **Some Linear Systems with a Special Structure**



## Chapter 6

# Examples of Linear Systems

Many problems in computational science lead to linear systems where the coefficient matrix has a special structure. In this chapter we present some problem that lead to a linear system with a tridiagonal, or almost tridiagonal coefficient matrix. Such linear systems can be solved by a version of Gaussian elimination adapted to the special structure. We also consider block multiplication and some useful facts about triangular matrices.

### 6.1 Cubic Spline Interpolation

Given  $n \geq 3$  interpolation sites  $\mathbf{x} = [x_1, \dots, x_n]$  with  $a := x_1 < \dots < x_n =: b$ , real  $y$  values  $\mathbf{y} = [y_1, \dots, y_n]$ , and derivative values  $\sigma_a, \sigma_b$ . We seek a function  $g : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$g(x_i) = y_i, \text{ for } i = 1, \dots, n, \quad g'(a) = \sigma_a, \quad g'(b) = \sigma_b. \quad (6.1)$$

The conditions  $g'(a) = \sigma_a, g'(b) = \sigma_b$  are called 1. derivative boundary conditions.

#### 6.1.1 Cubic $C^2$ Splines

The interpolant  $g$  should be a piecewise polynomial of the form

$$g(x) := \begin{cases} p_1(x), & \text{if } x < x_2, \\ p_2(x), & \text{if } x_2 \leq x < x_3, \\ \vdots & \\ p_{n-2}(x), & \text{if } x_{n-2} \leq x < x_{n-1}, \\ p_{n-1}(x), & \text{if } x_{n-1} \leq x, \end{cases} \quad (6.2)$$

where each  $p_i$  is a polynomial of degree  $\leq 3$ . In addition  $g \in C^2 := C^2(\mathbb{R})$ , i.e.,  $g$  should be continuous and have continuous first and second derivatives on  $\mathbb{R}$ . Clearly

$g \in C^2$  if and only if

$$p_{i-1}(x_i) = p_i(x_i), \quad p'_{i-1}(x_i) = p'_i(x_i), \quad p''_{i-1}(x_i) = p''_i(x_i), \quad i = 2, \dots, n-1. \quad (6.3)$$

We call  $g$  a **cubic  $C^2$  spline** with **knots  $\mathbf{x}$** . The sites  $x_2, \dots, x_{n-1}$  are called interior knots. The name spline is inherited from its "physical uncle", i.e., an elastic ruler that is used to draw smooth curves. Heavy weights, called **ducks**, are used to force the physical spline to pass through, or near given locations. (Cf. Figure 6.1).



**Figure 6.1.** A physical spline with ducks.

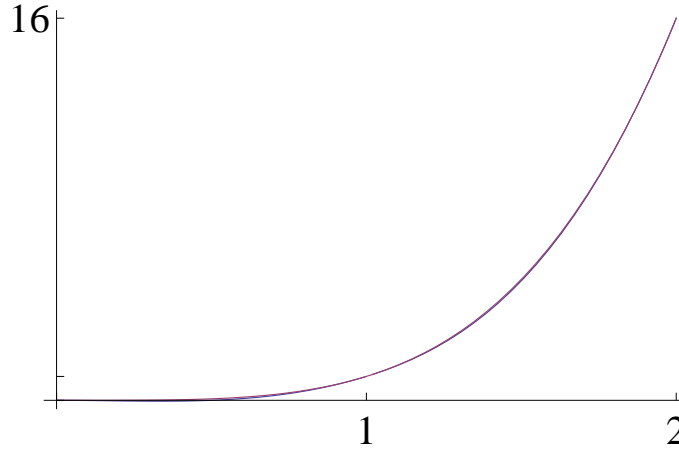
**Example 6.1** Show that  $g$  given by

$$g(x) := \begin{cases} p_1(x) = -x^2 + 2x^3, & \text{if } x < 1, \\ p_2(x) = -4 + 12x - 13x^2 + 6x^3, & \text{if } 1 \leq x, \end{cases} \quad (6.4)$$

is a cubic  $C^2$  spline interpolating the data

$$\mathbf{x} = [0, 1, 2], \quad \mathbf{y} = [0, 1, 4], \quad \sigma_a = 0, \quad \sigma_b = 32.$$

*Discussion:* Clearly  $g$  is in the form (6.2) with one interior knot  $x_2 = 1$ . Since  $p_1(1) = 1 = p_2(1)$ ,  $p'_1(1) = 4 = p'_2(1)$ ,  $p''_1(1) = 32 = p''_2(1)$  we see that  $g$  is a cubic  $C^2$  spline.  $g$  interpolates the data since  $g(x_1) = p_1(0) = 0 = y_1$ ,  $g(x_2) = p_2(1) = 1 = y_2$ ,  $g(x_3) = p_2(2) = 16 = y_3$ ,  $g'(0) = p'_1(0) = 0 = \sigma_a$ ,  $g'(2) = p'_2(2) = 32 = \sigma_b$ . The data is sampled from the function given by the rule  $f(x) = x^4$ . A plot of  $f$  and  $g$  is shown in Figure 6.2. It is hard to distinguish the two curves.



**Figure 6.2.** A two piece cubic spline interpolant to  $f(x) = x^4$ .

We will represent each  $p_i$  in shifted power form

$$p_i(x) = c_{1i} + c_{2i}(x - x_i) + c_{3i}(x - x_i)^2 + c_{4i}(x - x_i)^3. \quad (6.5)$$

A cubic spline is completely determined by the shifts  $\mathbf{x}_s$  and the coefficients  $\mathbf{C}$

$$\mathbf{x}_s := [x_1, \dots, x_k], \quad \mathbf{C} := \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1,k} \\ c_{21} & c_{22} & \cdots & c_{2,k} \\ c_{31} & c_{32} & \cdots & c_{3,k} \\ c_{41} & c_{42} & \cdots & c_{4,k} \end{bmatrix} \in \mathbb{R}^{4,k}, \quad k = n - 1. \quad (6.6)$$

We call (6.6) the **pp representation** of  $g$ . The shifted power form of (6.4) is

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

and the pp representation is

$$\mathbf{x}_s = [0, 1], \quad \mathbf{C} = \begin{bmatrix} 0 & 1 \\ 0 & 4 \\ -1 & 5 \\ 2 & 6 \end{bmatrix}. \quad (6.7)$$

To plot a cubic spline  $g$  we need to compute  $y$  values  $q_j = g(r_j)$  at a number of  $x$  values  $\mathbf{r} = [r_1, \dots, r_m] \in \mathbb{R}^m$  for some reasonably large integer  $m$ . For each  $r_j$  we have  $g(r_j) = p_{i_j}(r_j)$  for some integer  $i_j$ . The following Matlab function determines  $\mathbf{i} = [i_1, \dots, i_m]$ . It uses the built in Matlab functions **length**, **min**, **max**, **sort**, **find**.

**Algorithm 6.2 (findsubintervals)** Given shifts  $\mathbf{x}_s = [x_1, \dots, x_k]$  and a real number  $r$ , an integer  $i$  is computed so that  $i = 1$  if  $r < x(2)$ ,  $i = k$  if  $r \geq x_k$ , and  $x_i \leq r < x_{i+1}$  otherwise. If  $\mathbf{r}$  is a vector then a vector  $\mathbf{i}$  is computed, such that the  $j$ th component of  $\mathbf{i}$  gives the location of the  $j$ th component of  $\mathbf{r}$ .

```
function i=findsubintervals(xs,r)
    k=length(xs); m=length(r);
    xs(1)=min(r)-1;
    [sorted,j] = sort([xs(:)' r(:)']);
    i = find(j>k)-(1:m);
```

Here is the algorithm that was used to compute points for the plot in Figure 6.2. It uses Algorithm 6.2.

**Algorithm 6.3 (cubppeval)** Given a pp representation  $(\mathbf{x}_s, \mathbf{C})$  of a cubic spline  $g$  together with  $x$  values  $\mathbf{r} \in \mathbb{R}^m$ . The vector  $\mathbf{q} = g(\mathbf{r})$  is computed.

```
function q=cubppeval(xs,C,r)
i=findsubintervals(xs,r); q=r;
for j=1:length(r)
    k=i(j); t=r(j)-xs(k);
    q(j)=[1 t t^2 t^3]*C(:,k);
end
```

### 6.1.2 Determining the Interpolant

To determine a cubic  $C^2$  spline  $g$  interpolating the data (6.1), we first introduce the unknown derivatives  $s_i := g'(x_i)$ ,  $i = 2, \dots, n-1$  as free parameters in the shifted power representation of  $g$  and then determine these parameters so that  $g \in C^2[a, b]$ . The following lemma gives the details.

**Lemma 6.4** Assume that  $g$  is a piecewise polynomial given by (6.2) interpolating the data (6.1) and let  $\mathbf{s} = [s_1, \dots, s_n]$ , where  $s_1 = \sigma_a$ ,  $s_n = \sigma_b$ , and  $s_2, \dots, s_{n-1}$  are any real numbers. Suppose the coefficients in the pp representation of  $g$  are given for  $i = 1, 2, \dots, n-1$  by

$$\begin{aligned} c_{1i} &:= y_i, & c_{2i} &:= s_i, \\ c_{3i} &:= (3\delta_i - 2s_i - s_{i+1})/h_i, \\ c_{4i} &:= (-2\delta_i + s_i + s_{i+1})/h_i^2, \\ h_i &:= x_{i+1} - x_i, & \delta_i &:= (y_{i+1} - y_i)/h_i. \end{aligned} \tag{6.8}$$

Then

$$g(x_i) = y_i, \quad g'(x_i) = s_i, \quad i = 1, \dots, n, \tag{6.9}$$

and  $g \in C^1(\mathbb{R})$  for any choice of  $\mathbf{s}$ . Moreover,  $g \in C^2$  if and only if the  $n-2$  unknown slopes  $s_2, \dots, s_{n-1}$  satisfy the  $n-2$  relations

$$h_i s_{i-1} + d_i s_i + h_{i-1} s_{i+1} = \beta_i, \quad i = 2, \dots, n-1, \tag{6.10}$$

where for  $i = 2, \dots, n-1$

$$\begin{aligned} d_i &= 2(h_{i-1} + h_i), \\ \beta_i &= 3(h_i\delta_{i-1} + h_{i-1}\delta_i). \end{aligned} \quad (6.11)$$

**Proof.** Taking derivatives in (6.5) gives

$$\begin{aligned} p_i(x) &= c_{1i} + c_{2i}(x - x_i) + c_{3i}(x - x_i)^2 + c_{4i}(x - x_i)^3, \\ p'_i(x) &= c_{2i} + 2c_{3i}(x - x_i) + 3c_{4i}(x - x_i)^2, \\ p''_i(x) &= 2(c_{3i} + 3c_{4i}(x - x_i)), \end{aligned} \quad (6.12)$$

where the  $c_{j,i}$  are given by (6.8). Clearly  $p_i(x_i) = c_{1i} = y_i$ ,  $p'_i(x_i) = c_{2i} = s_i$ , and

$$\begin{aligned} p_i(x_{i+1}) &= c_{1i} + c_{2i}h_i + c_{3i}h_i^2 + c_{4i}h_i^3 \\ &= y_i + h_i(s_i + 3\delta_i - 2s_i - s_{i+1} - 2\delta_i + s_i + s_{i+1}) \\ &= y_i + h_i\delta_i = y_{i+1}, \\ p'_i(x_{i+1}) &= c_{2i} + 2c_{3i}h_i + 3c_{4i}h_i^2 \\ &= s_i + 2(3\delta_i - 2s_i - s_{i+1}) + 3(-2\delta_i + s_i + s_{i+1}) = s_{i+1} \end{aligned}$$

so that (6.9) holds. For  $2 \leq i \leq n-1$ , since  $p_{i-1}(x_i) = y_i = p_i(x_i)$  and  $p'_{i-1}(x_i) = s_i = p'_i(x_i)$  we see that  $g \in C^1$  and  $g \in C^2$  if and only if  $p''_{i-1}(x_i) = p''_i(x_i)$ . By (6.12)  $g \in C^2$  if and only if for  $i = 2, \dots, n-1$

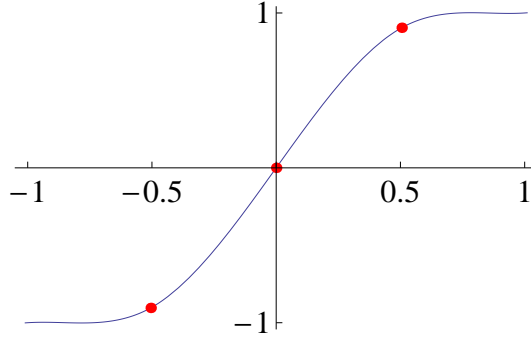
$$\begin{aligned} 0 &= h_{i-1}h_i\left(\frac{1}{2}p''_{i-1}(x_i) - \frac{1}{2}p''_i(x_i)\right) \\ &= h_{i-1}h_i(c_{3,i-1} + 3c_{4,i-1}h_{i-1} - c_{3,i}) \\ &= h_i(3\delta_{i-1} - 2s_{i-1} - s_i) + 3h_i(-2\delta_{i-1} + s_{i-1} + s_i) - h_{i-1}(3\delta_i - 2s_i - s_{i+1}) \\ &= h_is_{i-1} + 2(h_{i-1} + h_i)s_i + h_{i-1}s_{i+1} - 3(h_i\delta_{i-1} + h_{i-1}\delta_i). \end{aligned}$$

Rearranging we obtain (6.10).  $\square$

Adding the trivial equations  $s_1 = \sigma_a$  and  $s_n = \sigma_b$  to (6.10) we obtain for  $n \geq 3$  the following  $n \times n$  system

$$\mathbf{N}_1 \mathbf{s} = \begin{bmatrix} 1 & 0 & & & \\ h_2 & d_2 & h_1 & & \\ & \ddots & \ddots & \ddots & \\ & & h_{n-1} & d_{n-1} & h_{n-2} \\ & & & 0 & 1 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_{n-1} \\ s_n \end{bmatrix} = \begin{bmatrix} \sigma_a \\ \beta_2 \\ \vdots \\ \beta_{n-1} \\ \sigma_b \end{bmatrix} =: \mathbf{b}_1. \quad (6.13)$$

In many cases the knots are uniformly spaced, i.e.,  $h_i = h$  for all  $i$ . In that case  $\beta_i = 3(h\delta_{i-1} + h\delta_i) = 3(y_{i+1} - y_{i-1})$  and dividing both sides of (6.10) by  $h$ , (6.13)



**Figure 6.3.** Cubic spline interpolation to the data in Example 6.6. The breakpoints  $(x_i, y_i)$ ,  $i = 2, 3, 4$  are marked with dots on the curve.

takes the form

$$\begin{bmatrix} 1 & 0 & & & \\ 1 & 4 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & 4 & 1 \\ & & & 0 & 1 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_{n-1} \\ s_n \end{bmatrix} = \begin{bmatrix} \sigma_a \\ 6(y_3 - y_1)/(2h) \\ \vdots \\ 6(y_n - y_{n-2})/(2h) \\ \sigma_b \end{bmatrix}. \quad (6.14)$$

**Exercise 6.5** Show that

$$\frac{f(x+h) - f(x-h)}{2h} = f'(x) + \frac{h^2}{6} f^{(3)}(\eta), \quad x-h < \eta < x+h.$$

This is known as the central difference approximation to the first derivative.

**Example 6.6** Consider the data  $\mathbf{x} = [-1, -1/2, 0, 1/2, 1]$ ,  $\mathbf{y} = [-1, -0.9, 0, 0.9, 1]$  and  $\sigma_a = \sigma_b = 1/10$ . By (6.14)

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 1 & 0 & 0 \\ 0 & 1 & 4 & 1 & 0 \\ 0 & 0 & 1 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{bmatrix} = \frac{1}{10} \begin{bmatrix} 1 \\ 60 \\ 108 \\ 60 \\ 1 \end{bmatrix}.$$

with solution  $\mathbf{s} = [7, 64, 157, 64, 7]/70$ . A lengthy calculation gives the coefficients in the pp representation

$$\mathbf{C} = \begin{bmatrix} -70 & -63 & 0 & 63 \\ 7 & 64 & 157 & 64 \\ -72 & 186 & 0 & -186 \\ 172 & -124 & -124 & 172 \end{bmatrix} / 70.$$

The cubic spline interpolant is shown in Figure 6.3. Here Algorithm 6.3 was used with 200 uniform plot points.



**Exercise 6.7** Derive the  $pp$  representation (6.4) of  $g$  in Example 6.1.

### 6.1.3 Strictly Diagonally Dominant Matrices

The matrix  $\mathbf{N}_1$  in (6.13) is strictly diagonally dominant, i.e., in each row the absolute value of the diagonal element is strictly greater than the sum of the absolute values of the off diagonal elements in that row. In general a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  is said to be **strictly diagonally dominant** if  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$  for  $i = 1, \dots, n$ .

Recall that a linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a unique solution for any  $\mathbf{b}$  if and only if  $\mathbf{A}$  is nonsingular, that is  $\mathbf{A}\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{0}$ .

**Lemma 6.8** A strictly diagonally dominant matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  is nonsingular.

**Proof.** Let  $\mathbf{x}$  be any solution of  $\mathbf{A}\mathbf{x} = \mathbf{0}$  and let  $i$  be such that  $|x_i| = \max_j |x_j|$ . Then

$$\begin{aligned} 0 &= |(\mathbf{A}\mathbf{x})_i| = \left| \sum_j a_{ij}x_j \right| = |a_{ii}x_i + \sum_{j \neq i} a_{ij}x_j| \geq |a_{ii}x_i| - \left| \sum_{j \neq i} a_{ij}x_j \right| \\ &\geq |a_{ii}||x_i| - \sum_{j \neq i} |a_{ij}||x_j| \geq |a_{ii}||x_i| - \sum_{j \neq i} |a_{ij}||x_i| = |x_i|(|a_{ii}| - \sum_{j \neq i} |a_{ij}|). \end{aligned}$$

Since  $\mathbf{A}$  is strictly diagonally dominant the inequality  $|x_i|(|a_{ii}| - \sum_{j \neq i} |a_{ij}|) \leq 0$  implies  $|x_i| = 0$ . But then  $\mathbf{x} = \mathbf{0}$  and  $\mathbf{A}$  is nonsingular.  $\square$

We conclude that the matrix  $\mathbf{N}_1$  is nonsingular and we can find a  $C^2$  cubic spline interpolant  $g$  for any choice of  $x_1 < x_2 < \dots < x_n$  and  $n \geq 3$ . Moreover we have uniqueness. By derivation and nonsingularity of  $\mathbf{N}_1$  there is only one cubic  $C^2$  spline satisfying (6.1).

**Exercise 6.9** What kind of interpolant do we obtain for  $n = 2$ ?

**Exercise 6.10 Not-a-knot boundary condition.** Suppose  $n \geq 5$ . We can replace the 1. derivative boundary conditions  $g'(a) = \sigma_a$ ,  $g'(b) = \sigma_b$  by requiring  $p_1 = p_2$  and  $p_{n-1} = p_{n-2}$ . Since  $x_2$  and  $x_{n-1}$  are no longer interior knots we refer to  $p_1 = p_2$  and  $p_{n-1} = p_{n-2}$  as the not-a-knot conditions. The  $C^2$  spline  $g$  now consists of only  $n - 3$  pieces

$$g(x) := \begin{cases} p_2(x), & \text{if } x < x_3, \\ p_3(x), & \text{if } x_3 \leq x < x_4, \\ \vdots & \\ p_{n-3}(x), & \text{if } x_{n-3} \leq x < x_{n-2}, \\ p_{n-2}(x), & \text{if } x_{n-2} \leq x, \end{cases} \quad (6.15)$$

and we solve the interpolation problem

$$g(x_i) = y_i, \text{ for } i = 1, \dots, n. \quad (6.16)$$

With the shifted power form (6.5) we obtain the  $pp$  representation

$$\mathbf{x}_s = [x_2, \dots, x_{n-2}], \quad \mathbf{C} := \begin{bmatrix} c_{12} & c_{13} & \cdots & c_{1,n-2} \\ c_{22} & c_{23} & \cdots & c_{2,n-2} \\ c_{32} & c_{33} & \cdots & c_{3,n-2} \\ c_{42} & c_{43} & \cdots & c_{4,n-2} \end{bmatrix} \in \mathbb{R}^{4,n-3}.$$

Using  $p_2(a) = y_1$  and  $p_{n-2}(b) = y_n$  in addition to (6.10) for  $i = 3, \dots, n-2$  we obtain the linear system

$$\mathbf{N}_3 \tilde{\mathbf{s}} := \begin{bmatrix} d_2 & 2h_1 & & & \\ h_3 & d_3 & h_2 & & \\ & \ddots & \ddots & \ddots & \\ & & h_{n-2} & d_{n-2} & h_{n-3} \\ & & & 2h_{n-1} & d_{n-1} \end{bmatrix} \begin{bmatrix} s_2 \\ s_3 \\ \vdots \\ s_{n-2} \\ s_{n-1} \end{bmatrix} = \begin{bmatrix} \nu_2 \\ \beta_3 \\ \vdots \\ \beta_{n-2} \\ \nu_{n-1} \end{bmatrix}, \quad (6.17)$$

where

$$\begin{aligned} \nu_2 &:= \frac{2h_2^2}{h_1 + h_2} \delta_1 + \frac{2h_1(2h_1 + 3h_2)}{h_1 + h_2} \delta_2, \\ \nu_{n-1} &:= \frac{2h_{n-2}^2}{h_{n-1} + h_{n-2}} \delta_{n-1} + \frac{2h_{n-1}(2h_{n-1} + 3h_{n-2})}{h_{n-1} + h_{n-2}} \delta_{n-2}. \end{aligned} \quad (6.18)$$

Show that the requirements  $p_2(a) = y_1$  and  $p_{n-2}(b) = y_n$  give the first and last equation in (6.17). Explain why  $\mathbf{N}_3$  is nonsingular. Note that  $s_{n-1}$  is not needed for the  $pp$  representation of  $g$ .

**Exercise 6.11 2. derivative boundary conditions.** Let  $m_a, m_b \in \mathbb{R}$ . Consider finding a cubic  $C^2$  spline  $g$  satisfying

$$g(x_i) = y_i, \text{ for } i = 1, \dots, n, \quad g''(a) = m_a, \quad g''(b) = m_b. \quad (6.19)$$

We call  $g''(a) = m_a, g''(b) = m_b$  for 2. derivative boundary conditions. Show, using (6.12) for the first and last equation that we obtain the linear system

$$\mathbf{N}_2 \mathbf{s} = \begin{bmatrix} 2 & 1 & & & \\ h_2 & d_2 & h_1 & & \\ & \ddots & \ddots & \ddots & \\ & & h_{n-1} & d_{n-1} & h_{n-2} \\ & & & 1 & 2 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_{n-1} \\ s_n \end{bmatrix} = \begin{bmatrix} \gamma_1 \\ \beta_2 \\ \vdots \\ \beta_{n-1} \\ \gamma_n \end{bmatrix}, \quad (6.20)$$

where

$$\gamma_1 := 3\delta_1 - h_1 m_a / 2, \quad \gamma_n := 3\delta_{n-1} + h_{n-1} m_b / 2. \quad (6.21)$$

The matrix  $\mathbf{N}_2$  is strictly diagonally dominant and therefore nonsingular.

## 6.2 The Second Derivative Matrix

Before discussing algorithms for solving the linear systems in the previous section we present briefly another problem that leads to a linear system with a tridiagonal coefficient matrix.

Consider the simple two point boundary value problem

$$-u''(x) = f(x), \quad x \in [0, 1], \quad u(0) = 0, \quad u(1) = 0, \quad (6.22)$$

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 (6.22) 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.

(6.22) can be solved approximately using the **finite difference method**. For this we choose a positive integer  $m$ , define the discretization parameter  $h := 1/(m+1)$ , and replace the interval  $[0, 1]$  by grid points  $x_j := jh$  for  $j = 0, 1, \dots, m+1$ . We then use the following finite difference approximation of the second derivative:

$$\frac{u(x-h) - 2u(x) + u(x+h)}{h^2} = u''(x) + \frac{h^2}{12}u^{(4)}(\xi),$$

for some  $\xi \in (x-h, x+h)$ . This can be shown by a Taylor expansion provided  $u \in C^4[x-h, x+h]$ .

We 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

$$\mathbf{T}\mathbf{v} = \begin{bmatrix} 2 & -1 & 0 & & & \\ -1 & 2 & -1 & & & \\ & 0 & \ddots & \ddots & \ddots & \\ & & & -1 & 2 & -1 \\ & & & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_{m-1} \\ v_m \end{bmatrix} = h^2 \begin{bmatrix} f(h) \\ f(2h) \\ \vdots \\ f((m-1)h) \\ f(mh) \end{bmatrix} =: \mathbf{b}. \quad (6.23)$$

The matrix  $\mathbf{T}$  is called the **second derivative matrix**. This matrix is diagonally dominant (cf. Definition 6.14), but not strictly diagonally dominant.

## 6.3 LU Factorization of a Tridiagonal System

Given a linear system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A} = \text{tridiag}(a_i, d_i, c_i) \in \mathbb{R}^{n,n}$  is nonsingular and tridiagonal. In order to solve  $\mathbf{Ax} = \mathbf{b}$  we try to construct triangular matrices

$\mathbf{L}$  and  $\mathbf{R}$  such that the product  $\mathbf{A} = \mathbf{L}\mathbf{R}$  has the form

$$\begin{bmatrix} d_1 & c_1 & & & \\ a_2 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-1} & d_{n-1} & c_{n-1} \\ & & & a_n & d_n \end{bmatrix} = \begin{bmatrix} 1 & & & & \\ l_2 & 1 & & & \\ & \ddots & \ddots & & \\ & & l_n & 1 & \end{bmatrix} \begin{bmatrix} r_1 & c_1 & & & \\ & \ddots & \ddots & & \\ & & r_{n-1} & c_{n-1} & \\ & & & r_n & \end{bmatrix}. \quad (6.24)$$

Note that  $\mathbf{L}$  has ones on the diagonal, and that we can use the same  $c_i$  elements on the super-diagonals of  $\mathbf{A}$  and  $\mathbf{R}$ . By equating elements in (6.24) we find

$$d_1 = r_1, \quad a_k = l_k r_{k-1}, \quad d_k = l_k c_{k-1} + r_k, \quad k = 2, 3, \dots, n.$$

Solving for  $l_k$  and  $r_k$  we see that the  $l$ 's and the  $r$ 's can be determined recursively

$$r_1 = d_1, \quad l_k = \frac{a_k}{r_{k-1}}, \quad r_k = d_k - l_k c_{k-1}, \quad k = 2, 3, \dots, n. \quad (6.25)$$

The solution of the system  $\mathbf{A}\mathbf{x} = \mathbf{L}(\mathbf{R}\mathbf{x}) = \mathbf{b}$  is found by solving  $\mathbf{L}\mathbf{y} = \mathbf{b}$  and  $\mathbf{R}\mathbf{x} = \mathbf{y}$  as follows:

$$\begin{aligned} y_1 &= b_1, & y_k &= b_k - l_k y_{k-1}, & k &= 2, 3, \dots, n, \\ x_n &= y_n / r_n, & x_k &= (y_k - c_k x_{k+1}) / r_k, & k &= n-1, \dots, 2, 1. \end{aligned} \quad (6.26)$$

The following algorithms factor and solve, if possible, a tridiagonal system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} = \text{tridiag}(\mathbf{a}, \mathbf{d}, \mathbf{c})$  is given by (6.24).

**Algorithm 6.12 (trifactor)** Vectors  $\mathbf{l}, \mathbf{r} \in \mathbb{C}^n$  are computed from  $\mathbf{a}, \mathbf{c}, \mathbf{d} \in \mathbb{C}^n$ . This implements the LU factorization of a tridiagonal matrix. The first (dummy) component in  $\mathbf{a}$  and last component of  $\mathbf{c}$  are not used.

```
function [l,r]=trifactor(a,d,c)
r=d; l=d;
for k=2:length(d)
    l(k)=a(k)/r(k-1);
    r(k)=d(k)-l(k)*c(k-1);
end
```

**Algorithm 6.13 (trisolve)** The solution  $x$  of the tridiagonal system  $\mathbf{LR}x = \mathbf{b}$  is found from (6.26). Here  $\mathbf{l}, \mathbf{r}, \mathbf{b} \in \mathbb{C}^n$ . The vectors  $\mathbf{l}, \mathbf{r}$  are typically output from **trifactor**.

```
function x=trisolve(l,r,c,b)
x=b; n=length(b);
for k=2:n
    x(k)=b(k)-l(k)*x(k-1);
end
x(n)=x(n)/r(n);
for k=n-1:-1:1
    x(k)=(x(k)-c(k)*x(k+1))/r(k);
end
```

### 6.3.1 Diagonal Dominance

We show that Algorithms 6.12, 6.13 are well defined for a class of tridiagonal matrices containing  $\mathbf{N}_1$  and  $\mathbf{T}$ .

**Definition 6.14** *The tridiagonal matrix*

$$\mathbf{A} := \text{tridiag}(a_i, d_i, c_i) = \begin{bmatrix} d_1 & c_1 & & & \\ a_2 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-1} & d_{n-1} & c_{n-1} \\ & & & a_n & d_n \end{bmatrix} \quad (6.27)$$

is **diagonally dominant** if

$$|d_1| > |c_1|, \quad |d_k| \geq |a_k| + |c_k|, \quad k = 2, \dots, n, \quad (6.28)$$

where  $c_n := 0$ .

Clearly the matrices  $\mathbf{N}_1$  and  $\mathbf{T}$  are diagonally dominant. The matrix  $\mathbf{A}_1 = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 2 \end{bmatrix}$  is also diagonally dominant, but it is singular. So diagonal dominance is not enough to guarantee nonsingularity.

**Theorem 6.15** *Suppose  $\mathbf{A}$  given by (6.27) is tridiagonal and diagonally dominant. If in addition  $a_i \neq 0$  for  $i = 2, \dots, n$  then  $\mathbf{A}$  is nonsingular and has a unique LU factorization given by (6.25).*

**Proof.** We show by induction on  $k$  that  $|r_k| > |c_k|$  for  $k = 1, \dots, n$ . Clearly  $|r_1| = |d_1| > |c_1|$ . Suppose  $|c_{k-1}|/|r_{k-1}| < 1$  for some  $2 \leq k \leq n$ . By (6.25) and diagonal dominance

$$|r_k| = |d_k - l_k c_{k-1}| = |d_k - \frac{a_k c_{k-1}}{r_{k-1}}| \geq |d_k| - \frac{|a_k| |c_{k-1}|}{|r_{k-1}|} > |d_k| - |a_k| \geq |c_k|.$$

Thus  $r_k \neq 0$  for  $k = 1, \dots, n$  and an LU factorization exists. It is unique since any LU factorization must satisfy (6.25). Since both  $\mathbf{L}$  and  $\mathbf{R}$  have nonzero diagonal elements they are nonsingular (cf. Lemma 6.32). Thus  $\mathbf{A} = \mathbf{LR}$  is nonsingular.  $\square$

The requirement  $r_n \neq 0$  is necessary for nonsingularity, but not for the LU factorization.

**Corollary 6.16** *Suppose that the conditions in Theorem 6.15 hold except that we assume  $a_n = 0$ . Then  $\mathbf{A}$  has a unique LU factorization. Moreover,  $\mathbf{A}$  is nonsingular if  $d_n \neq 0$ .*

**Proof.** By the proof of Theorem 6.15 we have  $|r_k| > |c_k|$  for  $k < n$  and  $r_n = d_n$ . So the formulas (6.25) are well defined,  $\mathbf{A}$  has a unique LU factorization, and  $r_n \neq 0$  if  $d_n \neq 0$ .  $\square$

The number of floating point operations (flops) to compute the LU factorization of a tridiagonal matrix using Algorithm 6.12 is only  $3n - 3$ , while the number of flops for Algorithm 6.13 is  $5n - 4$ . This means that the number of flops (the complexity) to solve a tridiagonal system is  $8n - 7$ , and this number only grows linearly with  $n$ . This should be compared to Gaussian elimination on a full  $n \times n$  system which is an  $O(n^3)$  process.

We could solve the system  $\mathbf{T}\mathbf{v} = \mathbf{b}$  by using the inverse  $\mathbf{T}^{-1}$  of  $\mathbf{T}$  and simply compute the matrix vector product  $\mathbf{v} = \mathbf{T}^{-1}\mathbf{b}$ . However this is not a good idea. In fact, all elements in  $\mathbf{T}^{-1}$  are nonzero. See Exercise 6.18.

**Exercise 6.17** *Show that the LU factorization of  $\mathbf{T}$  is  $\mathbf{T} = \mathbf{LR}$ , where*

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ -\frac{1}{2} & 1 & \ddots & & \vdots \\ 0 & -\frac{2}{3} & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & -\frac{m-1}{m} & 1 \end{bmatrix}, \mathbf{R} = \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ 0 & \frac{3}{2} & -1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \frac{m}{m-1} & -1 \\ 0 & \cdots & \cdots & 0 & \frac{m+1}{m} \end{bmatrix}. \quad (6.29)$$

**Exercise 6.18** *Show that the inverse of  $\mathbf{T}$  is*

$$(\mathbf{T}^{-1})_{i,j} = (\mathbf{T}^{-1})_{j,i} = \left(1 - \frac{i}{m+1}\right) * j, \quad 1 \leq j \leq i \leq m. \quad (6.30)$$

### 6.3.2 Periodic Boundary Conditions

Given  $a = x_1 < \dots < x_n = b$  and  $(y_j)_{j=1}^n$  with the periodic condition  $y_1 = y_n$ . For all  $j \in \mathbb{Z}$  we define the periodic extensions  $x_{j+n-1} = x_j + b - a$  and  $y_{j+n-1} = y_j$ . Consider the interpolation problem

$$g(x_j) = y_j, \quad j = 1, \dots, n, \quad g'(a) = g'(b), \quad g''(a) = g''(b). \quad (6.31)$$

We extend  $g$  by periodicity so that  $g : \mathbb{R} \rightarrow \mathbb{R}$  is  $b - a$  periodic, i.e.,  $g(x + b - a) = g(x)$  for all  $x \in \mathbb{R}$ . With these extensions we can use the  $C^2$  equations (6.10) for  $i = 2, \dots, n$ . Since  $s_1 = s_n$  and  $s_{n+1} = s_2$  we obtain the system

$$\mathbf{N}_p \mathbf{s} = \begin{bmatrix} d_2 & h_1 & 0 & h_2 \\ h_3 & d_3 & h_2 & 0 \\ & \ddots & \ddots & \ddots \\ 0 & & h_{n-1} & d_{n-1} & h_{n-2} \\ h_{n-1} & & 0 & h_n & d_n \end{bmatrix} \begin{bmatrix} s_2 \\ s_3 \\ \vdots \\ s_{n-1} \\ s_n \end{bmatrix} = \begin{bmatrix} \beta_2 \\ \beta_3 \\ \vdots \\ \beta_{n-1} \\ \beta_n \end{bmatrix}, \quad (6.32)$$

where  $h_n := h_1$  and  $\delta_n := \delta_1$ . This system is strictly diagonally dominant, but not tridiagonal. For an LU factorization see Exercise 6.20.

Periodic boundary conditions are useful for modeling a closed parametric plane curve  $\gamma(u) = [x(u), y(u)]$  or a space curve  $\gamma(u) = [x(u), y(u), z(u)]$ ,  $u \in [a, b]$ . First pick knots  $a = u_1 < u_2 < \dots < u_n = b$  in the parameter interval  $[a, b]$ . Then compute a cubic  $C^2$  spline vector  $\mathbf{g} = [g_x, g_y]$  or  $\mathbf{g} = [g_x, g_y, g_z]$  interpolating the points  $(u_i, \gamma(u_i))_{i=1}^n$  with boundary conditions  $\mathbf{g}'(a) = \mathbf{g}'(b)$  and  $\mathbf{g}''(a) = \mathbf{g}''(b)$ .

**Example 6.19** Consider the unit circle

$$\gamma(u) := (\cos(2\pi u), \sin(2\pi u)), \quad u \in [0, 1].$$

To obtain a  $C^2$  cubic spline approximation to this circle we use the uniform data sites

$$\mathbf{u} = [u_1, \dots, u_n], \quad u_j = (j-1)/(n-1), \quad j = 1, \dots, n.$$

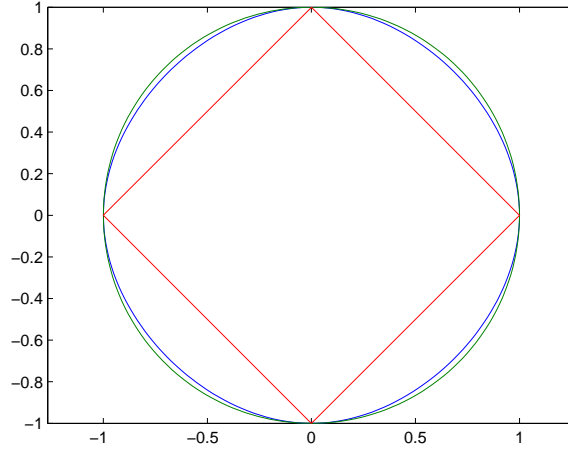
Let  $\mathbf{g} = (g_x, g_y)$  be a vector of cubic  $C^2$ -splines satisfying  $\mathbf{g}(u_j) = \gamma(u_j)$ , for  $j = 1, \dots, n$ . Solving (6.32) for each component  $g_x, g_y$  we obtain the pp representation of  $\mathbf{g}$  as  $(\mathbf{u}_s, \mathbf{C}_x, \mathbf{C}_y)$ . For  $n = 5$  we find  $\mathbf{u}_s = [0, 1/4, 1/2, 3/4]$  and

$$\mathbf{C}_x = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & -6 & 0 & 6 \\ -24 & 0 & 24 & 0 \\ 32 & 32 & -32 & -32 \end{bmatrix}, \quad \mathbf{C}_y = \begin{bmatrix} 0 & 1 & 0 & -1 \\ 6 & 0 & -6 & 0 \\ 0 & -24 & 0 & 24 \\ -32 & 32 & 32 & -32 \end{bmatrix}.$$

We show  $\mathbf{g}$  together with the exact circle  $\gamma$  and the piecewise linear interpolant to the data in Figure 6.4. The computations to obtain the data for this plot are outlined in the following series of exercises.

**Exercise 6.20** Given a strictly diagonally dominant matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  of the form (6.32) illustrated here for  $n = 5$

$$\mathbf{A} = \begin{bmatrix} d_1 & c_1 & 0 & 0 & a_1 \\ a_2 & d_2 & c_2 & 0 & 0 \\ 0 & a_3 & d_3 & c_3 & 0 \\ 0 & 0 & a_4 & d_4 & c_4 \\ c_5 & 0 & 0 & a_5 & d_5 \end{bmatrix}$$



**Figure 6.4.** A 4 piece periodic cubic spline interpolant to the unit circle.

Such a matrix has an  $LU$  factorization of the form

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ l_2 & 1 & 0 & 0 & 0 \\ 0 & l_3 & 1 & 0 & 0 \\ 0 & 0 & l_4 & 1 & 0 \\ g_1 & g_2 & g_3 & g_4 & 1 \end{bmatrix} \begin{bmatrix} r_1 & c_1 & 0 & 0 & e_1 \\ 0 & r_2 & c_2 & 0 & e_2 \\ 0 & 0 & r_3 & c_3 & e_3 \\ 0 & 0 & 0 & r_4 & e_4 \\ 0 & 0 & 0 & 0 & r_5 \end{bmatrix},$$

where (you do not have to prove this)

$$\begin{aligned} r_1 &= d_1; \quad e_1 = a_1; \quad g_1 = c_1/r_1; \\ l_i &= a_i/r_{i-1}, \quad i = 2, \dots, n-1, \\ r_i &= d_i - c_{i-1}l_i, \quad i = 2, \dots, n-1, \\ e_i &= -e_{i-1}l_i, \quad i = 2, \dots, n-2, \\ g_i &= -c_{i-1}g_{i-1}/r_i, \quad i = 2, \dots, n-2, \\ e_{n-1} &= c_{n-1} - e_{n-2}l_{n-1}, \\ g_{n-1} &= (a_n - c_{n-2}g_{n-2})/r_{n-1}, \\ r_n &= d_n - \sum_{i=1}^{n-1} g_i e_i. \end{aligned} \tag{6.33}$$

Since  $\mathbf{A}$  is strictly diagonally dominant it can be shown that  $|r_i| > |c_i|$  for  $i = 1, \dots, n-1$  and hence the  $LU$  factorization exists.



We can solve the system  $\mathbf{LU} = \mathbf{b}$  as follows

$$\begin{aligned}
 y_1 &= b_1, \\
 y_i &= b_i - l_i y_{i-1}, \quad i = 2, \dots, n-1, \\
 y_n &= b_n - \sum_{i=1}^{n-1} g_i y_i, \\
 x_n &= y_n / c_n, \quad x_{n-1} = (y_{n-1} - e_{n-1} x_n) / c_{n-1}, \\
 x_i &= (y_i - c_i x_{i+1} - e_i x_n) / r_i, \quad i = n-2, n-3, \dots, 1.
 \end{aligned} \tag{6.34}$$

It follows that we can solve the system  $\mathbf{Ax} = \mathbf{b}$  in  $O(n)$  flops.

- (a) Write a function `[l,r,e,g]=trifactorp(a,d,c)` that returns the vectors given by (6.33). Verify your code by computing  $\mathbf{A} - \mathbf{L} * \mathbf{U} \in \mathbb{R}^{5,5}$  for  $\mathbf{a} = \mathbf{c} = \mathbf{d}/4 = [1, 1, 1, 1, 1]^T$ . When this works try it on random vectors.
- (b) Write a function `x=trisolvep(l,r,c,e,g,b)` that uses the output from `trifactorp` to solve  $\mathbf{LRx} = \mathbf{b}$  using (6.34). Verify the code by first using the integer data from (a) and  $\mathbf{b} = [6, 6, 6, 6, 6]^T$  and then with  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}$  random vectors.

**Exercise 6.21** Given data  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $a = x_1 < \dots < x_n = b$  and  $y_1 = y_n$ . Write a function `C=cubsplinep(x,y)` that finds the coefficients  $\mathbf{C} \in \mathbb{R}^{4,n-1}$  of a periodic cubic spline satisfying (6.2), (6.5), and (6.31). Use the linear solvers in Exercise 6.20 for the system (6.32). Test the program on  $\mathbf{x} = [0, 1/4, 2/4, 3/4, 1]$  and  $y = \cos(2\pi x)$ .

**Exercise 6.22** In this exercise we want to fit a periodic cubic spline curve  $\mathbf{g} = [g_x, g_y]$  to periodic data. Given data  $\mathbf{u}, \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $a = u_1 < \dots < u_n = b$ ,  $x_1 = x_n$  and  $y_1 = y_n$ . We can think of the data as representing a closed parametric curve  $\gamma(u) = [x(u), y(u)]$ . The goal is to compute points  $[\mathbf{q}_x, \mathbf{q}_y] = \mathbf{g}(\mathbf{r})$  on  $\mathbf{g}$ , where  $\mathbf{r}$  is a vector of  $u$ -values. Write a function `[qx,qy]=cubsplinecurvep(u,x,y,r)` that first finds the coefficients  $\mathbf{C}_x, \mathbf{C}_y \in \mathbb{R}^{4,n-1}$  interpolating the two data sets  $(\mathbf{u}, \mathbf{x})$  and  $(\mathbf{u}, \mathbf{y})$ . For simplicity you can call the function in Exercise 6.21 twice. Then you can use Algorithm 6.3 twice to find  $\mathbf{q}_x$  and  $\mathbf{q}_y$ .

**Exercise 6.23** To do this exercise you should first complete Exercises 6.20, 6.21, and 6.22.

- (a) Write a function `circle(n)` that interpolates to a circle with a periodic spline curve  $\mathbf{g}$  as in Example 6.19.
- (b) Make a plot for  $n = 5$  as in Figure 6.4. Make another plot for  $n = 9$ .

**Exercise 6.24** We consider a finite difference method for the two point boundary value problem

$$\begin{aligned} -u''(x) + r(x)u'(x) + q(x)u(x) &= f(x), \text{ for } x \in [a, b], \\ u(a) &= g_0, \quad u(b) = g_1. \end{aligned} \quad (6.35)$$

We assume that the given functions  $f, g$  and  $r$  are continuous on  $[a, b]$  and that  $q(x) \geq 0$  for  $x \in [a, b]$ . It can then be shown that (6.35) has a unique solution  $u$ .

To solve (6.35) 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

$$\frac{-v_{j-1} + 2v_j - v_{j+1}}{h^2} + r(x_j)\frac{v_{j+1} - v_{j-1}}{2h} + q(x_j)v_j = f(x_j), \quad j = 1, \dots, m, \quad (6.36)$$

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

- (a) Show that (6.36) leads to a tridiagonal linear system  $\mathbf{A}\mathbf{v} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{m,m}$  given by (6.27) has elements

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

and

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

- (b) Show that the linear system satisfies the conditions in Theorem 6.15 if the spacing  $h$  is so small that  $\frac{h}{2}|r(x)| < 1$  for all  $x \in [a, b]$ .
- (c) Propose a method based on Algorithms 6.12, 6.13 to find  $v_1, \dots, v_m$ .
- (d) Consider the problem (6.35) 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 (6.36) 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  $h$ .
- (e) Make a combined plot of the solution  $u$  and the computed points  $v_j$ ,  $j = 0, \dots, m+1$  for  $h = 0.1$ .
- (f) One can show that the error is proportional to  $h^p$  for some integer  $p$ . Estimate  $p$  based on the error for  $h = 0.1, 0.05, 0.025, 0.0125$ .

## 6.4 Block Multiplication and Triangular Matrices

### 6.4.1 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

$$(i) \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \left[ \begin{array}{c|cc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right], \quad (ii) [\mathbf{a}_{.1}, \mathbf{a}_{.2}, \mathbf{a}_{.3}] = \left[ \begin{array}{c|cc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right],$$

$$(iii) \begin{bmatrix} \mathbf{a}_{1.}^T \\ \mathbf{a}_{2.}^T \\ \mathbf{a}_{3.}^T \end{bmatrix} = \left[ \begin{array}{ccc|c} 1 & 2 & 3 & \\ \hline 4 & 5 & 6 & \\ 7 & 8 & 9 & \end{array} \right], \quad (iv) [\mathbf{A}_{11}, \mathbf{A}_{12}] = \left[ \begin{array}{cc|c} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right].$$

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  $\mathbf{A} \in \mathbb{C}^{m,p}$  and  $\mathbf{B} \in \mathbb{C}^{p,n}$ . 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  $i = 1, \dots, p$ .

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

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

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

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

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  $\mathbf{A}_1 \in \mathbb{C}^{k,p}$  and  $\mathbf{A}_2 \in \mathbb{C}^{m-k,p}$  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  $\mathbf{A}_1 \in \mathbb{C}^{m,s}$ ,  $\mathbf{A}_2 \in \mathbb{C}^{m,p-s}$ ,  $\mathbf{B}_1 \in \mathbb{C}^{s,n}$  and  $\mathbf{B}_2 \in \mathbb{C}^{p-s,n}$  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,  $(\mathbf{A}\mathbf{B})_{ij} = \sum_{k=1}^p a_{ik} b_{kj} = \sum_{k=1}^s a_{ik} b_{kj} + \sum_{k=s+1}^p a_{ik} b_{kj} = (\mathbf{A}_1 \mathbf{B}_1)_{ij} + (\mathbf{A}_2 \mathbf{B}_2)_{ij} = (\mathbf{A}_1 \mathbf{B}_1 + \mathbf{A}_2 \mathbf{B}_2)_{ij}$ .

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 similar for the other blocks. To show this we use Rule 4. to obtain

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

We complete the proof using Rules 5. and 6.

8. For the general case see Section 3.1.1.

**Exercise 6.25** For any matrix  $\mathbf{A}$  show that  $a_{ij} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_j$  for all  $i, j$ .

**Exercise 6.26** Let  $\mathbf{B} = \mathbf{A}^T$ . Explain why this product is defined for any matrix. Show that  $b_{ij} = \langle \mathbf{a}_i, \mathbf{a}_j \rangle := \mathbf{a}_i^T \mathbf{a}_j$  for all  $i, j$ .

**Exercise 6.27** For  $\mathbf{A} \in \mathbb{R}^{m,n}$  and  $\mathbf{B} \in \mathbb{R}^{p,n}$  show that

$$\mathbf{A}\mathbf{B}^T = \mathbf{a}_{.1} \mathbf{b}_{.1}^T + \mathbf{a}_{.2} \mathbf{b}_{.2}^T + \cdots + \mathbf{a}_{.n} \mathbf{b}_{.n}^T.$$

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

**Exercise 6.28** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{B} \in \mathbb{R}^{m,p}$ , and  $\mathbf{X} \in \mathbb{R}^{n,p}$ . Show that

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

**Exercise 6.29** 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 6.30** Suppose  $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n,n}$  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  $\mathbf{A}_1, \mathbf{B}_1, \mathbf{C}_1 \in \mathbb{R}^{n-1, n-1}$ . Show that

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

### 6.4.2 Triangular matrices

Recall that a matrix  $\mathbf{R}$  is upper- or right triangular if  $r_{ij} = 0$  for  $i > j$ , and a matrix  $\mathbf{L}$  is lower- or left triangular if  $l_{ij} = 0$  for  $i < j$ . If  $\mathbf{R}$  is upper triangular then  $\mathbf{R}^T$  is lower triangular.

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

**Lemma 6.31** *Suppose*

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

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

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{A}_{22}^{-1} \end{bmatrix} \quad (6.37)$$

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

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

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

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

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 6.32** *An upper (lower) triangular matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{n,n}$  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  $\mathbf{A} \in \mathbb{C}^{k+1,k+1}$  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  $\mathbf{A}_k \in \mathbb{C}^{k,k}$  is upper triangular. By Lemma 1.1  $\mathbf{A}$  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 6.33** *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 6.34** *For a unit upper (lower) triangular matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$ :*

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 1.2, while Lemma 1.3 implies 2.  $\square$

## Chapter 7

# LU Factorizations

In Chapter 6 we saw how an LU factorization of the coefficient matrix can be used to solve certain tridiagonal systems efficiently. In this chapter we consider the general theory of LU factorizations<sup>1</sup>. We consider some related factorizations called block LU, PLU, symmetric LU, and Cholesky.

### 7.1 The LU Factorization

We say that  $\mathbf{A} = \mathbf{L}\mathbf{R}$  is an **LU factorization** of  $\mathbf{A} \in \mathbb{C}^{n,n}$  if  $\mathbf{L} \in \mathbb{C}^{n,n}$  is lower triangular (**left triangular**) and  $\mathbf{R} \in \mathbb{C}^{n,n}$  is upper triangular (**right triangular**). In addition we will assume that  $\mathbf{L}$  is unit triangular, i.e., it has ones on the diagonal. The LU factorization of the 2. derivative matrix  $\mathbf{T}$  was given in (6.29). But not every nonsingular matrix has an LU factorization.

**Example 7.1** An LU factorization of  $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$  must satisfy the equations

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_1 & 1 \end{bmatrix} \begin{bmatrix} r_1 & r_3 \\ 0 & r_2 \end{bmatrix} = \begin{bmatrix} r_1 & r_3 \\ l_1 r_1 & l_1 r_3 + r_2 \end{bmatrix}$$

for the unknowns  $l_1$  in  $\mathbf{L}$  and  $r_1, r_2, r_3$  in  $\mathbf{R}$ . Comparing  $(1, 1)$  elements we see that  $r_1 = 0$ , which makes it impossible to satisfy the condition  $1 = l_1 r_1$  for the  $(2, 1)$  element. We conclude that  $\mathbf{A}$  has no LU factorization.

We will make use of some special submatrices.

**Definition 7.2** For  $k = 1, \dots, n$  the matrices  $\mathbf{A}_k \in \mathbb{C}^{k,k}$  given by

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

---

<sup>1</sup>In the literature an upper triangular matrix is denoted by  $\mathbf{U}$  in an LU factorization and  $\mathbf{R}$  in a QR factorization. (see Chapter 16). We have chosen to use  $\mathbf{R}$  to denote an upper triangular matrix both for LU and QR factorizations.

are called the **leading principal submatrices** of  $\mathbf{A} = \mathbf{A}_n \in \mathbb{C}^{n,n}$ . More generally, a matrix  $\mathbf{B} \in \mathbb{C}^{k,k}$  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 example, for  $k = 1$  the only principal submatrices are the diagonal elements of  $\mathbf{A}$ .

**Example 7.3** The principal submatrices of  $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$  are

$$[1], [5], [9], \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix}, \begin{bmatrix} 1 & 3 \\ 7 & 9 \end{bmatrix}, \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}, \mathbf{A}.$$

The leading principal submatrices are

$$[1], \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix}, \mathbf{A}.$$

**Theorem 7.4** Suppose the leading principal submatrices  $\mathbf{A}_k$  of  $\mathbf{A} \in \mathbb{C}^{n,n}$  are nonsingular for  $k = 1, \dots, n-1$ . Then  $\mathbf{A}$  has a unique LU factorization.

**Proof.** 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  $\mathbf{A}_{n-1} = \mathbf{L}_{n-1}\mathbf{R}_{n-1}$ , and that  $\mathbf{A}_1, \dots, \mathbf{A}_{n-1}$  are nonsingular. Since  $\mathbf{A}_{n-1}$  is nonsingular it follows that  $\mathbf{L}_{n-1}$  and  $\mathbf{R}_{n-1}$  are nonsingular. But then

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{n-1} & \mathbf{b} \\ \mathbf{c}^T & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{n-1} & \mathbf{0} \\ \mathbf{c}^T \mathbf{R}_{n-1}^{-1} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{R}_{n-1} & \mathbf{v} \\ 0 & a_{nn} - \mathbf{c}^T \mathbf{R}_{n-1}^{-1} \mathbf{v} \end{bmatrix} = \mathbf{L}\mathbf{R}, \quad (7.1)$$

where  $\mathbf{v} = \mathbf{L}_{n-1}^{-1} \mathbf{b}$ , and this is an LU factorization of  $\mathbf{A}$ . Since  $\mathbf{L}_{n-1}$  and  $\mathbf{R}_{n-1}$  are nonsingular the block (2,1) element in  $\mathbf{L}$  and the block (1,2) element in  $\mathbf{R}$  are uniquely given in (7.1), and then  $r_{nn}$  is also determined uniquely. Thus the LU factorization is unique by construction.  $\square$

The following observation is useful.

**Lemma 7.5** Suppose  $\mathbf{A} = \mathbf{L}\mathbf{R}$  is an LU factorization of  $\mathbf{A} \in \mathbb{C}^{n,n}$ . For  $k = 1, \dots, n$  let  $\mathbf{A}_k, \mathbf{L}_k, \mathbf{R}_k$  be the leading principal submatrices of  $\mathbf{A}, \mathbf{L}, \mathbf{R}$ , respectively. Then  $\mathbf{A}_k = \mathbf{L}_k \mathbf{R}_k$  is an LU factorization of  $\mathbf{A}_k$  for  $k = 1, \dots, n$ .

**Proof.** For  $k = 1, \dots, n-1$  we partition  $\mathbf{A} = \mathbf{L}\mathbf{R}$  as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_k & \mathbf{B}_k \\ \mathbf{C}_k & \mathbf{D}_k \end{bmatrix} = \begin{bmatrix} \mathbf{L}_k & \mathbf{0} \\ \mathbf{M}_k & \mathbf{N}_k \end{bmatrix} \begin{bmatrix} \mathbf{R}_k & \mathbf{S}_k \\ 0 & \mathbf{T}_k \end{bmatrix} = \mathbf{L}\mathbf{R}, \quad (7.2)$$



where  $\mathbf{D}_k, \mathbf{N}_k, \mathbf{T}_k \in \mathbb{C}^{n-k, n-k}$ . Using block multiplication we find  $\mathbf{A}_k = \mathbf{L}_k \mathbf{R}_k$ . Since  $\mathbf{L}_k$  is unit lower triangular and  $\mathbf{R}_k$  is upper triangular we see that this gives an LU factorization of  $\mathbf{A}_k$ .  $\square$

There is a converse of Theorem 7.4.

**Theorem 7.6** *Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  has an LU factorization. If  $\mathbf{A}$  is nonsingular then the leading principal submatrices  $\mathbf{A}_k$  are nonsingular for  $k = 1, \dots, n-1$ .*

**Proof.** Suppose  $\mathbf{A}$  is nonsingular with the LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{R}$ . Since  $\mathbf{A}$  is nonsingular it follows that  $\mathbf{L}$  and  $\mathbf{R}$  are nonsingular. Let  $1 \leq k \leq n$ . By Lemma 7.5 it follows that  $\mathbf{A}_k = \mathbf{L}_k \mathbf{R}_k$ . Since  $\mathbf{L}_k$  is unit lower triangular it is nonsingular. Moreover  $\mathbf{R}_k$  is nonsingular since its diagonal elements are among the nonzero diagonal elements of  $\mathbf{R}$ . But then  $\mathbf{A}_k$  is nonsingular.  $\square$

The following lemma shows that the LU factorization of a nonsingular matrix is unique.

**Corollary 7.7** *The LU factorization of a nonsingular matrix is unique whenever it exists.*

**Proof.** By Theorem 7.6 the leading principal submatrices are nonsingular for  $k = 1, \dots, n-1$ . But then by Theorem 7.4 the LU factorization is unique.  $\square$

**Remark 7.8** *Theorem 7.6 is not true in general if  $\mathbf{A}$  is singular. An LU factorization of an upper triangular matrix  $\mathbf{A}$  is  $\mathbf{A} = \mathbf{I}\mathbf{A}$ , and if  $\mathbf{A}$  is singular it can have zeros anywhere on the diagonal. By Lemma 6.32, if some  $a_{kk}$  is zero then  $\mathbf{A}_k$  is singular.*

**Remark 7.9** *The LU factorization of a singular matrix need not be unique. In particular, for the zero matrix any unit lower triangular matrix can be used as  $\mathbf{L}$  in an LU factorization.*

**Remark 7.10** *We have shown that a nonsingular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  has an 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  $\mathbf{P}^T \mathbf{A} = \mathbf{L}\mathbf{R}$  or equivalently  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{R}$ , where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is unit lower triangular, and  $\mathbf{R}$  is upper triangular. We call this a **PLU factorization** of  $\mathbf{A}$ . (Cf. Section 7.7 and Appendix A.)*

**Exercise 7.11** *Show that  $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  has an LU factorization. Note that we have only interchanged rows in Example 7.1*

**Exercise 7.12** Find an LU factorization of the singular matrix  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ . Is it unique?

**Exercise 7.13** Suppose  $\mathbf{A}$  has an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{R}$ . Show that  $\det(\mathbf{A}_k) = r_{11}r_{22}\cdots r_{kk}$  for  $k = 1, \dots, n$ .

**Exercise 7.14** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  and  $\mathbf{A}_k$  is nonsingular for  $k = 1, \dots, n-1$ . Use Exercise 7.13 to show that the diagonal elements  $r_{kk}$  in the LU factorization are

$$r_{11} = a_{11}, \quad r_{kk} = \frac{\det(\mathbf{A}_k)}{\det(\mathbf{A}_{k-1})}, \quad \text{for } k = 2, \dots, n. \quad (7.3)$$

## 7.2 Block LU Factorization

Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is a block matrix of the form

$$\mathbf{A} := \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1m} \\ \vdots & & \vdots \\ \mathbf{A}_{m1} & \cdots & \mathbf{A}_{mm} \end{bmatrix}, \quad (7.4)$$

where each (diagonal) block  $\mathbf{A}_{ii}$  is square. We call the factorization

$$\mathbf{A} = \mathbf{L}\mathbf{R} = \begin{bmatrix} \mathbf{I} & & & \\ \mathbf{L}_{21} & \mathbf{I} & & \\ \vdots & & \ddots & \\ \mathbf{L}_{m1} & \cdots & \mathbf{L}_{m,m-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{R}_{11} & \cdots & \mathbf{R}_{1m} \\ & \mathbf{R}_{21} & \cdots & \mathbf{R}_{2m} \\ & & \ddots & \vdots \\ & & & \mathbf{R}_{mm} \end{bmatrix} \quad (7.5)$$

a **block LU factorization of  $\mathbf{A}$** . Here the  $i$ th diagonal blocks  $\mathbf{I}$  and  $\mathbf{R}_{ii}$  in  $\mathbf{L}$  and  $\mathbf{R}$  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 7.15** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is a block matrix of the form (7.4), and 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  $k = 1, \dots, m-1$ . Then  $\mathbf{A}$  has a unique block LU factorization (7.5). Conversely, if  $\mathbf{A}$  is nonsingular and has a block LU factorization then  $\mathbf{A}_k$  is 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 7.4 suppose  $\mathbf{A}_{m-1}$  has a unique LU factorization  $\mathbf{A}_{m-1} = \mathbf{L}_{m-1}\mathbf{R}_{m-1}$ ,

and that  $\mathbf{A}_1, \dots, \mathbf{A}_{m-1}$  are nonsingular. Then  $\mathbf{L}_{m-1}$  and  $\mathbf{R}_{m-1}$  are nonsingular and

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{m-1} & \mathbf{B} \\ \mathbf{C}^T & \mathbf{A}_{mm} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{m-1} & \mathbf{0} \\ \mathbf{C}^T \mathbf{R}_{m-1}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{R}_{m-1} & \mathbf{L}_{m-1}^{-1} \mathbf{B} \\ 0 & \mathbf{A}_{mm} - \mathbf{C}^T \mathbf{R}_{m-1}^{-1} \mathbf{L}_{m-1}^{-1} \mathbf{B} \end{bmatrix}, \quad (7.6)$$

is a block LU factorization of  $\mathbf{A}$ . It is unique by derivation. Conversely, suppose  $\mathbf{A}$  is nonsingular and has a block LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{R}$ . Then as in Lemma 7.5 it is easily seen that  $\mathbf{A}_k = \mathbf{L}_k \mathbf{R}_k$  is a block LU factorization of  $\mathbf{A}_k$  for  $k = 1, \dots, m$ . By Lemma 6.31 and induction a block triangular matrix is nonsingular if and only if the diagonal blocks are nonsingular and we see that  $\mathbf{L}_k$  and  $\mathbf{R}_k$  are nonsingular, and hence  $\mathbf{A}_k$  is nonsingular for  $k = 1, \dots, m-1$ .  $\square$

**Remark 7.16** *The number of flops 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 7.17** *Note that (7.5) is not an LU factorization of  $\mathbf{A}$  since the  $\mathbf{R}_{ii}$ 's are not upper triangular in general. To relate the block LU factorization to the usual LU factorization we assume that each  $\mathbf{R}_{ii}$  has an LU factorization  $\mathbf{R}_{ii} = \tilde{\mathbf{L}}_{ii} \tilde{\mathbf{R}}_{ii}$ . Then  $\mathbf{A} = \hat{\mathbf{L}} \hat{\mathbf{R}}$ , where  $\hat{\mathbf{L}} := \mathbf{L} \text{diag}(\tilde{\mathbf{L}}_{ii})$  and  $\hat{\mathbf{R}} := \text{diag}(\tilde{\mathbf{L}}_{ii}^{-1}) \mathbf{R}$ , and this is an ordinary LU factorization of  $\mathbf{A}$ .*

**Exercise 7.18** *Show that  $\hat{\mathbf{L}}$  is unit lower triangular and  $\hat{\mathbf{R}}$  is upper triangular.*

## 7.3 The Symmetric LU Factorization

We consider next LU factorization of a real symmetric matrix.

**Definition 7.19** *Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$ . A factorization  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$ , where  $\mathbf{L}$  is unit lower triangular and  $\mathbf{D}$  is diagonal is called a **symmetric LU factorization** of  $\mathbf{A}$ .*

A matrix which has a symmetric LU factorization must be symmetric since  $\mathbf{A}^T = (\mathbf{L}\mathbf{D}\mathbf{L}^T)^T = \mathbf{L}\mathbf{D}\mathbf{L}^T = \mathbf{A}$ .

**Theorem 7.20** *Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is nonsingular. Then  $\mathbf{A}$  has a symmetric LU factorization if and only if  $\mathbf{A} = \mathbf{A}^T$  and  $\mathbf{A}_k$  is nonsingular for  $k = 1, \dots, n-1$ . The symmetric LU factorization is unique.*

**Proof.** If  $\mathbf{A}_1, \dots, \mathbf{A}_{n-1}$  are nonsingular then Theorem 7.4 implies that  $\mathbf{A}$  has a unique LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{R}$ . Since  $\mathbf{A}$  is nonsingular it follows that  $\mathbf{R}$  is nonsingular and since  $\mathbf{R}$  is triangular the diagonal matrix  $\mathbf{D} := \text{diag}(r_{11}, \dots, r_{nn})$  is nonsingular (cf. Lemma 6.32). But then  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{M}^T$ , where  $\mathbf{M}^T = \mathbf{D}^{-1}\mathbf{R}$  is

unit upper triangular. By symmetry  $\mathbf{A} = \mathbf{L}(\mathbf{D}\mathbf{M}^T) = \mathbf{M}(\mathbf{D}\mathbf{L}^T) = \mathbf{A}^T$  are two LU factorizations of  $\mathbf{A}$ , and by uniqueness  $\mathbf{M} = \mathbf{L}$ . Thus  $\mathbf{A}$  has a unique symmetric LU factorization.

Conversely, if  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  is the symmetric LU factorization of  $\mathbf{A}$  then  $\mathbf{A}$  is symmetric since  $\mathbf{L}\mathbf{D}\mathbf{L}^T$  is symmetric, and  $\mathbf{A}$  has an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{R}$  with  $\mathbf{R} = \mathbf{D}\mathbf{L}^T$ . By Theorem 7.6 we conclude that  $\mathbf{A}_1, \dots, \mathbf{A}_{n-1}$  are nonsingular.  $\square$

## 7.4 Positive Definite- and Positive Semidefinite Matrices

Symmetric positive definite matrices occur often in scientific computing. For example, the second derivative matrix is symmetric positive definite, see Lemma 7.21 below. For symmetric positive definite and symmetric positive semidefinite matrices there is a special version of the symmetric LU factorization. Before considering this factorization we study some properties of positive (semi)definite matrices. We study only real matrices, but consider also the nonsymmetric case..

### 7.4.1 Definition and Examples

Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  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

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}. \quad (7.7)$$

- The zero-matrix is symmetric positive semidefinite, while the unit matrix is symmetric positive definite.

- 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 (7.7) implies that  $\mathbf{x} = \mathbf{0}$ .

**Lemma 7.21** *The second derivative matrix  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{n,n}$  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  $\mathbf{x} = \mathbf{0}$ . Hence  $\mathbf{T}$  is positive definite.  $\square$

**Example 7.22** *Consider (cf. (C.1)) the gradient  $\nabla f$  and hessian  $\nabla \nabla^T f$  of a function  $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$*

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n, \quad \nabla \nabla^T f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix} \in \mathbb{R}^{n,n}.$$

We assume that  $f$  has continuous first and second partial derivatives on  $\Omega$ .

Under suitable conditions on the domain  $\Omega$  it is shown in advanced calculus texts that if  $\nabla f(\mathbf{x}) = \mathbf{0}$  and  $\nabla \nabla^T f(\mathbf{x})$  is 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.

## 7.4.2 Some Criteria for the Nonsymmetric Case

We treat the positive definite and positive semidefinite cases in parallel.

**Theorem 7.23** *Let  $m, n$  be positive integers. If  $\mathbf{A} \in \mathbb{R}^{n,n}$  is positive semidefinite and  $\mathbf{X} \in \mathbb{R}^{n,m}$  then  $\mathbf{B} := \mathbf{X}^T \mathbf{A} \mathbf{X} \in \mathbb{R}^{m,m}$  is positive semidefinite. If in addition  $\mathbf{A}$  is positive definite and  $\mathbf{X}$  has linearly independent columns then  $\mathbf{B}$  is positive definite.*

**Proof.** Let  $\mathbf{y} \in \mathbb{R}^m$  and set  $\mathbf{x} := \mathbf{X}\mathbf{y}$ . Then  $\mathbf{y}^T \mathbf{B}\mathbf{y} = \mathbf{x}^T \mathbf{A}\mathbf{x} \geq 0$ . If  $\mathbf{A}$  is positive definite and  $\mathbf{X}$  has linearly independent columns then  $\mathbf{x}$  is nonzero if  $\mathbf{y}$  is nonzero and  $\mathbf{y}^T \mathbf{B}\mathbf{y} = \mathbf{x}^T \mathbf{A}\mathbf{x} > 0$ .  $\square$

Taking  $\mathbf{A} := \mathbf{I}$  and  $\mathbf{X} := \mathbf{A}$  we obtain

**Corollary 7.24** *Let  $m, n$  be positive integers. If  $\mathbf{A} \in \mathbb{R}^{m,n}$  then  $\mathbf{A}^T \mathbf{A}$  is positive semidefinite. If in addition  $\mathbf{A}$  has linearly independent columns then  $\mathbf{A}^T \mathbf{A}$  is positive definite.*

**Theorem 7.25** *Any principal submatrix of a positive (semi)definite matrix is positive (semi)definite.*

**Proof.** Suppose the submatrix  $\mathbf{B}$  is defined by the rows and columns  $r_1, \dots, r_k$  of  $\mathbf{A}$ . Then  $\mathbf{B} := \mathbf{X}^T \mathbf{A} \mathbf{X}$ , where  $\mathbf{X} = [\mathbf{e}_{r_1}, \dots, \mathbf{e}_{r_k}] \in \mathbb{R}^{n,k}$ , and  $\mathbf{B}$  is positive (semi)definite by Theorem 7.23.  $\square$

If  $\mathbf{A}$  is positive definite then the leading principal submatrices are nonsingular and we obtain:

**Corollary 7.26** *A positive definite matrix has a unique LU factorization.*

**Theorem 7.27** *A positive (semi)definite matrix  $\mathbf{A}$  has positive (nonnegative) eigenvalues. Conversely if  $\mathbf{A}$  has positive (nonnegative) eigenvalues and orthonormal eigenvectors then it is positive (semi)definite.*

**Proof.** Consider the positive definite case. Suppose  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with  $\mathbf{x} \neq \mathbf{0}$ . Multiplying both sides by  $\mathbf{x}^T$  and solving for  $\lambda$  we find  $\lambda = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} > 0$ . Suppose conversely that  $\mathbf{A} \in \mathbb{R}^{n,n}$  has eigenpairs  $(\lambda_j, \mathbf{u}_j)$ ,  $j = 1, \dots, n$ , where the eigenvalues are positive and the eigenvectors satisfy  $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$ ,  $i, j = 1, \dots, n$ . Let  $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{n,n}$  and  $\mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n)$ . Since  $\mathbf{A}\mathbf{u}_j = \lambda_j \mathbf{u}_j$  for all  $j$  we have  $\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{D}$  and since  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  we obtain  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{D}$ . 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  $\mathbf{c}$  is nonzero since  $\mathbf{U}^T$  is nonsingular. We find  $\mathbf{U}\mathbf{c} = \mathbf{U}\mathbf{U}^T \mathbf{x} = \mathbf{x}$  and so

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{U}\mathbf{c})^T \mathbf{A} \mathbf{U}\mathbf{c} = \mathbf{c}^T \mathbf{U}^T \mathbf{A} \mathbf{U} \mathbf{c} = \mathbf{c}^T \mathbf{D} \mathbf{c} = \sum_{j=1}^n \lambda_j c_j^2,$$

and it follows that  $\mathbf{A}$  is positive definite. The positive semidefinite case is similar.  $\square$

**Theorem 7.28** *If  $\mathbf{A}$  is positive (semi)definite then  $\det(\mathbf{A}) > 0$  ( $\det(\mathbf{A}) \geq 0$ ).*

**Proof.** Since the determinant of a matrix is equal to the product of its eigenvalues this follows from Theorem 7.27.  $\square$

## 7.5 The Symmetric Case and Cholesky Factorization

For symmetric positive definite matrices there is an alternative to the symmetric LU factorization known as the Cholesky factorization. We consider also a closely related factorization of symmetric positive semidefinite matrices.

We need the following necessary conditions for symmetric positive semidefinite matrices.

**Lemma 7.29** *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.  $a_{ii} = 0 \implies a_{ij} = a_{ji} = 0$ , fixed  $i$ , all  $j$ .

**Proof.** 3. follows from 1. and 4. from 2. We have

$$0 \leq (\alpha \mathbf{e}_i + \beta \mathbf{e}_j)^T \mathbf{A} (\alpha \mathbf{e}_i + \beta \mathbf{e}_j) = \alpha^2 a_{ii} + \beta^2 a_{jj} + 2\alpha\beta a_{ij}, \text{ all } i, j, \alpha, \beta \in \mathbb{R}. \quad (7.8)$$

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 positive. Taking  $\alpha = -a_{ij}, \beta = a_{ii}$  in (7.8) 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.

**Definition 7.30** *A factorization  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  where  $\mathbf{R}$  is upper triangular with positive diagonal elements is called a **Cholesky factorization**. A factorization  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  where  $\mathbf{R}$  is upper 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  $\mathbf{R}$  must be nonsingular. Thus the diagonal elements of  $\mathbf{R}$  cannot be zero.

**Exercise 7.31** *Show that a symmetric matrix has a Cholesky factorization if and only if it has a symmetric LU factorization with positive diagonal elements in  $\mathbf{D}$ .*

**Theorem 7.32** A matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  has a Cholesky factorization  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  if and only if it is symmetric positive definite.

**Proof.** If  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  is a Cholesky factorization then  $\mathbf{A}$  is symmetric. Since  $\mathbf{R}$  has positive diagonal elements it is nonsingular. Thus  $\mathbf{A}$  is symmetric positive definite by Corollary 7.24. The proof of the converse will lead to an algorithm. We use induction on  $n$ . A positive definite matrix of order one has a Cholesky factorization since the one and only element in  $\mathbf{A}$  is positive. Suppose any symmetric positive definite matrix of order  $n-1$  has a Cholesky factorization and suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$ . We partition  $\mathbf{A}$  as follows

$$\mathbf{A} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix}, \quad \alpha \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^{n-1}, \mathbf{B} \in \mathbb{R}^{n-1,n-1}. \quad (7.9)$$

Clearly  $\alpha = \mathbf{e}_1^T \mathbf{A} \mathbf{e}_1 > 0$ . We claim that  $\mathbf{C} := \mathbf{B} - \mathbf{v} \mathbf{v}^T / \alpha$  is symmetric positive definite.  $\mathbf{C}$  is symmetric. To show that  $\mathbf{C}$  is positive definite we let  $\mathbf{y} \in \mathbb{R}^{n-1}$  be nonzero and define  $\mathbf{x}^T := [-\mathbf{v}^T \mathbf{y} / \alpha, \mathbf{y}^T] \in \mathbb{R}^n$ . Then  $\mathbf{x} \neq \mathbf{0}$  and

$$0 < \mathbf{x}^T \mathbf{A} \mathbf{x} = [-\mathbf{v}^T \mathbf{y} / \alpha, \mathbf{y}^T] \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} \begin{bmatrix} -\mathbf{v}^T \mathbf{y} / \alpha \\ \mathbf{y} \end{bmatrix} = \mathbf{y}^T \mathbf{C} \mathbf{y}. \quad (7.10)$$

Here we used that  $(\mathbf{v}^T \mathbf{y}) \mathbf{v}^T \mathbf{y} = (\mathbf{v}^T \mathbf{y})^T \mathbf{v}^T \mathbf{y} = \mathbf{y}^T \mathbf{v} \mathbf{v}^T \mathbf{y}$ . So  $\mathbf{C} \in \mathbb{R}^{n-1,n-1}$  is symmetric positive definite and by the induction hypothesis it has a Cholesky factorization  $\mathbf{C} = \mathbf{R}_1^T \mathbf{R}_1$ . The matrix

$$\mathbf{R} := \begin{bmatrix} \beta & \mathbf{v}^T / \beta \\ \mathbf{0} & \mathbf{R}_1 \end{bmatrix}, \quad \beta := \sqrt{\alpha}, \quad (7.11)$$

is upper triangular with positive diagonal elements and

$$\mathbf{R}^T \mathbf{R} = \begin{bmatrix} \beta & \mathbf{0} \\ \mathbf{v} / \beta & \mathbf{R}_1^T \end{bmatrix} \begin{bmatrix} \beta & \mathbf{v}^T / \beta \\ \mathbf{0} & \mathbf{R}_1 \end{bmatrix} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} = \mathbf{A}$$

is a Cholesky factorization of  $\mathbf{A}$ .  $\square$

We can now show

**Theorem 7.33** The following is equivalent for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ .

1.  $\mathbf{A}$  is positive definite.
2.  $\mathbf{A}$  has only positive eigenvalues.
3. All leading principal minors are positive.
4.  $\mathbf{A} = \mathbf{B}^T \mathbf{B}$  for a nonsingular  $\mathbf{B} \in \mathbb{R}^{n,n}$ .

**Proof.** A symmetric matrix has a set of eigenvectors that form an orthonormal basis for  $\mathbb{R}^n$  (Cf. Theorem 10.7). Therefore, by Theorem 7.27 we know that  $1 \Leftrightarrow 2$ .



We show that  $1 \Rightarrow 3 \Rightarrow 4 \Rightarrow 1$ .

**1  $\Rightarrow$  3:** By Theorem 7.25 the leading principal submatrix  $\mathbf{A}_k$  of  $\mathbf{A}$  is positive definite, and has a positive determinant by Theorem 7.28.

**3  $\Rightarrow$  4:** Since all principal minors of  $\mathbf{A}$  are positive the principal submatrices  $\mathbf{A}_k$  are nonsingular for all  $k$  and therefore  $\mathbf{A}$  has a symmetric LU factorization. By Exercise 7.31  $\mathbf{A}$  has a Cholesky factorization and we can take  $\mathbf{B} = \mathbf{R}$ .

**4  $\Rightarrow$  1:** This follows from Corollary 7.24.  $\square$

Consider next the semi-Cholesky factorization.

**Theorem 7.34** *A matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  has a semi-Cholesky factorization  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  if and only if it is symmetric positive semidefinite.*

**Proof.** If  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  is a semi-Cholesky factorization then  $\mathbf{A}$  is symmetric and it is positive semidefinite by Corollary 7.24. Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is symmetric positive semidefinite. A symmetric positive semidefinite matrix of order one has a semi-Cholesky factorization since  $a_{11}$  is nonnegative. Suppose by induction on  $n$  that any symmetric positive semidefinite matrix  $\mathbf{C}$  of order  $n-1$  has a semi-Cholesky factorization. We partition  $\mathbf{A}$  as in (7.9). There are two cases. If  $\alpha > 0$  then we obtain a semi-Cholesky factorization of  $\mathbf{A}$  as in the proof of Theorem 7.32 since  $\mathbf{C}$  is symmetric positive semidefinite. This follows as in (7.10) since now  $0 \leq \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{y}^T \mathbf{C} \mathbf{y}$ . If  $\alpha = 0$  then it follows from 4. in Lemma 7.29 that  $\mathbf{v} = \mathbf{0}$ . Moreover,  $\mathbf{B} \in \mathbb{R}^{n-1,n-1}$  in (7.9) is positive semidefinite and therefore has a semi-Cholesky factorization  $\mathbf{R}_1$ . But then  $\mathbf{R} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{R}_1 \end{bmatrix}$  is a semi-Cholesky factorization of  $\mathbf{A}$ . Indeed,  $\mathbf{R}$  is upper triangular and

$$\mathbf{R}^T \mathbf{R} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{R}_1^T \end{bmatrix} \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{R}_1 \end{bmatrix} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B} \end{bmatrix} = \mathbf{A}.$$

$\square$

**Theorem 7.35** *The following is equivalent for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ .*

1.  $\mathbf{A}$  is positive semidefinite.
2.  $\mathbf{A}$  has only nonnegative eigenvalues.
3.  $\mathbf{A} = \mathbf{B}^T \mathbf{B}$  for some  $\mathbf{B} \in \mathbb{R}^{n,n}$ .
4. All principal minors are nonnegative.

**Proof.** The proof of  $1. \Leftrightarrow 2$  follows as in the proof of Theorem 7.33.  $1. \Rightarrow 3$ . follows from Theorem 7.34 while  $1. \Rightarrow 4$ . is a consequence of Theorem 7.25. 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 [15]). 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 7.35 we require nonnegativity of all principal minors, while only positivity of leading principal minors was required for positive definite matrices (cf. Theorem 7.33). 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.

## 7.6 An Algorithm for SemiCholesky Factorization of a Banded Matrix

Recall that a matrix  $\mathbf{A}$  has bandwidth  $d \geq 0$  if  $a_{ij} = 0$  for  $|i - j| > d$ . A (semi)Cholesky factorization preserves bandwidth.

**Theorem 7.36** *The Cholesky factor  $\mathbf{R}$  given by (7.11) has the same bandwidth as  $\mathbf{A}$ .*

**Proof.** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  has bandwidth  $d \geq 0$ . Then  $\mathbf{v}^T = [\mathbf{u}^T, \mathbf{0}^T]$  in (7.9), 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  $n$ ,  $\mathbf{C} = \mathbf{R}_1^T \mathbf{R}_1$ , where  $\mathbf{R}_1$  has the same bandwidth as  $\mathbf{C}$ . But then  $\mathbf{R}$  in (7.11) 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 upper part of  $\mathbf{A}$ . The first row of  $\mathbf{R}$  is  $[\beta, \mathbf{v}^T/\beta]$  if  $\alpha > 0$ . If  $\alpha = 0$  then by 4 in Lemma 7.29 the first row of  $\mathbf{A}$  is zero and this is also the first row of  $\mathbf{R}$ .

Suppose we store the first row of  $\mathbf{R}$  in the first row of  $\mathbf{A}$  and the upper part of  $\mathbf{C} = \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$  in the upper part of  $\mathbf{A}(2:n, 2:n)$ . The first row of  $\mathbf{R}$  and the upper part of  $\mathbf{C}$  can be computed as follows.

if $A(1, 1) > 0$ $A(1, 1) = \sqrt{A(1, 1)}$ $A(1, 2:n) = A(1, 2:n)/A(1, 1)$ for $i = 2:n$ $A(i, i:n) = A(i, i:n) - A(1, i) * A(1, i:n)$	(7.12)
---	--------

The code can be made more efficient when  $\mathbf{A}$  is a band matrix. If the bandwidth is  $d$  we simply replace all occurrences of  $n$  by  $\min(i + d, n)$ .

Continuing the reduction we arrive at the following algorithm.

**Algorithm 7.37 (bandcholesky)** Suppose  $\mathbf{A}$  is symmetric positive semidefinite. An upper triangular matrix  $\mathbf{R}$  is computed so that  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$ . 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 `triu`.

```
function R=bandcholesky(A,d)
n=length(A);
for k=1:n
    if A(k,k)>0
        kp=min(n,k+d);
        A(k,k)=sqrt(A(k,k));
        A(k,k+1:kp)=A(k,k+1:kp)/A(k,k);
        for i=k+1:kp
            A(i,i:kp)=A(i,i:kp)-A(k,i)*A(k,i:kp);
        end
    else
        A(k,k:kp)=zeros(1,kp-k+1);
    end
end
R=triu(A);
```

In the algorithm we overwrite the upper triangle of  $\mathbf{A}$  with the elements of  $\mathbf{R}$ . Row  $k$  of  $\mathbf{R}$  is zero for those  $k$  where  $r_{kk} = 0$ . 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 [3]. Algorithm 7.37 is based on outer products  $\mathbf{v}\mathbf{v}^T$ . An advantage of this formulation is that it can be extended to symmetric positive semidefinite matrices.

Consider next forward and backward substitution. Since  $\mathbf{R}^T$  is lower triangular and banded the  $k$ th component of  $\mathbf{R}^T \mathbf{y} = \mathbf{b}$  is  $\sum_{j=\max(1,k-d)}^{k-1} r_{jk} y_j + r_{kk} y_k = b_k$ , and solving for  $y_k$

$$y_k = (b_k - \sum_{j=\max(1,k-d)}^{k-1} r_{jk} y_j) / r_{kk}, \text{ for } k = 1, \dots, n, \quad (7.13)$$

Similarly the  $k$ th component of  $\mathbf{R}\mathbf{x} = \mathbf{y}$  is  $r_{kk} x_k + \sum_{i=k+1}^{\min(n,k+d)} r_{ki} x_i = y_k$ , and solving for  $x_k$

$$x_k = (y_k - \sum_{i=k+1}^{\min(n,k+d)} r_{ki} x_i) / r_{kk}, \text{ for } k = n, n-1, \dots, 1. \quad (7.14)$$

This give the following algorithms

**Algorithm 7.38 (bandforwardsolve)** Solves the lower triangular system  $\mathbf{R}^T \mathbf{y} = \mathbf{b}$ .  $\mathbf{R}$  is upper triangular and banded with  $r_{kj} = 0$  for  $j - k > d$ .

```
function y=bandforwardsolve(R,b,d)
    n=length(b); y=b(:);
    for k=1:n
        km=max(1,k-d);
        y(k)=(y(k)-R(km:k-1,k)'*y(km:k-1))/R(k,k);
    end
```

**Algorithm 7.39 (bandbacksolve)** Solves the upper triangular system  $\mathbf{R}\mathbf{x} = \mathbf{y}$ .  $\mathbf{R}$  is upper triangular and banded with  $r_{kj} = 0$  for  $j - k > d$ .

```
function x=bandbacksolve(R,y,d)
    n=length(y); x=y(:);
    for i=n:-1:1
        kp=min(n,k+d);
        x(k)=(x(k)-R(k,k+1:kp)*x(k+1:kp))/R(k,k);
    end
```

For a full matrix ( $d = n$ ) the number of flops needed for the Cholesky factorization including  $n$  square roots is given by

$$\sum_{k=1}^n \sum_{i=k+1}^n (1 + \sum_{j=i}^n 2) + n = \frac{1}{3}n(n + \frac{1}{2})(n + 1) \approx n^3/3.$$

The number  $n^3/3$  is half the number of flops needed for Gaussian elimination of an arbitrary matrix. We obtain this reduction since the Cholesky factorization takes advantage of the symmetry of  $\mathbf{A}$ .

The number of flops for the banded algorithms is given approximately by

$$\sum_{k=1}^n \sum_{i=k+1}^{k+d} (1 + \sum_{j=i}^{k+d} 2) + n = O(nd^2)$$

for Algorithm 7.37 and  $O(2nd)$  for each of Algorithms 7.38 and 7.39. When  $d$  is small compared to  $n$  we see that these numbers are considerably smaller than the  $O(n^3/3)$  and  $O(2n^2)$  counts for the factorization of a full symmetric matrix.

There is also a banded version of the symmetric LU factorization which requires approximately the same number of flops as the Cholesky factorization. The choice between using a symmetric LU factorization or an  $\mathbf{R}^T \mathbf{R}$  factorization depends on several factors. Usually an LU or a symmetric LU factorization is preferred for matrices with small bandwidth (tridiagonal, pentadiagonal), while the  $\mathbf{R}^T \mathbf{R}$  factorization is restricted to symmetric positive semidefinite matrices and is often used when the bandwidth is larger.

## 7.7 The PLU Factorization

Suppose  $\mathbf{A}$  is nonsingular. We show existence of a factorization  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{R}$ , where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is a unit lower triangular, and  $\mathbf{R}$  is upper triangular. Recall that a **permutation matrix** is a matrix of the form

$$\mathbf{P} = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_n}] \in \mathbb{R}^{n,n},$$

where  $\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}$  is a permutation of the unit vectors  $\mathbf{e}_1, \dots, \mathbf{e}_n \in \mathbb{R}^n$ . Since  $\mathbf{P}^T \mathbf{P} = \mathbf{I}$  the inverse of  $\mathbf{P}$  is equal to its transpose,  $\mathbf{P}^{-1} = \mathbf{P}^T$  and  $\mathbf{P}\mathbf{P}^T = \mathbf{I}$  as well. A special case is an **(j,k)-Exchange Matrix**  $\mathbf{I}_{jk}$  obtained by exchanging 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 exchange matrix is symmetric and equal to its own inverse. Pre-multiplying a matrix by an exchange matrix interchanges two rows of the matrix, while post-multiplication interchanges two columns.

**Theorem 7.40 (The PLU theorem)** *A nonsingular matrix  $\mathbf{A}$  has a factorization  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{R}$ , where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is unit lower triangular, and  $\mathbf{R}$  is upper triangular.*

**Proof.** We use induction on  $n$ . The result is obvious for  $n = 1$ . Suppose any nonsingular matrix of order  $n-1$  has a PLU factorization and consider a nonsingular matrix  $\mathbf{A}$  of order  $n$ . Since  $\mathbf{A}$  is nonsingular one of the elements, say  $a_{r1}$ , in the first column of  $\mathbf{A}$  must be nonzero. Let  $\mathbf{B} := \mathbf{I}_{r1}\mathbf{A}$  and set

$$\mathbf{M}_1 := \mathbf{I} - \mathbf{m}\mathbf{e}_1^T, \quad \mathbf{m} = [0, \frac{b_{21}}{b_{11}}, \dots, \frac{b_{n1}}{b_{11}}]^T.$$

Note that  $\mathbf{M}_1$  is unit lower triangular and therefore nonsingular. We have  $\mathbf{M}_1^{-1} = \mathbf{I} + \mathbf{m}\mathbf{e}_1^T$  since

$$(\mathbf{I} + \mathbf{m}\mathbf{e}_1^T)(\mathbf{I} - \mathbf{m}\mathbf{e}_1^T) = \mathbf{I} - \mathbf{m}\mathbf{e}_1^T + \mathbf{m}\mathbf{e}_1^T - \mathbf{m}(\mathbf{e}_1^T \mathbf{m})\mathbf{e}_1^T = \mathbf{I}.$$

The first column of  $\mathbf{M}_1\mathbf{B}$  is

$$\mathbf{M}_1\mathbf{B}\mathbf{e}_1 = \mathbf{B}\mathbf{e}_1 - \mathbf{m}\mathbf{e}_1^T\mathbf{B}\mathbf{e}_1 = \mathbf{B}\mathbf{e}_1 - b_{11}\mathbf{m} = [b_{11}, 0, \dots, 0]^T$$

and we can write

$$\mathbf{M}_1\mathbf{B} = \mathbf{M}_1\mathbf{I}_{r1}\mathbf{A} = \begin{bmatrix} b_{11} & \mathbf{c}_2^T \\ 0 & \mathbf{D}_2 \end{bmatrix}, \quad \text{with } \mathbf{D}_2 \in \mathbb{R}^{n-1, n-1}. \quad (7.15)$$

The matrix  $\mathbf{M}_1\mathbf{I}_{r1}\mathbf{A}$  is a product of nonsingular matrices and therefore nonsingular. By Lemma 6.31 the matrix  $\mathbf{D}_2$  is nonsingular and by the induction hypothesis we have  $\mathbf{D}_2 = \mathbf{P}_2\mathbf{L}_2\mathbf{R}_2$  or  $\mathbf{P}_2^T\mathbf{D}_2 = \mathbf{L}_2\mathbf{R}_2$ , where  $\mathbf{P}_2 \in \mathbb{R}^{n-1, n-1}$  is a permutation matrix,  $\mathbf{L}_2$  is unit lower triangular and  $\mathbf{R}_2$  is upper triangular. Define matrices  $\mathbf{Q}_2, \mathbf{M}_2, \mathbf{R}$  of order  $n$  by

$$\mathbf{Q}_2 = \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{P}_2 \end{bmatrix}, \quad \mathbf{M}_2 = \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{L}_2 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} b_{11} & \mathbf{c}_2^T \\ 0 & \mathbf{R}_2 \end{bmatrix}.$$

Then

$$\begin{aligned} Q_2^T M_1 I_{r_1} A &= \begin{bmatrix} 1 & 0 \\ 0 & P_2^T \end{bmatrix} \begin{bmatrix} b_{11} & c_2^T \\ 0 & D_2 \end{bmatrix} = \begin{bmatrix} b_{11} & c_2^T \\ 0 & P_2^T D_2 \end{bmatrix} \\ &= \begin{bmatrix} b_{11} & c_2^T \\ 0 & L_2 R_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & L_2 \end{bmatrix} \begin{bmatrix} b_{11} & c_2^T \\ 0 & R_2 \end{bmatrix} = M_2 R, \end{aligned}$$

and hence

$$A = I_{r_1} M_1^{-1} Q_2 M_2 R = (I_{r_1} Q_2) (Q_2^T M_1^{-1} Q_2) M_2 R.$$

Now

$$\begin{aligned} Q_2^T M_1^{-1} Q_2 &= \begin{bmatrix} 1 & 0 \\ 0 & P_2^T \end{bmatrix} (I + m e_1^T) \begin{bmatrix} 1 & 0 \\ 0 & P_2 \end{bmatrix} = I + \begin{bmatrix} 1 & 0 \\ 0 & P_2^T \end{bmatrix} m e_1^T \begin{bmatrix} 1 & 0 \\ 0 & P_2 \end{bmatrix} \\ &= I + \begin{bmatrix} 0 \\ P_2^T m(2:n) \end{bmatrix} e_1^T. \end{aligned}$$

Thus  $Q_2^T M_1^{-1} Q_2$  is unit lower triangular and we have  $A = PLR$ , where  $P = I_{r_1} Q_2$  is a permutation matrix,  $L = Q_2^T M_1^{-1} Q_2 M_2$  is unit lower triangular, and  $R$  is upper triangular.  $\square$

To find the PLU factorization of a matrix we can use Gaussian elimination with row interchanges (pivoting). See Appendix A for details.

## Chapter 8

# 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.

### 8.1 Test Matrices

In this section we introduce some matrices which we will use to compare various algorithms in later chapters.

#### 8.1.1 The 2D Poisson Problem

Consider the problem

$$-\nabla^2 u := -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f \text{ on } \Omega := (0, 1)^2 = \{(x, y) : 0 < x < 1, 0 < y < 1\}, \quad (8.1)$$

$$u := 0 \text{ on } \partial\Omega.$$

Here  $\Omega$  is the open unit square while  $\partial\Omega$  is the boundary of  $\Omega$ . The function  $f$  is given and continuous on  $\Omega$  and we seek a function  $u = u(x, y)$  such that (8.1) 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

$$\{(jh, kh) : j, k = 0, 1, \dots, m+1\}, \quad \text{where } h = 1/(m+1).$$

The points  $\Omega_h := \{(jh, kh) : j, k = 1, \dots, m\}$  are the interior points, while the remaining points are the boundary points. The solution is zero at the boundary points. For an interior point we insert the difference approximations

$$\frac{\partial^2 u(jh, kh)}{\partial x^2} \approx \frac{v_{j-1,k} - 2v_{j,k} + v_{j+1,k}}{h^2}, \quad \frac{\partial^2 u(jh, kh)}{\partial y^2} \approx \frac{v_{j,k-1} - 2v_{j,k} + v_{j,k+1}}{h^2}$$

in (8.1) and multiply both sides by  $h^2$  to obtain

$$(-v_{j-1,k} + 2v_{j,k} - v_{j+1,k}) + (-v_{j,k-1} + 2v_{j,k} - v_{j,k+1}) = h^2 f_{j,k} \quad (8.2)$$

or

$$4v_{j,k} - v_{j-1,k} - v_{j+1,k} - v_{j,k-1} - v_{j,k+1} = h^2 f_{j,k} := h^2 f(jh, kh). \quad (8.3)$$

From the boundary conditions we have in addition

$$v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0, \quad j, k = 0, 1, \dots, m+1. \quad (8.4)$$

The equations (8.3) and (8.4) define a linear set of equations for the unknowns  $\mathbf{V} = [v_{jk}] \in \mathbb{R}^{m,m}$ .

Observe that (8.2) can be written as a matrix equation in the form

$$\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F} \quad \text{with} \quad h = 1/(m+1), \quad (8.5)$$

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m,m}$  is the second derivative matrix given by (6.23),  $\mathbf{V} = (v_{jk}) \in \mathbb{R}^{m,m}$ , and  $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m,m}$ . 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 (8.2).

To write (8.3) and (8.4) 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 8.1** For any  $\mathbf{B} \in \mathbb{R}^{m,n}$  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 8.1. We call this the **natural ordering**. The location of the elements in (8.3) form a 5-point stencil, as shown in Figure 8.2.

To find the matrix  $\mathbf{A}$  we note that for values of  $j, k$  where the 5-point stencil does not touch the boundary, (8.3) 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

$$\mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n,n}, \quad \mathbf{b} \in \mathbb{R}^n, \quad n = m^2, \quad (8.6)$$



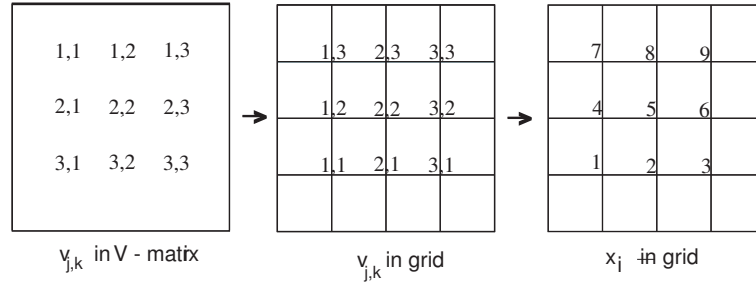


Figure 8.1. Numbering of grid points

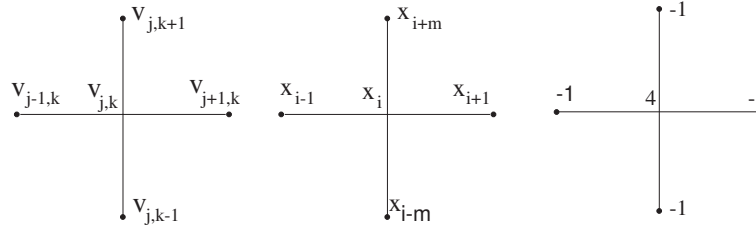


Figure 8.2. The 5-point stencil

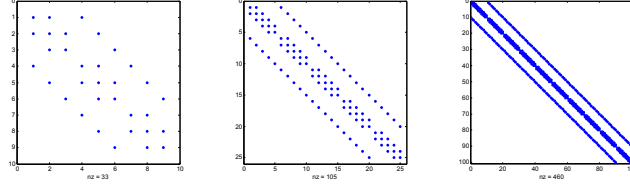
where  $\mathbf{x} = \text{vec}(\mathbf{V})$ ,  $\mathbf{b} = h^2 \text{vec}(\mathbf{F})$  with  $\mathbf{F} = (f_{jk}) \in \mathbb{R}^{m,m}$  and  $\mathbf{A}$  is the **Poisson matrix** given by

$$\begin{aligned}
 a_{ii} &= 4, & i &= 1, \dots, n \\
 a_{i+1,i} &= a_{i,i+1} = -1, & i &= 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m \\
 a_{i+m,i} &= a_{i,i+m} = -1, & i &= 1, \dots, n-m \\
 a_{ij} &= 0, & & \text{otherwise.}
 \end{aligned} \tag{8.7}$$

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 8.2** Write down the Poisson matrix for  $m = 2$  and show that it is strictly diagonally dominant.



**Figure 8.3.** Band structure of the 2D test matrix,  $n = 9$ ,  $n = 25$ ,  $n = 100$

### 8.1.2 The test Matrices

The second derivative matrix  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$  is a special case of the tridiagonal matrix

$$\mathbf{T}_1 := \begin{bmatrix} d & a & 0 & & & \\ a & d & a & & & \\ 0 & \ddots & \ddots & \ddots & & \\ & & & & 0 & \\ & & & & a & d & a \\ & & & & 0 & a & d \end{bmatrix}, \quad (8.8)$$

where  $a, d \in \mathbb{R}$ . We call this the **1D test matrix**.

The (2 dimensional) Poisson matrix is a special case of the matrix  $\mathbf{T}_2 = \mathbf{A} \in \mathbb{R}^{n,n}$  with elements

$$\begin{aligned} a_{i,i+1} = a_{i+1,i} &= a, & i = 1, \dots, n-1, & \quad i \neq m, 2m, \dots, (m-1)m, \\ a_{i,i+m} = a_{i+m,i} &= a, & i = 1, \dots, n-m, \\ a_{i,i} &= 2d, & i = 1, \dots, n, \\ a_{i,j} &= 0, & \text{otherwise,} \end{aligned} \quad (8.9)$$

and where  $a, d$  are real numbers. We will refer to this matrix as simply the **2D test matrix**. The 2D test matrix is a symmetric, banded matrix with bandwidth  $m = \sqrt{n}$ . The position of the nonzero elements is shown in Figure 8.3. With  $a = -1$  and  $d = 2$ , we obtain the Poisson matrix given by (8.7). This matrix is strictly diagonally dominant for  $m = 2, n = 4$ , but only diagonally dominant for  $m > 2$ . In this chapter we show that the Poisson matrix is symmetric positive definite and therefore nonsingular for any  $m \in \mathbb{N}$ . If  $|d| > 2|a|$  then  $\mathbf{T}_1$  is strictly diagonally dominant. In particular, with  $a = 1/9$  and  $d = 5/18$ , we obtain a version of  $\mathbf{T}_2$  which we refer to as the **averaging matrix**.

## 8.2 The Kronecker Product

**Definition 8.3** For any positive integers  $p, q, r, s$  we define the Kronecker product of two matrices  $\mathbf{A} \in \mathbb{R}^{p,q}$  and  $\mathbf{B} \in \mathbb{R}^{r,s}$  as a matrix  $\mathbf{C} \in \mathbb{R}^{pr,qs}$  given in block form as

$$\mathbf{C} = \begin{bmatrix} \mathbf{A}b_{1,1} & \mathbf{A}b_{1,2} & \cdots & \mathbf{A}b_{1,s} \\ \mathbf{A}b_{2,1} & \mathbf{A}b_{2,2} & \cdots & \mathbf{A}b_{2,s} \\ \vdots & \vdots & & \vdots \\ \mathbf{A}b_{r,1} & \mathbf{A}b_{r,2} & \cdots & \mathbf{A}b_{r,s} \end{bmatrix}.$$

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

$$\mathbf{T}_1 \otimes \mathbf{I} = \left[ \begin{array}{cc|cc} d & a & 0 & 0 \\ a & d & 0 & 0 \\ \hline 0 & 0 & d & a \\ 0 & 0 & a & d \end{array} \right] \quad \text{and} \quad \mathbf{I} \otimes \mathbf{T}_1 = \left[ \begin{array}{cc|cc} d & 0 & a & 0 \\ 0 & d & 0 & a \\ \hline a & 0 & d & 0 \\ 0 & a & 0 & d \end{array} \right].$$

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 8.4** Let for positive integers  $r, s, k$ ,  $\mathbf{A} \in \mathbb{R}^{r,r}$ ,  $\mathbf{B} \in \mathbb{R}^{s,s}$  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  $\lambda, \mu$  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

$$\begin{aligned}
 (\lambda \mathbf{A}) \otimes (\mu \mathbf{B}) &= \lambda \mu (\mathbf{A} \otimes \mathbf{B}), \\
 (\mathbf{A}_1 + \mathbf{A}_2) \otimes \mathbf{B} &= \mathbf{A}_1 \otimes \mathbf{B} + \mathbf{A}_2 \otimes \mathbf{B}, \\
 \mathbf{A} \otimes (\mathbf{B}_1 + \mathbf{B}_2) &= \mathbf{A} \otimes \mathbf{B}_1 + \mathbf{A} \otimes \mathbf{B}_2, \\
 (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} &= \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}), \\
 (\mathbf{A} \otimes \mathbf{B})^T &= \mathbf{A}^T \otimes \mathbf{B}^T.
 \end{aligned} \tag{8.10}$$

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 [10].

**Exercise 8.5** Prove (8.10).

The following *mixed product rule* is an essential tool for dealing with Kronecker products and sums.

**Lemma 8.6** 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

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}). \tag{8.11}$$

**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 8.7** 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  $\mathbf{A} \in \mathbb{R}^{r,r}$  and  $\mathbf{B} \in \mathbb{R}^{s,s}$  and

$$\mathbf{A}\mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad i = 1, \dots, r, \quad \mathbf{B}\mathbf{v}_j = \mu_j \mathbf{v}_j, \quad j = 1, \dots, s,$$

then

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = \lambda_i \mu_j (\mathbf{u}_i \otimes \mathbf{v}_j), \quad i = 1, \dots, r, \quad j = 1, \dots, s. \tag{8.12}$$

**Proof.** Using (8.11) the proof is a one liner. For all  $i, j$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\mathbf{A}\mathbf{u}_i) \otimes (\mathbf{B}\mathbf{v}_j) = (\lambda_i \mathbf{u}_i) \otimes (\mu_j \mathbf{v}_j) = (\lambda_i \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j).$$

□

Consider next a Kronecker sum.

**Lemma 8.8** *For positive integers  $r, s$  let  $\mathbf{A} \in \mathbb{R}^{r,r}$  and  $\mathbf{B} \in \mathbb{R}^{s,s}$ . 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*

$$\mathbf{A}\mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad i = 1, \dots, r, \quad \mathbf{B}\mathbf{v}_j = \mu_j \mathbf{v}_j, \quad j = 1, \dots, s,$$

then

$$(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\lambda_i + \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j), \quad i = 1, \dots, r, \quad j = 1, \dots, s. \quad (8.13)$$

**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 8.7 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. □

In many cases the Kronecker product and sum inherit properties of their factors.

**Lemma 8.9**

1. If  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular then  $\mathbf{A} \otimes \mathbf{B}$  is nonsingular. Moreover  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ .
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  $\mathbf{A} \in \mathbb{R}^{r,r}$  and  $\mathbf{B} \in \mathbb{R}^{s,s}$ . 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 (8.10). Suppose  $\mathbf{A}$  is positive definite and  $\mathbf{B}$  is positive semidefinite. Then  $\mathbf{A}$  has positive eigenvalues and  $\mathbf{B}$  has nonnegative eigenvalues. By Lemma 8.8 the eigenvalues of  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  are all positive and 3. follows. □

In (8.5) we derived the matrix equation  $\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}$  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  $\mathbf{A} = \mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}$ . This conversion could have been carried out with less effort using the following result.

**Lemma 8.10** Suppose  $\mathbf{A} \in \mathbb{R}^{r,r}$ ,  $\mathbf{B} \in \mathbb{R}^{s,s}$ , and  $\mathbf{F}, \mathbf{V} \in \mathbb{R}^{r,s}$ . Then we have

$$(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F}) \Leftrightarrow \mathbf{AVB}^T = \mathbf{F}, \quad (8.14)$$

$$(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F}) \Leftrightarrow \mathbf{AV} + \mathbf{VB}^T = \mathbf{F}. \quad (8.15)$$

**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

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\ \Leftrightarrow \begin{bmatrix} \mathbf{Ab}_{11} & \cdots & \mathbf{Ab}_{1s} \\ \vdots & & \vdots \\ \mathbf{Ab}_{s1} & \cdots & \mathbf{Ab}_{ss} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_s \end{bmatrix} &= \begin{bmatrix} \mathbf{f}_1 \\ \vdots \\ \mathbf{f}_s \end{bmatrix} \\ \Leftrightarrow \mathbf{A} \left[ \sum_j b_{1j} \mathbf{v}_j, \dots, \sum_j b_{sj} \mathbf{v}_j \right] &= [\mathbf{f}_1, \dots, \mathbf{f}_s] \\ \Leftrightarrow \mathbf{A}[\mathbf{Vb}_1, \dots, \mathbf{Vb}_s] = \mathbf{F} &\Leftrightarrow \mathbf{AVB}^T = \mathbf{F}. \end{aligned}$$

This proves (8.14). (8.15) follows immediately from (8.14) as follows

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\ \Leftrightarrow (\mathbf{AVI}_s^T + \mathbf{I}_r \mathbf{VB}^T) &= \mathbf{F} \Leftrightarrow \mathbf{AV} + \mathbf{VB}^T = \mathbf{F}. \end{aligned}$$

□

For more on Kronecker products see [10].

### 8.3 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

$$\mathbf{S} = \left[ \sin \frac{jk\pi}{m+1} \right]_{j,k=1}^m \in \mathbb{R}^{m,m}. \quad (8.16)$$

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 8.11** Suppose  $\mathbf{T}_1 = (t_{kj})_{k,j} = \text{tridiag}(a, d, a) \in \mathbb{R}^{m,m}$  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

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (8.17)$$

$$\lambda_j = d + 2a \cos(j\pi h). \quad (8.18)$$

2. The eigenvalues are distinct and the eigenvectors are orthogonal

$$\mathbf{s}_j^T \mathbf{s}_k = \frac{1}{2h} \delta_{j,k}, \quad j, k = 1, \dots, m. \quad (8.19)$$

**Proof.** Let  $s_{k,j} = \sin(jk\pi h)$  be the  $k$ th element in  $\mathbf{s}_j$  (the  $(k, j)$ -element in  $\mathbf{S}$ ). Observe that  $s_{k,0} = s_{k,m+1} = 0$ . Consider the  $k$ th component  $(\mathbf{T}_1 \mathbf{s}_j)_k$  of  $\mathbf{T}_1 \mathbf{s}_j$ . Note that  $t_{k,j} = 0$  except for  $t_{k,k-1} = t_{k,k+1} = a$  and  $t_{k,k} = d$ . With  $A := kj\pi h$  and  $B := j\pi h$  we find

$$\begin{aligned} (\mathbf{T}_1 \mathbf{s}_j)_k &= \sum_{l=1}^m t_{k,l} \sin(lj\pi h) = \sum_{l=k-1}^{k+1} t_{k,l} \sin(lj\pi h) \\ &= a \sin((k-1)j\pi h) + d \sin(kj\pi h) + a \sin((k+1)j\pi h) \\ &= a \sin(A-B) + d \sin A + a \sin(A+B) \\ &= 2a \cos B \sin A + d \sin A = (2a \cos B + d) \sin A = \lambda_j s_{k,j}, \end{aligned}$$

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 8.12 below that the eigenvectors  $\mathbf{s}_j$  are orthogonal. To finish the proof of (8.19) 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

$$\sum_{k=0}^m \cos(2kj\pi h) + i \sum_{k=0}^m \sin(2kj\pi h) = \sum_{k=0}^m e^{2ikj\pi h} = \frac{e^{2i(m+1)j\pi h} - 1}{e^{2ij\pi h} - 1} = \frac{e^{2ij\pi} - 1}{e^{2ij\pi h} - 1} = 0,$$

and (8.19) follows.  $\square$

**Lemma 8.12** The eigenvalues of a Hermitian matrix are real. Moreover, eigenvectors corresponding to distinct eigenvalues are orthogonal.

**Proof.** Suppose  $\mathbf{A}^H = \mathbf{A}$  and  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with  $\mathbf{x} \neq 0$ . We multiply both sides of  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  by  $\mathbf{x}^H$  and divide by  $\mathbf{x}^H\mathbf{x}$  to obtain  $\lambda = \frac{\mathbf{x}^H\mathbf{A}\mathbf{x}}{\mathbf{x}^H\mathbf{x}}$ . Taking complex conjugates we find  $\bar{\lambda} = \lambda^H = \frac{(\mathbf{x}^H\mathbf{A}\mathbf{x})^H}{(\mathbf{x}^H\mathbf{x})^H} = \frac{\mathbf{x}^H\mathbf{A}^H\mathbf{x}}{\mathbf{x}^H\mathbf{x}} = \frac{\mathbf{x}^H\mathbf{A}\mathbf{x}}{\mathbf{x}^H\mathbf{x}} = \lambda$ , and  $\lambda$  is real.

Suppose in addition that  $(\mu, \mathbf{y})$  is another eigenpair for  $\mathbf{A}$  with  $\mu \neq \lambda$ . Multiplying  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  by  $\mathbf{y}^H$  gives

$$\lambda\mathbf{y}^H\mathbf{x} = \mathbf{y}^H\mathbf{A}\mathbf{x} = (\mathbf{x}^H\mathbf{A}^H\mathbf{y})^H = (\mathbf{x}^H\mathbf{A}\mathbf{y})^H = (\mu\mathbf{x}^H\mathbf{y})^H = \mu\mathbf{y}^H\mathbf{x},$$

using that  $\mu$  is real. Since  $\lambda \neq \mu$  it follows that  $\mathbf{y}^H\mathbf{x} = 0$  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 8.13** For fixed  $m \geq 2$  let  $\mathbf{T}_2$  be the matrix given by (8.9) 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

$$\mathbf{x}_{j,k} = \mathbf{s}_j \otimes \mathbf{s}_k, \quad (8.20)$$

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (8.21)$$

$$\lambda_{j,k} = 2d + 2a \cos(j\pi h) + 2a \cos(k\pi h). \quad (8.22)$$

2. The eigenvectors are orthogonal

$$\mathbf{x}_{j,k}^T \mathbf{x}_{p,q} = \frac{1}{4h^2} \delta_{j,p} \delta_{k,q}, \quad j, k, p, q = 1, \dots, m. \quad (8.23)$$

3.  $\mathbf{T}_2$  is symmetric positive definite if  $d > 0$  and  $d \geq 2|a|$ .

4. The Poisson and averaging matrix are symmetric positive definite.

**Proof.**

1. follows from Lemma 8.11 and Lemma 8.8 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 (8.19) 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  $d \geq 2|a|$  and this holds both for both choices  $a = -1$ ,  $d = 2$  and  $a = 1/5$ ,  $d = 5/18$ . Thus the matrices in 4. are positive definite.  $\square$

**Exercise 8.14** Write down the eigenvalues of  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$  using Lemma 8.11 and conclude that  $\mathbf{T}$  is symmetric positive definite.



**Exercise 8.15** Use Lemma 8.11 to show that the matrix  $\mathbf{T}_1 := \text{tridiag}(a, d, a) \in \mathbb{R}^{n,n}$  is symmetric positive definite if  $d > 0$  and  $d \geq 2|a|$ .

**Exercise 8.16** For  $m = 2$  the matrix (8.9) 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  $x = [1, 1, 1, 1]^T$ . Verify that apart from a scaling of the eigenvector this agrees with (8.22) and (8.21) for  $j = k = 1$  and  $m = 2$ .

**Exercise 8.17** 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. (8.1))

$$\begin{aligned} \text{(a)} \quad & -(\square_h v)_{j,k} = (\mu f)_{j,k} \quad j, k = 1, \dots, m \\ \text{(b)} \quad & v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0, \quad j, k = 0, 1, \dots, m+1, \\ \text{(c)} \quad & -(\square_h v)_{j,k} = [20v_{j,k} - 4v_{j-1,k} - 4v_{j,k-1} - 4v_{j+1,k} - 4v_{j,k+1} \\ & \quad - v_{j-1,k-1} - v_{j+1,k-1} - v_{j-1,k+1} - v_{j+1,k+1}]/(6h^2), \\ \text{(d)} \quad & (\mu f)_{j,k} = [8f_{j,k} + f_{j-1,k} + f_{j,k-1} + f_{j+1,k} + f_{j,k+1}]/12. \end{aligned} \tag{8.24}$$

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 (8.24) defines an  $O(h^4)$  approximation to (8.1).

**Exercise 8.18** Consider the nine point difference approximation to (8.1) given by (8.24) in Problem 8.17.

a) Show that (8.24) is equivalent to the matrix equation

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} - \frac{1}{6}\mathbf{T}\mathbf{V}\mathbf{T} = h^2\mu\mathbf{F}. \tag{8.25}$$

Here  $\mu\mathbf{F}$  has elements  $(\mu f)_{j,k}$  given by (8.24d).

b) Show that the standard form of the matrix equation (8.25) 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 8.19** Consider the biharmonic equation

$$\begin{aligned} \nabla^4 u(s, t) &= \nabla^2(\nabla^2 u(s, t)) = f(s, t) & (s, t) \in \Omega, \\ u(s, t) &= 0, \quad \nabla^2 u(s, t) = 0 & (s, t) \in \partial\Omega. \end{aligned} \tag{8.26}$$

Here  $\Omega$  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  $v = -\nabla^2 u$ . Then (8.26) can be written as a system

$$\begin{aligned} -\nabla^2 v(s, t) &= f(s, t) & (s, t) \in \Omega \\ -\nabla^2 u(s, t) &= v(s, t) & (s, t) \in \Omega \\ u(s, t) &= v(s, t) = 0 & (s, t) \in \partial\Omega. \end{aligned} \quad (8.27)$$

- b) Discretizing, using (8.2), with  $\mathbf{T} = \text{diag}(-1, 2, -1) \in \mathbb{R}^{m,m}$ ,  $h = 1/(m+1)$ , 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

$$\mathbf{T}^2\mathbf{U} + 2\mathbf{T}\mathbf{U}\mathbf{T} + \mathbf{U}\mathbf{T}^2 = h^4\mathbf{F} \quad \text{and} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (8.28)$$

where  $\mathbf{A} = \mathbf{T}^2 \otimes \mathbf{I} + 2\mathbf{T} \otimes \mathbf{T} + \mathbf{I} \otimes \mathbf{T}^2 \in \mathbb{R}^{n,n}$ ,  $\mathbf{x} = \text{vec}(\mathbf{U})$ , 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 a basis of an algorithm? Why?

## Chapter 9

# Fast Direct Solution of a Large Linear System

### 9.1 Algorithms for a Banded Positive Definite System

In this chapter we present a fast method for solving  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is the Poisson matrix (8.7). Thus for  $n = 3$

$$\mathbf{A} = \left[ \begin{array}{ccc|ccc|ccc} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ \hline -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ \hline 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{array} \right]$$

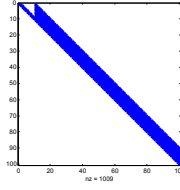
$$= \left[ \begin{array}{ccc|ccc} \mathbf{T} + 2\mathbf{I} & -\mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{T} + 2\mathbf{I} & -\mathbf{I} \\ \mathbf{0} & -\mathbf{I} & \mathbf{T} + 2\mathbf{I} \end{array} \right],$$

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$ . For this matrix we know by now that

1. It is symmetric positive definite.
2. It is banded.
3. It is block-tridiagonal.
4. We know the eigenvalues and eigenvectors of  $\mathbf{A}$ .

#### 9.1.1 Cholesky Factorization

Since  $\mathbf{A}$  is symmetric positive definite we can use the Cholesky factorization Algorithm 7.37. Since  $\mathbf{A}$  is banded with bandwidth  $d = \sqrt{n}$  the complexity of this



**Figure 9.1.** Fill-inn in the Cholesky factor of the Poisson matrix ( $n = 100$ ).

factorization is  $O(nd^2) = O(n^2)$ . We need to store  $\mathbf{A}$  possibly in sparse form.

The nonzero elements in  $\mathbf{R}$  are shown in Figure 9.1. Note that the zeros between the diagonals in  $\mathbf{A}$  have become nonzero in  $\mathbf{R}$ . This is known as **fill-inn**.

### 9.1.2 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

$$\begin{bmatrix} D_1 & C_1 & & & \\ A_2 & D_2 & C_2 & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-1} & D_{m-1} & C_{m-1} \\ & & & A_m & D_m \end{bmatrix} = \begin{bmatrix} I & & & & \\ L_2 & I & & & \\ & \ddots & \ddots & \ddots & \\ & & L_m & I & \end{bmatrix} \begin{bmatrix} R_1 & C_1 & & & \\ & \ddots & \ddots & \ddots & \\ & & R_{m-1} & C_{m-1} & \\ & & & R_m \end{bmatrix}. \quad (9.1)$$

Here  $D_1, \dots, D_m$  and  $R_1, \dots, R_m$  are square matrices while  $A_2, \dots, A_m$  and  $C_1, \dots, C_{m-1}$  can be rectangular.

Using block multiplication the formulas (6.25) generalize to

$$R_1 = D_1, \quad L_k = A_k R_{k-1}^{-1}, \quad R_k = D_k - L_k C_{k-1}, \quad k = 2, 3, \dots, m. \quad (9.2)$$

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  $\mathbf{Rx} = \mathbf{y}$  are as follows:

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{b}_1, & \mathbf{y}_k &= \mathbf{b}_k - \mathbf{L}_k \mathbf{y}_{k-1}, & k &= 2, 3, \dots, m, \\ \mathbf{x}_m &= \mathbf{R}_m^{-1} \mathbf{y}_m, & \mathbf{x}_k &= \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_{k+1}), & k &= m-1, \dots, 2, 1. \end{aligned} \quad (9.3)$$

The solution is then  $\mathbf{x}^T = [\mathbf{x}_1^T, \dots, \mathbf{x}_m^T]$ . To find  $\mathbf{L}_k$  in (9.2) we solve the linear systems  $\mathbf{L}_k \mathbf{R}_{k-1} = \mathbf{A}_k$ . Similarly we need to solve a linear system to find  $\mathbf{x}_k$  in (9.3).

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.

### 9.1.3 Other Methods

Other methods include

- Iterative methods. We study this in Chapters 13, 14, 15.
- Multigrid. See [5].
- Fast solvers based on diagonalization and the Fast Fourier Transform. See Sections 9.2, 9.3.

## 9.2 A Fast Poisson Solver based on Diagonalization

The algorithm we now derive will only require  $O(n^{3/2})$  flops and we only need to work with matrices of order  $m$ . Using the Fast Fourier Transform the number of flops can be reduced further to  $O(n \log n)$ .

To start we recall that  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written as a matrix equation in the form (cf. (8.5))

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F} \quad \text{with} \quad h = 1/(m+1),$$

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m,m}$  is the second derivative matrix,  $\mathbf{V} = (v_{jk}) \in \mathbb{R}^{m,m}$  are the unknowns, and  $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m,m}$  contains function values.

Recall that the eigenpairs of  $\mathbf{T}$  are given by

$$\begin{aligned} \mathbf{T}\mathbf{s}_j &= \lambda_j\mathbf{s}_j, \quad j = 1, \dots, m, \\ \mathbf{s}_j &= [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \\ \lambda_j &= 2 - 2\cos(j\pi h) = 4\sin^2(j\pi h/2), \quad h = 1/(m+1), \\ \mathbf{s}_j^T \mathbf{s}_k &= \delta_{jk}/(2h) \quad \text{for all } j, k. \end{aligned}$$

Let

$$\mathbf{S} := [\mathbf{s}_1, \dots, \mathbf{s}_m] = [\sin(jk\pi h)]_{j,k=1}^m \in \mathbb{R}^{m,m}, \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_m). \quad (9.4)$$

Then  $\mathbf{T}\mathbf{S} = \mathbf{S}\mathbf{D}$  and  $\mathbf{S}^T\mathbf{S} = \mathbf{S}^2 = \mathbf{I}/(2h)$ . Define  $\mathbf{X} \in \mathbb{R}^{m,m}$  by  $\mathbf{V} = \mathbf{S}\mathbf{X}\mathbf{S}$ , where  $\mathbf{V}$  is the solution of  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ . Then

$$\begin{aligned} \mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} &= h^2\mathbf{F} \\ \mathbf{V} &\stackrel{\mathbf{S}\mathbf{X}\mathbf{S}}{\Longleftrightarrow} \mathbf{T}\mathbf{S}\mathbf{X}\mathbf{S} + \mathbf{S}\mathbf{X}\mathbf{S}\mathbf{T} = h^2\mathbf{F} \\ &\stackrel{\mathbf{S}(\cdot)\mathbf{S}}{\Longleftrightarrow} \mathbf{S}\mathbf{T}\mathbf{S}\mathbf{X}\mathbf{S}^2 + \mathbf{S}^2\mathbf{X}\mathbf{S}\mathbf{T}\mathbf{S} = h^2\mathbf{S}\mathbf{F}\mathbf{S} \\ &\stackrel{\mathbf{T}\mathbf{S}=\mathbf{S}\mathbf{D}}{\Longleftrightarrow} \mathbf{S}^2\mathbf{D}\mathbf{X}\mathbf{S}^2 + \mathbf{S}^2\mathbf{X}\mathbf{S}^2\mathbf{D} = h^2\mathbf{S}\mathbf{F}\mathbf{S} \\ &\stackrel{\mathbf{S}^2=\mathbf{I}/(2h)}{\Longleftrightarrow} \mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} = 4h^4\mathbf{S}\mathbf{F}\mathbf{S}. \end{aligned}$$

An equation of the form  $\mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} = \mathbf{B}$ , 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  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ .

**Algorithm 9.1 (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  $\mathbf{F} = (f(jh, kh)) \in \mathbb{R}^{m,m}$ . We compute  $\mathbf{V} \in \mathbb{R}^{m,m}$ , where  $v_{jk} \approx u(jh, kh)$  by solving the equation  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$  using diagonalization of  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m,m}$ .

```
function V=fastpoisson(F)
m=length(F); h=1/(m+1); hv=pi*h*(1:m)';
sigma=sin(hv/2).^2;
S=sin(hv*(1:m));
G=S*F*S;
X=h^4*G./(sigma*ones(1,m)+ones(m,1)*sigma');
V=zeros(m+2,m+2);
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  $\lambda_j$ . Since the statement "X=h<sup>4</sup>\*G./(sigma\*ones(1,m)+ones(m,1)\*sigma)" only requires  $O(m^2)$  flops the complexity of this algorithm is for large  $m$  determined by the 4  $m$ -by- $m$  matrix multiplications and is given by  $O(4 \times 2m^3) = O(8n^{3/2})$ .<sup>2</sup>

### 9.3 A Fast Poisson Solver based on the Discrete Sine and Fourier Transforms

In Algorithm 9.1 we need to compute the product of the sine matrix  $\mathbf{S} \in \mathbb{R}^{m,m}$  given by (9.4) and a matrix  $\mathbf{A} \in \mathbb{R}^{m,m}$ . 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. 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.

#### 9.3.1 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  $\mathbf{w} = \mathbf{S}\mathbf{v}$ , where  $\mathbf{S}$  is the sine matrix given by (9.4). 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,

<sup>2</sup>It 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 9.4.

since  $\mathbf{S}$  is symmetric we have  $\mathbf{B} = (\mathbf{S}\mathbf{A}^T)^T$  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 9.1 by carrying out Discrete Sine Transforms on 4  $m$ -by- $m$  matrices in addition to the computation of  $\mathbf{X}$ .

### 9.3.2 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  $N \in \mathbb{N}$

$$\omega_N = \exp^{-2\pi i/N} = \cos(2\pi/N) - i \sin(2\pi/N), \quad (9.5)$$

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

$$\mathbf{F}_N = \left( \omega_N^{(j-1)(k-1)} \right)_{j,k=1}^N \in \mathbb{C}^{N,N}. \quad (9.6)$$

This matrix is known as the **Fourier matrix**. If  $\mathbf{A} \in \mathbb{R}^{N,m}$  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

$$\mathbf{F}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4 & \omega_4^2 & \omega_4^3 \\ 1 & \omega_4^2 & \omega_4^4 & \omega_4^6 \\ 1 & \omega_4^3 & \omega_4^6 & \omega_4^9 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -i & -1 & i \\ 1 & -1 & 1 & -1 \\ 1 & i & -1 & -i \end{bmatrix}. \quad (9.7)$$

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 9.2** *Given a positive integer  $m$  and a vector  $\mathbf{x} \in \mathbb{R}^m$ . Component  $k$  of  $\mathbf{S}\mathbf{x}$  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*

$$(\mathbf{S}\mathbf{x})_k = \frac{i}{2} (\mathbf{F}_{2m+2}\mathbf{z})_{k+1}, \quad k = 1, \dots, m.$$

**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  $N = 2m + 2$ .

### 9.3.3 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 flops 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  $\mathbf{P}_N \in \mathbb{R}^{N,N}$  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 (9.7)

$$\mathbf{P}_4 = [\mathbf{e}_1 \ \mathbf{e}_3 \ \mathbf{e}_2 \ \mathbf{e}_4] = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{F}_4 \mathbf{P}_4 = \left[ \begin{array}{cc|cc} 1 & 1 & 1 & 1 \\ 1 & -1 & -i & i \\ \hline 1 & 1 & -1 & -1 \\ 1 & -1 & i & -i \end{array} \right],$$

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 9.3** *If  $N = 2m$  is even then*

$$\mathbf{F}_{2m} \mathbf{P}_{2m} = \left[ \begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m \mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m \mathbf{F}_m \end{array} \right], \quad (9.8)$$

where

$$\mathbf{D}_m = \text{diag}(1, \omega_N, \omega_N^2, \dots, \omega_N^{m-1}). \quad (9.9)$$

**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

$$\begin{aligned} (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p,q} &= \omega_N^{j(2k)} &= \omega_m^{jk} &= (\mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p+m,q} &= \omega_N^{(j+m)(2k)} &= \omega_m^{(j+m)k} &= (\mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p,q+m} &= \omega_N^{j(2k+1)} &= \omega_N^j \omega_m^{jk} &= (\mathbf{D}_m \mathbf{F}_m)_{p,q}, \\ (\mathbf{F}_{2m} \mathbf{P}_{2m})_{p+m,q+m} &= \omega_N^{(j+m)(2k+1)} &= -\omega_N^j \omega_m^{jk} &= (-\mathbf{D}_m \mathbf{F}_m)_{p,q}. \end{aligned}$$

It follows that the four  $m$ -by- $m$  blocks of  $\mathbf{F}_{2m} \mathbf{P}_{2m}$  have the required structure.  $\square$

Using Theorem 9.3 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

$$\begin{aligned} \mathbf{F}_{2m} \mathbf{y} &= \mathbf{F}_{2m} \mathbf{P}_{2m} \mathbf{P}_{2m}^T \mathbf{y} = \mathbf{F}_{2m} \mathbf{P}_{2m} \mathbf{w} \\ &= \left[ \begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m \mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m \mathbf{F}_m \end{array} \right] \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 + \mathbf{q}_2 \\ \mathbf{q}_1 - \mathbf{q}_2 \end{bmatrix}, \end{aligned}$$

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 9.4 (Recursive FFT)** For  $\mathbf{y} \in \mathbb{C}^n$  we compute the Fourier transform  $\mathbf{z} = \mathbf{F}_n \mathbf{y}$ .

```
function z=fftrec(y)
n=length(y);
if n==1
    z=y;
else
    q1=fftrec(y(1:2:n-1));
    q2=exp(-2*pi*i/n).^(0:n/2-1).*fftrec(y(2:2:n));
    z=[q1+q2 q1-q2];
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 [23] for further details.

The complexity of the FFT is given by  $\gamma N \log_2 N$  for some constant  $\gamma$  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 flops) 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  $\gamma$  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 [23].

The efficiency improvement using the FFT to compute the DFT is spectacular for large  $N$ . The direct multiplication  $\mathbf{F}_N \mathbf{y}$  requires  $O(8n^2)$  flops since complex arithmetic is involved. Assuming that the FFT uses  $5N \log_2 N$  flops 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 flops then the direct multiplication would take something like 84000 seconds or 23 hours.

### 9.3.4 A Poisson Solver based on the FFT

We now have all the ingredients to compute the matrix products  $\mathbf{SA}$  and  $\mathbf{AS}$  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 9.1. We first compute  $\mathbf{H} = \mathbf{SF}$  using Lemma 9.2 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  $O(\gamma(2m+2) \log_2(2m+2))$  flops 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  $\mathbf{A}\mathbf{x} = \mathbf{b}$  we would be solving if Cholesky factorization was used. This should be compared to the  $O(8n^{3/2})$  flops used in Algorithm 9.1 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 [23] for other efficient ways to implement the DST.

## 9.4 Problems

**Exercise 9.1** Show that the Fourier matrix  $\mathbf{F}_4$  is symmetric, but not Hermitian.

**Exercise 9.2** Verify Lemma 9.2 directly when  $m = 1$ .

**Exercise 9.3** Show that the exact solution of the discrete Poisson equation (8.3) and (8.4) 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 9.4** Algorithm 9.1 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  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ .

(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

$$(\mathbf{T} + \lambda_j \mathbf{I})\mathbf{x}_j = \mathbf{c}_j \quad j = 1, \dots, m, \quad (9.10)$$

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 (6.25) and (6.26) in  $8m - 7$  flops. Give an algorithm to find  $\mathbf{X}$  which only requires  $O(\delta m^2)$  flops for some constant  $\delta$  independent of  $m$ .

(c) Describe a method to compute  $\mathbf{V}$  which only requires  $O(4m^3) = O(4n^{3/2})$  flops.

(d) Describe a method based on the Fast Fourier Transform which requires  $O(\gamma n \log_2 n)$  where  $\gamma$  is the same constant as mentioned at the end of the last section.

**Exercise 9.5** 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 8.18 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 (8.25). 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 8.18 b) is symmetric positive definite and that (8.24) always has a solution  $\mathbf{V}$ .

**Exercise 9.6** Derive an algorithm for solving (8.24) which for large  $m$  requires essentially the same number of operations as in Algorithm 9.1. (We assume that  $\mu\mathbf{F}$  already has been formed).

**Exercise 9.7** For the biharmonic problem we derived in Exercise 8.19 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  $\mathbf{X}$  by  $\mathbf{U} = \mathbf{S}\mathbf{X}\mathbf{S} = (x_{j,k})$  where  $\mathbf{U}$  is the solution of (8.28). 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 9.8** Use Exercise 9.7 to derive an algorithm

```
function U=simplefastbiharmonic(F)
```

which requires only  $O(\delta n^{3/2})$  operations to find  $\mathbf{U}$  in Problem 8.19. Here  $\delta$  is some constant independent of  $n$ .

**Exercise 9.9** In Exercise 9.8 compute the solution  $\mathbf{U}$  corresponding to  $\mathbf{F} = \mathbf{ones}(m, m)$ . For some small  $m$ 's check that you get the same solution obtained by solving the standard form  $\mathbf{A}\mathbf{x} = \mathbf{b}$  in (8.28). You can use  $\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$  for solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Use `F(:)` to vectorize a matrix and `reshape(x,m,m)` to turn a vector  $\mathbf{x} \in \mathbb{R}^{m^2}$  into an  $m \times m$  matrix. Make a plot of  $\mathbf{U}$  for say  $m = 50$ .

**Exercise 9.10** Repeat Exercises 8.19, 9.8 and 9.9 using the nine point rule (8.24) to solve the system (8.27).

**Part III**

**Some Matrix Theory**



## Chapter 10

# Orthonormal Eigenpairs and the Schur Form

A matrix is said to have **orthonormal (orthogonal) eigenpairs** if the eigenvectors are orthonormal (orthogonal). Two examples are the 2. derivative matrix  $T$  in Lemma 8.11 and the discrete Poisson matrix, cf. Lemma 8.13. In this chapter we characterize the family of matrices that have orthonormal eigenpairs. These matrices are called **normal matrices** and they contain the symmetric, Hermitian, and unitary matrices among their members.

If  $B = S^{-1}AS$  and  $S = U \in \mathbb{C}^{n,n}$  is unitary, then  $S^{-1} = U^H$  and  $B = U^H AU$ . In this case we say that  $B$  is **unitary similar** to  $A$ . In the real case where  $A$  and  $U$  are real matrices and  $U$  is orthonormal, we have  $S^{-1} = U^T$  and  $B = U^T AU$ . Unitary and orthonormal transformations are important in numerical algorithms since they are insensitive to noise in the elements of the matrix.

If  $B = U^H AU$  then  $AU = UB$ . If  $B = \text{diag}(\lambda_j)$  is diagonal and  $U = [u_1, \dots, u_n]$ , then  $Au_j = \lambda_j u_j$  for  $j = 1, \dots, n$ , and it follows that the columns of  $U$  are orthonormal eigenvectors of  $A$ . Conversely, if  $A$  has orthonormal eigenvectors  $u_1, \dots, u_n$ , then  $AU = UB$  or  $B = U^H AU$ , where the columns of  $U$  are the eigenvectors of  $A$  and  $B$  is diagonal. Thus  $A$  is unitary similar to a diagonal matrix if and only if  $A$  has a set of orthonormal eigenvectors.

## 10.1 The Schur Form

Not every matrix can be diagonalized by a similarity transformation, see Theorems 5.19, 5.20, and 5.27. But it can be triangularized, even by a unitary similarity transformation.

**Theorem 10.1 (Schur Triangularization)** *For each  $A \in \mathbb{C}^{n,n}$  there exists a unitary matrix  $U \in \mathbb{C}^{n,n}$  such that  $R := U^H AU$  is upper triangular.*

**Proof.** We use induction on  $n$ . For  $n = 1$  the matrix  $U$  is the  $1 \times 1$  identity matrix. Assume that the theorem is true for matrices of order  $k$  and suppose  $A \in \mathbb{C}^{n,n}$ , where  $n := k + 1$ . Let  $(\lambda_1, v_1)$  be an eigenpair for  $A$  with  $\|v_1\|_2 = 1$ . By

Theorem 2.62 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  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n,n}$  is unitary, and the first column of the product  $\mathbf{V}^H \mathbf{A} \mathbf{V}$  is

$$\mathbf{V}^H \mathbf{A} \mathbf{V} \mathbf{e}_1 = \mathbf{V}^H \mathbf{A} \mathbf{v}_1 = \lambda_1 \mathbf{V}^H \mathbf{v}_1 = \lambda_1 \mathbf{e}_1.$$

It follows that

$$\mathbf{V}^H \mathbf{A} \mathbf{V} = \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^H \\ \hline \mathbf{0} & \mathbf{M} \end{array} \right], \text{ for some } \mathbf{M} \in \mathbb{C}^{k,k} \text{ and } \mathbf{x} \in \mathbb{C}^k. \quad (10.1)$$

By the induction hypothesis there is a unitary matrix  $\mathbf{W}_1 \in \mathbb{C}^{k,k}$  such that  $\mathbf{W}_1^H \mathbf{M} \mathbf{W}_1$  is upper triangular. Define

$$\mathbf{W} = \left[ \begin{array}{c|c} 1 & \mathbf{0}^H \\ \hline \mathbf{0} & \mathbf{W}_1 \end{array} \right] \text{ and } \mathbf{U} = \mathbf{V} \mathbf{W}.$$

Then  $\mathbf{W}$  and  $\mathbf{U}$  (cf. Theorem 3.26) are unitary and

$$\begin{aligned} \mathbf{U}^H \mathbf{A} \mathbf{U} &= \mathbf{W}^H (\mathbf{V}^H \mathbf{A} \mathbf{V}) \mathbf{W} = \left[ \begin{array}{c|c} 1 & \mathbf{0}^H \\ \hline \mathbf{0} & \mathbf{W}_1^H \end{array} \right] \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^H \\ \hline \mathbf{0} & \mathbf{M} \end{array} \right] \left[ \begin{array}{c|c} 1 & \mathbf{0}^H \\ \hline \mathbf{0} & \mathbf{W}_1 \end{array} \right] \\ &= \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^H \mathbf{W}_1 \\ \hline \mathbf{0} & \mathbf{W}_1^H \mathbf{M} \mathbf{W}_1 \end{array} \right], \end{aligned}$$

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  $\lambda$ . Thus we can find another eigenvalue of  $\mathbf{A}$  by working with a smaller matrix  $\mathbf{M}$  and where one occurrence of  $\lambda$  has been eliminated. This is an example of a **deflation** technique which is very useful in numerical work.

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 10.1.

**Theorem 10.2** *For each  $\mathbf{A} \in \mathbb{R}^{n,n}$  with real eigenvalues there exists an orthonormal matrix  $\mathbf{U} \in \mathbb{R}^{n,n}$  such that  $\mathbf{U}^T \mathbf{A} \mathbf{U}$  is upper triangular.*

**Proof.** Consider the proof of Theorem 10.1. Since  $\mathbf{A}$  and  $\lambda_1$  are real the eigenvector  $\mathbf{v}_1$  is real and the matrix  $\mathbf{W}$  is real and orthonormal. By the induction hypothesis  $\mathbf{V}$  is real and orthonormal. But then also  $\mathbf{U} = \mathbf{V} \mathbf{W}$  is real and orthonormal.  $\square$

**Exercise 10.3** *Show that the Schur triangulation of  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 2 \end{bmatrix}$  is  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \begin{bmatrix} -1 & -2 \\ 0 & 4 \end{bmatrix}$ , where  $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$ .*

From the Schur triangulation  $\mathbf{R} = \mathbf{U}^H \mathbf{A} \mathbf{U}$  we obtain the **Schur factorization**  $\mathbf{A} = \mathbf{U} \mathbf{R} \mathbf{U}^H$ . The matrices  $\mathbf{U}$  and  $\mathbf{R}$  are called the **Schur factors**.

A real matrix with complex eigenvalues cannot be reduced to triangular form by an orthonormal similarity transformation. Indeed, if  $\mathbf{R} = \mathbf{U}^T \mathbf{A} \mathbf{U}$  is triangular,



one of the diagonal elements of  $\mathbf{R}$  (one of the eigenvalues of  $\mathbf{A}$ ) will be complex. But then  $\mathbf{U}$  cannot be real. How far can we reduce a real matrix  $\mathbf{A}$  by an orthonormal 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  $\mu, \nu$  are real. The real  $2 \times 2$  matrix

$$\mathbf{M} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix} \quad (10.2)$$

has eigenvalues  $\lambda$  and  $\bar{\lambda}$ . We say that a matrix is quasi-triangular if it is block triangular with only  $1 \times 1$  and  $2 \times 2$  blocks on the diagonal. Moreover, no  $2 \times 2$  block should have real eigenvalues. As an example consider

$$\mathbf{R} = \left[ \begin{array}{cc|cc|cc} 2 & 1 & 3 & 4 & 5 \\ -1 & 2 & 4 & 3 & 2 \\ \hline 0 & 0 & 1 & 2 & 3 \\ 0 & 0 & 0 & 3 & 2 \\ \hline 0 & 0 & 0 & -1 & 1 \end{array} \right].$$

$\mathbf{R}$  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 5.3 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  $1 \times 1$  block and complex eigenvalues  $2+i$ ,  $2-i$  with multiplicity 2 corresponding to the two  $2 \times 2$  blocks.

For a proof that any  $\mathbf{A} \in \mathbb{R}^{n,n}$  can be brought to quasi-triangular form by a real orthonormal similarity transformation see Section 10.4.

## 10.2 Hermitian and Normal Matrices

For some matrices the Schur factor  $\mathbf{R}$  will be diagonal.

**Definition 10.4 (Normal Matrix)** A matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  is said to be **normal** if  $\mathbf{A}\mathbf{A}^H = \mathbf{A}^H\mathbf{A}$ .

Examples of normal matrices are

1.  $\mathbf{A}^H = \mathbf{A}$ , (Hermitian)
2.  $\mathbf{A}^H = -\mathbf{A}$ , (Skew-Hermitian)
3.  $\mathbf{A}^H = \mathbf{A}^{-1}$ , (Unitary)
4.  $\mathbf{A} = \mathbf{D}$ . (Diagonal)

For real matrices "Hermitian" and "symmetric" are synonyms.

The following theorem says that a matrix has orthonormal eigenpairs if and only if it is normal.

**Theorem 10.5** *A matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  is unitary similar with a diagonal matrix if and only if it is normal.*

**Proof.** If  $\mathbf{B} = \mathbf{U}^H \mathbf{A} \mathbf{U}$ , with  $\mathbf{B}$  diagonal, and  $\mathbf{U}^H \mathbf{U} = \mathbf{I}$ , then

$$\begin{aligned}\mathbf{A} \mathbf{A}^H &= (\mathbf{U} \mathbf{B} \mathbf{U}^H)(\mathbf{U} \mathbf{B}^H \mathbf{U}^H) = \mathbf{U} \mathbf{B} \mathbf{B}^H \mathbf{U}^H \text{ and} \\ \mathbf{A}^H \mathbf{A} &= (\mathbf{U} \mathbf{B}^H \mathbf{U}^H)(\mathbf{U} \mathbf{B} \mathbf{U}^H) = \mathbf{U} \mathbf{B}^H \mathbf{B} \mathbf{U}^H.\end{aligned}$$

Now  $\mathbf{B} \mathbf{B}^H = \mathbf{B}^H \mathbf{B}$  since  $\mathbf{B}$  is diagonal, and  $\mathbf{A}$  is normal.

Suppose  $\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$ . By Theorem 10.1 we can find  $\mathbf{U}$  with  $\mathbf{U}^H \mathbf{U} = \mathbf{I}$  such that  $\mathbf{B} = \mathbf{U}^H \mathbf{A} \mathbf{U}$  is upper triangular. Since  $\mathbf{A}$  is normal  $\mathbf{B}$  is normal. Indeed,

$$\mathbf{B} \mathbf{B}^H = \mathbf{U}^H \mathbf{A} \mathbf{U} \mathbf{U}^H \mathbf{A}^H \mathbf{U} = \mathbf{U}^H \mathbf{A} \mathbf{A}^H \mathbf{U} = \mathbf{U}^H \mathbf{A}^H \mathbf{A} \mathbf{U} = \mathbf{B}^H \mathbf{B}.$$

The proof is complete if we can show that an upper triangular normal matrix  $\mathbf{B}$  must be diagonal. The diagonal elements in  $\mathbf{E} := \mathbf{B}^H \mathbf{B}$  and  $\mathbf{F} := \mathbf{B} \mathbf{B}^H$  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.  $\square$

### 10.2.1 The Spectral Theorem

The special cases where  $\mathbf{A}$  is Hermitian or real and symmetric deserve special attention.

**Theorem 10.6** *Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is Hermitian. Then  $\mathbf{A}$  has real eigenvalues  $\lambda_1, \dots, \lambda_n$ . Moreover, there is a unitary matrix  $\mathbf{U} \in \mathbb{C}^{n,n}$  such that  $\mathbf{U}^H \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n)$ . For the columns  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of  $\mathbf{U}$  we have  $\mathbf{A} \mathbf{u}_j = \lambda_j \mathbf{u}_j$  for  $j = 1, \dots, n$ . Thus  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  are orthonormal eigenvectors of  $\mathbf{A}$ .*

**Proof.** That the eigenvalues are real was shown in Lemma 8.12. Since a Hermitian matrix is normal it follows from Theorem 10.5 that  $\mathbf{U}^H \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n)$  for some unitary matrix  $\mathbf{U}$ . That  $\mathbf{A} \mathbf{u}_j = \lambda_j \mathbf{u}_j$  for  $j = 1, \dots, n$  follows from the initial discussion in this Chapter.  $\square$

The following real version is known as the Spectral Theorem.

**Theorem 10.7** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  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  $\mathbf{U} \in \mathbb{R}^{n,n}$  such that

$$\mathbf{U}^T \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n).$$

For the columns  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of  $\mathbf{U}$  we have  $\mathbf{A}\mathbf{u}_j = \lambda_j \mathbf{u}_j$  for  $j = 1, \dots, n$ . Thus  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  are orthonormal eigenvectors of  $\mathbf{A}$ .

**Proof.** Since  $\mathbf{A}^H = \mathbf{A}$  it follows from Theorem 10.6 that the eigenvalues are real. By Theorem 10.2 there is a matrix  $\mathbf{U} \in \mathbb{R}^{n,n}$  with  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  so that  $\mathbf{B} = \mathbf{U}^T \mathbf{A} \mathbf{U}$  is upper triangular. Since  $\mathbf{A}^T = \mathbf{A}$ , we have  $\mathbf{B}^T = \mathbf{B}$ . But then  $\mathbf{B}$  must be diagonal. The columns  $\mathbf{u}_1, \dots, \mathbf{u}_n$  of  $\mathbf{U}$  satisfies  $\mathbf{A}\mathbf{u}_j = \lambda_j \mathbf{u}_j$  for all  $j$  and are orthonormal eigenvectors of  $\mathbf{A}$ .  $\square$

**Example 10.8** 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  $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ .

**Exercise 10.9** Suppose  $\mathbf{C} = \mathbf{A} + i\mathbf{B}$ , where  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,n}$ . Show that  $\mathbf{C}$  is skew-Hermitian if and only if  $\mathbf{A}^T = -\mathbf{A}$  and  $\mathbf{B}^T = \mathbf{B}$ .

**Exercise 10.10** Show that any eigenvalue of a skew-Hermitian matrix is purely imaginary.

## 10.3 The Rayleigh Quotient and Minmax Theorems

### 10.3.1 The Rayleigh Quotient

The Rayleigh quotient is an important tool when studying eigenvalues.

**Definition 10.11** For  $\mathbf{A} \in \mathbb{C}^{n,n}$  and any  $\mathbf{x} \in \mathbb{C}^n$  the quantity  $R(\mathbf{x}) = R_{\mathbf{A}}(\mathbf{x}) := \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}}$  is called a **Rayleigh quotient** for  $\mathbf{A}$ .

If  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$  then  $R(\mathbf{x}) = \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} = \lambda$ .

**Exercise 10.12** More generally for  $\mathbf{A} \in \mathbb{C}^{n,n}$  and any  $\mathbf{y}, \mathbf{x} \in \mathbb{C}^n$  with  $\mathbf{y}^H \mathbf{x} \neq 0$  the quantity  $R(\mathbf{y}, \mathbf{x}) = R_{\mathbf{A}}(\mathbf{y}, \mathbf{x}) := \frac{\mathbf{y}^H \mathbf{A} \mathbf{x}}{\mathbf{y}^H \mathbf{x}}$  is also called a **Rayleigh quotient** for  $\mathbf{A}$ . Show that if  $(\lambda, \mathbf{x})$  is a (right) eigenpair for  $\mathbf{A}$  then  $R(\mathbf{y}, \mathbf{x}) = \lambda$  for any  $\mathbf{y}$  with  $\mathbf{y}^H \mathbf{x} \neq 0$ . Also show that if  $(\lambda, \mathbf{y})$  is a left eigenpair for  $\mathbf{A}$  then  $R(\mathbf{y}, \mathbf{x}) = \lambda$  for any  $\mathbf{x}$  with  $\mathbf{y}^H \mathbf{x} \neq 0$ .

The following lemma gives some useful formulas.

**Lemma 10.13** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  and let  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  be an orthonormal basis for a subspace  $\mathcal{S} \subset \mathbb{C}^n$ . If  $\mathbf{x}, \mathbf{y} \in \mathcal{S}$  with  $\mathbf{x} = \sum_{j=1}^k c_j \mathbf{u}_j$  and  $\mathbf{y} = \sum_{j=1}^k d_j \mathbf{u}_j$ , then

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^H \mathbf{y} = \sum_{j=1}^k \bar{c}_j d_j, \quad (10.3)$$

$$R(\mathbf{x}) = \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} = \frac{\sum_{j=1}^k \lambda_j |c_j|^2}{\sum_{j=1}^k |c_j|^2} = \sum_{j=1}^k \lambda_j |c_j|^2, \text{ if } \|\mathbf{x}\|_2 = 1. \quad (10.4)$$

**Proof.** We have

$$\langle \mathbf{x}, \mathbf{y} \rangle = \left\langle \sum_{i=1}^k c_i \mathbf{u}_i, \sum_{j=1}^k d_j \mathbf{u}_j \right\rangle = \sum_{i=1}^k \sum_{j=1}^k \bar{c}_i d_j \langle \mathbf{u}_i, \mathbf{u}_j \rangle = \sum_{j=1}^k \bar{c}_j d_j$$

and (10.3) follows. Since  $\mathbf{A} \mathbf{x} = \sum_{j=1}^k c_j \mathbf{A} \mathbf{u}_j = \sum_{j=1}^k c_j \lambda_j \mathbf{u}_j$  we obtain from (10.3) that  $\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle = \sum_{j=1}^k \lambda_j |c_j|^2$  and also  $\|\mathbf{x}\|_2^2 = \langle \mathbf{x}, \mathbf{x} \rangle = \sum_{j=1}^k |c_j|^2$ . This shows both equalities in (10.4).  $\square$

The Rayleigh quotient is especially useful when the matrix  $\mathbf{A}$  is Hermitian. Since  $\mathbf{A}$  is normal it has orthonormal eigenpairs  $\{(\lambda_1, \mathbf{u}_1), \dots, (\lambda_n, \mathbf{u}_n)\}$  and the eigenvalues are real and can be ordered, say  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ . In this case we have for any  $i \leq k$  and  $c_i, \dots, c_k$  not all zero

$$\lambda_k \leq \frac{\sum_{j=i}^k \lambda_j |c_j|^2}{\sum_{j=i}^k |c_j|^2} \leq \lambda_i, \text{ if } \lambda_i \geq \lambda_{i+1} \geq \dots \geq \lambda_k. \quad (10.5)$$

Indeed, to show the lower (upper) bound we replace all  $\lambda$ 's in the numerator  $\sum_{j=i}^k \lambda_j |c_j|^2$  by  $\lambda_k$  ( $\lambda_i$ ). From (10.5) it follows that the value of the Rayleigh quotient for a Hermitian matrix must lie between the smallest and largest eigenvalue. Since  $R(\mathbf{u}_1) = \lambda_1$  and  $R(\mathbf{u}_n) = \lambda_n$  we can express the smallest and largest eigenvalue in terms of extrema of the Rayleigh quotient.

$$\lambda_n = \min_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) = \lambda_1 \quad (10.6)$$

### 10.3.2 Minmax and Maxmin Theorems

More generally we have a minmax and maxmin characterization of the eigenvalues of a Hermitian matrix. In the following theorem  $\mathcal{S}$  is a subspace of  $\mathbb{C}^n$  of the indicated dimension.

**Theorem 10.14 (The Courant-Fischer Theorem)** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is Hermitian with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$  ordered so that  $\lambda_1 \geq \dots \geq \lambda_n$ . Then for  $k = 1, \dots, n$

$$\lambda_k = \min_{\dim(\mathcal{S})=n-k+1} \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) = \max_{\dim(\mathcal{S})=k} \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}). \quad (10.7)$$

**Proof.** We prove the maxmin version and leave the minmax version as an exercise. Let  $\{(\lambda_1, \mathbf{u}_1), \dots, (\lambda_n, \mathbf{u}_n)\}$  be orthonormal eigenpairs for  $\mathbf{A}$ . Fix  $k$ . We will show that  $\max \min R \leq \lambda_k$  and  $\max \min R \geq \lambda_k$ , where  $\max \min R$  is shorthand for the expression after the second equality in (10.7). Let  $\mathcal{S}$  be any subspace of  $\mathbb{C}^n$  of dimension  $k$  and define  $\mathcal{S}' = \text{span}\{\mathbf{u}_k, \dots, \mathbf{u}_n\}$ . Since  $\mathcal{S} + \mathcal{S}' \subset \mathbb{C}^n$  we have  $\dim(\mathcal{S} + \mathcal{S}') \leq n$  and we can use (2.6) to find

$$\dim(\mathcal{S} \cap \mathcal{S}') = \dim(\mathcal{S}) + \dim(\mathcal{S}') - \dim(\mathcal{S} + \mathcal{S}') \geq k + (n - k + 1) - n = 1,$$

and it follows that  $\mathcal{S} \cap \mathcal{S}'$  is nonempty. Let  $\mathbf{y} \in \mathcal{S} \cap \mathcal{S}' = \sum_{j=k}^n d_j \mathbf{u}_j$ . By (10.4) applied to  $\mathcal{S}'$  we find

$$R(\mathbf{y}) = \frac{\sum_{j=k}^n \lambda_j |d_j|^2}{\sum_{j=k}^n |d_j|^2} \leq \lambda_k.$$

This implies that  $\min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \leq \lambda_k$  and therefore, since  $\mathcal{S}$  is arbitrary,  $\max \min R \leq \lambda_k$ . To show the inequality in the opposite direction we use the subspace  $\mathcal{S}_k := \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ . Suppose  $\mathbf{x} = \sum_{j=1}^k c_j \mathbf{u}_j$  is any nonzero element in  $\mathcal{S}_k$ . Then

$$R(\mathbf{x}) = \frac{\sum_{j=1}^k \lambda_j |c_j|^2}{\sum_{j=1}^k |c_j|^2} \geq \lambda_k.$$

Since  $\mathbf{x}$  is arbitrary we obtain  $\min_{\substack{\mathbf{x} \in \mathcal{S}_k \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \geq \lambda_k$  and therefore  $\max \min R \geq \lambda_k$ .  $\square$

**Exercise 10.15** *Modify the proof of the maxmin version of the Courant-Fischer theorem to prove the minmax version.*

Using Theorem 10.14 we can prove inequalities of eigenvalues without knowing the eigenvectors and we can get both upper and lower bounds.

**Corollary 10.16** *Let  $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{C}^{n,n}$  be Hermitian with eigenvalues  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ ,  $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$ , and  $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_n$ , respectively. If  $\mathbf{C} = \mathbf{A} + \mathbf{B}$  then*

$$\alpha_i + \beta_n \leq \gamma_i \leq \alpha_i + \beta_1, \text{ for } i = 1, 2, \dots, n. \quad (10.8)$$

**Proof.** Let  $\mathbf{u}_1, \dots, \mathbf{u}_n$  be orthonormal eigenvectors for  $\mathbf{A}$  and let for fixed  $i$ ,  $\mathcal{S} := \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_{n-i+1}\}$ . By Theorem 10.14 and (10.5) we obtain

$$\gamma_i \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{C}}(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{A}}(\mathbf{x}) + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{B}}(\mathbf{x}) = \alpha_i + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{B}}(\mathbf{x}) \leq \alpha_i + \beta_1,$$

and this proves the upper inequality. For the lower one we define  $\mathbf{D} := -\mathbf{B}$  and observe that  $-\beta_n$  is the largest eigenvalue of  $\mathbf{D}$ . Since  $\mathbf{A} = \mathbf{C} + \mathbf{D}$  it follows from the result just proved that  $\alpha_i \leq \gamma_i - \beta_n$ , which is the same as the lower inequality.  $\square$

In many applications of this result the matrix  $\mathbf{B}$  will be small and then the theorem states that the eigenvalues of  $\mathbf{C}$  are close to those of  $\mathbf{A}$ . Moreover, it associates a unique eigenvalue of  $\mathbf{A}$  with each eigenvalue of  $\mathbf{C}$ .

**Exercise 10.17** Show that in Corollary 10.16, if  $\mathbf{B}$  is symmetric positive semidefinite then  $\gamma_i \geq \alpha_i$ .

### 10.3.3 The Hoffman-Wielandt Theorem

We can also give a bound involving all eigenvalues.

**Theorem 10.18 (Hoffman-Wielandt Theorem)** Suppose  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$  are both normal matrices with eigenvalues  $\lambda_1, \dots, \lambda_n$  and  $\mu_1, \dots, \mu_n$ , respectively. Then there is a permutation  $i_1, \dots, i_n$  of  $1, 2, \dots, n$  such that

$$\sum_{j=1}^n |\mu_{i_j} - \lambda_j|^2 \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2. \quad (10.9)$$

Taking  $\mathbf{B} = \mathbf{A} + \mathbf{E}$  this shows, in terms of absolute error, that as long as  $\mathbf{A} + \mathbf{E}$  is normal, i. e., we perturb in a "normal way", then the eigenvalue problem for a normal matrix is well conditioned. Small perturbation in the elements of  $\mathbf{A}$  lead to small changes in the eigenvalues.

For a proof of this theorem see [[19], 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 10.19** Show that (10.9) 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?

## 10.4 Proof of the Real Schur Form

In this section we prove the following theorem.

**Theorem 10.20** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$ . Then we can find  $\mathbf{U} \in \mathbb{R}^{n,n}$  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 10.2 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  $\lambda$ . 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

$$\mathbf{A}\mathbf{x} = \mu\mathbf{x} - \nu\mathbf{y}, \quad \mathbf{A}\mathbf{y} = \nu\mathbf{x} + \mu\mathbf{y}. \quad (10.10)$$

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 (10.10) 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  $\alpha$ . Replacing  $\mathbf{y}$  by  $\alpha\mathbf{x}$  in (10.10), we find  $\mathbf{Ax} = (\mu - \alpha\nu)\mathbf{x}$  and  $\mathbf{Ax} = \mathbf{Ay}/\alpha = (\mu + \nu/\alpha)\mathbf{x}$ . 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.

(10.10) can be written in matrix form as

$$\mathbf{AX}_1 = \mathbf{X}_1\mathbf{M}, \quad \mathbf{X}_1 = (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n,2}, \quad (10.11)$$

where  $\mathbf{M}$  is given by (10.2). By Theorem 16.3 we can find an orthonormal matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$  such that

$$\mathbf{QX}_1 = \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix}$$

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

$$\mathbf{QX} = [\mathbf{QX}_1, \mathbf{Qq}_3, \dots, \mathbf{Qq}_n] = \begin{bmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix}.$$

Since  $\mathbf{R}$  is nonsingular,  $\mathbf{QX}$  and  $\mathbf{X}$  are nonsingular. Moreover, using (10.11)

$$\mathbf{X}^{-1}\mathbf{AX} = [\mathbf{X}^{-1}\mathbf{AX}_1, \mathbf{X}^{-1}\mathbf{AX}_2] = [\mathbf{X}^{-1}\mathbf{X}_1\mathbf{M}, \mathbf{X}^{-1}\mathbf{AX}_2] = \begin{bmatrix} \mathbf{M} & \mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix}$$

for some matrices  $\mathbf{B} \in \mathbb{R}^{2,n-2}$ ,  $\mathbf{C} \in \mathbb{R}^{n-2,n-2}$ . Now

$$\mathbf{QAQ}^T = (\mathbf{QX})\mathbf{X}^{-1}\mathbf{AX}(\mathbf{QX})^{-1} = \begin{bmatrix} \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{R}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-2} \end{bmatrix},$$

or

$$\mathbf{QAQ}^T = \begin{bmatrix} \mathbf{RMR}^{-1} & \mathbf{RB} \\ \mathbf{0} & \mathbf{C} \end{bmatrix}. \quad (10.12)$$

By Theorem 5.15 the  $2 \times 2$  matrix  $\mathbf{RMR}^{-1}$  has the same eigenvalues  $\lambda$  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  $\mathbf{U} = \mathbf{QV}$ . Then  $\mathbf{U} \in \mathbb{R}^{n,n}$ ,  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  and  $\mathbf{U}^T \mathbf{AU}$  is quasi-triangular.  $\square$





## Chapter 11

# 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.

### 11.1 Singular Values and Singular Vectors

We know from Theorem 10.5 that a square matrix  $\mathbf{A}$  can be diagonalized by a unitary similarity transformation if and only if it is normal. In particular, if  $\mathbf{A} \in \mathbb{C}^{n,n}$  is normal with eigenvalues  $\lambda_1, \dots, \lambda_n$  then

$$\mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n) \text{ or } \mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^H, \text{ where } \mathbf{U}^H \mathbf{U} = \mathbf{I}. \quad (11.1)$$

In this section we show that any matrix, even a rectangular one, can be diagonalized provided we allow two different unitary matrices. Thus

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^H, \text{ where } \mathbf{\Sigma} \text{ is a diagonal matrix, } \mathbf{U}^H \mathbf{U} = \mathbf{I}, \text{ and } \mathbf{V}^H \mathbf{V} = \mathbf{I}. \quad (11.2)$$

The diagonal entries of  $\mathbf{\Sigma}$ , are called singular values and the columns of  $\mathbf{U}$  and  $\mathbf{V}$  are singular vectors. The formula  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^H$  is known as the singular value decomposition of  $\mathbf{A}$ .

#### 11.1.1 SVD and SVF

Every matrix has a singular value decomposition (SVD) and a reduced form called the singular value factorization (SVF). To derive these we start with a lemma and a theorem.

**Lemma 11.1** *Suppose  $m, n \in \mathbb{N}$  and  $\mathbf{A} \in \mathbb{C}^{m,n}$ . The matrix  $\mathbf{A}^H \mathbf{A}$  has eigenpairs  $(\lambda_j, \mathbf{v}_j)$  for  $j = 1, \dots, n$ , where  $\mathbf{v}_j^H \mathbf{v}_k = \delta_{jk}$  and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ . Moreover*

$$\sigma_j := \sqrt{\lambda_j} = \|\mathbf{A} \mathbf{v}_j\|_2, \text{ for } j = 1, \dots, n. \quad (11.3)$$

**Proof.** The matrix  $\mathbf{A}^H \mathbf{A} \in \mathbb{C}^{n,n}$  is Hermitian, and by Theorem 10.6 it has real eigenvalues  $\lambda_j$  and orthonormal eigenvectors  $\mathbf{v}_j$  for  $j = 1, \dots, n$ . For each  $j$   $\|\mathbf{A}\mathbf{v}_j\|_2^2 = (\mathbf{A}\mathbf{v}_j)^H \mathbf{A}\mathbf{v}_j = \mathbf{v}_j^H \mathbf{A}^H \mathbf{A} \mathbf{v}_j = \mathbf{v}_j^H \lambda_j \mathbf{v}_j = \lambda_j$ , since  $\mathbf{v}_j^H \mathbf{v}_j = 1$ , and (11.3) follows.  $\square$

The nonnegative square roots of the  $n$  eigenvalues of  $\mathbf{A}^H \mathbf{A}$  are called the **singular values** of  $\mathbf{A} \in \mathbb{C}^{m,n}$ . They are usually ordered so that

$$\sigma_1 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots = \sigma_n. \quad (11.4)$$

We will show that the number  $r$  of positive singular values equals the rank of  $\mathbf{A}$ . Moreover, the eigenvectors of  $\mathbf{A}^H \mathbf{A}$  determine orthonormal bases for the column space  $\text{span}(\mathbf{A})$  and null space  $\ker(\mathbf{A})$  of  $\mathbf{A}$ .

**Theorem 11.2** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  and let  $(\sigma_j^2, \mathbf{v}_j)$  for  $j = 1, \dots, n$  be orthonormal eigenpairs for  $\mathbf{A}^H \mathbf{A}$  with  $\sigma_1, \dots, \sigma_n$  ordered as in (11.4). Then  $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$  is an orthogonal basis for the column space  $\text{span}(\mathbf{A})$  of  $\mathbf{A}$  and  $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$  is an orthonormal basis for the nullspace  $\ker(\mathbf{A})$  of  $\mathbf{A}$ .

**Proof.** The proof will be complete if we can show

1.  $\mathbf{A}\mathbf{v}_j \neq \mathbf{0}$  if and only if  $1 \leq j \leq r$ .
2.  $(\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r)$  are orthogonal and nonzero.
3.  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{v}_j \Rightarrow \mathbf{A}\mathbf{x} = \sum_{j=1}^r c_j \mathbf{A}\mathbf{v}_j$ .
4.  $\text{span}\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\} \subset \text{span}(\mathbf{A})$ .
5.  $\text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\} \subset \ker(\mathbf{A})$ .
6.  $\text{span}(\mathbf{A}) \subset \text{span}\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ .
7.  $\ker(\mathbf{A}) \subset \text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ .

1 follows by combining (11.3) and (11.4). 1 and the calculation

$$(\mathbf{A}\mathbf{v}_j)^H \mathbf{A}\mathbf{v}_k = \mathbf{v}_j^H \mathbf{A}^H \mathbf{A} \mathbf{v}_k = \mathbf{v}_j^H \sigma_k^2 \mathbf{v}_k = 0, j \neq k$$

implies 2. 3 is a consequence of 1. Clearly  $\mathbf{A}\mathbf{v}_j \in \text{span}(\mathbf{A})$  for  $j = 1, \dots, r$  and  $\mathbf{v}_j \in \ker(\mathbf{A})$  by 1. Since  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A})$  are vector spaces 4 and 5 follow. 3 implies 6. Finally, by 2  $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$  are linearly independent. So if  $\mathbf{x} \in \ker(\mathbf{A})$  then by 3  $c_1 = \dots = c_r = 0$  and 7 follows.  $\square$

Every matrix has a singular value decomposition.

**Theorem 11.3 (SVD)** Let  $m, n \in \mathbb{N}$  and suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  has rank  $r$ . Then  $\mathbf{A}$  has exactly  $r$  positive singular values  $\sigma_1 \geq \dots \geq \sigma_r > 0$ . Moreover,  $\mathbf{A}$  has the singular value decomposition

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^H, \quad \mathbf{U} \in \mathbb{C}^{m,m}, \quad \mathbf{\Sigma} \in \mathbb{R}^{m,n}, \quad \mathbf{V} \in \mathbb{C}^{n,n},$$

where  $U$  and  $V$  are unitary and

$$\Sigma := \begin{bmatrix} \Sigma_1 & \mathbf{0}_{r,n-r} \\ \mathbf{0}_{m-r,r} & \mathbf{0}_{m-r,n-r} \end{bmatrix} \in \mathbb{R}^{m,n}, \text{ where } \Sigma_1 := \text{diag}(\sigma_1, \dots, \sigma_r). \quad (11.5)$$

Here  $\mathbf{0}_{k,l} \in \mathbb{R}^{k,l}$  is the zero matrix and  $\mathbf{0}_{k,l} = [\ ]$  is the empty matrix, if  $k = 0$  or  $l = 0$ .

If  $A$  is real then  $A = U\Sigma V^T$ , where  $U \in \mathbb{R}^{m,m}$  and  $V \in \mathbb{R}^{n,n}$  are orthonormal, and  $\Sigma$  is again given by (11.5).

**Proof.** Suppose  $(\lambda_j, \mathbf{v}_j)$  for  $j = 1, \dots, n$  are orthonormal eigenpairs for  $A^H A$  and define  $\Sigma$  by (11.5), where  $\sigma_j = \sqrt{\lambda_j}$  for all  $j$ . By Theorem 11.2 the set  $\{A\mathbf{v}_1, \dots, A\mathbf{v}_r\}$  is an orthogonal basis for the column space of  $A$  and it follows that  $r$  is the number of positive singular values. We turn  $\{A\mathbf{v}_1, \dots, A\mathbf{v}_r\}$  into an orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  for  $\text{span}(A)$  by setting (cf. (11.3))

$$\mathbf{u}_j := \frac{A\mathbf{v}_j}{\|A\mathbf{v}_j\|_2} = \frac{1}{\sigma_j} A\mathbf{v}_j, \text{ for } j = 1, \dots, r.$$

By Theorem 11.2

$$A\mathbf{v}_j = \sigma_j \mathbf{u}_j, \quad j = 1, \dots, r \text{ and } A\mathbf{v}_j = \mathbf{0}, \quad j = r+1, \dots, n. \quad (11.6)$$

We extend  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  to an orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$  for  $\mathbb{C}^m$  and define

$$U := [\mathbf{u}_1, \dots, \mathbf{u}_m] \in \mathbb{C}^{m,m} \text{ and } V := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n,n}.$$

Since  $U$  and  $V$  have orthonormal columns they are unitary matrices, and from (11.5) and (11.6)

$$U\Sigma = U[\sigma_1 \mathbf{e}_1, \dots, \sigma_r \mathbf{e}_r, \mathbf{0}, \dots, \mathbf{0}] = [\sigma_1 \mathbf{u}_1, \dots, \sigma_r \mathbf{u}_r, \mathbf{0}, \dots, \mathbf{0}] = [A\mathbf{v}_1, \dots, A\mathbf{v}_n].$$

Thus  $U\Sigma = AV$  and since  $V$  is unitary we find  $U\Sigma V^H = AVV^H = A$ .

For a matrix with real entries the eigenvectors of  $A^T A$  are real and the singular value decomposition takes the stated form.  $\square$

From the singular value decomposition we obtain a reduced factorization called the singular value factorization and an outer product form of this factorization.

**Corollary 11.4 (SVF)** Suppose  $A = U\Sigma V^H$  is a singular value decomposition of a rank  $r$  matrix  $A \in \mathbb{C}^{m,n}$ . Then  $A$  has the singular value factorization

$$A = U_1 \Sigma_1 V_1^H = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^H, \quad U_1 \in \mathbb{C}^{m,r}, \quad \Sigma_1 \in \mathbb{R}^{r,r}, \quad V_1 \in \mathbb{C}^{n,r},$$

where

$$\begin{aligned} \Sigma_1 &= \text{diag}(\sigma_1, \dots, \sigma_r), \\ U &= [\mathbf{u}_1, \dots, \mathbf{u}_m] = [U_1, U_2], \quad U_1 \in \mathbb{C}^{m,r}, \quad U_2 \in \mathbb{C}^{m,m-r}, \\ V &= [\mathbf{v}_1, \dots, \mathbf{v}_n] = [V_1, V_2], \quad V_1 \in \mathbb{C}^{n,r}, \quad V_2 \in \mathbb{C}^{n,n-r}, \end{aligned} \quad (11.7)$$

**Proof.** We find

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^H = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^H.$$

□

SVF and SVD are not unique. The singular values are unique since they are the nonnegative square roots of the eigenvalues of  $\mathbf{A}^H \mathbf{A}$ . However the matrices  $\mathbf{U}$  and  $\mathbf{V}$  are in general not uniquely given.

### 11.1.2 Examples

**Example 11.5 (Nonsingular matrix)** Derive the following SVD.

$$\mathbf{A} := \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}. \quad (11.8)$$

*Discussion:* The eigenpairs of  $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 97 & 96 \\ 96 & 153 \end{bmatrix} / 25$  are

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 9 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Taking square roots and normalizing we find  $\sigma_1 = 3$ ,  $\sigma_2 = 1$ ,  $\mathbf{v}_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix} / 5$ ,  $\mathbf{v}_2 = \begin{bmatrix} 4 \\ -3 \end{bmatrix} / 5$ . Thus  $\mathbf{u}_1 := \mathbf{A}\mathbf{v}_1 / \sigma_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix} / 5$  and  $\mathbf{u}_2 := \mathbf{A}\mathbf{v}_2 / \sigma_2 = \begin{bmatrix} -4 \\ 3 \end{bmatrix} / 5$  and this shows (11.8). Since  $m = n = r$  we have  $\mathbf{U}_1 = \mathbf{U}$ ,  $\mathbf{\Sigma}_1 = \mathbf{\Sigma}$  and  $\mathbf{V}_1 = \mathbf{V}$ . In general the SVD and SVF are the same for a nonsingular matrix. See also Example 11.11 for some further discussion.

**Example 11.6 (Full row rank)** Find the singular value decomposition of

$$\mathbf{A} := \frac{1}{15} \begin{bmatrix} 14 & 4 & 16 \\ 2 & 22 & 13 \end{bmatrix} \in \mathbb{R}^{2,3}.$$

*Discussion:* The eigenpairs of  $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 8 & 4 & 10 \\ 4 & 20 & 14 \\ 10 & 14 & 17 \end{bmatrix} / 9$  are

$$\mathbf{B} \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} = 0 \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix}.$$

Thus  $r = 2$  and

$$\mathbf{\Sigma} := \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{V} := \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix}.$$

From (11.6) we find  $\mathbf{u}_1 = \mathbf{A}\mathbf{v}_1 / \sigma_1 = [3, 4]^T / 5$ , and  $\mathbf{u}_2 = \mathbf{A}\mathbf{v}_2 / \sigma_2 = [4, -3]^T / 5$  and

$$\mathbf{U} = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}.$$

Since  $r = 2$  it follows that  $\text{rank}(\mathbf{A}) = 2$ ,  $\{\mathbf{u}_1, \mathbf{u}_2\}$  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

$$\mathbf{A} = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \end{bmatrix} = 2 \frac{1}{15} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \end{bmatrix} + 1 \frac{1}{15} \begin{bmatrix} -4 \\ -3 \end{bmatrix} \begin{bmatrix} 2 & -2 & 1 \end{bmatrix}.$$

**Example 11.7 (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 11.6 we can simply transpose the SVD of  $\mathbf{A}$  in that example. Thus

$$\mathbf{A}_1 = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = \mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T = \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}. \quad (11.9)$$

Alternatively we can follow the recipe from the previous examples. The eigenpairs of

$$\mathbf{B}_1 = \mathbf{A}_1^T \mathbf{A}_1 = \frac{1}{25} \begin{bmatrix} 52 & 36 \\ 36 & 73 \end{bmatrix}$$

are

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 4 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = 1 \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Thus  $\sigma_1 = 2$ ,  $\sigma_2 = 1$ . Now

$$\mathbf{u}_1 = \mathbf{A}_1 \mathbf{v}_1 / \sigma_1 = [1, 2, 2]^T / 3, \quad \mathbf{u}_2 = \mathbf{A}_1 \mathbf{v}_2 / \sigma_1 = [2, -2, 1]^T / 3.$$

Since  $m = 3$  we also need  $\mathbf{u}_3$  which should be orthogonal to  $\mathbf{u}_1$  and  $\mathbf{u}_2$ .  $\mathbf{u}_3 = [2, 1, -2]^T$  is such a vector and we obtain (11.9).

**Example 11.8** ( $r < n < m$ ) Consider

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}.$$

For this matrix all the zero matrices in (11.5) are nonempty. The eigenpairs of

$$\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

are

$$\mathbf{B} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0 \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

and we find  $\sigma_1 = 2$ ,  $\sigma_2 = 0$ , Thus  $r = 1$ ,  $m = 3$ ,  $n = 2$  and

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \Sigma_1 = [2], \quad V = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Now (11.6) implies  $\mathbf{u}_1 = \mathbf{A}\mathbf{v}_1/\sigma_1 = \mathbf{s}_1/\sqrt{2}$ , where  $\mathbf{s}_1 = [1, 1, 0]^T$ . To find the other columns of  $\mathbf{U}$  we 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

$$\mathbf{s}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{s}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{s}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

we find from (2.23)

$$\mathbf{w}_1 = \mathbf{s}_1, \quad \mathbf{w}_2 = \mathbf{s}_2 - \frac{\mathbf{s}_2^T \mathbf{w}_1}{\mathbf{w}_1^T \mathbf{w}_1} \mathbf{w}_1 = \begin{bmatrix} -1/2 \\ 1/2 \\ 0 \end{bmatrix}, \quad \mathbf{w}_3 = \mathbf{s}_3 - \frac{\mathbf{s}_3^T \mathbf{w}_1}{\mathbf{w}_1^T \mathbf{w}_1} \mathbf{w}_1 - \frac{\mathbf{s}_3^T \mathbf{w}_2}{\mathbf{w}_2^T \mathbf{w}_2} \mathbf{w}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Normalizing the  $\mathbf{w}_i$ 's we obtain  $\mathbf{u}_1 = \mathbf{s}_1/\|\mathbf{s}_1\|_2 = [1/\sqrt{2}, 1/\sqrt{2}, 0]^T$ ,  $\mathbf{u}_2 = \mathbf{s}_2/\|\mathbf{s}_2\|_2 = [-1/\sqrt{2}, 1/\sqrt{2}, 0]^T$ , and  $\mathbf{u}_3 = \mathbf{s}_3/\|\mathbf{s}_3\|_2 = [0, 0, 1]^T$ . Therefore,  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ , where

$$\mathbf{U} := \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{3,3}, \quad \Sigma := \begin{bmatrix} 2 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{3,2}, \quad \mathbf{V} := \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \in \mathbb{R}^{2,2}.$$

**Exercise 11.9** Find the singular value decomposition of the following matrices

(a)  $A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}.$

(b)  $A = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 2 & 2 \end{bmatrix}.$

**Exercise 11.10** Find the singular value decomposition of the following matrices

(a)  $A = \mathbf{e}_1$  the first unit vector in  $\mathbb{R}^m$ .

(b)  $A = \mathbf{e}_n^T$  the last unit vector in  $\mathbb{R}^n$ .

(c)  $A = \begin{bmatrix} -1 & 0 \\ 0 & 3 \end{bmatrix}.$

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  $\mathbf{A}^H \mathbf{A}$  can lead to extra errors in the computation. State of the art computer implementations of the singular value decomposition use an adapted version of the QR algorithm where the matrix  $\mathbf{A}^H \mathbf{A}$  is not formed. The QR algorithm is discussed in Chapter 19.

### 11.1.3 Singular Values of Normal and Positive Semidefinite Matrices

The singular values of a normal matrix are the absolute values of its eigenvalues. For if  $\mathbf{A} \in \mathbb{C}^{n,n}$  is normal with eigenvalues  $|\lambda_1| \geq \dots \geq |\lambda_n|$ , then it follows from Theorem 10.5 that  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^H$ , where  $\mathbf{U}^H\mathbf{U} = \mathbf{I}$ , and  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is a diagonal matrix. We find  $\mathbf{A}^H\mathbf{A} = \mathbf{U}\mathbf{D}^H\mathbf{D}\mathbf{U}^H$ , where  $\mathbf{D}^H\mathbf{D} = \text{diag}(|\lambda_1|^2, \dots, |\lambda_n|^2)$ . It follows that  $\sigma_i^2 = |\lambda_i|^2$  or  $\sigma_i = |\lambda_i|$  for  $i = 1, \dots, n$ .

For a symmetric positive semi-definite matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  the singular values are identical to the eigenvalues. The factorization  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^H$  above is both a SVD decomposition and factorization provided we have sorted the nonnegative eigenvalues in nondecreasing order.

**Example 11.11** *The matrix  $\mathbf{A}$  in Example 11.5 is normal so that the singular values of  $\mathbf{A}$  are equal to the absolute value of the eigenvalues of  $\mathbf{A}$ . The eigenvalues of  $\mathbf{A}$  are  $\lambda_1 = 3$  and  $\lambda_2 = -1$ . Thus  $\lambda_2 \neq \sigma_2$ .*

### 11.1.4 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  $\mathbf{z} \rightarrow \mathbf{A}\mathbf{z}$ . 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  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  is the singular value decomposition of  $\mathbf{A}$ . Since  $\mathbf{A}$  has rank  $n$  we have  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$ , with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$  and  $\mathbf{A}^{-1} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$ . Let  $\mathcal{E} := \mathbf{A}\mathcal{S} \subset \mathbb{C}^n$  be the image of  $\mathcal{S}$  under the transformation  $\mathbf{T}$ . If  $\mathbf{x} \in \mathcal{E}$  then  $\mathbf{x} = \mathbf{A}\mathbf{z}$  for some  $\mathbf{z} \in \mathcal{S}$  and we find

$$\begin{aligned} 1 &= \|\mathbf{z}\|_2^2 = \|\mathbf{A}^{-1}\mathbf{A}\mathbf{z}\|_2^2 = \|\mathbf{A}^{-1}\mathbf{x}\|_2^2 = \|\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{x}\|_2^2 \\ &= \|\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{x}\|_2^2 = \|\mathbf{\Sigma}^{-1}\mathbf{y}\|_2^2 = \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2}, \end{aligned}$$

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  $\sigma_j$  along the unit vectors  $\mathbf{e}_j$  for  $j = 1, \dots, n$ . Since the orthogonal transformation  $\mathbf{U}\mathbf{y} \rightarrow \mathbf{x}$  preserves length, the image  $\mathcal{E} = \mathbf{A}\mathcal{S}$  is an ellipsoid with semiaxes along the left singular vectors  $\mathbf{u}_j = \mathbf{U}\mathbf{e}_j$  of length  $\sigma_j$ . Since  $\mathbf{A}\mathbf{v}_j = \sigma_j\mathbf{u}_j$ , the right singular vectors are orthogonal points in  $\mathcal{S}$  that are mapped onto the semiaxes of  $\mathcal{E}$ .

**Example 11.12** *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 11.5. Recall that  $\sigma_1 = 3$ ,  $\sigma_2 = 1$ ,  $\mathbf{u}_1 = [3, 4]^T/5$  and  $\mathbf{u}_2 = [-4, 3]^T/5$ . The ellipsoids  $y_1^2/\sigma_1^2 + y_2^2/\sigma_2^2 = 1$  and  $\mathcal{E} = \mathbf{A}\mathcal{S}$  are shown in Figure 11.1. 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 ellipsoid on the right is

$$\frac{(\frac{3}{5}x_1 + \frac{4}{5}x_2)^2}{9} + \frac{(-\frac{4}{5}x_1 + \frac{3}{5}x_2)^2}{1} = 1,$$

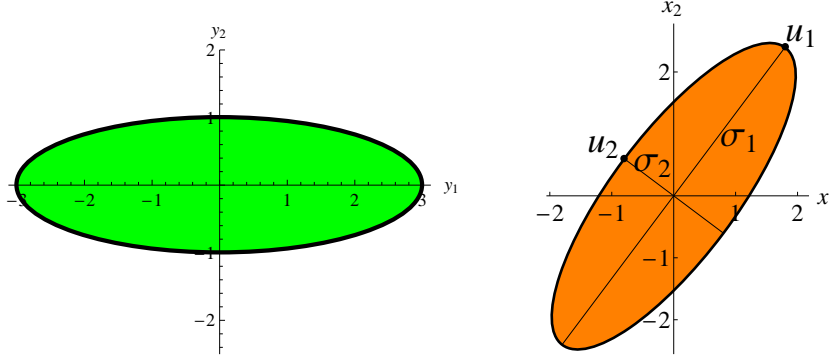


Figure 11.1. The ellipse  $y_1^2/9 + y_2^2 = 1$  (left) and the rotated ellipse  $\mathbf{A}\mathbf{S}$  (right).

## 11.2 Singular Vectors

The columns  $\mathbf{u}_1, \dots, \mathbf{u}_m$  of  $\mathbf{U}$  are called **left singular vectors**, and the columns  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of  $\mathbf{V}$  are called **right singular vectors**. These vectors satisfy the following relations.

**Theorem 11.13** If  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  is the singular value decomposition of  $\mathbf{A}$  then

$$\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma} \text{ and } \mathbf{A}^H\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^H. \quad (11.10)$$

If  $\mathbf{U}$  and  $\mathbf{V}$  are partitioned as in (11.7) then

$$\begin{aligned} 1. \mathbf{A}\mathbf{V}_1 &= \mathbf{U}_1\mathbf{\Sigma}_1, & \text{or } \mathbf{A}\mathbf{v}_i &= \sigma_i\mathbf{u}_i \text{ for } i = 1, \dots, r, \\ 2. \mathbf{A}\mathbf{V}_2 &= \mathbf{0}, & \text{or } \mathbf{A}\mathbf{v}_i &= \mathbf{0} \text{ for } i = r+1, \dots, n, \\ 3. \mathbf{A}^H\mathbf{U}_1 &= \mathbf{V}_1\mathbf{\Sigma}_1, & \text{or } \mathbf{A}^H\mathbf{u}_i &= \sigma_i\mathbf{v}_i \text{ for } i = 1, \dots, r, \\ 4. \mathbf{A}^H\mathbf{U}_2 &= \mathbf{0}, & \text{or } \mathbf{A}^H\mathbf{u}_i &= \mathbf{0} \text{ for } i = r+1, \dots, m. \end{aligned} \quad (11.11)$$

**Proof.** Since  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$  the first equation in (11.10) follows. Taking conjugate transposes and multiplying by  $\mathbf{U}$  we have  $\mathbf{A}^H\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^H\mathbf{U}^H\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^H$  and the second relation follows. In terms of partitioned matrices, (11.10) gives

$$\mathbf{A}[\mathbf{v}_1, \mathbf{v}_2] = [\mathbf{u}_1, \mathbf{u}_2] \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{A}^H[\mathbf{u}_1, \mathbf{u}_2] = [\mathbf{v}_1, \mathbf{v}_2] \begin{bmatrix} \mathbf{\Sigma}_1^H & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},$$

and this leads to the equations in (11.11).  $\square$



**Theorem 11.14** *The singular vectors of  $\mathbf{A} \in \mathbb{C}^{m,n}$  are orthonormal bases for the four fundamental subspaces of  $\mathbf{A}$ . In particular*

1.  $\mathbf{U}_1$  is an orthonormal basis for  $\text{span}(\mathbf{A})$ ,
  2.  $\mathbf{V}_2$  is an orthonormal basis for  $\ker(\mathbf{A})$ ,
  3.  $\mathbf{V}_1$  is an orthonormal basis for  $\text{span}(\mathbf{A}^H)$ ,
  4.  $\mathbf{U}_2$  is an orthonormal basis for  $\ker(\mathbf{A}^H)$ .
- (11.12)

**Proof.** Since  $\mathbf{u}_j = \mathbf{A}\mathbf{v}_j/\sigma_j$  for  $j = 1, \dots, r$  it follows from Theorem 11.2 that the first  $r$  left singular vectors of  $\mathbf{A}$  form an orthonormal basis for  $\text{span}(\mathbf{A})$  and the last  $n - r$  right singular vectors of  $\mathbf{A}$  form an orthonormal basis for  $\ker(\mathbf{A})$ . The same holds for  $\mathbf{A}^H$  and we have seen that the left and right singular vectors for  $\mathbf{A}^H$  are the columns of  $\mathbf{V}$  and  $\mathbf{U}$ , respectively.  $\square$

By counting the number of columns in the four submatrices  $\mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1, \mathbf{V}_2$ , we obtain from Theorem 11.14 a new proof of the following fundamental result (Cf. Theorem 3.16).

**Corollary 11.15** *Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$ . Then*

1.  $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}) = n$ ,
2.  $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}^H) = m$ ,
3.  $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^H)$ .

**Exercise 11.16** *Let  $\mathbf{A}$  and  $\mathbf{B}$  be as in Example 11.7. Give orthonormal bases for  $\text{span}(\mathbf{B})$  and  $\ker(\mathbf{B})$  and explain why  $\text{span}(\mathbf{B}) \oplus \ker(\mathbf{B})$  is an orthogonal decomposition of  $\mathbb{R}^3$ .*

### 11.2.1 The SVD of $\mathbf{A}^H\mathbf{A}$ and $\mathbf{A}\mathbf{A}^H$

The singular value decomposition of  $\mathbf{A}^H\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^H$  is related to the spectral decomposition of these matrices.

**Theorem 11.17** *Suppose  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H = \mathbf{U}_1\Sigma_1\mathbf{V}_1^H$  is the singular value decomposition and factorization of  $\mathbf{A}$ . Then a singular value decomposition and factorization of the matrices  $\mathbf{A}^H\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^H$  are given by*

$$\mathbf{A}^H\mathbf{A} = \mathbf{V}\Sigma^H\Sigma\mathbf{V}^H = \mathbf{V}_1\Sigma_1^2\mathbf{V}_1^H \text{ and } \mathbf{A}\mathbf{A}^H = \mathbf{U}\Sigma\Sigma^H\mathbf{U}^H = \mathbf{U}_1\Sigma_1^2\mathbf{U}_1^H. \quad (11.13)$$

Moreover,

$$\mathbf{A}^H\mathbf{A}\mathbf{V}_1 = \mathbf{V}_1\Sigma_1^2, \quad \mathbf{A}^H\mathbf{A}\mathbf{V}_2 = \mathbf{0}, \quad (11.14)$$

and

$$\mathbf{A}\mathbf{A}^H\mathbf{U}_1 = \mathbf{U}_1\Sigma_1^2, \quad \mathbf{A}\mathbf{A}^H\mathbf{U}_2 = \mathbf{0}, \quad (11.15)$$

**Proof.** We compute  $\mathbf{A}^H \mathbf{A} = \mathbf{V} \boldsymbol{\Sigma}^H \mathbf{U}^H \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^H = \mathbf{V} \boldsymbol{\Sigma}^H \boldsymbol{\Sigma} \mathbf{V}^H$  and

$$\mathbf{V} \boldsymbol{\Sigma}^H \boldsymbol{\Sigma} \mathbf{V}^H = [\mathbf{V}_1, \mathbf{V}_2] \begin{bmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} = \mathbf{V}_1 \boldsymbol{\Sigma}_1^2 \mathbf{V}_1^H,$$

with an analogous computation for  $\mathbf{A} \mathbf{A}^H$ . The equation (11.15) follows from the computation  $\mathbf{A} \mathbf{A}^H \mathbf{U}_1 = \mathbf{U}_1 \boldsymbol{\Sigma}_1^2 \mathbf{U}_1^H \mathbf{U}_1 = \mathbf{U}_1 \boldsymbol{\Sigma}_1^2$ ,  $\mathbf{A} \mathbf{A}^H \mathbf{U}_2 = \mathbf{U}_1 \boldsymbol{\Sigma}_1^2 \mathbf{U}_1^H \mathbf{U}_2 = \mathbf{0}$ . The proof of (11.14) is analogous.  $\square$

Theorem 11.17 leads to

**Theorem 11.18** For any  $\mathbf{A} \in \mathbb{C}^{m,n}$  we have

1.  $\text{rank } \mathbf{A} = \text{rank}(\mathbf{A}^H \mathbf{A}) = \text{rank}(\mathbf{A} \mathbf{A}^H)$ ,
2.  $\text{null}(\mathbf{A}^H \mathbf{A}) = \text{null } \mathbf{A}$ , and  $\text{null}(\mathbf{A} \mathbf{A}^H) = \text{null}(\mathbf{A}^H)$ ,
3.  $\text{span}(\mathbf{A}^H \mathbf{A}) = \text{span}(\mathbf{A}^H)$  and  $\ker(\mathbf{A}^H \mathbf{A}) = \ker(\mathbf{A})$ .

**Proof.** The three matrices  $\mathbf{A}$ ,  $\mathbf{A}^H \mathbf{A}$ , and  $\mathbf{A} \mathbf{A}^H$  have the same number of nonzero singular values and we obtain 1. Moreover, 2. and 3. follow from Corollary 11.15 and (11.12), respectively, applied to  $\mathbf{A}^H \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^H$ .  $\square$

**Exercise 11.19** Let  $\mathbf{A} \in \mathbb{C}^{m,n}$  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

$$\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^H & \mathbf{0} \end{bmatrix}$$

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

$$\mathbf{p}_i = \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix}, \quad \mathbf{q}_i = \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix}, \quad \mathbf{r}_j = \begin{bmatrix} \mathbf{u}_j \\ \mathbf{0} \end{bmatrix}, \quad \text{for } i = 1, \dots, n \text{ and } j = n+1, \dots, m.$$

## 11.3 The Pseudo-Inverse and Orthogonal Projections

### 11.3.1 The Pseudo-Inverse

Suppose  $\mathbf{A} = \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^H$  is a singular value factorization of  $\mathbf{A} \in \mathbb{C}^{m,n}$ . The matrix  $\mathbf{A}^\dagger \in \mathbb{C}^{n,m}$  given by

$$\mathbf{A}^\dagger := \mathbf{V}_1 \boldsymbol{\Sigma}_1^{-1} \mathbf{U}_1^H \tag{11.16}$$

is called the **pseudo-inverse** of  $\mathbf{A}$ . It is independent of the particular factorization used to define it. We show this in Exercises 11.20, 11.21. In terms of the singular value decomposition we have

$$\mathbf{A}^\dagger = \mathbf{V} \boldsymbol{\Sigma}^\dagger \mathbf{U}^H, \quad \text{where } \boldsymbol{\Sigma}^\dagger := \begin{bmatrix} \boldsymbol{\Sigma}_1^{-1} & \mathbf{0}_{r, m-r} \\ \mathbf{0}_{n-r, r} & \mathbf{0}_{n-r, m-r} \end{bmatrix}.$$

If  $\mathbf{A}$  is square and nonsingular then  $\mathbf{A}^\dagger \mathbf{A} = \mathbf{A} \mathbf{A}^\dagger = \mathbf{I}$  and  $\mathbf{A}^\dagger$  is the usual inverse of  $\mathbf{A}$ . Any matrix has a pseudoinverse, and so  $\mathbf{A}^\dagger$  is a generalization of the usual inverse.

**Exercise 11.20** Show that  $\mathbf{B} := \mathbf{A}^\dagger$  satisfies (1)  $\mathbf{A} \mathbf{B} \mathbf{A} = \mathbf{A}$ , (2)  $\mathbf{B} \mathbf{A} \mathbf{B} = \mathbf{B}$ , (3)  $(\mathbf{B} \mathbf{A})^H = \mathbf{B} \mathbf{A}$ , and (4)  $(\mathbf{A} \mathbf{B})^H = \mathbf{A} \mathbf{B}$ .

Conversely, it can be shown that if  $\mathbf{B} \in \mathbb{C}^{n,m}$  satisfies the four equations in Exercise 11.20 then  $\mathbf{B} = \mathbf{A}^\dagger$ . Thus  $\mathbf{A}^\dagger$  is uniquely defined by these axioms and is independent of the particular singular value factorization used to define it.

**Exercise 11.21** Given  $\mathbf{A} \in \mathbb{C}^{m,n}$ , and suppose  $\mathbf{B}, \mathbf{C} \in \mathbb{C}^{n,m}$  satisfy

$$\begin{aligned} \mathbf{A} \mathbf{B} \mathbf{A} &= \mathbf{A} & (1) & \quad \mathbf{A} \mathbf{C} \mathbf{A} = \mathbf{A}, \\ \mathbf{B} \mathbf{A} \mathbf{B} &= \mathbf{B} & (2) & \quad \mathbf{C} \mathbf{A} \mathbf{C} = \mathbf{C}, \\ (\mathbf{A} \mathbf{B})^H &= \mathbf{A} \mathbf{B} & (3) & \quad (\mathbf{A} \mathbf{C})^H = \mathbf{A} \mathbf{C}, \\ (\mathbf{B} \mathbf{A})^H &= \mathbf{B} \mathbf{A} & (4) & \quad (\mathbf{C} \mathbf{A})^H = \mathbf{C} \mathbf{A}. \end{aligned}$$

Verify the following proof that  $\mathbf{B} = \mathbf{C}$ .

$$\begin{aligned} \mathbf{B} &= (\mathbf{B} \mathbf{A}) \mathbf{B} = (\mathbf{A}^H)^H \mathbf{B}^H \mathbf{B} = (\mathbf{A}^H \mathbf{C}^H) \mathbf{A}^H \mathbf{B}^H \mathbf{B} = \mathbf{C} \mathbf{A} (\mathbf{A}^H \mathbf{B}^H) \mathbf{B} \\ &= \mathbf{C} \mathbf{A} (\mathbf{B} \mathbf{A} \mathbf{B}) = (\mathbf{C}) \mathbf{A} \mathbf{B} = \mathbf{C} (\mathbf{A} \mathbf{C}) \mathbf{A} \mathbf{B} = \mathbf{C} \mathbf{C}^H \mathbf{A}^H (\mathbf{A} \mathbf{B}) \\ &= \mathbf{C} \mathbf{C}^H (\mathbf{A}^H \mathbf{B}^H \mathbf{A}^H) = \mathbf{C} (\mathbf{C}^H \mathbf{A}^H) = \mathbf{C} \mathbf{A} \mathbf{C} = \mathbf{C}. \end{aligned}$$

**Exercise 11.22** Show that the matrices  $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$  and  $\mathbf{B} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$  in Example 17.12 satisfy the axioms in Exercise 11.20. Thus we can conclude that  $\mathbf{B} = \mathbf{A}^\dagger$  without computing the singular value decomposition of  $\mathbf{A}$ .

**Exercise 11.23** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  has linearly independent columns. Show that  $\mathbf{A}^H \mathbf{A}$  is nonsingular and  $\mathbf{A}^\dagger = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H$ . If  $\mathbf{A}$  has linearly independent rows, then show that  $\mathbf{A} \mathbf{A}^H$  is nonsingular and  $\mathbf{A}^\dagger = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1}$ .

**Exercise 11.24** Show that  $\mathbf{u}^\dagger = (\mathbf{u}^H \mathbf{u})^{-1} \mathbf{u}^H$  if  $\mathbf{u} \in \mathbb{C}^{n,1}$  is nonzero.

**Exercise 11.25** If  $\mathbf{A} = \mathbf{u} \mathbf{v}^H$  where  $\mathbf{u} \in \mathbb{C}^m$ ,  $\mathbf{v} \in \mathbb{C}^n$  are nonzero, show that

$$\mathbf{A}^\dagger = \frac{1}{\alpha} \mathbf{A}^H, \quad \alpha = \|\mathbf{u}\|_2^2 \|\mathbf{v}\|_2^2.$$

**Exercise 11.26** 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 11.27** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$ . Show that

- a)  $(\mathbf{A}^H)^\dagger = (\mathbf{A}^\dagger)^H$ .
- b)  $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$ .
- c)  $(\alpha \mathbf{A})^\dagger = \frac{1}{\alpha} \mathbf{A}^\dagger$ ,  $\alpha \neq 0$ .

**Exercise 11.28** Suppose  $k, m, n \in \mathbb{N}$ ,  $\mathbf{A} \in \mathbb{C}^{m,n}$ ,  $\mathbf{B} \in \mathbb{C}^{n,k}$ . Suppose  $\mathbf{A}$  has linearly independent columns and  $\mathbf{B}$  has linearly independent rows.

- a) 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  $\mathbf{A}^\dagger \mathbf{A} = \mathbf{B} \mathbf{B}^\dagger = \mathbf{I}$  that  $\mathbf{F}$  is the pseudo-inverse of  $\mathbf{E}$ .
- b) 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 11.29** Show that  $\mathbf{A}^H = \mathbf{A}^\dagger$  if and only if all singular values of  $\mathbf{A}$  are either zero or one.

### 11.3.2 Orthogonal Projections

The singular value decomposition and the pseudo-inverse can be used to compute orthogonal projections into the subspaces  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A}^H)$ .

We start by recalling some facts about sums, direct sums and orthogonal sums of subspaces (Cf. Chapter 2). Suppose  $\mathcal{S}$  and  $\mathcal{T}$  are subspaces of a vector space  $(\mathcal{V}, \mathbb{F})$ , where typically  $\mathbb{F} = \mathbb{C}$  or  $\mathbb{F} = \mathbb{R}$ . We define

- **Sum:**  $\mathcal{X} := \mathcal{S} + \mathcal{T} := \{\mathbf{s} + \mathbf{t} : \mathbf{s} \in \mathcal{S} \text{ and } \mathbf{t} \in \mathcal{T}\}$ .
- **Direct Sum:** If  $\mathcal{S} \cap \mathcal{T} = \{\mathbf{0}\}$ , then  $\mathcal{S} \oplus \mathcal{T} := \mathcal{S} + \mathcal{T}$ .
- **Orthogonal Sum:** Suppose  $(\mathcal{V}, \mathbb{F}, \langle \cdot, \cdot \rangle)$  is an inner product space. Then  $\mathcal{S} \oplus \mathcal{T}$  is an orthogonal sum if  $\langle \mathbf{s}, \mathbf{t} \rangle = 0$  for all  $\mathbf{s} \in \mathcal{S}$  and all  $\mathbf{t} \in \mathcal{T}$ .
- **Orthogonal Complement:**  $\mathcal{T} = \mathcal{S}^\perp := \{\mathbf{x} \in \mathcal{X} : \langle \mathbf{s}, \mathbf{x} \rangle = 0 \text{ for all } \mathbf{s} \in \mathcal{S}\}$ .
- If  $\mathcal{S} \oplus \mathcal{T}$  is an orthogonal sum and  $\mathbf{v} = \mathbf{s} + \mathbf{t} \in \mathcal{S} \oplus \mathcal{T}$  with  $\mathbf{s} \in \mathcal{S}$  and  $\mathbf{t} \in \mathcal{T}$  then  $\mathbf{s}$  is called the **orthogonal projection** of  $\mathbf{v}$  into  $\mathcal{S}$ .

We recall that

- $\mathcal{S} + \mathcal{T} = \mathcal{T} + \mathcal{S}$  and  $\mathcal{S} + \mathcal{T}$  is a subspace of  $\mathcal{V}$ .
- $\dim(\mathcal{S} + \mathcal{T}) = \dim \mathcal{S} + \dim \mathcal{T} - \dim(\mathcal{S} \cap \mathcal{T})$ .
- $\dim(\mathcal{S} \oplus \mathcal{T}) = \dim \mathcal{S} + \dim \mathcal{T}$ .
- Every  $\mathbf{v} \in \mathcal{S} \oplus \mathcal{T}$  can be decomposed uniquely as  $\mathbf{v} = \mathbf{s} + \mathbf{t}$ , where  $\mathbf{s} \in \mathcal{S}$  and  $\mathbf{t} \in \mathcal{T}$ .

Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  has a singular value factorization  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  and let  $\mathbf{b} \in \mathbb{C}^m$ . Then

$$\mathbf{b} = \mathbf{U}\mathbf{U}^H\mathbf{b} = [\mathbf{U}_1 \ \mathbf{U}_2] \begin{bmatrix} \mathbf{U}_1^H \\ \mathbf{U}_2^H \end{bmatrix} \mathbf{b} = \mathbf{U}_1\mathbf{U}_1^H\mathbf{b} + \mathbf{U}_2\mathbf{U}_2^H\mathbf{b} =: \mathbf{b}_1 + \mathbf{b}_2,$$

where  $\mathbf{b}_j = \mathbf{U}_j\mathbf{c}_j$  and  $\mathbf{c}_j := \mathbf{U}_j^H\mathbf{b}$  for  $j = 1, 2$ . Since  $\mathbf{U}_1(\mathbf{U}_2)$  is an orthonormal basis for  $\text{span}(\mathbf{A})$  ( $\ker(\mathbf{A}^H)$ ), we have  $\mathbf{b}_1(\mathbf{b}_2) \in \text{span}(\mathbf{A})$  ( $\ker(\mathbf{A}^H)$ ). Since

$$\text{span}(\mathbf{A}) \cap \ker(\mathbf{A}^H) = \text{span}(\mathbf{U}_1) \cap \ker(\mathbf{U}_2) = \{\mathbf{0}\}$$

the decomposition is a direct sum decomposition and since  $\mathbf{b}_1^H\mathbf{b}_2 = \mathbf{c}_1^H\mathbf{U}_1^H\mathbf{U}_2\mathbf{c}_2 = 0$ , it is in fact an orthogonal decomposition of  $\mathbb{C}^m$ . This shows in particular that  $\mathbb{C}^m = \text{span}(\mathbf{A}) \oplus \ker(\mathbf{A}^H)$ .

**Theorem 11.30** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{b} \in \mathbb{C}^m$ . Then

$$\mathbf{b}_1 := \mathbf{A}\mathbf{A}^\dagger\mathbf{b} \tag{11.17}$$

is the orthogonal projection of  $\mathbf{b}$  into  $\text{span}(\mathbf{A})$ , and

$$\mathbf{b}_2 := (\mathbf{I} - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b} \tag{11.18}$$

is the orthogonal projection of  $\mathbf{b}$  into the orthogonal complement  $\ker(\mathbf{A}^H)$  of  $\text{span}(\mathbf{A})$ .

**Proof.** We have already shown that  $\mathbf{b}_1 = \mathbf{U}_1\mathbf{U}_1^H\mathbf{b}$  and  $\mathbf{b}_2 = \mathbf{U}_2\mathbf{U}_2^H\mathbf{b}$ . It is verified directly that  $\mathbf{A}\mathbf{A}^\dagger = \mathbf{U}_1\mathbf{U}_1^H$ , and then  $\mathbf{b}_2 = \mathbf{b} - \mathbf{b}_1 = (\mathbf{I} - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b}$ .  $\square$

**Example 11.31** 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  $\mathbf{b}_2 = (\mathbf{I}_3 - \mathbf{A}\mathbf{A}^\dagger)\mathbf{b} = \mathbf{U}_2\mathbf{U}_2^T\mathbf{b} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ b_3 \end{bmatrix}$ .

**Exercise 11.32** Show that if  $\mathbf{A}$  has rank  $n$  then  $\mathbf{A}(\mathbf{A}^H\mathbf{A})^{-1}\mathbf{A}^H\mathbf{b}$  is the projection of  $\mathbf{b}$  into  $\text{span}(\mathbf{A})$ . (Cf. Exercise 11.23.)

**Exercise 11.33** Consider the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{C}^{n,n}$  has rank  $r > 0$  and  $\mathbf{b} \in \mathbb{C}^n$ . Let

$$\mathbf{U}^H\mathbf{A}\mathbf{V} = \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

represent the singular value decomposition of  $\mathbf{A}$ .

- a) Let  $\mathbf{c} = [c_1, \dots, c_n]^T = \mathbf{U}^H \mathbf{b}$  and  $\mathbf{y} = [y_1, \dots, y_n]^T = \mathbf{V}^H \mathbf{x}$ . Show that  $\mathbf{Ax} = \mathbf{b}$  if and only if

$$\begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{y} = \mathbf{c}.$$

- b) Show that  $\mathbf{Ax} = \mathbf{b}$  has a solution  $\mathbf{x}$  if and only if  $c_{r+1} = \dots = c_n = 0$ .
- c) Deduce that a linear system  $\mathbf{Ax} = \mathbf{b}$  has either no solution, one solution or infinitely many solutions.

**Exercise 11.34** For any  $\mathbf{A} \in \mathbb{C}^{m,n}$ ,  $\mathbf{b} \in \mathbb{C}^n$  show that one and only one of the following systems has a solution

$$(1) \quad \mathbf{Ax} = \mathbf{b}, \quad (2) \quad \mathbf{A}^H \mathbf{y} = \mathbf{0}, \mathbf{y}^H \mathbf{b} \neq 0.$$

In other words either  $\mathbf{b} \in \text{span}(\mathbf{A})$ , or we can find  $\mathbf{y} \in \ker(\mathbf{A}^H)$  such that  $\mathbf{y}^H \mathbf{b} \neq 0$ . This is called **Fredholms alternative**.

## 11.4 The Minmax Theorem for Singular Values and the Hoffman-Wielandt Theorem

We have a minmax and maxmin characterization for singular values.

**Theorem 11.35 (The Courant-Fischer Theorem for Singular Values)** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  has singular values  $\sigma_1, \sigma_2, \dots, \sigma_n$  ordered so that  $\sigma_1 \geq \dots \geq \sigma_n$ . Then for  $k = 1, \dots, n$

$$\sigma_k = \min_{\dim(S)=n-k+1} \max_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} = \max_{\dim(S)=k} \min_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}. \quad (11.19)$$

**Proof.** Since

$$\frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} = \frac{\langle \mathbf{Ax}, \mathbf{Ax} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \frac{\langle \mathbf{x}, \mathbf{A}^H \mathbf{Ax} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}$$

denotes the Rayleigh quotient  $R(\mathbf{A}^H \mathbf{A})$  of  $\mathbf{A}^H \mathbf{A}$ , and since the singular values of  $\mathbf{A}$  are the nonnegative square roots of the eigenvalues of  $\mathbf{A}^H \mathbf{A}$ , the results follow from the Courant-Fischer Theorem for eigenvalues, see Theorem 10.14.  $\square$

By taking  $k = 1$  and  $k = n$  in (11.19) we obtain for any  $\mathbf{A} \in \mathbb{C}^{m,n}$

$$\sigma_1 = \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}, \quad \sigma_n = \min_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}. \quad (11.20)$$

This follows since the only subspace of  $\mathbb{C}^n$  of dimension  $n$  is  $\mathbb{C}^n$  itself.

The Hoffman-Wielandt Theorem for eigenvalues of Hermitian matrices, Theorem 10.18 can be written

$$\sum_{j=1}^n |\mu_j - \lambda_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2 := \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2, \quad (11.21)$$

where  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$  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.

**Theorem 11.36 (Hoffman-Wielandt Theorem for singular values)** *For any  $m, n \in \mathbb{N}$  and  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m,n}$  we have*

$$\sum_{j=1}^n |\beta_j - \alpha_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2. \quad (11.22)$$

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.

**Proof.** Suppose first  $m \geq n$ . The block matrix  $\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^H & \mathbf{0} \end{bmatrix} \in \mathbb{C}^{m+n, m+n}$  has  $2n$  eigenvalues  $\alpha_1, -\alpha_1, \dots, \alpha_n, -\alpha_n$  and  $m - n$  additional zero eigenvalues (cf. Exercise 11.19). Moreover,  $\mathbf{C}$  is Hermitian. By the Hoffman-Wielandt Theorem for eigenvalues we obtain

$$\begin{aligned} 2 \sum_{j=1}^n |\beta_j - \alpha_j|^2 &\leq \left\| \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^H & \mathbf{0} \end{bmatrix} \right\|_F^2 \\ &= \|\mathbf{B} - \mathbf{A}\|_F^2 + \|\mathbf{B}^H - \mathbf{A}^H\|_F^2 = 2\|\mathbf{B} - \mathbf{A}\|_F^2, \end{aligned}$$

and (11.22) follows for  $m \geq n$ . Since  $\mathbf{A}$  and  $\mathbf{A}^H$  have the same nonzero singular values, the result also holds for  $n > m$ .  $\square$





## Chapter 12

# Matrix Norms

To measure the size of a matrix we can use a matrix norm. In this chapter we give a systematic study of matrix norms.

### 12.1 Matrix Norms

For simplicity we consider only matrix norms on the vector space  $(\mathbb{C}^{m,n}, \mathbb{C})$ . All results also holds for  $(\mathbb{R}^{m,n}, \mathbb{R})$ .

**Definition 12.1 (Matrix Norms)** *Suppose  $m, n$  are positive integers. A function  $\|\cdot\|: \mathbb{C}^{m,n} \rightarrow \mathbb{R}$  is called a **matrix norm** on  $\mathbb{C}^{m,n}$  if for all  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m,n}$  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  $(\mathbb{C}^{m,n}, \mathbb{C})$  of  $m \times n$  matrices. Adapting Theorem 2.34 to this special situation gives

**Theorem 12.2** *All matrix norms are equivalent. Thus, if  $\|\cdot\|$  and  $\|\cdot\|'$  are two matrix norms on  $\mathbb{C}^{m,n}$  then there are positive constants  $\mu$  and  $M$  such that*

$$\mu \|\mathbf{A}\| \leq \|\mathbf{A}\|' \leq M \|\mathbf{A}\|$$

*holds for all  $\mathbf{A} \in \mathbb{C}^{m,n}$ . Moreover, a matrix norm is a continuous function.*

#### 12.1.1 The Frobenius Norm

From any vector norm  $\|\cdot\|_V$  on  $\mathbb{C}^{mn}$  we can define a matrix norm on  $\mathbb{C}^{m,n}$  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

$$\|\mathbf{A}\|_S := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|, \quad \|\mathbf{A}\|_F := \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}, \quad \|\mathbf{A}\|_M := \max_{i,j} |a_{ij}|. \quad (12.1)$$

Of these norms the Frobenius norm is the most useful. It satisfies the following properties.

**Lemma 12.3** *For any matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  we have*

1.  $\|\mathbf{A}^H\|_F = \|\mathbf{A}\|_F$ ,
2.  $\|\mathbf{A}\|_F^2 = \sum_{j=1}^n \|\mathbf{a}_{\cdot j}\|_2^2 = \sum_{i=1}^m \|\mathbf{a}_i\|_2^2$ ,
3.  $\|\mathbf{U}\mathbf{A}\|_F = \|\mathbf{A}\mathbf{V}\|_F = \|\mathbf{A}\|_F$  for any unitary matrices  $\mathbf{U} \in \mathbb{C}^{m,m}$  and  $\mathbf{V} \in \mathbb{C}^{n,n}$ ,
4.  $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$  for any  $\mathbf{B} \in \mathbb{C}^{n,k}$ ,
5.  $\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$ , for all  $\mathbf{x} \in \mathbb{C}^n$ .

**Proof.**

1.  $\|\mathbf{A}^H\|_F^2 = \sum_{j=1}^n \sum_{i=1}^m |\bar{a}_{ij}|^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 = \|\mathbf{A}\|_F^2$ .
2. Obvious.
3. Recall that  $\|\mathbf{U}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$  for all  $\mathbf{x} \in \mathbb{C}^n$  if  $\mathbf{U}^H \mathbf{U} = \mathbf{I}$ . Applying this to each column  $\mathbf{a}_{\cdot j}$  of  $\mathbf{A}$  we find  $\|\mathbf{U}\mathbf{A}\|_F^2 \stackrel{2.}{=} \sum_{j=1}^n \|\mathbf{U}\mathbf{a}_{\cdot j}\|_2^2 = \sum_{j=1}^n \|\mathbf{a}_{\cdot j}\|_2^2 \stackrel{2.}{=} \|\mathbf{A}\|_F^2$ . Similarly, since  $\mathbf{V}\mathbf{V}^H = \mathbf{I}$  we find  $\|\mathbf{A}\mathbf{V}\|_F \stackrel{1.}{=} \|\mathbf{V}^H \mathbf{A}^H\|_F = \|\mathbf{A}^H\|_F \stackrel{1.}{=} \|\mathbf{A}\|_F$ .
4. Using Cauchy-Schwarz' inequality and 2. we obtain

$$\|\mathbf{A}\mathbf{B}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^k (\mathbf{a}_i^T \mathbf{b}_{\cdot j})^2 \leq \sum_{i=1}^n \sum_{j=1}^k \|\mathbf{a}_i\|_2^2 \|\mathbf{b}_{\cdot j}\|_2^2 = \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.$$

5. Since  $\|\mathbf{v}\|_F = \|\mathbf{v}\|_2$  for a vector this follows by taking  $k = 1$  and  $\mathbf{B} = \mathbf{x}$  in 4.

□

There is a relation between the Frobenius norm and the singular values.

**Theorem 12.4** *We have  $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_n^2}$ , where  $\sigma_1, \dots, \sigma_n$  are the singular values of  $\mathbf{A}$ .*

**Proof.** Using Lemma 12.3 we find  $\|\mathbf{A}\|_F \stackrel{3.}{=} \|\mathbf{U}^T \mathbf{A} \mathbf{V}\|_F = \|\mathbf{\Sigma}\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_n^2}$ .  
□

### 12.1.2 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  $\mathbf{A}\mathbf{B}$  is defined in a fixed space  $\mathbb{C}^{n,n}$ , while in the rectangular case matrix multiplication combines matrices in different spaces. The following definition captures this distinction.

**Definition 12.5 (Consistent Matrix Norms)** *A matrix norm is called **consistent** on  $\mathbb{C}^{n,n}$  if*

$$4. \quad \|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (\text{submultiplicativity})$$

*holds for all  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$ . A matrix norm is **consistent** if it is defined on  $\mathbb{C}^{m,n}$  for all  $m, n \in \mathbb{N}$ , and 4. holds for all matrices  $\mathbf{A}, \mathbf{B}$  for which the product  $\mathbf{A}\mathbf{B}$  is defined.*

Clearly the three norms in (12.1) are defined for all  $m, n \in \mathbb{N}$ . From Lemma 12.3 it follows that the Frobenius norm is consistent.

**Exercise 12.6** *Show that the sum norm is consistent.*

**Exercise 12.7** *Show that the max norm is not consistent by considering  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ .*

**Exercise 12.8**

(a) *Show that the norm*

$$\|\mathbf{A}\| := \sqrt{mn} \|\mathbf{A}\|_M, \quad \mathbf{A} \in \mathbb{C}^{m,n}$$

*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  $\mathbb{C}^{n,n}$  we have the inequality

$$\|\mathbf{A}^k\| \leq \|\mathbf{A}\|^k \text{ for } k \in \mathbb{N}. \quad (12.2)$$

When working with norms we often have to bound the vector norm of a matrix times a vector by the norm of the matrix times the norm of the vector. We have the following definition.

**Definition 12.9 (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  $\mathbb{C}^{m,n}$ . We say that the matrix norm  $\|\cdot\|$  is **subordinate** to the vector norms  $\|\cdot\|_\alpha$  and  $\|\cdot\|_\beta$  if  $\|\mathbf{A}\mathbf{x}\|_\alpha \leq \|\mathbf{A}\| \|\mathbf{x}\|_\beta$  for all  $\mathbf{A} \in \mathbb{C}^{m,n}$  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 12.3 we have  $\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$ , for all  $\mathbf{x} \in \mathbb{C}^n$ . Thus the Frobenius norm is subordinate to the Euclidian vector norm.

**Exercise 12.10** Show that the sum norm is subordinate to the  $l_1$ -norm.

**Exercise 12.11** (a) Show that the max norm is subordinate to the  $\infty$  and 1 norm, i. e.,  $\|\mathbf{Ax}\|_\infty \leq \|\mathbf{A}\|_M \|\mathbf{x}\|_1$  holds for all  $\mathbf{A} \in \mathbb{C}^{m,n}$  and all  $\mathbf{x} \in \mathbb{C}^n$ .

(b) Show that  $\|\mathbf{Ae}_l\|_\infty = \|\mathbf{A}\|_M \|\mathbf{e}_l\|_1$ , where  $\|\mathbf{A}\|_M = |a_{kl}|$ .

(c) Show that  $\|\mathbf{A}\|_M = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_\infty}{\|\mathbf{x}\|_1}$ .

### 12.1.3 Operator Norms

Corresponding to vector norms on  $\mathbb{C}^n$  and  $\mathbb{C}^m$  there is an induced matrix norm on  $\mathbb{C}^{m,n}$  which we call the **operator norm**.

**Definition 12.12 (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  $\mathbf{A} \in \mathbb{C}^{m,n}$  we define

$$\|\mathbf{A}\| := \|\mathbf{A}\|_{\alpha,\beta} := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta}. \quad (12.3)$$

We call this the  $(\alpha, \beta)$  **operator norm**, the  $(\alpha, \beta)$ -norm, or simply the  $\alpha$ -norm if  $\alpha = \beta$ .

Before we show that the  $(\alpha, \beta)$ -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

$$\|\mathbf{A}\|_{\alpha,\beta} = \max_{\mathbf{x} \notin \ker(\mathbf{A})} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta} = \max_{\|\mathbf{x}\|_\beta=1} \|\mathbf{Ax}\|_\alpha. \quad (12.4)$$

That we only need to consider  $\mathbf{x}$ 's outside the null space  $\ker(\mathbf{A})$  of  $\mathbf{A}$  is obvious. We can take the max over the  $\beta$ -norm unit sphere in  $\mathbb{C}^n$  since

$$\max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_\alpha}{\|\mathbf{x}\|_\beta} = \max_{\mathbf{x} \neq \mathbf{0}} \left\| \mathbf{A} \left( \frac{\mathbf{x}}{\|\mathbf{x}\|_\beta} \right) \right\|_\alpha = \max_{\|\mathbf{x}\|_\beta=1} \|\mathbf{Ax}\|_\alpha.$$

2. The operator norm  $\|\mathbf{A}\|$  is subordinate to the vector norms  $\|\cdot\|_\alpha$  and  $\|\cdot\|_\beta$ . Thus

$$\|\mathbf{Ax}\|_\alpha \leq \|\mathbf{A}\| \|\mathbf{x}\|_\beta \text{ for all } \mathbf{A} \in \mathbb{C}^{m,n} \text{ and } \mathbf{x} \in \mathbb{C}^n. \quad (12.5)$$

3. We can use max instead of sup in (12.3). This follows by the following compactness argument. Since all vector norms on  $\mathbb{C}^n$  are equivalent the unit sphere  $\mathcal{S}_\beta = \{\mathbf{x} \in \mathbb{C}^n : \|\mathbf{x}\|_\beta = 1\}$  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  $f : \mathcal{S}_\beta \rightarrow \mathbb{R}$  given by  $f(\mathbf{x}) = \|\mathbf{Ax}\|_\alpha$  is continuous. But then  $f$  attains its max and min and we have

$$\|\mathbf{A}\|_{\alpha,\beta} = \|\mathbf{Ax}^*\|_\alpha \text{ for some } \mathbf{x}^* \in \mathbb{C}^n \text{ with } \|\mathbf{x}^*\|_\beta = 1. \quad (12.6)$$

**Lemma 12.13** *The operator norm given by (12.3) is a matrix norm on  $\mathbb{C}^{m,n}$ . 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 (12.4). In 2. and 3. below we take the max over the unit sphere  $S_\beta$ .

1. Nonnegativity is obvious. If  $\|\mathbf{A}\| = 0$  then  $\|\mathbf{A}\mathbf{y}\|_\beta = 0$  for each  $\mathbf{y} \in \mathbb{C}^n$ . In particular, each column  $\mathbf{A}\mathbf{e}_j$  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.  $\|\mathbf{AB}\| = \max_{\mathbf{B}\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{AB}\mathbf{x}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \max_{\mathbf{B}\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{AB}\mathbf{x}\|_\alpha}{\|\mathbf{B}\mathbf{x}\|_\alpha} \frac{\|\mathbf{B}\mathbf{x}\|_\alpha}{\|\mathbf{x}\|_\alpha} \leq \max_{\mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{y}\|_\alpha}{\|\mathbf{y}\|_\alpha} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{B}\mathbf{x}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \|\mathbf{A}\| \|\mathbf{B}\|$ .

□

For any  $\alpha$ -norm of the  $n \times n$  identity matrix we find

$$\|\mathbf{I}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{I}\mathbf{x}\|_\alpha}{\|\mathbf{x}\|_\alpha} = \max_{\mathbf{x} \neq \mathbf{0}} 1 = 1.$$

For the Frobenius norm we find  $\|\mathbf{I}\|_F = \sqrt{n}$ , and this shows that the Frobenius norm is not an operator norm for  $n > 1$ .

### 12.1.4 The $p$ -Norms

Recall that the  $p$  or  $\ell_p$  vector norms (2.10) 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  $1 \leq p \leq \infty$

$$\|\mathbf{A}\|_p := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p} = \max_{\|\mathbf{y}\|_p=1} \|\mathbf{A}\mathbf{y}\|_p. \quad (12.7)$$

In the most important cases we have explicit expressions for these norms.

**Theorem 12.14** *For  $\mathbf{A} \in \mathbb{C}^{m,n}$  we have*

$$\begin{aligned} \|\mathbf{A}\|_1 &:= \max_{1 \leq j \leq n} \sum_{k=1}^m |a_{k,j}|, & (\text{max column sum}) \\ \|\mathbf{A}\|_2 &:= \sigma_1, & (\text{largest singular value of } \mathbf{A}) \\ \|\mathbf{A}\|_\infty &:= \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{k,j}|, & (\text{max row sum}). \end{aligned} \quad (12.8)$$

The expression  $\|\mathbf{A}\|_2$  is called the **two-norm** or the **spectral norm** of  $\mathbf{A}$ .

**Proof.** The result for  $p = 2$  follows from the minmax theorem for singular values. Indeed, by (11.20) 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  $\|\mathbf{Ax}\|_p \leq K_p$  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  $\|\mathbf{Ay}^*\|_p = K_p$ .

It then follows from (12.7) that  $\|\mathbf{A}\|_p = \|\mathbf{Ay}^*\|_p = K_p$ .

**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)

$$\|\mathbf{Ax}\|_1 = \sum_{k=1}^m \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \sum_{k=1}^m \sum_{j=1}^n |a_{kj}| |x_j| = \sum_{j=1}^n \left( \sum_{k=1}^m |a_{kj}| \right) |x_j| \leq K_1.$$

(b)  $\|\mathbf{Ay}^*\|_1 = K_1$ .

**$\infty$ -norm:** Define  $K_\infty$ ,  $r$  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)  $\|\mathbf{Ax}\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{kj}| |x_j| \leq K_\infty$ .

(b)  $\|\mathbf{Ay}^*\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} e^{-i\theta_j} \right| = K_\infty$ .

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 12.15** In Example 11.6 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 12.16** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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

$$\|\mathbf{A}\|_2 = \sigma_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\sigma_n}, \quad (12.9)$$

$$\|\mathbf{A}\|_2 = \lambda_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\lambda_n}, \quad \text{if } \mathbf{A} \text{ is symmetric positive definite,} \quad (12.10)$$

$$\|\mathbf{A}\|_2 = |\lambda_1| \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{|\lambda_n|}, \quad \text{if } \mathbf{A} \text{ is normal.} \quad (12.11)$$

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  $\mathbf{A}^{-1}$ , (12.9) follows. As shown in Section 11.1.3 the singular values of a symmetric positive definite matrix (normal matrix) are equal to the eigenvalues (absolute value of the eigenvalues). This implies (12.10) and (12.11).  $\square$

**Exercise 12.17** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is nonsingular. Use (12.9) and (11.20) to show that

$$\|\mathbf{A}^{-1}\|_2 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{x}\|_2}{\|\mathbf{A}\mathbf{x}\|_2}.$$

**Exercise 12.18** 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 12.19** For any  $\mathbf{A} \in \mathbb{C}^{m,n}$  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  $\mathbf{A}^H \mathbf{A}$  corresponding to the largest singular value of  $\mathbf{A}$ . Then

$$\|\mathbf{A}\|_2^2 \|\mathbf{v}_1\|_1 = \sigma_1^2 \|\mathbf{v}_1\|_1 = \|\sigma_1^2 \mathbf{v}_1\|_1 = \|\mathbf{A}^H \mathbf{A} \mathbf{v}_1\|_1 \leq \|\mathbf{A}^H\|_1 \|\mathbf{A}\|_1 \|\mathbf{v}_1\|_1.$$

Observing that  $\|\mathbf{A}^H\|_1 = \|\mathbf{A}\|_\infty$  by Theorem 12.14 and canceling  $\|\mathbf{v}_1\|_1$  proves the result.  $\square$

### 12.1.5 Unitary Invariant Matrix Norms

**Definition 12.20** A matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{m,n}$  is called **unitary invariant** if  $\|\mathbf{U}\mathbf{A}\mathbf{V}\| = \|\mathbf{A}\|$  for any  $\mathbf{A} \in \mathbb{C}^{m,n}$  and any unitary matrices  $\mathbf{U} \in \mathbb{C}^{m,m}$  and  $\mathbf{V} \in \mathbb{C}^{n,n}$ .

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 12.3 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 [9] p. 308.

**Theorem 12.21** The Frobenius norm and the spectral norm are unitary invariant. Moreover  $\|\mathbf{A}^H\|_F = \|\mathbf{A}\|_F$  and  $\|\mathbf{A}^H\|_2 = \|\mathbf{A}\|_2$ .

**Proof.** The results for the Frobenius norm follow from Lemma 12.3. Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  and let  $\mathbf{U} \in \mathbb{C}^{m,m}$  and  $\mathbf{V} \in \mathbb{C}^{n,n}$  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  $\mathbf{A}^H$  have the same nonzero singular values, and it follows from Theorem 12.14 that  $\|\mathbf{A}^H\|_2 = \|\mathbf{A}\|_2$ . Moreover  $\mathbf{V}^H$  is unitary. Using these facts we find

$$\|\mathbf{A}\mathbf{V}\|_2 = \|(\mathbf{A}\mathbf{V})^H\|_2 = \|\mathbf{V}^H\mathbf{A}^H\|_2 = \|\mathbf{A}^H\|_2 = \|\mathbf{A}\|_2.$$

□

**Exercise 12.22** Show that  $\|\mathbf{V}\mathbf{A}\|_2 = \|\mathbf{A}\|_2$  holds even for a rectangular  $\mathbf{V}$  as long as  $\mathbf{V}^H\mathbf{V} = \mathbf{I}$ .

**Exercise 12.23** Find  $\mathbf{A} \in \mathbb{R}^{2,2}$  and  $\mathbf{U} \in \mathbb{R}^{2,1}$  with  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$  such that  $\|\mathbf{A}\mathbf{U}\|_2 < \|\mathbf{A}\|_2$ . Thus, in general,  $\|\mathbf{A}\mathbf{U}\|_2 = \|\mathbf{A}\|_2$  does not hold for a rectangular  $\mathbf{U}$  even if  $\mathbf{U}^H\mathbf{U} = \mathbf{I}$ .

**Exercise 12.24** 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 12.25** 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 12.26** If  $\mathbf{A} \in \mathbb{C}^{m,n}$  has elements  $a_{ij}$ , let  $|\mathbf{A}| \in \mathbb{C}^{m,n}$  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  $\mathbf{A} \in \mathbb{C}^{m,n}$   $\|\mathbf{A}\|_F = \||\mathbf{A}|\|_F$ ,  $\|\mathbf{A}\|_p = \||\mathbf{A}|\|_p$  for  $p = 1, \infty$ .

(c) Show that for any  $\mathbf{A} \in \mathbb{C}^{m,n}$   $\|\mathbf{A}\|_2 \leq \||\mathbf{A}|\|_2$ .

(d) Find a real symmetric  $2 \times 2$  matrix  $\mathbf{A}$  such that  $\|\mathbf{A}\|_2 < \||\mathbf{A}|\|_2$ .

**Exercise 12.27** Let  $m, n \in \mathbb{N}$  and  $\mathbf{A} \in \mathbb{C}^{m,n}$ . Show that

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=\|\mathbf{y}\|_2=1} |\mathbf{y}^H \mathbf{A} \mathbf{x}|.$$



### 12.1.6 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 [9]) that a vector norm on  $\mathbb{C}^n$  is an absolute norm if and only if it is a **monotone norm**, i.e.,

$$|x_i| \leq |y_i|, i = 1, \dots, n \implies \|\mathbf{x}\| \leq \|\mathbf{y}\|, \text{ for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

Absolute and monotone matrix norms are defined as for vector norms.

**Exercise 12.28** Show that the Frobenius norm and the  $1, \infty$  operator norms are absolute norms.

**Exercise 12.29** Show that the spectral norm is not an absolute norm.

The study of matrix norms will be continued in Chapter 13.

## 12.2 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  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . 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  $\mathbb{C}^{n,n}$  which in addition is subordinate to the vector norm. Thus for any  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$  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 12.3.

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 12.30** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (12.12)$$

**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 (12.12). 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 (12.12).  $\|\mathbf{e}\|/\|\mathbf{b}\|$  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  $\|\mathbf{e}\|/\|\mathbf{b}\|$ .  $K(\mathbf{A})$  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  $\|\mathbf{x}\| = \|\mathbf{Ix}\| \leq \|\mathbf{I}\| \|\mathbf{x}\|$  for any  $\mathbf{x}$ , by subordination we have  $\|\mathbf{I}\| \geq 1$  and therefore by submultiplicativity  $\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \geq \|\mathbf{AA}^{-1}\| = \|\mathbf{I}\| \geq 1$ .

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  $\ell_1$ ,  $\ell_\infty$ , or Frobenius norm when one wishes to compute it or estimate it.

Explicit expressions for the 2-norm condition number follow from Theorem 12.16.

**Theorem 12.31** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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,

$$K_2(\mathbf{A}) = \begin{cases} \lambda_1/\lambda_n, & \text{if } \mathbf{A} \text{ is symmetric positive definite,} \\ |\lambda_1|/|\lambda_n|, & \text{if } \mathbf{A} \text{ is normal.} \end{cases} \quad (12.13)$$

It follows that  $\mathbf{A}$  is ill-conditioned with respect to inversion if and only if  $\sigma_1/\sigma_n$  is large, or  $\lambda_1/\lambda_n$  is large when  $\mathbf{A}$  is symmetric positive definite.

**Exercise 12.32** The upper and lower bounds for  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$  given by (12.12) 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 (12.12) 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  $\mathbf{e} = \mathbf{A}\mathbf{y}_{\mathbf{A}}$ .

We consider next a perturbation  $\mathbf{E}$  in a nonsingular matrix  $\mathbf{A}$ . The following result shows that  $\mathbf{A} + \mathbf{E}$  is nonsingular if  $\mathbf{E}$  is sufficiently small and that small changes in  $\mathbf{A}$  give small changes in the inverse if  $\mathbf{A}$  is well conditioned.

**Theorem 12.33** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is nonsingular and let  $\|\cdot\|$  be a consistent matrix norm on  $\mathbb{C}^{n,n}$ . If  $\mathbf{E} \in \mathbb{C}^{n,n}$  is so small that  $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1$  then  $\mathbf{A} + \mathbf{E}$  is nonsingular and

$$\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - r}. \quad (12.14)$$

If  $r < 1/2$  then

$$\frac{\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\|}{\|\mathbf{A}^{-1}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (12.15)$$

**Proof.** We show in (12.24) in Section 12.4 that if  $\mathbf{B} \in \mathbb{C}^{n,n}$  and  $\|\mathbf{B}\| < 1$  then  $\mathbf{I} - \mathbf{B}$  is nonsingular and

$$\|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}. \quad (12.16)$$

Since  $r < 1$  the matrix  $\mathbf{I} - \mathbf{B} := \mathbf{I} + \mathbf{A}^{-1}\mathbf{E}$  is nonsingular. Since  $(\mathbf{I} - \mathbf{B})^{-1}\mathbf{A}^{-1}(\mathbf{A} + \mathbf{E}) = \mathbf{I}$  we see that  $\mathbf{A} + \mathbf{E}$  is nonsingular with inverse  $(\mathbf{I} - \mathbf{B})^{-1}\mathbf{A}^{-1}$ . Hence,  $\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \|(\mathbf{I} - \mathbf{B})^{-1}\|\|\mathbf{A}^{-1}\|$  and (12.14) follows from (12.16). From the identity

$$(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1} = -\mathbf{A}^{-1}\mathbf{E}(\mathbf{A} + \mathbf{E})^{-1}$$

we obtain by (12.14)

$$\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\| \leq \|\mathbf{A}^{-1}\|\|\mathbf{E}\|\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|} \frac{\|\mathbf{A}^{-1}\|}{1 - r}.$$

Dividing by  $\|\mathbf{A}^{-1}\|$  and setting  $r = 1/2$  proves (12.15).  $\square$

We can now show the following upper bounds.

**Theorem 12.34** Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n,n}$ ,  $\mathbf{b} \in \mathbb{C}^n$  with  $\mathbf{A}$  invertible and  $\mathbf{b} \neq \mathbf{0}$ . If  $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1/2$  for some operator norm then  $\mathbf{A} + \mathbf{E}$  is invertible. If  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$  then

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{y}\|} \leq \|\mathbf{A}^{-1}\mathbf{E}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}, \quad (12.17)$$

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (12.18)$$

**Proof.** That the matrix  $\mathbf{A} + \mathbf{E}$  is invertible follows from Theorem 12.33. (12.17) follows easily by taking norms in the equation  $\mathbf{x} - \mathbf{y} = \mathbf{A}^{-1}\mathbf{E}\mathbf{y}$  and dividing by  $\|\mathbf{y}\|$ . From the identity  $\mathbf{y} - \mathbf{x} = ((\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1})\mathbf{A}\mathbf{x}$  we obtain  $\|\mathbf{y} - \mathbf{x}\| \leq \|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\| \|\mathbf{A}\| \|\mathbf{x}\|$  and (12.18) follows from (12.14).  $\square$

In Theorem 12.34 we gave a bound for the relative error in  $\mathbf{x}$  as an approximation to  $\mathbf{y}$ , (12.17), and the relative error in  $\mathbf{y}$  as an approximation to  $\mathbf{x}$ , (12.18).  $\|\mathbf{E}\|/\|\mathbf{A}\|$  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 (12.17) can give a fairly realistic estimate for  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$ .

Suppose we have computed an approximate solution  $\mathbf{y}$  to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . The vector  $\mathbf{r}(\mathbf{y}) := \mathbf{A}\mathbf{y} - \mathbf{b}$  is called the **residual vector**, or just the residual. We can bound  $\mathbf{x} - \mathbf{y}$  in term of  $\mathbf{r}$ .

**Theorem 12.35** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$ ,  $\mathbf{b} \in \mathbb{C}^n$ ,  $\mathbf{A}$  is nonsingular and  $\mathbf{b} \neq \mathbf{0}$ . Let  $\mathbf{r}(\mathbf{y}) = \mathbf{A}\mathbf{y} - \mathbf{b}$  for each  $\mathbf{y} \in \mathbb{C}^n$ . If  $\mathbf{A}\mathbf{x} = \mathbf{b}$  then

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|}. \quad (12.19)$$

**Proof.** We simply take  $\mathbf{e} = \mathbf{r}(\mathbf{y})$  in Theorem 12.30.  $\square$

If  $\mathbf{A}$  is well-conditioned, (12.19) 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.

**Exercise 12.36** Let  $\|\cdot\|_p$  be the  $l_p$  vector norm and let  $\text{cond}_p(\mathbf{T}) = \|\mathbf{T}\|_p \|\mathbf{T}^{-1}\|_p$ , where  $\|\mathbf{T}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{T}\mathbf{x}\|_p / \|\mathbf{x}\|_p$  be the  $p$ -condition number of  $\mathbf{T} \in \mathbb{R}^{m,m}$ . In this exercise we find the  $p$ -condition numbers for the matrix  $\mathbf{T} := \text{tridiag}(-1, 2, -1)$  in terms of  $h := 1/(m+1)$ . You will need the explicit inverse of  $\mathbf{T}$  given by (6.30) and the eigenvalues given in Lemma 8.11.

a) Show that

$$\text{cond}_1(\mathbf{T}) = \text{cond}_\infty(\mathbf{T}) = \frac{1}{2} \begin{cases} h^{-2}, & m \text{ odd}, m > 1, \\ h^{-2} - 1, & m \text{ even}. \end{cases} \quad (12.20)$$

b) Show that for  $p = 2$  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

$$\frac{4}{\pi^2}h^{-2} - \frac{2}{3} < \text{cond}_2(\mathbf{T}) < \frac{4}{\pi^2}h^{-2}. \quad (12.21)$$

*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.*

## 12.3 Determining the Rank of a Matrix

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 11.3 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.

The following discussion can be used to decide how many of the computed singular values one can set equal to zero. Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  with  $m \geq n$  has singular value decomposition  $\mathbf{A} = \mathbf{U} \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^H$ , where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$  with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ . 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  $\mathbf{A}' := \mathbf{U} \begin{bmatrix} \boldsymbol{\Sigma}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}^H$ , where  $\boldsymbol{\Sigma}' := \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{n,n}$ . By Theorem 12.4

$$\|\mathbf{A} - \mathbf{A}'\|_F = \left\| \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\Sigma}' \\ \mathbf{0} \end{bmatrix} \right\|_F = \sqrt{\sigma_{r+1}^2 + \dots + \sigma_n^2} < \epsilon.$$

Thus  $\mathbf{A}$  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 12.37 (Best low rank approximation)** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$  has singular values  $\sigma_1 \geq \dots \geq \sigma_n \geq 0$ . For any  $r \leq \text{rank}(\mathbf{A})$  we have

$$\|\mathbf{A} - \mathbf{A}'\|_F = \min_{\substack{\mathbf{B} \in \mathbb{R}^{m,n} \\ \text{rank}(\mathbf{B})=r}} \|\mathbf{A} - \mathbf{B}\|_F = \sqrt{\sigma_{r+1}^2 + \dots + \sigma_n^2}.$$

For the proof of this theorem we refer to p. 322 of [20].

**Exercise 12.38** Consider the singular value decomposition

$$A := \begin{bmatrix} 0 & 3 & 3 \\ 4 & 1 & -1 \\ 4 & 1 & -1 \\ 0 & 3 & 3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 6 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & -\frac{1}{3} \\ \frac{1}{3} & -\frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \end{bmatrix}$$

- (a) Give orthonormal bases for  $\text{span}(A)$ ,  $\text{span}(A^T)$ ,  $\ker(A)$ ,  $\ker(A^T)$  and  $\text{span}(A)^\perp$ .
- (b) Explain why for all matrices  $B \in \mathbb{R}^{4,3}$  of rank one we have  $\|A - B\|_F \geq 6$ .
- (c) Give a matrix  $A_1$  of rank one such that  $\|A - A_1\|_F = 6$ .

**Exercise 12.39** 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  $\mathbf{B}\mathbf{x} = \mathbf{0}$ , 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 the matrix  $\mathbf{A}$  is very close to being singular for large  $n$ .
- (b) Use Theorem 12.37 to show that the smallest singular value  $\sigma_n$  of  $\mathbf{A}$  is bounded above by  $2^{2-n}$ .

## 12.4 Convergence and Spectral Radius

We start with some basic notions that we need.

### 12.4.1 Convergence in $\mathbb{R}^{m,n}$ and $\mathbb{C}^{m,n}$

**Definition 12.40** Consider an infinite sequence of matrices  $\{\mathbf{A}_k\} = \mathbf{A}_0, \mathbf{A}_1, \mathbf{A}_2, \dots$  in  $\mathbb{C}^{m,n}$ .

1.  $\{\mathbf{A}_k\}$  is said to converge to the limit  $\mathbf{A}$  in  $\mathbb{C}^{m,n}$  if each element sequence  $\{\mathbf{A}_k(ij)\}_k$  converges to the corresponding element  $\mathbf{A}(ij)$  for  $i = 1, \dots, m$  and  $j = 1, \dots, n$ .
2.  $\{\mathbf{A}_k\}$  is a **Cauchy sequence** if for all  $\epsilon > 0$  there is an integer  $N \in \mathbb{N}$  such that for each  $k, l \geq N$  and all  $i, j$  we have  $|\mathbf{A}_k(ij) - \mathbf{A}_l(ij)| \leq \epsilon$ .
3.  $\{\mathbf{A}_k\}$  is bounded if there is a constant  $M$  such that  $|\mathbf{A}_k(ij)| \leq M$  for all  $i, j, k$ .

By stacking the columns of  $\mathbf{A}$  into a vector in  $\mathbb{C}^{mn}$  we can use the results in Section 2.5 and obtain

- Theorem 12.41** 1. A sequence  $\{\mathbf{A}_k\}$  in  $\mathbb{C}^{m,n}$  converges to a matrix  $\mathbf{A} \in \mathbb{C}^{m,n}$  if and only if  $\lim_{k \rightarrow \infty} \|\mathbf{A}_k - \mathbf{A}\| = 0$  for any matrix norm  $\|\cdot\|$ .
2. A sequence  $\{\mathbf{A}_k\}$  in  $\mathbb{C}^{m,n}$  is convergent if and only if it is a Cauchy sequence.
3. Every bounded sequence  $\{\mathbf{A}_k\}$  in  $\mathbb{C}^{m,n}$  has a convergent subsequence.

### 12.4.2 The Spectral Radius

We define the **spectral radius** of a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  as the maximum absolute value of its eigenvalues.

$$\rho(\mathbf{A}) := \max_{\lambda \in \sigma(\mathbf{A})} |\lambda|. \quad (12.22)$$

**Theorem 12.42** For any matrix norm  $\|\cdot\|$  which is consistent on  $\mathbb{C}^{n,n}$  and any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$ .

**Proof.** Let  $(\lambda, \mathbf{x})$  be an eigenpair for  $\mathbf{A}$  and define  $\mathbf{X} := [\mathbf{x}, \dots, \mathbf{x}] \in \mathbb{C}^{n,n}$ . Then  $\lambda \mathbf{X} = \mathbf{A} \mathbf{X}$ , which implies  $|\lambda| \|\mathbf{X}\| = \|\lambda \mathbf{X}\| = \|\mathbf{A} \mathbf{X}\| \leq \|\mathbf{A}\| \|\mathbf{X}\|$ . Since  $\|\mathbf{X}\| \neq 0$  we obtain  $|\lambda| \leq \|\mathbf{A}\|$ .  $\square$

The inequality  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$  can almost be made into an equality by choosing the norm carefully.

**Theorem 12.43** Let  $\mathbf{A} \in \mathbb{C}^{n,n}$  and  $\epsilon > 0$  be given. There is a consistent matrix norm  $\|\cdot\|'$  on  $\mathbb{C}^{n,n}$  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 10.1 there is a unitary matrix  $\mathbf{U}$  and an upper triangular matrix  $\mathbf{R} = [r_{ij}]$  such that  $\mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{R}$ . For  $t > 0$  we define  $\mathbf{D}_t := \text{diag}(t, t^2, \dots, t^n) \in \mathbb{R}^{n,n}$ , 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  $\mathbf{B} \in \mathbb{C}^{n,n}$  and  $t > 0$  we define  $\|\mathbf{B}\|_t := \|\mathbf{D}_t \mathbf{U}^H \mathbf{B} \mathbf{U} \mathbf{D}_t^{-1}\|_1$ . We leave it as an exercise to show that this is a consistent matrix norm on  $\mathbb{C}^{n,n}$ . 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  $\epsilon$ . Then

$$\begin{aligned} \|\mathbf{A}\|' &= \|\mathbf{D}_t \mathbf{U}^H \mathbf{A} \mathbf{U} \mathbf{D}_t^{-1}\|_1 = \|\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |(D_t \mathbf{R} \mathbf{D}_t^{-1})_{ij}| \\ &\leq \max_{1 \leq j \leq n} (|\lambda_j| + \epsilon) = \rho(\mathbf{A}) + \epsilon. \end{aligned}$$

$\square$

**Theorem 12.44** For any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have

$$\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0} \iff \rho(\mathbf{A}) < 1.$$

**Proof.** Suppose  $\rho(\mathbf{A}) < 1$ . By Theorem 12.43 there is a consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n,n}$  such that  $\|\mathbf{A}\| < 1$ . But then  $\|\mathbf{A}^k\| \leq \|\mathbf{A}\|^k \rightarrow 0$  as  $k \rightarrow \infty$ . Hence  $\mathbf{A}^k \rightarrow \mathbf{0}$ . Conversely, suppose  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{A}$  with  $|\lambda| \geq 1$ . Since  $\mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}$ , by Theorem 5.3 it follows that  $\mathbf{A}^k \mathbf{x}$  does not tend to zero. But then we cannot have  $\mathbf{A}^k \rightarrow \mathbf{0}$ .  $\square$

**Theorem 12.45** For any consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n,n}$  and any  $\mathbf{A} \in \mathbb{C}^{n,n}$  we have

$$\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A}). \quad (12.23)$$

**Proof.** By Theorems 5.3 and 12.42 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 12.44 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  $\epsilon$  is arbitrary the result follows.  $\square$

**Exercise 12.46** The convergence  $\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A})$  can be quite slow. Consider

$$\mathbf{A} := \begin{bmatrix} \lambda & a & 0 & \cdots & 0 & 0 \\ 0 & \lambda & a & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \vdots & \\ 0 & 0 & 0 & \cdots & \lambda & a \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix} \in \mathbb{R}^{n,n}.$$

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  $k \geq n-1$ .

- Make a plot of  $f(k)$  for  $\lambda = 0.9$ ,  $a = 10$ , and  $k \leq 200$ . Your program should also compute  $\max_k f(k)$ . Use your program to determine how large  $k$  must be before  $f(k) < 10^{-8}$ .
- 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  $\mathbf{E}^n = \mathbf{0}$ .
- We have  $\mathbf{A}^k = (a\mathbf{E} + \lambda \mathbf{I})^k = \sum_{j=0}^{\min\{k, n-1\}} \binom{k}{j} \lambda^{k-j} a^j \mathbf{E}_n^j$  and conclude that the  $(1, n)$  element is given by  $f(k)$  for  $k \geq n-1$ .



### 12.4.3 Neumann Series

A geometric series of matrices is known as a Neumann Series.

**Theorem 12.47 (Neumann Series)** Suppose  $B \in \mathbb{C}^{n,n}$ . Then

1. The series  $\sum_{k=0}^{\infty} B^k$  converges if and only if  $\rho(B) < 1$ .
2. If  $\rho(B) < 1$  then  $(I - B)$  is nonsingular and  $(I - B)^{-1} = \sum_{k=0}^{\infty} B^k$ .
3. If  $\|B\| < 1$  for some consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n,n}$  then

$$\|(I - B)^{-1}\| \leq \frac{1}{1 - \|B\|}. \quad (12.24)$$

**Proof.**

1. Suppose  $\rho(B) < 1$ . We use Theorem 12.41 and show that the sequence  $\{A_m\}$  of partial sums  $A_m := \sum_{k=0}^m B^k$  is a Cauchy sequence. Let  $\epsilon > 0$ . By Theorem 12.43 there is a consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n,n}$  such that  $\|B\| < 1$ . Then for  $l > m$

$$\|A_l - A_m\| = \left\| \sum_{k=m+1}^l B^k \right\| \leq \sum_{k=m+1}^l \|B\|^k \leq \frac{\|B\|^{m+1}}{1 - \|B\|} \leq \epsilon$$

provided  $m \geq N$  and  $N$  is such that  $\frac{\|B\|^{N+1}}{1 - \|B\|} \leq \epsilon$ . Thus  $\{A_m\}$  is a Cauchy sequence and hence convergent.

Conversely, suppose  $(\lambda, x)$  is an eigenpair for  $B$  with  $\lambda \geq 1$ . Now for  $l > m$

$$\|(A_l - A_m)x\| = \left\| \sum_{k=m+1}^l B^k x \right\| = \left\| \sum_{k=m+1}^l \lambda^k x \right\| = \|x\| \sum_{k=m+1}^l |\lambda|^k \geq |\lambda|^{m+1} \|x\|.$$

But then  $\{A_m\}$  cannot be a Cauchy sequence and hence not convergent.

2. By induction on  $m$  it follows that

$$\left( \sum_{k=0}^m B^k \right) (I - B) = I - B^{m+1}. \quad (12.25)$$

For if  $(\sum_{k=0}^{m-1} B^k)(I - B) = I - B^m$  then

$$\left( \sum_{k=0}^m B^k \right) (I - B) = \left( \sum_{k=0}^{m-1} B^k + B^m \right) (I - B) = I - B^m + B^m - B^{m+1} = I - B^{m+1}.$$

Since  $\rho(B) < 1$  we conclude that  $B^{m+1} \rightarrow 0$  and hence taking limits in (12.25) we obtain  $(\sum_{k=0}^{\infty} B^k)(I - B) = I$  which completes the proof of 2.

3. By 1:  $\|(I - B)^{-1}\| = \|\sum_{k=0}^{\infty} B^k\| \leq \sum_{k=0}^{\infty} \|B\|^k = \frac{1}{1 - \|B\|}$ .

□

**Exercise 12.48** Show that  $\|B\|_t := \|D_t U^H B U D_t^{-1}\|_1$  defined in the proof of Theorem 12.43 is a consistent matrix norm on  $\mathbb{C}^{n,n}$ .

**Exercise 12.49** Suppose  $A \in \mathbb{C}^{n,n}$  is nonsingular and  $E \in \mathbb{C}^{n,n}$ . Show that  $A + E$  is nonsingular if and only if  $\rho(A^{-1}E) < 1$ .

## **Part IV**

# **Iterative Methods for Large Linear Systems**



## Chapter 13

# 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 [26]. 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.

### 13.1 Classical Iterative Methods; Component Form

Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $x_i$ , we obtain a fixed-point form of  $\mathbf{Ax} = \mathbf{b}$

$$x_i = \left( - \sum_{j=1}^{i-1} a_{ij}x_j - \sum_{j=i+1}^n a_{ij}x_j + b_i \right) / a_{ii}, \quad i = 1, 2, \dots, n. \quad (13.1)$$

Suppose we know an approximation  $\mathbf{x}^{(k)} = [x_1^{(k)}, \dots, x_n^{(k)}]^T$  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 (13.1) and compute a new approximation by

$$x_i^{(k+1)} = \left( - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} + b_i \right) / a_{ii}, \text{ for } i = 1, 2, \dots, n. \quad (13.2)$$

2. **Gauss-Seidel's method (GS method)** is a modification of Jacobi's method, where we use the new  $x_i^{(k+1)}$  immediately after it has been computed.

$$x_i^{(k+1)} = \left( - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} + b_i \right) / a_{ii}, \text{ for } i = 1, 2, \dots, n. \quad (13.3)$$

3. The **Successive Over Relaxation method (SOR method)** is obtained by introducing an acceleration parameter  $0 < \omega < 2$  in the GS method. We write  $x_i = \omega x_i + (1 - \omega)x_i$  and this leads to the method

$$x_i^{(k+1)} = \omega \left( - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} + b_i \right) / a_{ii} + (1 - \omega)x_i^{(k)}. \quad (13.4)$$

The SOR method reduces to the Gauss-Seidel method for  $\omega = 1$ . Denoting the right hand side of (13.3) by  $\mathbf{x}_{gs}^{(k+1)}$  we can write (13.4) as  $\mathbf{x}^{(k+1)} = \omega \mathbf{x}_{gs}^{(k+1)} + (1 - \omega)\mathbf{x}^{(k)}$ , and we see that  $\mathbf{x}^{(k+1)}$  is located on the straight line passing through the two points  $\mathbf{x}_{gs}^{(k+1)}$  and  $\mathbf{x}^{(k)}$ . The restriction  $0 < \omega < 2$  is necessary for convergence (cf. Theorem 13.22). Normally we choose the relaxation parameter  $\omega$  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  $\mathbf{x}_{gs}^{(k+1)}$  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 (13.4), computing an approximation denoted  $\mathbf{x}^{(k+1/2)}$  instead of  $\mathbf{x}^{(k+1)}$ , is followed by a back SOR sweep computing

$$x_i^{(k+1)} = \omega \left( - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1/2)} - \sum_{j=i+1}^n a_{ij}x_j^{(k+1)} + b_i \right) / a_{ii} + (1 - \omega)x_i^{(k+1/2)} \quad (13.5)$$

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 15.2 for an alternative example of a preconditioner.

We will refer to the J,GS, and SOR methods as the **classical (iteration) methods**.

## 13.2 The Discrete Poisson System

Consider the classical methods applied to the discrete Poisson matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  given by (8.7). Let  $n = m^2$  and set  $h = 1/(m+1)$ . In component form the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written (cf. (8.3))

$$4u_{i,j} - u_{i-1,j} - u_{i+1,j} - u_{i,j-1} - u_{i,j+1} = h^2 f_{i,j}, \quad i, j = 1, \dots, m,$$

with homogenous boundary conditions (8.4). Solving for  $u_{i,j}$  we obtain

$$u_{i,j} = (u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1} + h^2 f_{i,j})/4, \quad i, j = 1, \dots, m. \quad (13.6)$$

The J, GS, and SOR methods can now be written

$$\begin{aligned} J: v_{ij}^{(k+1)} &= (v_{i-1,j}^{(k)} + v_{i,j-1}^{(k)} + v_{i+1,j}^{(k)} + v_{i,j+1}^{(k)} + h^2 f_{ij})/4 \\ GS: v_{ij}^{(k+1)} &= (v_{i-1,j}^{(k+1)} + v_{i,j-1}^{(k+1)} + v_{i+1,j}^{(k)} + v_{i,j+1}^{(k)} + h^2 f_{ij})/4 \\ SOR: v_{ij}^{(k+1)} &= \omega(v_{i-1,j}^{(k+1)} + v_{i,j-1}^{(k+1)} + v_{i+1,j}^{(k)} + v_{i,j+1}^{(k)} + h^2 f_{ij})/4 + (1-\omega)v_{ij}^{(k)}. \end{aligned} \quad (13.7)$$

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.

Here is a Matlab program to test the convergence of Jacobi's method on the discrete Poisson problem.

**Algorithm 13.1 (Jacobi)** We carry out Jacobi iterations on the linear system (13.6) with  $\mathbf{F} = (f_{ij}) \in \mathbb{R}^{m,m}$ , starting with  $\mathbf{V}^{(0)} = \mathbf{0} \in \mathbb{R}^{m+2,m+2}$ . The output is the number of iterations  $k$ , to obtain  $\|\mathbf{V}^{(k)} - \mathbf{U}\|_M := \max_{i,j} |v_{ij}^{(k)} - u_{ij}| < tol$ . Here  $(u_{ij}) \in \mathbb{R}^{m+2,m+2}$  is the "exact" solution of (13.6) computed using the fast Poisson solver in Algorithm 9.1. We set  $k = K + 1$  if convergence is not obtained in  $K$  iterations.

```
function k=jdp(F,K,tol)
m=length(F);
U=fastpoisson(F);
V=zeros(m+2,m+2); W=V;
E=F/(m+1)^2;
for k=1:K
    for i=2:m+1
        for j=2:m+1
            W(i,j)=(V(i-1,j)+V(i+1,j)+V(i,j-1)+...
                    +V(i,j+1)+E(i-1,j-1))/4;
        end
    end
    if max(max(abs(W-U)))<tol, return
end
V=W;
end
k=K+1;
```

	$k_{100}$	$k_{2500}$	$k_{10\,000}$	$k_{40\,000}$	$k_{160\,000}$
J	385	8386			
GS	194	4194			
SOR	35	164	324	645	1286

**Table 13.1.** 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  $10^{-8}$ .

In Table 13.1 we show the output  $k = k_n$  from this algorithm using  $\mathbf{F} = \mathbf{ones}(m, m)$  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  $\omega$  known as the optimal acceleration parameter  $\omega = 2/(1 + \sin(\pi/(m+1)))$ . We will derive this value later. For the GS and SOR methods we have used Algorithm 13.2.

**Algorithm 13.2 (SOR)** This is the analog of Algorithm 13.1 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. The optimal value for the discrete Poisson problem is  $w = 2/(1 + \sin(\pi/(m+1)))$ .

```
function k=sordp(F,K,w,tol)
m=length(F);
U=fastpoisson(F);
V=zeros(m+2,m+2);
E=F/(m+1)^2;
for k=1:K
    for i=2:m+1
        for j=2:m+1
            V(i,j)=w*(V(i-1,j)+V(i+1,j)+V(i,j-1)...
                    +V(i,j+1)+E(i-1,j-1))/4+(1-w)*V(i,j);
        end
    end
    if max(max(abs(V-U))) < tol, return
end
k=K+1;
```

We make several remarks about these programs and the results in Table 13.1.

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  $\omega$  is rather spectacular.
2. We show in Section 13.5.1 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



- $\omega$ . The choice of  $tol$  will only influence the constants multiplying  $n$  or  $\sqrt{n}$ .
3. From (13.7) it follows that each iteration requires  $O(n)$  flops. Thus the number of flops to achieve a given tolerance is  $O(k_n \times n)$ . Therefore the number of flops for the J and GS method is  $O(n^2)$ , while it is only  $O(n^{3/2})$  for the SOR method with optimal  $\omega$ . 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  $v_{ij}^{(k+1)}$  in the same location as  $v_{ij}^{(k)}$ . For Jacobi's method we need an extra array. (W in Algorithm 13.1).
  5. Jacobi's method has the advantage that it can be easily parallelized.

### 13.3 Matrix Formulations of the Classical Methods

To study convergence we need matrix formulations of the classical methods. In general we can construct an iterative method by choosing a nonsingular matrix  $M$  and write  $Ax = b$  in the equivalent form  $Bx = c$ , where  $B = M^{-1}A$  and  $c = M^{-1}b$ . The system  $Bx = c$  can be written  $x = x - Bx + c = (I - B)x + c$ , and this defines the iterative method

$$x^{(k+1)} := Gx^{(k)} + c, \quad G = I - B = I - M^{-1}A, \quad c = M^{-1}b. \quad (13.8)$$

Different choices of  $M$  leads to different iterative methods. The matrix  $M$  can be interpreted in two ways. It is a **preconditioning matrix** since a good choice of  $M$  will lead to a system  $Bx = c$  with smaller condition number. It is also known as a **splitting matrix**, since if we split  $A$  in the form  $A = M + (A - M)$  then  $Ax = b$  can be written  $Mx = (M - A)x + b$  and this leads to the iterative method

$$Mx^{(k+1)} = (M - A)x^{(k)} + b \quad (13.9)$$

which is equivalent to (13.8).

The matrix  $M$  should be chosen so that  $G$  has small spectral radius and such that the linear system (13.9) is easy to solve for  $x^{(k+1)}$ . These are conflicting demands.  $M$  should be an approximation to  $A$  to obtain a  $B$  with small condition number, but then (13.9) might not be easy to solve for  $x^{(k+1)}$ .

#### 13.3.1 The Splitting Matrices for the Classical Methods

To describe  $M$  for the classical methods we write  $A$  as a sum of three matrices,  $A = D - A_L - A_R$ , where  $-A_L$ ,  $D$ , and  $-A_R$  are the lower, diagonal, and upper

part of  $\mathbf{A}$ , respectively. Thus  $\mathbf{D} := \text{diag}(a_{11}, \dots, a_{nn})$ ,

$$\mathbf{A}_L := \begin{bmatrix} 0 & & & \\ -a_{21} & 0 & & \\ \vdots & \ddots & \ddots & \\ -a_{n,1} & \cdots & -a_{n,n-1} & 0 \end{bmatrix}, \quad \mathbf{A}_R := \begin{bmatrix} 0 & -a_{12} & \cdots & -a_{1n} \\ & \ddots & \ddots & \vdots \\ & & 0 & -a_{n-1,n} \\ & & & 0 \end{bmatrix}. \quad (13.10)$$

**Proposition 13.3** *The splitting matrices  $\mathbf{M}_J, \mathbf{M}_1, \mathbf{M}_\omega$  for the J, GS, and SOR method are given by*

$$\mathbf{M}_J = \mathbf{D}, \quad \mathbf{M}_1 = \mathbf{D} - \mathbf{A}_L, \quad \mathbf{M}_\omega = \omega^{-1}\mathbf{D} - \mathbf{A}_L. \quad (13.11)$$

**Proof.** The equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written  $\mathbf{D}\mathbf{x} - \mathbf{A}_L\mathbf{x} - \mathbf{A}_R\mathbf{x} = \mathbf{b}$  or  $\mathbf{D}\mathbf{x} = \mathbf{A}_L\mathbf{x} + \mathbf{A}_R\mathbf{x} + \mathbf{b}$ . This leads to

$$\begin{aligned} J: \quad \mathbf{D}\mathbf{x}^{(k+1)} &= \mathbf{A}_L\mathbf{x}^{(k)} + \mathbf{A}_R\mathbf{x}^{(k)} + \mathbf{b}, \\ GS: \quad \mathbf{D}\mathbf{x}^{(k+1)} &= \mathbf{A}_L\mathbf{x}^{(k+1)} + \mathbf{A}_R\mathbf{x}^{(k)} + \mathbf{b}, \\ SOR: \quad \mathbf{D}\mathbf{x}^{(k+1)} &= \omega(\mathbf{A}_L\mathbf{x}^{(k+1)} + \mathbf{A}_R\mathbf{x}^{(k)} + \mathbf{b}) + (1 - \omega)\mathbf{D}\mathbf{x}^{(k)}. \end{aligned} \quad (13.12)$$

Writing these equations in the form (13.9) we obtain (13.11).  $\square$

**Example 13.4** *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*

$$\mathbf{G}_\omega = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} \omega/2 & 0 \\ \omega^2/4 & \omega/2 \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 1 - \omega & \omega/2 \\ \omega(1 - \omega)/2 & 1 - \omega + \omega^2/4 \end{bmatrix}. \quad (13.13)$$

*For the J and GS method we have*

$$\mathbf{G}_J = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix}, \quad \mathbf{G}_1 = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix}. \quad (13.14)$$

*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*

$$x_1^{(k+1)} = \frac{1}{2}x_2^{(k)} + \frac{1}{2}, \quad x_2^{(k+1)} = \frac{1}{2}x_1^{(k+1)} + \frac{1}{2}.$$

Substituting the value of  $x_1^{(k+1)}$  from the first equation into the second equation we find

$$x_2^{(k+1)} = \frac{1}{2} \left( \frac{1}{2} x_2^{(k)} + \frac{1}{2} \right) + \frac{1}{2} = \frac{1}{4} x_2^{(k)} + \frac{3}{4}.$$

Thus

$$\mathbf{x}^{(k+1)} = \begin{bmatrix} x_1^{(k+1)} \\ x_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix} \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \end{bmatrix} + \begin{bmatrix} 1/2 \\ 3/4 \end{bmatrix} = \mathbf{G}_1 \mathbf{x}^{(k)} + \mathbf{c}.$$

## 13.4 Convergence of Fixed-point Iteration

We have seen that the classical methods can be written in the form (13.8) 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  $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$  since

$$\mathbf{x} = \lim_{k \rightarrow \infty} \mathbf{x}^{(k+1)} = \lim_{k \rightarrow \infty} (\mathbf{G}\mathbf{x}^{(k)} + \mathbf{c}) = \mathbf{G} \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} + \mathbf{c} = \mathbf{G}\mathbf{x} + \mathbf{c}.$$

For a general  $\mathbf{G} \in \mathbb{C}^{n,n}$  and  $\mathbf{c} \in \mathbb{C}^n$  a solution of  $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$  is called a **fixed-point** and the iteration  $\mathbf{x}^{(k+1)} = \mathbf{G}\mathbf{x}^{(k)} + \mathbf{c}$  a **fixed-point iteration**. The fixed-point is unique if  $\mathbf{I} - \mathbf{G}$  is nonsingular.

Consider next convergence of fixed-point iteration.

**Definition 13.5** 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  $\mathbf{x}^{(0)}$ .

To study convergence we consider for  $k \geq 0$  the error

$$\boldsymbol{\epsilon}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}.$$

**Lemma 13.6** 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.** Subtraction of  $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$  from  $\mathbf{x}^{(k+1)} = \mathbf{G}\mathbf{x}^{(k)} + \mathbf{c}$  leads to cancellation of  $\mathbf{c}$  and  $\boldsymbol{\epsilon}^{(k+1)} = \mathbf{G}\boldsymbol{\epsilon}^{(k)}$ . By induction  $\boldsymbol{\epsilon}^{(k)} = \mathbf{G}^k \boldsymbol{\epsilon}^{(0)}$  for  $k = 0, 1, 2, \dots$ . It follows that  $\boldsymbol{\epsilon}^{(k)} \rightarrow \mathbf{0}$  for all  $\boldsymbol{\epsilon}^{(0)}$  if and only if  $\mathbf{G}^k \rightarrow \mathbf{0}$ .  $\square$

Recall that the spectral radius of a matrix  $\mathbf{G} \in \mathbb{C}^{n,n}$  with eigenvalues  $\lambda_1, \dots, \lambda_n$  is defined as  $\rho(\mathbf{G}) = \max_j |\lambda_j|$ . Using Theorem 12.42 we obtain the following theorem:

**Theorem 13.7** Suppose  $\mathbf{G} \in \mathbb{C}^{n,n}$  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  $\mathbb{C}^{n,n}$  (cf. Theorem 12.42) we obtain

**Corollary 13.8** 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 13.9** Show that both Jacobi's method and Gauss-Seidel's method diverge for  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ .

**Exercise 13.10** Explain why  $J$  and  $GS$  converge for the cubic spline matrices  $\mathbf{N}_1$ ,  $\mathbf{N}_2$ , and  $\mathbf{N}_3$ , in Chapter 6. (This is mainly of academic interest since for tridiagonal strictly diagonally dominant matrices Gaussian elimination has complexity  $O(n)$  and is preferable for such systems.)

**Exercise 13.11** Show that the  $J$  method converges if  $\mathbf{A}$  is strictly diagonally dominant, i. e.,  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$  for  $i = 1, \dots, n$ .

**Exercise 13.12** 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  $|\epsilon_j^{(k+1)}| \leq r \|\epsilon^{(k)}\|_\infty$  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  $\epsilon^{(k)} = \mathbf{G}^k \epsilon^{(0)}$  we obtain

$$\|\epsilon^{(k)}\| = \|\mathbf{G}^k \epsilon^{(0)}\| \leq \|\mathbf{G}^k\| \|\epsilon^{(0)}\| \approx \rho(\mathbf{G})^k \|\epsilon^{(0)}\|.$$

For the last formula we apply Theorem 12.45 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 13.13** Suppose  $\rho(\mathbf{G}) = 1 - \eta$  for some  $0 < \eta < 1$ ,  $\|\cdot\|$  a consistent matrix norm, and let  $s \in \mathbb{N}$ . Then

$$\tilde{k} := \frac{\log(10)s}{\eta} \tag{13.15}$$

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 13.14** Consider the iteration in Example 13.4. Show that  $\rho(\mathbf{G}_J) = 1/2$ . Then show that  $x_1^{(k)} = x_2^{(k)} = 1 - 2^{-k}$  for  $k \geq 0$ . Thus the estimate in Lemma 13.13 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 12.46).

### 13.4.1 Stopping the Iteration

In Algorithms 13.1 and 13.2 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 13.15** *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*

$$\|\mathbf{x}^{(k)} - \mathbf{x}\| \leq \frac{\|\mathbf{G}\|}{1 - \|\mathbf{G}\|} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|, \quad k \geq 1. \quad (13.16)$$

**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 (13.16).  $\square$

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 (13.8) in an alternative form. If  $\mathbf{M}$  is the splitting matrix of the method then by (13.9) we have  $\mathbf{M}\mathbf{x}^{(k+1)} = \mathbf{M}\mathbf{x}^{(k)} - \mathbf{A}\mathbf{x}^{(k)} + \mathbf{b}$ . This leads to

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{M}^{-1}\mathbf{r}^{(k)}, \quad \mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}. \quad (13.17)$$

Testing on  $\mathbf{r}^{(k)}$  works fine if  $\mathbf{A}$  is well conditioned, but Theorem 12.35 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.

### 13.4.2 Richardson's Method (R method)

This method is based on the simple splitting  $\mathbf{M}_R := \alpha\mathbf{I}$ , where  $\alpha$  is a nonzero scalar. By (13.17) we obtain Richardson's method in the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{-1}\mathbf{r}^{(k)}, \quad \mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}. \quad (13.18)$$

If all eigenvalues of  $\mathbf{A}$  have positive real parts then the R method converges provided  $\alpha$  is sufficiently large.

**Proposition 13.16** *Suppose all eigenvalues of  $\mathbf{A}$  have positive real parts and that  $\alpha$  is real. Then there is an  $\alpha_0$  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

$$\alpha^* := \frac{\lambda_1 + \lambda_n}{2}, \text{ and } \rho(\mathbf{G}(\alpha^*)) = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}. \quad (13.19)$$

**Proof.** The eigenvalues of  $\mathbf{G}(\alpha)$  are

$$\mu_j(\alpha) = 1 - \lambda_j/\alpha, \quad j = 1, \dots, n,$$

and if  $u_j := \operatorname{Re} \lambda_j > 0$  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  $\rho(\mathbf{G}(\alpha)) > \rho(\mathbf{G}(\alpha^*))$ .  $\square$

### 13.5 Convergence of the Classical Methods for the Discrete Poisson Matrix

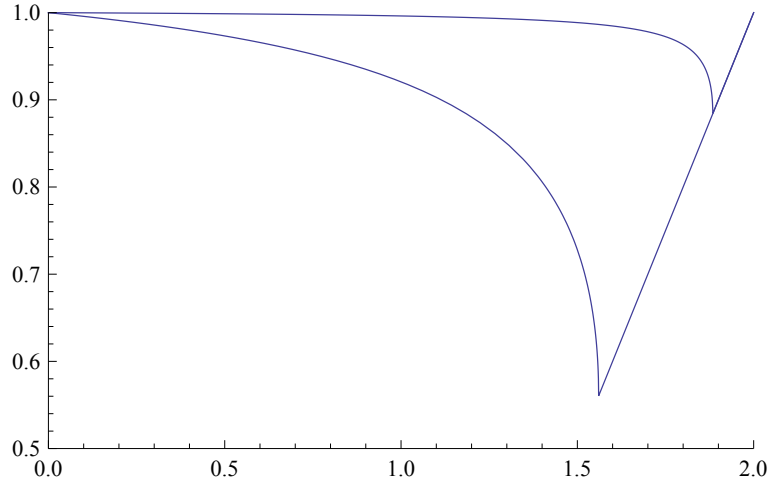
The matrix  $\mathbf{A}$  in (8.7) is symmetric positive definite (cf. Theorem 8.13). We show in Theorem 13.23 that the SOR method converges for all  $0 < \omega < 2$  if  $\mathbf{A}$  is symmetric positive definite. So the GS method converges, but the J method does not converge for all symmetric positive definite matrices.

**Exercise 13.17** 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.



**Figure 13.1.**  $\rho(\mathbf{G}_\omega)$  with  $\omega \in [0, 2]$  for  $n = 100$ , (lower curve) and  $n = 2500$  (upper curve).

Recall that by (8.22) the eigenvalues  $\lambda_{j,k}$  of  $\mathbf{A}$  given by (8.7) are

$$\lambda_{j,k} = 4 - 2\cos(j\pi h) - 2\cos(k\pi h), \quad j, k = 1, \dots, m, \quad h = 1/(m+1).$$

Consider first Jacobi's method. The matrix  $\mathbf{G}_J = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{I} - \mathbf{A}/4$  has eigenvalues

$$\mu_{j,k} = 1 - \frac{1}{4}\lambda_{j,k} = \frac{1}{2}\cos(j\pi h) + \frac{1}{2}\cos(k\pi h), \quad j, k = 1, \dots, m. \quad (13.20)$$

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 13.6.

**Theorem 13.18** *Consider the SOR iteration (13.7), with the natural ordering. The spectral radius of  $\mathbf{G}_\omega$  is*

$$\rho(\mathbf{G}_\omega) = \begin{cases} \frac{1}{4} \left( \omega\beta + \sqrt{(\omega\beta)^2 - 4(\omega - 1)} \right)^2, & \text{for } 0 < \omega \leq \omega^*, \\ \omega - 1, & \text{for } \omega^* < \omega < 2, \end{cases} \quad (13.21)$$

where  $\beta := \rho(\mathbf{G}_J)$  and

$$\omega^* := \frac{2}{1 + \sqrt{1 - \beta^2}} > 1. \quad (13.22)$$

Moreover,

$$\rho(\mathbf{G}_\omega) > \rho(\mathbf{G}_{\omega^*}) \text{ for } \omega \in (0, 2) \setminus \{\omega^*\}. \quad (13.23)$$

	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 13.2.** Spectral radii for  $\mathbf{G}_J$ ,  $\mathbf{G}_1$ ,  $\mathbf{G}_{\omega^*}$  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 13.1 for  $n = 100$  (lower curve) and  $n = 2500$  (upper curve). As  $\omega$  increases the spectral radius of  $\mathbf{G}_{\omega}$  decreases monotonically to the minimum  $\omega^*$ . Then it increases linearly to the value one for  $\omega = 2$ . We call  $\omega^*$  the **optimal relaxation parameter**.

For the discrete Poisson problem we have  $\beta = \cos(\pi h)$  and it follows from (13.21), (13.22) that

$$\omega^* = \frac{2}{1 + \sin(\pi h)}, \quad \rho(\mathbf{G}_{\omega^*}) = \omega^* - 1 = \frac{1 - \sin(\pi h)}{1 + \sin(\pi h)}, \quad h = \frac{1}{m+1}. \quad (13.24)$$

Letting  $\omega = 1$  in (13.21) 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 13.2 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 13.1.

### 13.5.1 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  $\omega$ . We use Lemma 13.13 that provided the estimate

$$\tilde{k}_n = \frac{\log(10)s}{\eta}, \quad \rho(\mathbf{G}) = 1 - \eta.$$

Note that  $h = 1/(m+1) \approx n^{-1/2}$ .

- 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 13.20** 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

$$x_1^{(k)} = x_2^{(k)} = 1 - a^k, \quad k \geq 0,$$

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 (13.15) for the rate of convergence for  $a = 0.9$  and  $s = 16$  and compare with the true number of iterations determined from  $|a|^k \leq 10^{-16}$ .

## 13.6 Convergence Analysis for SOR

The iteration matrix  $\mathbf{G}_{\omega}$  for the SOR method can be written in two alternative forms that are both useful for the analysis.

**Lemma 13.21** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  and  $\mathbf{D} = \text{diag}(a_{11}, \dots, a_{nn})$  are both nonsingular. Then

$$\mathbf{G}_{\omega} = \mathbf{I} - (\omega^{-1}\mathbf{D} - \mathbf{A}_L)^{-1}\mathbf{A} = (\mathbf{I} - \omega\mathbf{L})^{-1}(\omega\mathbf{R} + (1 - \omega)\mathbf{I}), \quad (13.25)$$

where  $\mathbf{A}_L$  and  $\mathbf{A}_R$  are given by (13.10) and

$$\mathbf{L} := \mathbf{D}^{-1}\mathbf{A}_L, \quad \mathbf{R} := \mathbf{D}^{-1}\mathbf{A}_R, \quad \text{so that } \mathbf{D}^{-1}\mathbf{A} = \mathbf{I} - \mathbf{L} - \mathbf{R}. \quad (13.26)$$

**Proof.** For the first form see (13.8) and Proposition 13.3. Solving the SOR part of (13.12) 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  $\mathbf{G}_{\omega}$  is given by the second form in (13.25).  $\square$

We start with the following convergence result.

**Theorem 13.22** The SOR method diverges if  $\omega$  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 (13.25) 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 6.32 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  $\lambda$  of  $\mathbf{G}_\omega$ . We conclude that  $\rho(\mathbf{G}_\omega) \geq |1 - \omega|$ . But then  $\rho(\mathbf{G}_\omega) \geq 1$  if  $\omega$  is not in the interval  $(0, 2)$  and by Theorem 13.7 SOR diverges.  $\square$

We next show that SOR converges for all  $\omega \in (0, 2)$  if  $\mathbf{A}$  is symmetric positive definite.

**Theorem 13.23** *SOR converges for a symmetric positive definite matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  if and only if  $0 < \omega < 2$ . In particular, Gauss-Seidel's method converges for a symmetric positive definite matrix.*

**Proof.** By Theorem 13.22 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

$$\mathbf{A}\mathbf{x} = (\omega^{-1}\mathbf{D} - \mathbf{A}_L)\mathbf{y}, \quad \mathbf{y} := (1 - \lambda)\mathbf{x}. \quad (13.27)$$

Since  $\mathbf{A} = -\mathbf{A}_L + \mathbf{D} - \mathbf{A}_R$  we find

$$(\omega^{-1}\mathbf{D} - \mathbf{D} + \mathbf{A}_R)\mathbf{y} = (\omega^{-1}\mathbf{D} - \mathbf{A}_L - \mathbf{A})\mathbf{y} \stackrel{(13.27)}{=} \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{y} = \lambda\mathbf{A}\mathbf{x},$$

so that by taking inner products and replacing  $\mathbf{A}_R^H$  by  $\mathbf{A}_L$

$$\begin{aligned} \langle \mathbf{y}, \lambda\mathbf{A}\mathbf{x} \rangle &= \langle \mathbf{y}, (\omega^{-1}\mathbf{D} - \mathbf{D} + \mathbf{A}_R)\mathbf{y} \rangle = \langle (\omega^{-1}\mathbf{D} - \mathbf{D} + \mathbf{A}_R^H)\mathbf{y}, \mathbf{y} \rangle \\ &= \langle (\omega^{-1}\mathbf{D} - \mathbf{D} + \mathbf{A}_L)\mathbf{y}, \mathbf{y} \rangle. \end{aligned} \quad (13.28)$$

Taking inner product with  $\mathbf{y}$  in (13.27) and adding to (13.28) 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,

$$(2\omega^{-1} - 1)|1 - \lambda|^2\langle \mathbf{D}\mathbf{x}, \mathbf{x} \rangle = (1 - |\lambda|^2)\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle. \quad (13.29)$$

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 (13.27) implies that  $\mathbf{A}$  is singular. Since  $0 < \omega < 2$  implies  $\omega^{-1} > 1/2$  the left side of (13.29) 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.  $\square$

The following analysis holds both for the discrete Poisson matrix and the averaging matrix given by (8.9). A more general theory is presented in [26]. Consider first how the eigenvalues of  $\mathbf{G}_J$  and  $\mathbf{G}_\omega$  are related.

**Theorem 13.24** *Consider for  $a, d \in \mathbb{R}$  the SOR method applied to the matrix (8.9), where we use the natural ordering. Moreover, assume  $\omega \in (0, 2)$ .*

1. *If  $\lambda \neq 0$  is an eigenvalue of  $\mathbf{G}_\omega$  then*

$$\mu := \frac{\lambda + \omega - 1}{\omega \lambda^{1/2}} \quad (13.30)$$

*is an eigenvalue of  $\mathbf{G}_J$ .*

2. *If  $\mu$  is an eigenvalue of  $\mathbf{G}_J$  and  $\lambda$  satisfies the equation*

$$\mu \omega \lambda^{1/2} = \lambda + \omega - 1 \quad (13.31)$$

*then  $\lambda$  is an eigenvalue of  $\mathbf{G}_\omega$ .*

**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  $(\lambda, \mathbf{w})$  is an eigenpair for  $\mathbf{G}_\omega$ . By (13.25)  $(\mathbf{I} - \omega \mathbf{L})^{-1}(\omega \mathbf{R} + (1 - \omega)\mathbf{I})\mathbf{w} = \lambda \mathbf{w}$  or

$$(\omega \mathbf{R} + \lambda \omega \mathbf{L})\mathbf{w} = (\lambda + \omega - 1)\mathbf{w}. \quad (13.32)$$

Let  $\mathbf{w} = \text{vec}(\mathbf{W})$ , where  $\mathbf{W} \in \mathbb{C}^{m,m}$ . Then (13.32) can be written

$$\frac{\omega}{4}(\lambda w_{i-1,j} + \lambda w_{i,j-1} + w_{i+1,j} + w_{i,j+1}) = (\lambda + \omega - 1)w_{i,j}, \quad (13.33)$$

where  $w_{i,j} = 0$  if  $i \in \{0, m+1\}$  or  $j \in \{0, m+1\}$ . We claim that  $(\mu, \mathbf{v})$  is an eigenpair for  $\mathbf{G}_J$ , where  $\mu$  is given by (13.30) and  $\mathbf{v} = \text{vec}(\mathbf{V})$  with

$$v_{i,j} := \lambda^{-(i+j)/2} w_{i,j}. \quad (13.34)$$

Indeed, replacing  $w_{i,j}$  by  $\lambda^{(i+j)/2} v_{i,j}$  in (13.33) 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

$$\mathbf{G}_J \mathbf{v} = (\mathbf{L} + \mathbf{R})\mathbf{v} = \frac{\lambda + \omega - 1}{\omega \lambda^{1/2}} \mathbf{v} = \mu \mathbf{v}.$$

For the converse let  $(\mu, \mathbf{v})$  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  $\lambda$  be a solution of (13.31). 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 (13.33)  $(\lambda, \mathbf{w})$  is an eigenpair for  $\mathbf{G}_\omega$ .  $\square$

**Proof of Theorem 13.18** By (8.22) 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  $\mu$  is an eigenvalue then  $-\mu$  is also an eigenvalue. Thus it is enough to consider positive eigenvalues  $\mu$ . For simplicity of notation let again  $a = -1$  and  $d = 2$ . Solving (13.31) for  $\lambda$  gives

$$\lambda(\mu) := \frac{1}{4} \left( \omega \mu \pm \sqrt{(\omega \mu)^2 - 4(\omega - 1)} \right)^2. \quad (13.35)$$

Both roots  $\lambda(\mu)$  are eigenvalues of  $\mathbf{G}_\omega$ . 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  $\omega$ ,  $\lambda(\mu)$  changes from real to complex at

$$\omega = \tilde{\omega}(\mu) := \frac{2}{1 + \sqrt{1 - \mu^2}}. \quad (13.36)$$

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 (13.35) are positive and the larger one is

$$\lambda(\mu) = \frac{1}{4} \left( \omega \mu + \sqrt{(\omega \mu)^2 - 4(\omega - 1)} \right)^2, \quad 0 < \omega \leq \tilde{\omega}(\mu). \quad (13.37)$$

Both  $\lambda(\mu)$  and  $\tilde{\omega}(\mu)$  are strictly increasing as functions of  $\mu$ . It follows that  $|\lambda(\mu)|$  is maximized for  $\mu = \rho(\mathbf{G}_J) =: \beta$  and for this value of  $\mu$  we obtain (13.21) for  $0 < \omega \leq \tilde{\omega}(\beta) = \omega^*$ .

Evidently  $\rho(\mathbf{G}_\omega) = \omega - 1$  is strictly increasing in  $\omega^* < \omega < 2$ . Equation (13.23) 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.



## Chapter 14

# The Conjugate Gradient Method

The conjugate gradient method is an iterative method for solving large sparse linear systems  $\mathbf{Ax} = \mathbf{b}$  with a symmetric positive definite coefficient matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ . It can also be used to minimize a quadratic function  $Q : \mathbb{R}^n \rightarrow \mathbb{R}$  given by  $Q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{b}$ , see the following chapter. We compute a sequence of approximations to the exact solution. Each new approximation  $\mathbf{x}^{(k+1)}$  is computed from the previous  $\mathbf{x}^{(k)}$  by a formula of the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad (14.1)$$

where  $\mathbf{p}^{(k)}$  is a vector, the **search direction**, and  $\alpha_k$  is a scalar determining the **step length**. A characteristic of the method is that the residuals  $\mathbf{r}^{(k)} := \mathbf{b} - \mathbf{Ax}^{(k)}$  (the negative gradients of  $Q(\mathbf{x}^{(k)})$ ) are orthogonal (or conjugate), i. e.,  $(\mathbf{r}^{(i)}, \mathbf{r}^{(j)}) = 0$  for  $i \neq j$ , where  $(\mathbf{x}, \mathbf{y}) := \mathbf{x}^T \mathbf{y}$  is the usual inner product in  $\mathbb{R}^n$ . This orthogonality property has given the method its name. If  $\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(n-1)}$  are nonzero then  $\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(n)}$  are  $n+1$  orthogonal vectors in  $\mathbb{R}^n$  and  $\mathbf{r}^{(n)}$  must be zero. It follows that the conjugate gradient method is a direct method. The exact solution is found in a finite number of operations. In practice, however the method is used as an iterative method for large linear systems since the residuals become small quite rapidly. For the Poisson problem the method converges as fast as the SOR-method with optimal acceleration parameter and we do not have to estimate the parameter. The conjugate gradient method was first published as a direct method in [7]. It was only some 20-30 years later that the iterative nature was seriously appreciated.

The number of iterations to achieve a desired accuracy is essentially proportional to the square root of the 2-norm condition number of the coefficient matrix of the linear system. Thus the smaller the condition number the faster the method converges.

Before deriving the method we give the algorithm, discuss implementation and give numerical examples.

## 14.1 The Conjugate Gradient Algorithm

Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is symmetric positive definite and let  $\mathbf{b} \in \mathbb{R}^n$ . To solve the linear system  $\mathbf{Ax} = \mathbf{b}$  we choose an initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ , set  $\mathbf{p}^{(0)} = \mathbf{r}^{(0)} = \mathbf{b} - \mathbf{Ax}^{(0)}$  and generate a sequence of vectors  $\{\mathbf{x}^{(k)}\}$  as follows:

For  $k = 0, 1, 2, \dots$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad \alpha_k = \frac{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})}{(\mathbf{p}^{(k)}, \mathbf{Ap}^{(k)})}, \quad (14.2)$$

$$\mathbf{r}^{(k+1)} := \mathbf{b} - \mathbf{Ax}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k \mathbf{Ap}^{(k)}, \quad (14.3)$$

$$\mathbf{p}^{(k+1)} = \mathbf{r}^{(k+1)} + \beta_k \mathbf{p}^{(k)}, \quad \beta_k = \frac{(\mathbf{r}^{(k+1)}, \mathbf{r}^{(k+1)})}{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})}. \quad (14.4)$$

Here  $(\mathbf{u}, \mathbf{v}) := \mathbf{u}^T \mathbf{v}$  is the usual inner product of two vectors.

The middle equation (14.3) follows from (14.2). Indeed,

$$\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{Ax}^{(k+1)} = \mathbf{b} - \mathbf{A}(\mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}) = \mathbf{r}^{(k)} - \alpha_k \mathbf{Ap}^{(k)}.$$

is the residual corresponding to  $\mathbf{x}^{(k+1)}$ .

**Example 14.1** Consider the linear system

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Starting with  $\mathbf{x}^{(0)} = \mathbf{0}$  we set  $\mathbf{p}^{(0)} = \mathbf{r}^{(0)} = \mathbf{b} = [1, 0]^T$ . Using (14.2) we find  $\alpha_0 = \frac{(\mathbf{r}^{(0)}, \mathbf{r}^{(0)})}{(\mathbf{p}^{(0)}, \mathbf{Ap}^{(0)})} = \frac{1}{2}$ . Then  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}$  and from (14.3) we find  $\mathbf{r}^{(1)} = \mathbf{r}^{(0)} - \alpha_0 \mathbf{Ap}^{(0)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 2 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}$ . By (14.4) we find  $\beta_0 = \frac{(\mathbf{r}^{(1)}, \mathbf{r}^{(1)})}{(\mathbf{r}^{(0)}, \mathbf{r}^{(0)})} = \frac{1}{4}$  so that

$$\mathbf{p}^{(1)} = \mathbf{r}^{(1)} + \beta_0 \mathbf{p}^{(0)} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1/4 \\ 1/2 \end{bmatrix}.$$

Continuing with the next iteration we obtain  $\alpha_1 = \frac{(\mathbf{r}^{(1)}, \mathbf{r}^{(1)})}{(\mathbf{p}^{(1)}, \mathbf{Ap}^{(1)})} = \frac{2}{3}$  and  $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)} = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} + \frac{2}{3} \begin{bmatrix} 1/4 \\ 1/2 \end{bmatrix} = \begin{bmatrix} 2/3 \\ 1/3 \end{bmatrix}$ . Since  $\mathbf{r}^{(2)} = \mathbf{r}^{(1)} - \alpha_1 \mathbf{Ap}^{(1)} = \mathbf{0}$  this is the exact solution found in  $n = 2$  iterations.

**Exercise 14.2** Do one iteration with the conjugate gradient method when  $\mathbf{x}^{(0)} = \mathbf{0}$ . (Answer:  $\mathbf{x}^{(1)} = \frac{(\mathbf{b}, \mathbf{b})}{(\mathbf{b}, \mathbf{Ab})} \mathbf{b}$ .)

**Exercise 14.3** 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  $\mathbf{x}^{(0)} = \mathbf{0}$ .



The formulas in (14.2)-(14.4) and the previous discussion form a basis for an algorithm.

**Algorithm 14.4 (Conjugate Gradient Iteration)** The symmetric positive definite linear system  $\mathbf{Ax} = \mathbf{b}$  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  $k > \text{imax}$ .  $K$  is the number of iterations used.

```
function [x,K]=cg(A,b,x,tol,imax)
r=b-A*x; p=r; rho=r'*r;
rho0=rho; for k=0:imax
    if sqrt(rho/rho0)<= tol
        K=k; return
    end
    t=A*p; a=rho/(p'*t);
    x=x+a*p; r=r-a*t;
    rhos=rho; rho=r'*r;
    p=r+(rho/rhos)*p;
end
K=imax+1;
```

The work involved in each iteration is

1. one matrix times vector ( $\mathbf{t} = \mathbf{Ap}$ ),
2. two inner products ( $(\mathbf{p}, \mathbf{t})$  and  $(\mathbf{r}, \mathbf{r})$ ),
3. three vector-plus-scalar-times-vector ( $\mathbf{x} = \mathbf{x} + \mathbf{ap}$ ,  $\mathbf{r} = \mathbf{r} - \mathbf{at}$  and  $\mathbf{p} = \mathbf{r} + (\text{rho}/\text{rhos})\mathbf{p}$ ),

The dominating part is the computation of  $\mathbf{t} = \mathbf{Ap}$ .

## 14.2 Numerical Example

We test the method on the example used in Chapter 13. 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  $\mathbf{f} = [1, \dots, 1]^T \in \mathbb{R}^n$ .

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 14.4 in  $O(n)$  flops. With such an implementation the total number of flops 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  $\mathbf{V}, \mathbf{R}, \mathbf{P}, \mathbf{B}, \mathbf{T} \in \mathbb{R}^{m,m}$  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.

$n$	2 500	10 000	40 000	1 000 000	4 000 000
$K$	22	22	21	21	20

**Table 14.6.** 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
$K$	140	294	587	1168
$K/\sqrt{n}$	1.86	1.87	1.86	1.85

**Table 14.7.** The number of iterations  $K$  for the Poisson problem on a  $\sqrt{n} \times \sqrt{n}$  grid for various  $n$

**Algorithm 14.5 (Testing Conjugate Gradient )**

```

 $\mathbf{A} = \text{tridiag}_m(a, d, a) \otimes \mathbf{I}_m + \mathbf{I}_m \otimes \text{tridiag}_m(a, d, a) \in \mathbb{R}^{m^2, m^2}$ 

function [V,K]=cgtest(m,a,d,tol,itmax)
h=1/(m+1); R=h*h*ones(m);
D=sparse(tridiagonal(a,d,a,m));
V=zeros(m,m); P=R; rho=sum(sum(R.*R)); rho0=rho;
for k=1:itmax
    if sqrt(rho/rho0)<= tol
        K=k; return
    end
    T=D*P+P*D; a=rho/sum(sum(P.*T)); V=V+a*P; R=R-a*T;
    rhos=rho; rho=sum(sum(R.*R)); P=R+(rho/rhos)*P;
end;
K=itmax+1;

```

Consider first the averaging matrix given by  $a = 1/9$  and  $d = 5/18$ . Starting with  $\mathbf{x}^{(0)} = \mathbf{0}$  and  $\text{tol} = 10^{-8}$  we obtain the values in Table 14.6.

The convergence is quite rapid. Note that each iteration only requires  $O(n)$  flops and since it appears that the number of iterations can be bounded independently of  $n$ , 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 corresponding to  $a = -1$  and  $d = 2$ . Again starting with  $\mathbf{x}^{(0)} = \mathbf{0}$  and  $\text{tol} = 10^{-8}$  and using  $CG$  in the form of Algorithm 14.5 we list  $K$ , the required number of iterations, and  $K/\sqrt{n}$ . We obtain the values in Table 14.7.

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 any acceleration parameter.

We will show in Section 14.4 that the number of iterations to achieve  $\|\mathbf{r}\|_2/\|\mathbf{r}\|_0 \leq$

$tol$  is bounded by the square root of the 2-norm condition number of  $\mathbf{T}_2$ .

For the averaging problem it follows from (8.22) that the largest and smallest eigenvalue of  $\mathbf{T}_2$  are  $\lambda_{max} = \frac{5}{9} + \frac{4}{9} \cos(\pi h)$  and  $\lambda_{min} = \frac{5}{9} - \frac{4}{9} \cos(\pi h)$ . Thus

$$\text{cond}_2(\mathbf{T}_2) = \frac{\lambda_{max}}{\lambda_{min}} = \frac{5 + 4 \cos(\pi h)}{5 - 4 \cos(\pi h)} \leq 9.$$

Thus the condition number is independent of  $n$  and the number of iterations can be bounded independently of  $n$ .

For the Poisson problem we find

$$\text{cond}_2(\mathbf{T}_2) = \frac{\lambda_{max}}{\lambda_{min}} = \frac{1 + \cos(\pi h)}{1 - \cos(\pi h)} = \text{cond}_2(\mathbf{T}) = O(n)$$

and we solve the discrete Poisson problem in  $O(n^{3/2})$  flops. Again this is the same as for the SOR method and for the fast method without the FFT. In comparison the Cholesky Algorithm requires  $O(n^2)$  flops both for the averaging and the Poisson problem.

### 14.3 Derivation and Basic Properties

Let  $\mathbf{A} \in \mathbb{R}^{n,n}$  be symmetric positive definite. We will use two inner products on  $\mathbb{R}^n$

1.  $(\mathbf{x}, \mathbf{y}) := \mathbf{x}^T \mathbf{y}$
2.  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{A} \mathbf{y}$ .

The first product is the usual inner product corresponding to the Euclidian norm, while the second product, called the  $\mathbf{A}$ -product or the energy product, is an inner product since  $\mathbf{A}$  is symmetric positive definite.

**Exercise 14.8** Show that the  $\mathbf{A}$ -inner product is an inner product.

We note that

$$\langle \mathbf{x}, \mathbf{y} \rangle = (\mathbf{x}, \mathbf{A} \mathbf{y}) = (\mathbf{A} \mathbf{x}, \mathbf{y}).$$

The associated norm

$$\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

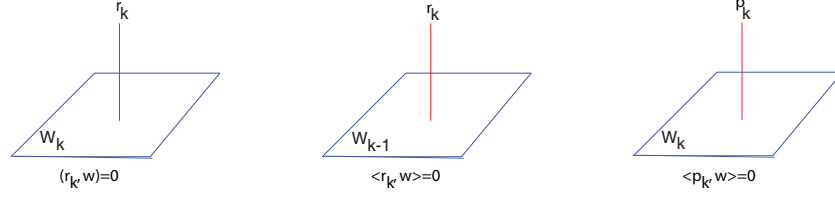
is called the  **$\mathbf{A}$ -norm** or **energy norm** of  $\mathbf{x}$ . Two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  are **orthogonal** if  $(\mathbf{x}, \mathbf{y}) = 0$  and  **$\mathbf{A}$ -orthogonal** if  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ .

Suppose  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  is an initial approximation to the solution of the linear system  $\mathbf{A} \mathbf{x} = \mathbf{b}$  and let  $\mathbf{r}^{(0)} := \mathbf{b} - \mathbf{A} \mathbf{x}^{(0)}$  be the corresponding residual. We consider the **Krylov subspaces**  $\mathbb{W}_k$  of  $\mathbb{R}^n$  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

$$\mathbb{W}_0 \subset \mathbb{W}_1 \subset \mathbb{W}_2 \subset \dots \subset \mathbb{W}_n \subset \mathbb{R}^n$$



**Figure 14.10.** Orthogonality in the conjugate gradient algorithm.

with  $\dim(\mathbb{W}_k) \leq k$  for all  $k \geq 0$ . We also note that if  $\mathbf{w} \in \mathbb{W}_k$  then  $\mathbf{A}\mathbf{w} \in \mathbb{W}_{k+1}$ . This implies

$$\mathbf{r}^{(k-1)}, \mathbf{p}^{(k-1)}, \mathbf{x}^{(k)} - \mathbf{x}^{(0)} \in \mathbb{W}_k, \quad k = 1, 2, \dots \quad (14.5)$$

For since  $\mathbf{p}^{(0)} = \mathbf{r}^{(0)}$  and  $\mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \alpha_0 \mathbf{p}^{(0)}$  this holds for  $k = 1$  and it holds for any  $k \geq 1$  by induction.

**Theorem 14.9** Suppose  $\mathbf{r}^{(j)} \neq 0$  for  $j = 0, 1, \dots, k$ . Then

1.  $\{\mathbf{r}^{(0)}, \mathbf{r}^{(1)}, \dots, \mathbf{r}^{(k)}\}$  is an orthogonal basis for  $\mathbb{W}_{k+1}$ .
2.  $\{\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}\}$  is an  $\mathbf{A}$ -orthogonal basis for  $\mathbb{W}_{k+1}$ .

**Proof.** We use induction on  $k$ . Suppose  $(\mathbf{r}^{(k)}, \mathbf{r}^{(i)}) = \langle \mathbf{p}^{(k)}, \mathbf{p}^{(i)} \rangle = 0$  for  $i < k$ . For  $j = 0, 1, \dots, k$

$$\begin{aligned} (\mathbf{r}^{(k+1)}, \mathbf{r}^{(j)}) &= (\mathbf{r}^{(k)} - \alpha_k \mathbf{A}\mathbf{p}^{(k)}, \mathbf{r}^{(j)}) \\ &= (\mathbf{r}^{(k)}, \mathbf{r}^{(j)}) - \alpha_k \langle \mathbf{p}^{(k)}, \mathbf{p}^{(j)} \rangle - \beta_{j-1} \mathbf{p}^{(j-1)} \\ &= (\mathbf{r}^{(k)}, \mathbf{r}^{(j)}) - \alpha_k \langle \mathbf{p}^{(k)}, \mathbf{p}^{(j)} \rangle. \end{aligned}$$

This is zero for  $j < k$  and vanishes for  $j = k$  by the formula for  $\alpha_k$ . Since  $\mathbf{r}^{(j)}$  is nonzero and  $\mathbf{r}^{(j)} \in \mathbb{W}_{j+1} \subset \mathbb{W}_{k+1}$  for  $j \leq k$  Claim 1. follows.

Since  $\mathbf{p}^{(j)} \in \mathbb{W}_{j+1} \subset \mathbb{W}_{k+1}$  Claim 2. will follow if we can show that

$$\mathbf{p}^{(j)} = \mathbf{r}^{(j)} - \sum_{i=0}^{j-1} \frac{\langle \mathbf{r}^{(j)}, \mathbf{p}^{(i)} \rangle}{\langle \mathbf{p}^{(i)}, \mathbf{p}^{(i)} \rangle} \mathbf{p}^{(i)}, \quad j = 0, 1, 2, \dots \quad (14.6)$$

For then  $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(k)}\}$  is the result of applying the Gram-Schmidt orthogonalization process to the linearly independent residuals  $\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(k)}\}$  using the inner product  $\langle \cdot, \cdot \rangle$  (cf. Theorem 2.57). Suppose (14.6) holds for  $j \leq k$ . Since  $\mathbf{A}\mathbf{p}^{(i)} \in \mathbb{W}_{i+2}$  and  $\mathbf{r}^{(k+1)}$  is orthogonal to  $\{\mathbf{r}^{(0)}, \mathbf{r}^{(1)}, \dots, \mathbf{r}^{(k)}\}$  and hence to any vector in  $\mathbb{W}_{k+1}$  we have  $\langle \mathbf{r}^{(k+1)}, \mathbf{p}^{(i)} \rangle = (\mathbf{r}^{(k+1)}, \mathbf{A}\mathbf{p}^{(i)}) = 0$  for  $i \leq k-1$ . Now

(14.6) follows for  $j = k + 1$  since

$$\begin{aligned}
 \mathbf{r}^{(k+1)} &= \sum_{i=0}^k \frac{\langle \mathbf{r}^{(k+1)}, \mathbf{p}^{(i)} \rangle}{\langle \mathbf{p}^{(i)}, \mathbf{p}^{(i)} \rangle} \mathbf{p}^{(i)} \\
 &= \mathbf{r}^{(k+1)} - \frac{(\mathbf{r}^{(k+1)}, \mathbf{A}\mathbf{p}^{(k)})}{\langle \mathbf{p}^{(k)}, \mathbf{p}^{(k)} \rangle} \mathbf{p}^{(k)} \\
 &= \mathbf{r}^{(k+1)} - \frac{(\mathbf{r}^{(k+1)}, \mathbf{r}^{(k)} - \mathbf{r}^{(k+1)})}{\alpha_k \langle \mathbf{p}^{(k)}, \mathbf{p}^{(k)} \rangle} \mathbf{p}^{(k)} \\
 &= \mathbf{r}^{(k+1)} + \frac{(\mathbf{r}^{(k+1)}, \mathbf{r}^{(k+1)})}{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})} \mathbf{p}^{(k)} \\
 &= \mathbf{r}^{(k+1)} + \beta_k \mathbf{p}^{(k)} = \mathbf{p}^{(k+1)}.
 \end{aligned}$$

□

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  $\{\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k-1)}\}$  Theorem 14.9 implies

$$(\mathbf{r}^{(k)}, \mathbf{w}) = \langle \mathbf{p}^{(k)}, \mathbf{w} \rangle = 0, \quad \mathbf{w} \in \mathbb{W}_k. \quad (14.7)$$

These orthogonal properties are illustrated in Figure 14.10.

**Corollary 14.11** Suppose  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n,n}$  is symmetric positive definite and  $\{\mathbf{x}^{(k)}\}$  is generated by the conjugate gradient algorithm. Then  $\mathbf{x}^{(k)} - \mathbf{x}^{(0)}$  is the best approximation to  $\mathbf{x} - \mathbf{x}^{(0)}$  in the  $\mathbf{A}$ -norm

$$\|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{A}} = \min_{\mathbf{w} \in \mathbb{W}_k} \|\mathbf{x} - \mathbf{x}^{(0)} - \mathbf{w}\|_{\mathbf{A}}. \quad (14.8)$$

**Proof.** Since

$$0 = (\mathbf{r}^{(k)}, \mathbf{w}) = (\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}^{(k)}, \mathbf{w}) = \langle \mathbf{x} - \mathbf{x}^{(0)} - (\mathbf{x}^{(k)} - \mathbf{x}^{(0)}), \mathbf{w} \rangle, \quad \mathbf{w} \in \mathbb{W}_k$$

it follows that  $\mathbf{x}^{(k)} - \mathbf{x}^{(0)}$  is the  $\mathbf{A}$ -orthogonal projection of  $\mathbf{x} - \mathbf{x}^{(0)}$  into  $\mathbb{W}_k$  and the result follows from Theorem 2.59. □

**Exercise 14.12** 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

$$\text{approximation } \mathbf{x} = \mathbf{0}. \text{ Answer: } [\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}] = \begin{bmatrix} 4 & 8 & 20 \\ 0 & -4 & -16 \\ 0 & 0 & 4 \end{bmatrix}.$$

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  $\mathbf{Ax} = \mathbf{b}$
- $\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(k-1)}$  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  $k = 1, 2, 3$
- $\{\|\mathbf{r}^{(k)}\|\}$  is monotonically decreasing

**Exercise 14.13** Study the proof of Lemma 14.21 which shows that for the Euclidean norm

$$\|\mathbf{r}^{(k+1)}\|_2 \leq \|\mathbf{r}^{(k)}\|_2, \quad k \geq 1.$$

**Exercise 14.14** Consider solving the least squares problem by using the conjugate gradient method on the normal equations  $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$ . Explain why only the following modifications in Algorithm 14.4 are necessary

1.  $r = \mathbf{A}'(\mathbf{b} - \mathbf{A}^*x)$ ;  $p = r$ ;
2.  $a = \text{rho}/(t'^*t)$ ;
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  $\mathbf{A}$ .

## 14.4 Convergence

The main result in this section is the following theorem.

**Theorem 14.15** *Suppose we apply the conjugate gradient method to a symmetric positive definite system  $\mathbf{Ax} = \mathbf{b}$ . Then the  $\mathbf{A}$ -norms of the errors satisfy*

$$\frac{\|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}^{(0)}\|_{\mathbf{A}}} = \frac{\|\mathbf{r}^{(k)}\|_2}{\|\mathbf{r}^{(0)}\|_2} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where  $\kappa = \text{cond}_2(\mathbf{A}) = \lambda_{\max}/\lambda_{\min}$  is the 2-norm condition number of  $\mathbf{A}$ .

This theorem explains what we observed in the previous section. Namely that the number of iterations is linked to  $\sqrt{\kappa}$ , the square root of the condition number of  $\mathbf{A}$ . Indeed, the following corollary gives an upper bound for the number of iterations in terms of  $\sqrt{\kappa}$ .

**Corollary 14.16** *If for some  $\epsilon > 0$  we have  $k \geq \frac{1}{2} \log(\frac{2}{\epsilon})\sqrt{\kappa}$  then  $\frac{\|\mathbf{r}^{(k)}\|_2}{\|\mathbf{r}^{(0)}\|_2} \leq \epsilon$ .*

We prove Theorem 14.15 for  $\mathbf{x}^{(0)} = \mathbf{0}$ . By Corollary 14.11  $\mathbf{x}^{(k)}$  is the best approximation to the solution  $\mathbf{x}$  in the  $\mathbf{A}$ -norm. We convert this best approximation property into a best approximation problem involving polynomials. In the following we let  $\Pi_k$  denote the class of univariate polynomials of degree  $\leq k$  with real coefficients.

**Theorem 14.17** *Suppose  $\mathbf{Ax} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{R}^{n,n}$  is symmetric positive definite with eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding orthonormal eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ . Then*

$$\|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{A}}^2 = \min_{\substack{Q \in \Pi_k \\ Q(0)=1}} \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j} Q(\lambda_j)^2, \quad (14.9)$$

where the  $\sigma_j$ 's are the coefficients when  $\mathbf{b}$  is expanded in terms of the basis of eigenvectors of  $\mathbf{A}$ ,  $\mathbf{b} = \sum_{j=1}^n \sigma_j \mathbf{u}_j$ .

**Proof.** If  $\mathbf{w} \in \mathbb{W}_k = \text{span}(\mathbf{b}, \mathbf{Ab}, \dots, \mathbf{A}^{k-1}\mathbf{b})$  then for some  $a_0, \dots, a_{k-1}$

$$\mathbf{w} = \sum_{j=0}^{k-1} a_j \mathbf{A}^j \mathbf{b} = P(\mathbf{A})\mathbf{b},$$

where

$$P(\mathbf{A}) = a_0 \mathbf{I} + a_1 \mathbf{A} + a_2 \mathbf{A}^2 + \dots + a_{k-1} \mathbf{A}^{k-1}$$

is a matrix polynomial corresponding to the ordinary polynomial  $P(t) = a_0 + a_1 t + \dots + a_{k-1} t^{k-1}$  of degree  $\leq k-1$ . Then

$$\begin{aligned} \|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 &= (\mathbf{x} - \mathbf{w}, \mathbf{A}(\mathbf{x} - \mathbf{w})) \\ &= (\mathbf{A}^{-1}(\mathbf{b} - \mathbf{Aw}), \mathbf{b} - \mathbf{Aw}) \\ &= (\mathbf{A}^{-1}(\mathbf{b} - \mathbf{AP}(\mathbf{A})\mathbf{b}), \mathbf{b} - \mathbf{AP}(\mathbf{A})\mathbf{b}) \\ &= (\mathbf{A}^{-1}Q(\mathbf{A})\mathbf{b}, Q(\mathbf{A})\mathbf{b}), \end{aligned} \quad (14.10)$$

where  $Q(\mathbf{A}) = I - \mathbf{A}P(\mathbf{A})$  is another matrix polynomial corresponding to the polynomial  $Q(t) = 1 - tP(t)$ . Observe that  $Q \in \Pi_k$  and  $Q(0) = 1$ . Using the eigenvector expansion for  $\mathbf{b}$  we obtain

$$Q(\mathbf{A})\mathbf{b} = \sum_{j=1}^n \sigma_j Q(\mathbf{A})\mathbf{u}_j = \sum_{j=1}^n \sigma_j Q(\lambda_j)\mathbf{u}_j. \quad (14.11)$$

The last equality follows from (5.3). We also have

$$\mathbf{A}^{-1}Q(\mathbf{A})\mathbf{u}_j = Q(\lambda_j)\mathbf{A}^{-1}\mathbf{u}_j = \frac{Q(\lambda_j)}{\lambda_j}\mathbf{u}_j. \quad (14.12)$$

Combining (14.10), (14.11), and (14.12) we find

$$\begin{aligned} \|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 &= (\mathbf{A}^{-1}Q(\mathbf{A})\mathbf{b}, Q(\mathbf{A})\mathbf{b}) \\ &= \left( \sum_{i=1}^n \sigma_i \frac{Q(\lambda_i)}{\lambda_i} \mathbf{u}_i, \sum_{j=1}^n \sigma_j Q(\lambda_j) \mathbf{u}_j \right) \\ &= \sum_{i,j} \sigma_i \sigma_j \frac{Q(\lambda_i)Q(\lambda_j)}{\lambda_i} (\mathbf{u}_i, \mathbf{u}_j) = \sum_{j=1}^n \sigma_j^2 \frac{Q(\lambda_j)^2}{\lambda_j}. \end{aligned}$$

Minimizing over  $\mathbf{w}$  is the same as minimizing over all  $Q \in \Pi_k$  with  $Q(0) = 1$  and the proof is complete.  $\square$

We will use the following weaker form of Theorem 14.17 to estimate the rate of convergence.

**Corollary 14.18** *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.** In the proof of Theorem 14.17 we showed that to each  $\mathbf{w} \in \mathbb{W}_k$  there corresponds a polynomial  $Q \in \Pi_k$  with  $Q(0) = 1$  such that

$$\|\mathbf{x} - \mathbf{w}\|_{\mathbf{A}}^2 = \sum_{j=1}^n \sigma_j^2 \frac{Q(\lambda_j)^2}{\lambda_j}.$$

Taking  $\mathbf{w} = \mathbf{x}^{(0)}$  we find  $\|\mathbf{x} - \mathbf{x}^{(0)}\|_{\mathbf{A}}^2 = \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j}$ . Therefore, by Theorem 14.17 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.  $\square$



We will apply Corollary 14.18 with  $Q(x)$  a suitably shifted and normalized version of the Chebyshev polynomial. Recall that the Chebyshev polynomial of degree  $n$  is defined recursively by

$$T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t), \quad n \geq 1$$

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 14.19** 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

$$x_{n+1} - 2tx_n + x_{n-1} = 0 \text{ for } n \geq 1, \text{ with } x_0 = 1, x_1 = t. \quad (14.13)$$

To find  $x_n$  we insert  $x_n = z^n$  into (14.13) and obtain  $z^{n+1} - 2tz^n + z^{n-1} = 0$  or  $z^2 - 2tz + 1 = 0$ . Let  $z_1$  and  $z_2$  be the roots of this quadratic equation. Then  $z_1^n, z_2^n$  and more generally  $c_1 z_1^n + c_2 z_2^n$  are solutions of (14.13) 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  $c_1 = c_2 = \frac{1}{2}$ . Solving the quadratic equation we find  $z_1 = \alpha := t + \sqrt{t^2 - 1}$  and  $z_2 = \alpha^{-1}$ . It follows that  $x_n = T_n(t) = \frac{1}{2}(\alpha^n + \alpha^{-n})$  which is the same as 2.  $\square$

**Exercise 14.20** Show that

$$T_n(t) = \cosh(n \operatorname{arccosh} t) \text{ for } |t| \geq 1,$$

where  $\operatorname{arccosh}$  is the inverse function of  $\cosh x := (e^x + e^{-x})/2$ .

**Proof of Theorem 14.15.**

**Proof.** Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of  $\mathbf{A}$  and let  $k \geq 0$ . We apply Corollary 14.18 with  $a = \min \lambda_j$ ,  $b = \max \lambda_j$ , and

$$Q(x) = T_k\left(\frac{b+a-2x}{b-a}\right) / T_k\left(\frac{b+a}{b-a}\right). \quad (14.14)$$

Note that  $Q$  is admissible since  $Q \in \Pi_k$  with  $Q(0) = 1$ . By Lemma 14.19

$$\max_{a \leq x \leq b} \left| T_k \left( \frac{b+a-2x}{b-a} \right) \right| = \max_{-1 \leq t \leq 1} |T_k(t)| = 1. \quad (14.15)$$

Moreover with  $t = (b+a)/(b-a)$  we have

$$t + \sqrt{t^2 - 1} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}, \quad \kappa = b/a.$$

Thus again by Lemma 14.19 we find

$$T_k \left( \frac{b+a}{b-a} \right) = \frac{1}{2} \left[ \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^k + \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \right] \geq \frac{1}{2} \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^k. \quad (14.16)$$

Using (14.15) and (14.16) in (14.14) completes the proof.  $\square$

#### Proof of Corollary 14.16.

**Proof.** The inequality

$$\frac{x-1}{x+1} < e^{-2/x} \quad \text{for } x > 1 \quad (14.17)$$

follows from the familiar series expansion of the exponential function. Indeed, with  $y = 1/x$  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 (14.17) follows. By Theorem 14.15 we then find

$$\frac{\|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}^{(0)}\|_{\mathbf{A}}} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k < 2e^{-2k/\sqrt{\kappa}}.$$

Solving the inequality  $2e^{-2k/\sqrt{\kappa}} < \epsilon$  leads immediately to the result.  $\square$

The Euclidian norm of the residuals  $\mathbf{b} - \mathbf{Ax}^{(k)}$  in the conjugate gradient iteration decreases monotonically (cf. Exercise 14.13). The following lemma shows that the Euclidian norm of the errors  $\mathbf{x} - \mathbf{x}^{(k)}$  are also monotonically decreasing.

**Lemma 14.21** *Let  $\mathbf{x}$  be the exact solution of  $\mathbf{Ax} = \mathbf{b}$ , define  $\epsilon_k = \mathbf{x} - \mathbf{x}^{(k)}$  for  $k \geq 0$  and let  $\|\cdot\|$  denote the Euclidian vector norm. If  $\mathbf{p}^{(j)} \neq 0$  for  $j \leq k$  then  $\|\epsilon_{k+1}\|_2 < \|\epsilon_k\|_2$ . More precisely,*

$$\|\epsilon_{k+1}\|_2^2 = \|\epsilon_k\|_2^2 - \frac{\|\mathbf{p}^{(k)}\|_2^2}{\|\mathbf{p}^{(k)}\|_{\mathbf{A}}^2} (\|\epsilon_{k+1}\|_{\mathbf{A}}^2 + \|\epsilon_k\|_{\mathbf{A}}^2). \quad (14.18)$$

**Proof.** Set

$$\rho_j := \|\mathbf{r}^{(j)}\|_2^2 \quad \text{and} \quad \pi_j := \|\mathbf{p}^{(j)}\|_{\mathbf{A}}^2, \quad j \geq 0$$

and let  $m$  be the smallest integer such that  $\|\boldsymbol{\epsilon}_m\|_2 = 0$ . Since  $\mathbf{p}^{(j)} \neq 0$  for  $j \leq k$  we have  $\dim \mathbb{W}_{k+1} = k+1$  which implies that  $\mathbf{r}^{(k)} \neq 0$  and hence  $m > k$ . For  $j < m$

$$\mathbf{x}^{(j+1)} = \mathbf{x}^{(j)} + \alpha_j \mathbf{p}^{(j)} = \mathbf{x}^{(j-1)} + \alpha_{j-1} \mathbf{p}^{(j-1)} + \alpha_j \mathbf{p}^{(j)} = \dots = \mathbf{x}^{(0)} + \sum_{i=0}^j \alpha_i \mathbf{p}^{(i)}$$

so that

$$\boldsymbol{\epsilon}_j = \mathbf{x}^{(m)} - \mathbf{x}^{(j)} = \sum_{i=j}^{m-1} \alpha_i \mathbf{p}^{(i)}, \quad \alpha_i = \frac{\rho_i}{\pi_i}. \quad (14.19)$$

For  $j > k$

$$(\mathbf{p}^{(j)}, \mathbf{p}^{(k)}) = (\mathbf{r}^{(j)} + \beta_{j-1} \mathbf{p}^{(j-1)}, \mathbf{p}^{(k)}) = \beta_{j-1} (\mathbf{p}^{(j-1)}, \mathbf{p}^{(k)}) = \dots = \beta_{j-1} \dots \beta_k (\mathbf{p}^{(k)}, \mathbf{p}^{(k)})$$

and since  $\beta_{j-1} \dots \beta_k = \rho_j / \rho_k$  we obtain

$$(\mathbf{p}^{(j)}, \mathbf{p}^{(k)}) = \frac{\rho_j}{\rho_k} (\mathbf{p}^{(k)}, \mathbf{p}^{(k)}), \quad j \geq k. \quad (14.20)$$

By  $\mathbf{A}$ -orthogonality and (14.19)

$$\|\boldsymbol{\epsilon}_j\|_{\mathbf{A}}^2 = \left\langle \sum_{i=j}^{m-1} \alpha_i \mathbf{p}^{(i)}, \sum_{i=j}^{m-1} \alpha_i \mathbf{p}^{(i)} \right\rangle = \sum_{i=j}^{m-1} \alpha_i^2 \pi_i = \sum_{i=j}^{m-1} \frac{\rho_i^2}{\pi_i}. \quad (14.21)$$

Now

$$\begin{aligned} \|\boldsymbol{\epsilon}_k\|_2^2 &= \|\boldsymbol{\epsilon}_{k+1} + \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2 = \|\boldsymbol{\epsilon}_{k+1} + \alpha_k \mathbf{p}^{(k)}\|_2^2 \\ &= \|\boldsymbol{\epsilon}_{k+1}\|_2^2 + \alpha_k (2(\mathbf{p}^{(k)}, \boldsymbol{\epsilon}_{k+1}) + \alpha_k \|\mathbf{p}^{(k)}\|_2^2). \end{aligned} \quad (14.22)$$

and moreover

$$\begin{aligned} \alpha_k (2(\mathbf{p}^{(k)}, \boldsymbol{\epsilon}_{k+1}) + \alpha_k \|\mathbf{p}^{(k)}\|_2^2) &\stackrel{(14.19)}{=} \alpha_k \left( 2 \sum_{j=k+1}^{m-1} \alpha_j (\mathbf{p}^{(j)}, \mathbf{p}^{(k)}) + \alpha_k \|\mathbf{p}^{(k)}\|_2^2 \right) \\ &\stackrel{(14.20)}{=} \alpha_k \left( 2 \sum_{j=k+1}^{m-1} \alpha_j \frac{\rho_j}{\rho_k} \|\mathbf{p}^{(k)}\|_2^2 + \alpha_k \|\mathbf{p}^{(k)}\|_2^2 \right) = \frac{\|\mathbf{p}^{(k)}\|_2^2}{\pi_k} \left( \sum_{j=k}^{m-1} \frac{\rho_j^2}{\pi_j} + \sum_{j=k+1}^{m-1} \frac{\rho_j^2}{\pi_j} \right) \\ &\stackrel{(14.21)}{=} \frac{\|\mathbf{p}^{(k)}\|_2^2}{\pi_k} (\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 + \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2). \end{aligned}$$

Inserting this in (14.22) proves the lemma.  $\square$



## Chapter 15

# Minimization and Preconditioning

We continue the study of the conjugate gradient method. Recall that the rate of convergence depends on the square root of the condition number of the coefficient matrix. For problems with a large condition number the convergence can be slow. For such problems a *preconditioned conjugate gradient method* is often used, and we consider this method here.

The conjugate gradient method can also be used as a minimization algorithm and we start discussing some aspects of minimization of quadratic functions.

### 15.1 Minimization

If  $\mathbf{A}$  is symmetric positive definite then the quadratic function

$$Q(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n$$

has a unique global minimum  $\mathbf{x}^* \in \mathbb{R}^n$  which is found by setting the gradient  $\mathbf{g}(\mathbf{x}) := \nabla Q(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}$  equal to zero (cf. Appendix C). So we find the minimum as a solution of the linear system  $\mathbf{A} \mathbf{x}^* = \mathbf{b}$ . We see also that the gradient of  $Q(\mathbf{x})$  is equal to the residual of  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , i.e.  $\mathbf{g}(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b} = \mathbf{r}(\mathbf{x})$ .

A general class of minimization algorithms for  $Q$  is given as follows:

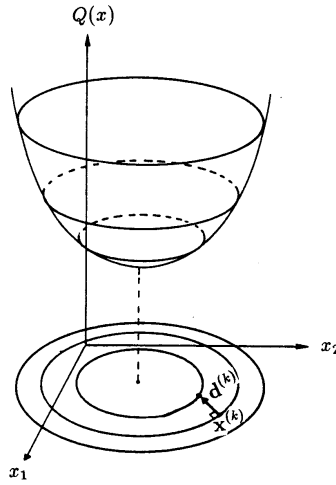
1. Choose  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ .
2. For  $k = 0, 1, 2, \dots$ 
  - (a) Choose a “search direction”  $\mathbf{d}^{(k)}$ .
  - (b) Choose a “step length”  $\sigma_k$ .
  - (c)  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \sigma_k \mathbf{d}^{(k)}$ .

We would like to generate a sequence  $\{\mathbf{x}^{(k)}\}$  of points such that  $\{\mathbf{x}^{(k)}\}$  converges quickly to the minimum  $\mathbf{x}$  of  $Q$ .

We can think of  $Q(\mathbf{x})$  as a paraboloid. To see this, let  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T$ , where  $\mathbf{U}$  is orthogonal and  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is diagonal, be the spectral decomposition of  $\mathbf{A}$  and change variables to  $\mathbf{v} = [v_1, \dots, v_n] := \mathbf{U}^T \mathbf{x}$  and  $\mathbf{c} := \mathbf{U}^T \mathbf{b} = [c_1, \dots, c_n]$ . Then

$$Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{U} \mathbf{D} \mathbf{U}^T \mathbf{x} - \mathbf{b}^T \mathbf{U} \mathbf{U}^T \mathbf{x} = \frac{1}{2} \mathbf{v}^T \mathbf{D} \mathbf{v} - \mathbf{c}^T \mathbf{v} = \frac{1}{2} \sum_{j=1}^n \lambda_j v_j^2 - \sum_{j=1}^n c_j v_j.$$

In particular for  $n = 2$  we have  $z := \frac{1}{2} \lambda_1 v_1^2 + \frac{1}{2} \lambda_2 v_2^2 - c_1 v_1 - c_2 v_2$  and since  $\lambda_1$  and  $\lambda_2$  are positive this is the equation for a paraboloid in  $(v_1, v_2, z)$  space as shown in the following figure.



Suppose  $\mathbf{x}^{(k)} \approx \mathbf{x}^*$ . To find a better approximation to the minimum we choose a search direction  $\mathbf{d}^{(k)}$  and go from  $\mathbf{x}^{(k)}$  along  $\mathbf{d}^{(k)}$  a certain distance determined by  $\sigma_k$ . To see how  $\sigma_k$  and  $\mathbf{d}^{(k)}$  should be chosen, we note that

$$Q(\mathbf{x}^{(k+1)}) = Q(\mathbf{x}^{(k)}) + \sigma_k \langle \mathbf{d}^{(k)}, \mathbf{r}^{(k)} \rangle + \frac{1}{2} \sigma_k^2 \langle \mathbf{d}^{(k)}, \mathbf{d}^{(k)} \rangle, \quad (15.1)$$

where  $\mathbf{r}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}$ . Since  $\mathbf{A}$  is symmetric positive definite, we have  $\sigma_k^2 \langle \mathbf{d}^{(k)}, \mathbf{d}^{(k)} \rangle > 0$  for all nonzero  $\sigma_k$  and  $\mathbf{d}^{(k)}$ . In order to make  $Q(\mathbf{x}^{(k+1)})$  smaller than  $Q(\mathbf{x}^{(k)})$ , we must at least pick  $\sigma_k$  and  $\mathbf{d}^{(k)}$  such that  $\sigma_k \langle \mathbf{d}^{(k)}, \mathbf{r}^{(k)} \rangle < 0$ . For such a direction we can determine the step length  $\sigma_k = \sigma_k^*$  such that  $Q(\mathbf{x}^{(k+1)})$  is as small as possible, i.e.

$$Q(\mathbf{x}^{(k+1)}) = \min_{\sigma \in \mathbb{R}} Q(\mathbf{x}^{(k)} + \sigma \mathbf{d}^{(k)}).$$

Differentiating with respect to  $\sigma_k$  in (15.1) and setting the right-hand side equal to zero, we find

$$\sigma_k^* := - \frac{\langle \mathbf{d}^{(k)}, \mathbf{r}^{(k)} \rangle}{\langle \mathbf{d}^{(k)}, \mathbf{d}^{(k)} \rangle}. \quad (15.2)$$

We find  $\frac{\partial^2 Q}{\partial \sigma_k^2} = \langle \mathbf{d}^{(k)}, \mathbf{d}^{(k)} \rangle > 0 \implies Q(\mathbf{x}^{(k)} + \sigma_k^* \mathbf{d}^{(k)}) = \min_{\sigma \in \mathbb{R}} Q(\mathbf{x}^{(k)} + \sigma \mathbf{d}^{(k)})$  and  $\sigma_k^*$  is called **optimal** with respect to  $\mathbf{d}^{(k)}$ .

In the method of **Steepest Descent** we choose  $\mathbf{d}^{(k)} = -\mathbf{r}^{(k)}$  and  $\sigma_k = \sigma_k^*$  so that

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \frac{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})}{\langle \mathbf{r}^{(k)}, \mathbf{r}^{(k)} \rangle} \mathbf{r}^{(k)}, \quad k = 0, 1, 2, \dots \quad (15.3)$$

The method of steepest descent will converge very slowly if  $\mathbf{A}$  is ill-conditioned. For then the ratio of the smallest and biggest eigenvalue becomes large and the paraboloid becomes very distorted. In this case the residuals need not point in the direction of the minimum. It can be shown that the number of iterations is proportional to the two-norm condition number  $\lambda_{\max}/\lambda_{\min}$  of  $\mathbf{A}$ .

Consider now the conjugate gradient method. Here we choose  $\mathbf{A}$ -orthogonal search directions  $\mathbf{d}^{(k)} = -\mathbf{p}^{(k)}$ . Since by (14.1)  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$  where  $\alpha_k = (\mathbf{p}^{(k)}, \mathbf{r}^{(k)})/(\mathbf{p}^{(k)}, \mathbf{A}\mathbf{p}^{(k)})$ , we see that the step length  $-\alpha_k$  is optimal with respect to  $-\mathbf{p}^{(k)}$ . Moreover the gradients  $\{\mathbf{r}^{(k)}\}$  are orthogonal. It can also be shown that

$$Q(\mathbf{x}^{(k+1)}) = \min_{w \in \mathcal{W}_{k+1}} Q(\mathbf{x}^{(0)} + w) \quad (15.4)$$

and in the next section we show that the number of iterations is proportional to the square root of the two-norm condition number of  $\mathbf{A}$ . So the conjugate gradient minimization algorithm converges much faster than the method of steepest descent for problems where the ratio  $\lambda_{\max}/\lambda_{\min}$  is large.

Conjugate gradient like algorithms can be used to minimize more general functions than  $Q$ , see [15].

**Exercise 15.1** Show that  $(\mathbf{r}^{(k)}, \mathbf{r}^{(k+1)}) = 0$  in the method of steepest descent. Does this mean that all the residuals are orthogonal?

**Exercise 15.2** Let  $Q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} - 2\mathbf{b}^T \mathbf{x}$  have a minimum at  $\mathbf{x}^* \in \mathbb{R}^n$ .

a) Show that  $Q(\mathbf{x}) = \|\mathbf{x}^* - \mathbf{x}\|_{\mathbf{A}}^2 - \|\mathbf{x}^*\|_{\mathbf{A}}^2$  for any  $\mathbf{x} \in \mathbb{R}^n$ .

b) Show (15.4).

## 15.2 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 **pre-conditioning**. 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  $\text{cond}_2(\mathbf{A})$ . 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  $\mathbf{B}$  is symmetric positive definite then we can apply CG to a symmetrized system and then transform the recurrence formulae to an iterative method for the original system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . This iterative method is

known as the **pre-conditioned conjugate gradient method**. We shall see that the convergence properties of this method is determined by the eigenvalues of  $\mathbf{BA}$ .

Suppose  $\mathbf{B}$  is symmetric positive definite. By Theorem 7.33 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 never be computed). Now

$$\mathbf{BAx} = \mathbf{Bb} \Leftrightarrow \mathbf{C}^T (\mathbf{CAC}^T) \mathbf{C}^{-T} \mathbf{x} = \mathbf{C}^T \mathbf{Cb} \Leftrightarrow (\mathbf{CAC}^T) \mathbf{y} = \mathbf{Cb}, \text{ \& } \mathbf{x} = \mathbf{C}^T \mathbf{y}.$$

We have 3 linear systems

$$\mathbf{Ax} = \mathbf{b} \tag{15.5}$$

$$\mathbf{BAx} = \mathbf{Bb} \tag{15.6}$$

$$(\mathbf{CAC}^T) \mathbf{y} = \mathbf{Cb}, \text{ \& } \mathbf{x} = \mathbf{C}^T \mathbf{y}. \tag{15.7}$$

Note that (15.5) and (15.7) 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 (15.7) 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  $\mathbf{z}^{(k)} = \mathbf{CAC}^T \mathbf{y}^{(k)} - \mathbf{Cb}$  we obtain the following from (14.2), (14.3), and (14.4).

$$\begin{aligned} \mathbf{y}^{(k+1)} &= \mathbf{y}^{(k)} + \alpha_k \mathbf{q}^{(k)}, \quad \alpha_k = (\mathbf{z}^{(k)}, \mathbf{z}^{(k)}) / (\mathbf{q}^{(k)}, (\mathbf{CAC}^T) \mathbf{q}^{(k)}), \\ \mathbf{z}^{(k+1)} &= \mathbf{z}^{(k)} + \alpha_k (\mathbf{CAC}^T) \mathbf{q}^{(k)}, \\ \mathbf{q}^{(k+1)} &= \mathbf{z}^{(k+1)} + \beta_k \mathbf{q}^{(k)}, \quad \beta_k = (\mathbf{z}^{(k+1)}, \mathbf{z}^{(k+1)}) / (\mathbf{z}^{(k)}, \mathbf{z}^{(k)}). \end{aligned}$$

With

$$\mathbf{x}^{(k)} := \mathbf{C}^T \mathbf{y}^{(k)}, \quad \mathbf{p}^{(k)} := \mathbf{C}^T \mathbf{q}^{(k)}, \quad \mathbf{s}^{(k)} := \mathbf{C}^T \mathbf{z}^{(k)}, \quad \mathbf{r}^{(k)} := \mathbf{C}^{-1} \mathbf{z}^{(k)} \tag{15.8}$$

this can be transformed into

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad \alpha_k = (\mathbf{s}^{(k)}, \mathbf{r}^{(k)}) / (\mathbf{p}^{(k)}, \mathbf{p}^{(k)}), \tag{15.9}$$

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} + \alpha_k \mathbf{Ap}^{(k)}, \tag{15.10}$$

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} + \alpha_k \mathbf{BAp}^{(k)}, \tag{15.11}$$

$$\mathbf{p}^{(k+1)} = \mathbf{s}^{(k+1)} + \beta_k \mathbf{p}^{(k)}, \quad \beta_k = (\mathbf{s}^{(k+1)}, \mathbf{r}^{(k+1)}) / (\mathbf{s}^{(k)}, \mathbf{r}^{(k)}). \tag{15.12}$$

Here  $\mathbf{x}^{(k)}$  will be an approximation to the solution  $\mathbf{x}$  of  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{r}^{(k)} = \mathbf{Ax}^{(k)} - \mathbf{b}$  is the residual in the original system and  $\mathbf{s}^{(k)} = \mathbf{BAx}^{(k)} - \mathbf{Bb}$  is the residual in the preconditioned system. This follows since by (15.8)

$$\mathbf{r}^{(k)} = \mathbf{C}^{-1} \mathbf{z}^{(k)} = \mathbf{C}^{-1} \mathbf{CAC}^T \mathbf{y}^{(k)} - \mathbf{b} = \mathbf{Ax}^{(k)} - \mathbf{b}$$

and  $\mathbf{s}^{(k)} = \mathbf{C}^T \mathbf{z}^{(k)} = \mathbf{C}^T \mathbf{C} \mathbf{r}^{(k)} = \mathbf{Br}^{(k)}$ . We now have the following preconditioned conjugate gradient algorithm for obtaining an approximation  $\mathbf{x}^{(k)}$  to the solution of a symmetric positive definite system  $\mathbf{Ax} = \mathbf{b}$ .



**Algorithm 15.3 (Preconditioned Conjugate Gradient Algorithm)**

1. Choose a starting vector  $\mathbf{x}^{(0)}$  (for example  $\mathbf{x}^{(0)} = \mathbf{0}$ )
2.  $\mathbf{r}_0 = \mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}$ ,  $\mathbf{p}_0 = \mathbf{s}_0 = \mathbf{B}\mathbf{r}_0$
3.  $\rho_0 = (\mathbf{s}_0, \mathbf{r}_0)$ ;  $k = 0$
4. while  $\sqrt{\rho_k/\rho_0} > \epsilon$  &  $k < kmax$ 
  - 4.1a  $\mathbf{t}_k = \mathbf{A}\mathbf{p}^{(k)}$
  - 4.1b  $\mathbf{w}_k = \mathbf{B}\mathbf{t}_k$
  - 4.2  $\alpha_k = \rho_k / (\mathbf{p}^{(k)}, \mathbf{t}_k)$
  - 4.3  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$
  - 4.4a  $\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} + \alpha_k \mathbf{t}_k$  ( $\mathbf{r}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}$ )
  - 4.4b  $\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} + \alpha_k \mathbf{w}_k$  ( $\mathbf{s}^{(k)} = \mathbf{B}\mathbf{A}\mathbf{x}^{(k)} - \mathbf{B}\mathbf{b}$ )
  - 4.5  $\rho_{k+1} = (\mathbf{s}^{(k+1)}, \mathbf{r}^{(k+1)})$
  - 4.6  $\mathbf{p}^{(k+1)} = \mathbf{s}^{(k+1)} + \frac{\rho_{k+1}}{\rho_k} \mathbf{p}^{(k)}$
  - 4.7  $k = k + 1$

This algorithm is quite similar to Algorithm 14.4. The main additional work is contained in statement 4.1b. We'll discuss this further in connection with an example.

We have the following convergence result for this algorithm.

**Theorem 15.4** *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 15.3 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 15.3 is equivalent to solving (15.7) by the conjugate gradient method Theorem 14.15 implies that

$$\frac{\|\mathbf{y} - \mathbf{y}^{(k)}\|_{\mathbf{C}\mathbf{A}\mathbf{C}^T}}{\|\mathbf{y} - \mathbf{y}_0\|_{\mathbf{C}\mathbf{A}\mathbf{C}^T}} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where  $\mathbf{y}^{(k)}$  is the conjugate gradient approximation to the solution  $\mathbf{y}$  of (15.7) and  $\kappa$  is the ratio of the largest and smallest eigenvalue of  $\mathbf{C}\mathbf{A}\mathbf{C}^T$ . Since  $\mathbf{B}\mathbf{A}$  and  $\mathbf{C}\mathbf{A}\mathbf{C}^T$  are similar this is the same as the  $\kappa$  in the theorem. By (15.8) we have

$$\begin{aligned} \|\mathbf{y} - \mathbf{y}^{(k)}\|_{\mathbf{C}\mathbf{A}\mathbf{C}^T}^2 &= (\mathbf{y} - \mathbf{y}^{(k)}, \mathbf{C}\mathbf{A}\mathbf{C}^T(\mathbf{y} - \mathbf{y}^{(k)})) \\ &= (\mathbf{C}^T(\mathbf{y} - \mathbf{y}^{(k)}), \mathbf{A}\mathbf{C}^T(\mathbf{y} - \mathbf{y}^{(k)})) = \|\mathbf{x} - \mathbf{x}^{(k)}\|_{\mathbf{A}}^2 \end{aligned}$$

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 flops, ideally  $O(n)$  for both.

### 15.3 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

$$-\nabla^2 u = -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f(x, y) \quad \text{for } (x, y) \in \Omega = (0, 1)^2 \quad (15.13)$$

$$u = 0 \text{ on } \partial\Omega,$$

where  $f$  is a given function,  $\Omega$  is the unit square in the plane, and  $\partial\Omega$  is the boundary of  $\Omega$ . For numerical solution we have the **discrete Poisson problem** which can either be written as a matrix equation

$$\begin{aligned} h^2 f_{j,k} &= 4v_{j,k} - v_{j-1,k} - v_{j+1,k} - v_{j,k-1} - v_{j,k+1}, & j, k &= 1, \dots, m \\ v_{0,k} &= v_{m+1,k} = v_{j,0} = v_{j,m+1} = 0, & j, k &= 0, 1, \dots, m+1, \end{aligned}$$

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

$$\begin{aligned} a_{ii} &= 4, & i &= 1, \dots, n \\ a_{i+1,i} = a_{i,i+1} &= -1, & i &= 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m \\ a_{i+m,i} = a_{i,i+m} &= -1, & i &= 1, \dots, n-m \\ a_{ij} &= 0, & & \text{otherwise.} \end{aligned}$$

#### 15.3.1 A Banded Matrix

Consider the problem

$$\begin{aligned} -\frac{\partial}{\partial x} \left( c(x, y) \frac{\partial u}{\partial x} \right) - \frac{\partial}{\partial y} \left( c(x, y) \frac{\partial u}{\partial y} \right) &= f(x, y) & (x, y) \in \Omega = (0, 1)^2 \\ u(x, y) &= 0 & (x, y) \in \partial\Omega. \end{aligned} \quad (15.14)$$

Here  $\Omega$  is the open unit square while  $\partial\Omega$  is the boundary of  $\Omega$ . The functions  $f$  and  $c$  are given and we seek a function  $u = u(x, y)$  such that (15.14) holds. We assume that  $c$  and  $f$  are defined and continuous on  $\Omega$  and that  $c(x, y) > 0$  for all  $(x, y) \in \Omega$ . The problem (15.14) 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 (15.14) 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

$$\begin{aligned} \frac{\partial}{\partial t} \left( f(t) \frac{\partial}{\partial t} g(t) \right) &\approx \left( f(t + \frac{h}{2}) \frac{\partial}{\partial t} g(t + h/2) - f(t - \frac{h}{2}) \frac{\partial}{\partial t} g(t - \frac{h}{2}) \right) / h \\ &\approx \left( f(t + \frac{h}{2}) (g(t + h) - g(t)) - f(t - \frac{h}{2}) (g(t) - g(t - h)) \right) / h^2 \end{aligned}$$

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 (15.14) turns out to be

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= h^2 f_{j,k} & j, k = 1, \dots, m \\ v_{j,k} &= 0 & j = 0, m+1 \text{ all } k \text{ or } k = 0, m+1 \text{ all } j, \end{aligned} \quad (15.15)$$

where

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= (c_{j,k-\frac{1}{2}} + c_{j-\frac{1}{2},k} + c_{j+\frac{1}{2},k} + c_{j,k+\frac{1}{2}}) v_{j,k} \\ &\quad - c_{j,k-\frac{1}{2}} v_{j,k-1} - c_{j-\frac{1}{2},k} v_{j-1,k} - c_{j+\frac{1}{2},k} v_{j+1,k} - c_{j,k+\frac{1}{2}} v_{j,k+1} \end{aligned} \quad (15.16)$$

and  $f_{j,k} = f(jh, kh)$ .

As before we let  $\mathbf{V} = (v_{j,k}) \in \mathbb{R}^{m,m}$  and  $\mathbf{F} = (f_{j,k}) \in \mathbb{R}^{m,m}$ . 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

$$\begin{aligned} a_{i,i} &= c_{j_i,k_i-\frac{1}{2}} + c_{j_i-\frac{1}{2},k_i} + c_{j_i+\frac{1}{2},k_i} + c_{j_i,k_i+\frac{1}{2}}, & i = 1, 2, \dots, n \\ a_{i+1,i} &= a_{i,i+1} = -c_{j_i+\frac{1}{2},k_i}, & i \bmod m \neq 0 \\ a_{i+m,i} &= a_{i,i+m} = -c_{j_i,k_i+\frac{1}{2}}, & i = 1, 2, \dots, n-m \\ a_{i,j} &= 0 & \text{otherwise,} \end{aligned} \quad (15.17)$$

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 (8.15). But we can show that  $\mathbf{A}$  is symmetric and it is positive definite as long as the function  $c$  is positive on  $\Omega$ . Recall that a matrix  $\mathbf{A}$  is positive definite if  $x^T \mathbf{A} x > 0$  for all  $x \neq 0$ .

**Theorem 15.1** *If  $c(x, y) > 0$  for  $(x, y) \in \Omega$  then the matrix  $\mathbf{A}$  given by (15.17) is symmetric positive definite.*

**Proof.**

To each  $x \in \mathbb{R}^n$  there corresponds a matrix  $\mathbf{V} \in \mathbb{R}^{m,m}$  such that  $x = \text{vec}(\mathbf{V})$ . We claim that

$$x^T \mathbf{A} x = \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} (v_{j,k+1} - v_{j,k})^2 + \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} (v_{j+1,k} - v_{j,k})^2, \quad (15.18)$$

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  $\Omega$  for the values of  $j, k$  in the sums it follows that they are positive and from (15.18) we see that  $x^T \mathbf{A} x \geq 0$  for all  $x \in \mathbb{R}^n$ . Moreover if  $x^T \mathbf{A} x = 0$  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 (15.18). From the connection between (15.16) and (15.17) we have

$$\begin{aligned} x^T \mathbf{A} x &= \sum_{j=1}^m \sum_{k=1}^m -(\mathbf{P}_h v)_{j,k} v_{j,k} \\ &= \sum_{j=1}^m \sum_{k=1}^m \left( c_{j,k-\frac{1}{2}} v_{j,k}^2 + c_{j-\frac{1}{2},k} v_{j,k}^2 + c_{j+\frac{1}{2},k} v_{j,k}^2 + c_{j,k+\frac{1}{2}} v_{j,k}^2 \right. \\ &\quad - c_{j,k-\frac{1}{2}} v_{j,k-1} v_{j,k} - c_{j,k+\frac{1}{2}} v_{j,k} v_{j,k+1} \\ &\quad \left. - c_{j-\frac{1}{2},k} v_{j-1,k} v_{j,k} - c_{j+\frac{1}{2},k} v_{j,k} v_{j+1,k} \right). \end{aligned}$$

Using the homogenous boundary conditions we have

$$\begin{aligned} \sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k}^2 &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1}^2, \\ \sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k-1} v_{j,k} &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1} v_{j,k}, \\ \sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j,k}^2 &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k}^2, \\ \sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j-1,k} v_{j,k} &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k} v_{j,k}. \end{aligned}$$

It follows that

$$\begin{aligned} x^T \mathbf{A} x &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} (v_{j,k}^2 + v_{j,k+1}^2 - 2v_{j,k} v_{j,k+1}) \\ &\quad + \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} (v_{j,k}^2 + v_{j+1,k}^2 - 2v_{j,k} v_{j+1,k}) \end{aligned}$$

$n$	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 15.2.** *The number of iterations  $K$  (no preconditioning) and  $K_{pre}$  (with preconditioning) for the problem (15.14) using the discrete Poisson problem as a preconditioner.*

and (15.18) follows.  $\square$   $\square$

### 15.3.2 Preconditioning

Consider solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A}$  is given by (15.17) 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 to find  $\mathbf{x}$ . Since the bandwidth of  $\mathbf{A}$  is  $m = \sqrt{n}$  both of these methods require  $O(n^2)$  flops for large  $n$ .

If we choose  $c(x, y) \equiv 1$  in (15.14), we get the Poisson problem (15.13). 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 (15.17). This suggests using  $\mathbf{B} = \mathbf{A}_p^{-1}$ , the inverse of the discrete Poisson matrix as a preconditioner for the system (15.15).

Consider Algorithm 15.3. With this preconditioner Statement 4.1b can be written  $\mathbf{A}_p \mathbf{w}_k = \mathbf{t}_k$ .

In Section 9.2 we developed a Simple fast Poisson Solver, Cf. Algorithm 9.1. This method can be utilized to solve  $\mathbf{A}_p \mathbf{w}_k = \mathbf{t}_k$ .

Consider the specific problem where

$$c(x, y) = e^{-x+y} \text{ and } f(x, y) = 1.$$

We have used Algorithm 14.4 (conjugate gradient without preconditioning), and Algorithm 15.3 (conjugate gradient with preconditioning) to solve the problem (15.14). We used  $\mathbf{x}^{(0)} = 0$  and  $\epsilon = 10^{-8}$ . The results are shown in Table 15.2.

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 flops 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 drastically.

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 (15.14) can be bounded independently of  $n$ .

**Theorem 15.3** *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  $\mathbf{BA} = \mathbf{A}_p^{-1}\mathbf{A}$  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  $\lambda$  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 (15.18) 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 (15.18) 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  $\lambda$  of  $\mathbf{BA} = \mathbf{A}_p^{-1}\mathbf{A}$ .  $\square$

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.

## **Part V**

# **Orthonormal Transformations and Least Squares**





## Chapter 16

# Orthonormal Transformations

Transformations by elementary lower triangular matrices can be used to reduce a matrix to triangular form (cf. Appendix A), Gaussian elimination. Lower triangular matrices are not the only kind of transformations which can be used for such a task. In this chapter we study how transformations by orthonormal matrices can be used to reduce a rectangular matrix to upper triangular (also called upper trapezoidal) form. This leads to a decomposition of the matrix known as a QR decomposition and a compact form which we refer to as a QR factorization. Orthonormal transformations have the advantage that they preserve the Euclidean norm of a vector, and the spectral norm and Frobenius norm of a matrix. Indeed, if  $Q \in \mathbb{R}^{m,m}$  is an orthonormal matrix then  $\|Qv\|_2 = \|v\|_2$ ,  $\|QA\|_2 = \|A\|_2$ , and  $\|QA\|_F = \|A\|_F$  for any vector  $v \in \mathbb{R}^m$  and any matrix  $A \in \mathbb{R}^{m,n}$ , (cf. Lemma 12.3 and Theorem 12.21). 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. The QR factorization can be used to solve least squares problems and linear equations. This will be considered in Chapter 17.

### 16.1 The QR Decomposition and QR Factorization.

**Definition 16.1** Let  $A \in \mathbb{C}^{m,n}$  with  $m \geq n \geq 1$ . We say that  $A = QR$  is a **QR decomposition** of  $A$  if  $Q \in \mathbb{C}^{m,m}$  is square and unitary and

$$R = \begin{bmatrix} R_1 \\ 0_{m-n,n} \end{bmatrix}$$

where  $R_1 \in \mathbb{C}^{n,n}$  is upper triangular and  $0_{m-n,n} \in \mathbb{C}^{m-n,n}$  is the zero matrix. We call  $A = QR$  a **QR factorization** of  $A$  if  $Q \in \mathbb{C}^{m,n}$  has orthonormal columns and  $R \in \mathbb{C}^{n,n}$  is upper triangular.

A QR factorization is obtained from a QR decomposition  $A = QR$  by simply using the first  $n$  columns of  $Q$ . Indeed, if we partition  $Q$  as  $[Q_1, Q_2]$  and  $R = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$ ,

where  $\mathbf{Q}_1 \in \mathbb{R}^{m,n}$  and  $\mathbf{R}_1 \in \mathbb{R}^{n,n}$  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 16.2** *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 now existence and uniqueness.

**Theorem 16.3** *Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  with  $m \geq n \geq 1$  has linearly independent columns. Then  $\mathbf{A}$  has a QR decomposition and a QR factorization. The QR factorization is unique if  $\mathbf{R}$  has positive diagonal elements.*

**Proof.** We prove only the real case. For existence it is enough to show that  $\mathbf{A}$  has a QR factorization. By Corollary ?? the matrix  $\mathbf{A}^T \mathbf{A}$  is symmetric positive definite, and by Theorem 7.32 it has a Cholesky factorization  $\mathbf{A}^T \mathbf{A} = \mathbf{R}^T \mathbf{R}$ , where  $\mathbf{R} \in \mathbb{R}^{n,n}$  is upper triangular and nonsingular. The matrix  $\mathbf{Q} := \mathbf{A} \mathbf{R}^{-1}$  has orthonormal columns since  $\mathbf{Q}^T \mathbf{Q} = \mathbf{R}^{-T} \mathbf{A}^T \mathbf{A} \mathbf{R}^{-1} = \mathbf{R}^{-T} \mathbf{R}^T \mathbf{R} \mathbf{R}^{-1} = \mathbf{I}$ . But then  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  is a QR factorization of  $\mathbf{A}$ . This shows existence. For uniqueness, if  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  is a QR factorization of  $\mathbf{A}$  and  $\mathbf{R}$  has positive diagonal elements then  $\mathbf{A}^T \mathbf{A} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} = \mathbf{R}^T \mathbf{R}$  is the Cholesky factorization of  $\mathbf{A}^T \mathbf{A}$ . Since the Cholesky factorization is unique it follows that  $\mathbf{R}$  is unique and hence  $\mathbf{Q} = \mathbf{A} \mathbf{R}^{-1}$  is unique.  $\square$

We show in Theorem 17.10 existence without the assumption of linearly independent columns.

The QR factorization can be used to prove a classical determinant inequality.

**Theorem 16.4 (Hadamard's Inequality)** *For any  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{n,n}$  we have*

$$|\det(\mathbf{A})| \leq \prod_{j=1}^n \|\mathbf{a}_j\|_2. \quad (16.1)$$

*Equality holds if and only if  $\mathbf{A}$  has a zero column or the columns of  $\mathbf{A}$  are orthogonal.*

**Proof.** Let  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  be a QR factorization of  $\mathbf{A}$ . Since

$$1 = \det(\mathbf{I}) = \det(\mathbf{Q}^H \mathbf{Q}) = \det(\mathbf{Q}^H) \det(\mathbf{Q}) = |\det(\mathbf{Q})|^2$$

we have  $|\det(\mathbf{Q})| = 1$ . Let  $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]$ . Then  $(\mathbf{A}^H \mathbf{A})_{jj} = \|\mathbf{a}_j\|_2^2 = (\mathbf{R}^H \mathbf{R})_{jj} = \|\mathbf{r}_j\|_2^2$ , and

$$|\det(\mathbf{A})| = |\det(\mathbf{Q}\mathbf{R})| = |\det(\mathbf{R})| = \prod_{j=1}^n |r_{jj}| \leq \prod_{j=1}^n \|\mathbf{r}_j\|_2 = \prod_{j=1}^n \|\mathbf{a}_j\|_2.$$

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  $\mathbf{A}^H \mathbf{A} = \mathbf{R}^H \mathbf{R}$  is diagonal, which means that the columns of  $\mathbf{A}$  are orthogonal.  $\square$

### Exercise 16.5

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{Q} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} 2 & 2 \\ 0 & 2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Show that  $\mathbf{Q}$  is orthonormal and that  $\mathbf{Q}\mathbf{R}$  is a QR decomposition of  $\mathbf{A}$ . Find a QR factorization of  $\mathbf{A}$ .

### 16.1.1 QR and Gram-Schmidt

The Gram-Schmidt orthogonalization of the columns of  $\mathbf{A}$  can be used to find the QR factorization of  $\mathbf{A}$ . If  $\mathbf{A} \in \mathbb{R}^{m,n}$  has rank  $n$ , then the set of columns  $\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$  forms a basis for  $\text{span}(\mathbf{A})$  and the Gram-Schmidt orthogonalization process (2.23) takes the form

$$\mathbf{v}_1 = \mathbf{a}_1, \quad \mathbf{v}_j = \mathbf{a}_j - \sum_{i=1}^{j-1} \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \mathbf{v}_i, \quad \text{for } j = 2, \dots, n. \quad (16.2)$$

By Theorem 2.57  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is an orthogonal basis for  $\text{span}(\mathbf{A})$ . If we write (16.2) in the form

$$\mathbf{a}_1 = \mathbf{v}_1, \quad \mathbf{a}_j = \sum_{i=1}^{j-1} \rho_{ij} \mathbf{v}_i + \mathbf{v}_j, \quad \text{where } \rho_{ij} = \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \text{ for } j = 2, \dots, n,$$

then  $\mathbf{A} = \mathbf{V}\hat{\mathbf{R}}$ , where  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{R}^{m,n}$  and

$$\hat{\mathbf{R}} := \begin{bmatrix} 1 & \rho_{12} & \rho_{13} & \rho_{14} & \cdots & \rho_{1n} \\ 0 & 1 & \rho_{23} & \rho_{24} & \cdots & \rho_{2n} \\ 0 & 0 & 1 & \rho_{34} & \cdots & \rho_{3n} \\ 0 & 0 & 0 & 1 & \cdots & \rho_{4n} \\ \vdots & & & & \ddots & \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{bmatrix}$$

is upper unit triangular. Since  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis for  $\text{span}(\mathbf{A})$  the matrix  $\mathbf{D} := \text{diag}(\|\mathbf{v}_1\|_2, \dots, \|\mathbf{v}_n\|_2)$  is nonsingular, and the matrix  $\mathbf{Q}_1 := \mathbf{V}\mathbf{D}^{-1} = [\frac{\mathbf{v}_1}{\|\mathbf{v}_1\|_2}, \dots, \frac{\mathbf{v}_n}{\|\mathbf{v}_n\|_2}]$  is orthonormal. Therefore,  $\mathbf{A} = \mathbf{Q}_1\mathbf{R}_1$ , with  $\mathbf{R}_1 := \mathbf{D}\hat{\mathbf{R}}$  is a QR factorization of  $\mathbf{A}$  with positive diagonal elements in  $\mathbf{R}_1$ .

**Exercise 16.6** Construct  $\mathbf{Q}_1$  and  $\mathbf{R}_1$  in Example 16.2 using Gram-Schmidt orthogonalization.

## 16.2 The Householder Transformation

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 is better numerically, but this will not be considered here. Instead we consider Householder transformations.

**Definition 16.7** A matrix  $\mathbf{H} \in \mathbb{R}^{n,n}$  of the form

$$\mathbf{H} := \mathbf{I} - \mathbf{u}\mathbf{u}^T, \text{ where } \mathbf{u} \in \mathbb{R}^n \text{ and } \mathbf{u}^T\mathbf{u} = 2$$

is called a **Householder transformation**. The name **elementary reflector** is also used.

For  $n = 2$  we find  $\mathbf{H} = \begin{bmatrix} 1-u_1^2 & -u_1u_2 \\ -u_2u_1 & 1-u_2^2 \end{bmatrix}$ . A Householder transformation is symmetric and orthonormal. In particular,

$$\mathbf{H}^T\mathbf{H} = \mathbf{H}^2 = (\mathbf{I} - \mathbf{u}\mathbf{u}^T)(\mathbf{I} - \mathbf{u}\mathbf{u}^T) = \mathbf{I} - 2\mathbf{u}\mathbf{u}^T + \mathbf{u}(\mathbf{u}^T\mathbf{u})\mathbf{u}^T = \mathbf{I}.$$

There are several ways to represent a Householder transformation. Householder used  $\mathbf{I} - 2\mathbf{u}\mathbf{u}^T$ , where  $\mathbf{u}^T\mathbf{u} = 1$ . For any nonzero  $\mathbf{v} \in \mathbb{R}^n$  the matrix

$$\mathbf{H} := \mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}} \quad (16.3)$$

is a Householder transformation. In fact  $\mathbf{H} = \mathbf{I} - \mathbf{u}\mathbf{u}^T$ , where  $\mathbf{u} := \sqrt{2}\frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ .

The main use of Householder transformations is to produce zeros in vectors. We start with

**Lemma 16.8** Suppose  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$  and  $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$ . Then  $(\mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{y}$ .

**Proof.** Since  $\mathbf{x}^T\mathbf{x} = \mathbf{y}^T\mathbf{y}$  we have

$$\mathbf{v}^T\mathbf{v} = (\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y}) = 2\mathbf{x}^T\mathbf{x} - 2\mathbf{y}^T\mathbf{x} = 2\mathbf{v}^T\mathbf{x}. \quad (16.4)$$

But then  $(\mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{x} - \frac{2\mathbf{v}^T\mathbf{x}}{\mathbf{v}^T\mathbf{v}}\mathbf{v} = \mathbf{x} - \mathbf{v} = \mathbf{y}$ .  $\square$

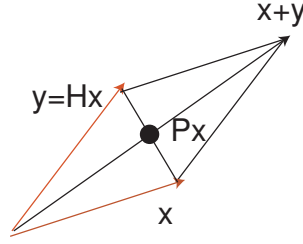


Figure 16.1. The Householder transformation

A geometric interpretation of this lemma is shown in Figure 16.1. We have

$$H = I - \frac{2vv^T}{v^T v} = P - \frac{vv^T}{v^T v}, \text{ where } P := I - \frac{vv^T}{v^T v},$$

and

$$Px = x - \frac{v^T x}{v^T v} v \stackrel{(16.4)}{=} x - \frac{1}{2}v = \frac{1}{2}(x + y).$$

It follows that  $Hx$  is the reflected image of  $x$ . The mirror contains the vector  $x + y$  and has normal  $x - y$ .

**Exercise 16.9** Show that  $\|x\|_2 = \|y\|_2$  implies that  $x - y$  is orthogonal to  $x + y$  and conclude that  $Px$  is the orthogonal projection of  $x$  into the subspace  $\text{span}(x + y)$ .

We can introduce zeros in a vector  $x$  by picking  $\alpha^2 = x^T x$  and  $y := \alpha e_1$  in Lemma 16.8. The equation  $\alpha^2 = x^T x$  has two solutions  $\alpha = +\|x\|_2$  and  $\alpha = -\|x\|_2$ . We want to develop an algorithm which defines a Householder transformation for any nonzero  $x$ . We achieve this by choosing  $\alpha$  to have opposite sign of  $x_1$ . Then  $v_1 = x_1 - \alpha \neq 0$  so  $v \neq 0$ . Another advantage of this choice is that we avoid cancelation in the subtraction in the first component of  $v = x - \alpha e_1$ . This leads to a numerically stable algorithm.

**Lemma 16.10** For a nonzero vector  $x \in \mathbb{R}^n$  we define

$$\alpha := \begin{cases} -\|x\|_2 & \text{if } x_1 > 0 \\ +\|x\|_2 & \text{otherwise,} \end{cases} \quad (16.5)$$

and

$$H := I - uu^T \text{ with } u = \frac{x/\alpha - e_1}{\sqrt{1 - x_1/\alpha}}. \quad (16.6)$$

Then  $H$  is a Householder transformation and  $Hx = \alpha e_1$ .

**Proof.** Let  $y := \alpha e_1$  and  $v := x - y$ . If  $x_1 > 0$  then  $y_1 = \alpha < 0$ , while if  $x_1 \leq 0$  then  $y_1 = \alpha > 0$ . It follows that  $x^T x = y^T y$  and  $v \neq 0$ . By Lemma 16.8 we have  $Hx = \alpha e_1$ , where  $H = I - 2\frac{vv^T}{v^T v}$  is a Householder transformation. Since

$$0 < v^T v = (x - \alpha e_1)^T (x - \alpha e_1) = x^T x - 2\alpha x_1 + \alpha^2 = 2\alpha(\alpha - x_1),$$

we find

$$\mathbf{H} = \mathbf{I} - \frac{2(\mathbf{x} - \alpha \mathbf{e}_1)(\mathbf{x} - \alpha \mathbf{e}_1)^T}{2\alpha(\alpha - x_1)} = \mathbf{I} - \frac{(\mathbf{x}/\alpha - \mathbf{e}_1)(\mathbf{x}/\alpha - \mathbf{e}_1)^T}{1 - x_1/\alpha} = \mathbf{I} - \mathbf{u}\mathbf{u}^T.$$

□

**Example 16.11** For  $\mathbf{x} := [1, 2, 2]^T$  we have  $\|\mathbf{x}\|_2 = 3$  and since  $x_1 > 0$  we choose  $\alpha = -3$ . We find  $\mathbf{u} = -[2, 1, 1]^T/\sqrt{3}$  and

$$\mathbf{H} = \mathbf{I} - \frac{1}{3} \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} -1 & -2 & -2 \\ -2 & 2 & -1 \\ -2 & -1 & 2 \end{bmatrix}.$$

The formulas in Lemma 16.10 are implemented in the following algorithm.

**Algorithm 16.12 (Generate a Householder transformation)** To given  $\mathbf{x} \in \mathbb{R}^n$  the following algorithm computes  $a = \alpha$  and the vector  $\mathbf{u}$  so that  $(\mathbf{I} - \mathbf{u}\mathbf{u}^T)\mathbf{x} = \alpha \mathbf{e}_1$ .

```
function [u, a]=housegen(x)
a=norm(x); u=x;
if a==0
    u(1)=sqrt(2); return;
end
if u(1)>0
    a=-a;
end
u=u/a; u(1)=u(1)-1;
u=u/sqrt(-u(1));
```

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.

**Exercise 16.13** Determine  $\mathbf{H}$  in Algorithm 16.12 when  $\mathbf{x} = \mathbf{e}_1$ .

Householder transformations can also be used to zero out only the lower part of a vector. Suppose  $\mathbf{y} \in \mathbb{R}^k$ ,  $\mathbf{z} \in \mathbb{R}^{n-k}$  and  $\alpha^2 = \mathbf{z}^T \mathbf{z}$ . Consider finding a Householder transformation  $\mathbf{H}$  such that  $\mathbf{H} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \alpha \mathbf{e}_1 \end{bmatrix}$ . Let  $\hat{\mathbf{u}}$  and  $\alpha$  be the output of Algorithm 16.12 called with  $\mathbf{x} = \mathbf{z}$ , i.e.,  $[\hat{\mathbf{u}}, \alpha] = \text{housegen}(\mathbf{z})$  and set  $\mathbf{u}^T = [\mathbf{0}^T, \hat{\mathbf{u}}^T]$ . Then

$$\mathbf{H} = \mathbf{I} - \mathbf{u}\mathbf{u}^T = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} - \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{u}} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{H}} \end{bmatrix},$$

where  $\hat{\mathbf{H}} = \mathbf{I} - \hat{\mathbf{u}}\hat{\mathbf{u}}^T$ . Since  $\mathbf{u}^T \mathbf{u} = \hat{\mathbf{u}}^T \hat{\mathbf{u}} = 2$  we see that  $\mathbf{H}$  and  $\hat{\mathbf{H}}$  are Householder transformations.

**Exercise 16.14** Construct an elementary reflector  $\mathbf{Q}$  such that  $\mathbf{Q}\mathbf{x} = \mathbf{y}$  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 16.15** Show that a  $2 \times 2$  Householder transformation can be written in the form

$$\mathbf{Q} = \begin{bmatrix} -\cos \phi & \sin \phi \\ \sin \phi & \cos \phi \end{bmatrix}.$$

Find  $\mathbf{Q}\mathbf{x}$  if  $\mathbf{x} = [\cos \phi, \sin \phi]^T$ .

**Exercise 16.16** a) Find Householder transformations  $\mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{R}^{3,3}$  such that

$$\mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \mathbf{Q}_2 \mathbf{Q}_1 \begin{bmatrix} 1 & 0 & 1 \\ -2 & -1 & 0 \\ 2 & 2 & 1 \end{bmatrix}$$

is upper triangular.

b) Find the QR factorization of  $\mathbf{A}$  where  $\mathbf{R}$  has positive diagonal elements.

## 16.3 Householder Triangulation

Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$ . 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 orthonormal matrices such that

$$\mathbf{H}_n \mathbf{H}_{n-1} \cdots \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix},$$

and where  $\mathbf{R}_1$  is upper triangular. Here each  $\mathbf{H}_k$  is a Householder transformation. Since the product of orthonormal matrices is orthonormal and each  $\mathbf{H}_k$  is symmetric we obtain the QR decomposition of  $\mathbf{A}$  in the form

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \text{ where } \mathbf{Q} := \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_n \text{ and } \mathbf{R} := \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}. \quad (16.7)$$

Define  $\mathbf{A}_1 = \mathbf{A}$  and suppose for  $k \geq 1$  that  $\mathbf{A}_k$  is upper triangular in its first  $k-1$  columns so that  $\mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \mathbf{D}_k \end{bmatrix}$ , where  $\mathbf{B}_k \in \mathbb{R}^{k-1,k-1}$  is upper triangular and  $\mathbf{D}_k \in \mathbb{R}^{n-k+1,n-k+1}$ . Let  $\hat{\mathbf{H}}_k = \mathbf{I} - \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^T$  be a Householder transformation which zero out the first column in  $\mathbf{D}_k$  under the diagonal, so that  $\hat{\mathbf{H}}_k(\mathbf{D}_k \mathbf{e}_1) = \alpha_k \mathbf{e}_1$ . 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  $\mathbf{B}_{k+1} \in \mathbb{R}^{k,k}$  is upper triangular and  $\mathbf{D}_{k+1} \in \mathbb{R}^{n-k,n-k}$ . 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  $\mathbf{R} := \mathbf{A}_{n+1} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ , where  $\mathbf{R}_1$  is upper triangular and  $\mathbf{R} = \mathbf{H}_n \cdots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A}$ . Thus  $\mathbf{A} = \mathbf{H}_1 \cdots \mathbf{H}_n \mathbf{R}$  and we obtain (16.7).

The process just described can be illustrated as follows when  $m = 4$  and  $n = 3$  using so called **Wilkinson diagrams**.

$$\begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{bmatrix} \xrightarrow{\mathbf{H}_1} \left[ \begin{array}{c|cc} r_{11} & r_{12} & r_{13} \\ \hline 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{array} \right] \xrightarrow{\mathbf{H}_2} \left[ \begin{array}{cc|c} r_{11} & r_{12} & r_{13} \\ \hline 0 & r_{22} & r_{23} \\ 0 & 0 & x \\ 0 & 0 & x \end{array} \right] \xrightarrow{\mathbf{H}_3} \left[ \begin{array}{ccc} r_{11} & r_{12} & r_{13} \\ \hline 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \\ \hline 0 & 0 & 0 \end{array} \right].$$

$$\mathbf{A}_1 = \mathbf{D}_1 \quad \mathbf{A}_2 = \begin{bmatrix} \mathbf{B}_2 & \mathbf{C}_2 \\ \mathbf{0} & \mathbf{D}_2 \end{bmatrix} \quad \mathbf{A}_3 = \begin{bmatrix} \mathbf{B}_3 & \mathbf{C}_3 \\ \mathbf{0} & \mathbf{D}_3 \end{bmatrix} \quad \mathbf{A}_4 = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$$

The transformation is applied to the lower right block.

The process can also be applied to  $\mathbf{A} \in \mathbb{R}^{m,n}$  if  $m \leq n$ . In this case  $m - 1$  Householder transformations will suffice and we obtain

$$\mathbf{H}_{m-1} \cdots \mathbf{H}_1 \mathbf{A} = [\mathbf{R}_1, \mathbf{S}_1] = \mathbf{R}, \quad (16.8)$$

where  $\mathbf{R}_1$  is upper triangular and  $\mathbf{S}_1 \in \mathbb{R}^{m,n-m}$ .

In an algorithm we can store most of the vectors  $\hat{\mathbf{u}}_k = [u_{kk}, \dots, u_{mk}]^T$  and  $\mathbf{R}_1$  in  $\mathbf{A}$ . However, the elements  $u_{kk}$  in  $\hat{\mathbf{u}}_k$  and  $r_{kk}$  in  $\mathbf{R}_1$  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 has to be stored in a separate vector.

## 16.4 Givens Rotations

In some applications, the matrix we want to triangulate has a special structure. Suppose for example that  $\mathbf{A} \in \mathbb{R}^{n,n}$  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 16.17** A **plane rotation** (also called a **Givens's rotation**) is a matrix of the form

$$\mathbf{P} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \quad \text{where } c^2 + s^2 = 1.$$



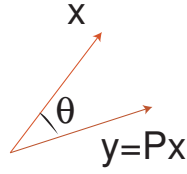


Figure 16.2. A plane rotation.

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 16.18** Show that if  $\mathbf{x} = \begin{bmatrix} r \cos \alpha \\ r \sin \alpha \end{bmatrix}$  then  $\mathbf{Px} = \begin{bmatrix} r \cos(\alpha - \theta) \\ r \sin(\alpha - \theta) \end{bmatrix}$ . Thus  $\mathbf{P}$  rotates a vector  $\mathbf{x}$  in the plane an angle  $\theta$  clockwise. See Figure 16.2.

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

$$\mathbf{Px} = \frac{1}{r} \begin{bmatrix} x_1 & x_2 \\ -x_2 & x_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{1}{r} \begin{bmatrix} x_1^2 + x_2^2 \\ 0 \end{bmatrix} = \begin{bmatrix} r \\ 0 \end{bmatrix},$$

and we have introduced a zero in  $\mathbf{x}$ . We can take  $\mathbf{P} = \mathbf{I}$  when  $\mathbf{x} = \mathbf{0}$ .

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  $\mathbf{P}_{ij} = (p_{kl}) \in \mathbb{R}^{n,n}$  by  $p_{kl} = \delta_{kl}$  except for positions  $ii, jj, ij, ji$ , which are given by

$$\begin{bmatrix} p_{ii} & p_{ij} \\ p_{ji} & p_{jj} \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \quad \text{where } c^2 + s^2 = 1.$$

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

$$\begin{bmatrix} \mathbf{B}(i, :) \\ \mathbf{B}(j, :) \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \mathbf{A}(i, :) \\ \mathbf{A}(j, :) \end{bmatrix}, \quad [\mathbf{C}(:, i) \quad \mathbf{C}(:, j)] = [\mathbf{A}(:, i) \quad \mathbf{A}(:, j)] \begin{bmatrix} c & s \\ -s & c \end{bmatrix}. \quad (16.9)$$

An upper Hessenberg matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  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.

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{12}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{23}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{34}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & 0 & r_{33} & r_{34} \\ 0 & 0 & 0 & r_{44} \end{bmatrix}.$$

For an algorithm see Exercise 16.19.

**Exercise 16.19** Let  $\mathbf{A} \in \mathbb{R}^{n,n}$  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 flops of this algorithm.

**Algorithm 16.20 (Upper Hessenberg linear system)** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  is nonsingular and upper Hessenberg and that  $\mathbf{b} \in \mathbb{R}^n$ . This algorithm uses Givens's rotations to solve the linear system  $\mathbf{Ax} = \mathbf{b}$ . It uses Algorithm A.7.

```
function x=rothesstri(A,b)
n=length(A); A=[A b];
for k=1:n-1
    r=norm([A(k,k),A(k+1,k)]);
    if r>0
        c=A(k,k)/r; s=A(k+1,k)/r;
        A([k k+1],k+1:n+1)=[c s;-s c]*A([k k+1],k+1:n+1);
    end
    A(k,k)=r; A(k+1,k)=0;
end
x=backsolve(A(:,1:n),A(:,n+1));
```

## Chapter 17

# Least Squares

Let  $\mathbf{A} \in \mathbb{C}^{m,n}$  and  $\mathbf{b} \in \mathbb{C}^m$  be given. Consider the linear system  $\mathbf{Ax} = \mathbf{b}$  of  $m$  equations in  $n$  unknowns. If  $m > n$ , we have more equations than unknowns and there might be no vector  $\mathbf{x}$  such that  $\mathbf{Ax} = \mathbf{b}$ . Let  $\mathbf{r}(\mathbf{x}) = \mathbf{Ax} - \mathbf{b} \in \mathbb{C}^m$ . We can then pick a vector norm  $\|\cdot\|$  and look for  $\mathbf{x}^* \in \mathbb{C}^n$  which minimizes  $\|\mathbf{r}(\mathbf{x})\|$ . The choice  $\|\cdot\| = \|\cdot\|_2$ , the Euclidean norm, is particularly convenient and will be studied here. We call an  $\mathbf{x}^* \in \mathbb{C}^n$  which minimizes  $\|\mathbf{r}(\mathbf{x})\|_2$  a least squares solution of  $\mathbf{Ax} = \mathbf{b}$ . Since the square root function is monotone, minimizing  $\|\mathbf{r}(\mathbf{x})\|_2^2$  and  $\|\mathbf{r}(\mathbf{x})\|_2$  are equivalent. We set

$$E(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \|\mathbf{r}(\mathbf{x})\|_2^2.$$

To find  $\mathbf{x}$  which minimizes  $E(\mathbf{x})$  or  $\sqrt{E(\mathbf{x})}$  is called the *least squares problem*.

**Theorem 17.1** *The least squares problem always has a solution. The solution is unique if and only if  $\mathbf{A}$  has linearly independent columns. Moreover, the following are equivalent.*

1.  $\mathbf{x}^*$  is a solution of the least squares problem.
2.  $\mathbf{A}^H \mathbf{Ax}^* = \mathbf{A}^H \mathbf{b}$
3.  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$ , where  $\mathbf{A}^\dagger$  is the pseudo-inverse of  $\mathbf{A}$  and  $\mathbf{z} \in \ker(\mathbf{A})$ .

We have  $\|\mathbf{x}^*\|_2 \geq \|\mathbf{A}^\dagger \mathbf{b}\|_2$  for all solutions  $\mathbf{x}^*$ .

**Proof.** Let  $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$ , where  $\mathbf{b}_1 \in \text{span}(\mathbf{A})$  and  $\mathbf{b}_2 \in \ker(\mathbf{A}^H)$  are the orthogonal projections into  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A}^H)$ , respectively (see Theorem 11.30). Since  $\mathbf{b}_2^H \mathbf{v} = 0$  for any  $\mathbf{v} \in \text{span}(\mathbf{A})$  we have  $\mathbf{b}_2^H (\mathbf{b}_1 - \mathbf{Ax}) = 0$  for any  $\mathbf{x} \in \mathbb{C}^n$ . Therefore, for  $\mathbf{x} \in \mathbb{C}^n$ ,

$$\|\mathbf{b} - \mathbf{Ax}\|_2^2 = \|(\mathbf{b}_1 - \mathbf{Ax}) + \mathbf{b}_2\|_2^2 = \|\mathbf{b}_1 - \mathbf{Ax}\|_2^2 + \|\mathbf{b}_2\|_2^2 \geq \|\mathbf{b}_2\|_2^2,$$

with equality if and only if  $\mathbf{Ax} = \mathbf{b}_1$ . Since  $\mathbf{b}_1 \in \text{span}(\mathbf{A})$  we can always find such an  $\mathbf{x}$  and existence follows.

1  $\iff$  2: By what we have shown  $\mathbf{x}^*$  solves the least squares problem if and only if  $\mathbf{A}\mathbf{x}^* = \mathbf{b}_1$  so that  $\mathbf{b} - \mathbf{A}\mathbf{x}^* = \mathbf{b}_1 + \mathbf{b}_2 - \mathbf{A}\mathbf{x}^* = \mathbf{b}_2 \in \ker(\mathbf{A}^H)$ , or  $\mathbf{A}^H(\mathbf{b} - \mathbf{A}\mathbf{x}^*) = \mathbf{0}$ .  
 1  $\implies$  3: We recall that the orthogonal projection of  $\mathbf{b}$  into  $\text{span}(\mathbf{A})$  can be written  $\mathbf{b}_1 = \mathbf{A}\mathbf{A}^\dagger\mathbf{b}$ , where  $\mathbf{A}^\dagger$  is the pseudo-inverse of  $\mathbf{A}$  (cf. Theorem 11.30). Let  $\mathbf{x}^*$  be a solution of the least squares problem so that  $\mathbf{A}\mathbf{x}^* = \mathbf{b}_1$ . We have to show that  $\mathbf{z} := \mathbf{x}^* - \mathbf{A}^\dagger\mathbf{b} \in \ker(\mathbf{A})$ . But we find  $\mathbf{A}\mathbf{z} = \mathbf{A}\mathbf{x}^* - \mathbf{A}\mathbf{A}^\dagger\mathbf{b} = \mathbf{b}_1 - \mathbf{A}\mathbf{A}^\dagger\mathbf{b} = \mathbf{0}$ .  
 3  $\Leftarrow$  1: If  $\mathbf{x}^* = \mathbf{A}^\dagger\mathbf{b} + \mathbf{z}$  for some  $\mathbf{z} \in \ker(\mathbf{A})$  then  $\mathbf{b}_1 - \mathbf{A}\mathbf{x}^* = \mathbf{b}_1 - \mathbf{A}\mathbf{A}^\dagger\mathbf{b} - \mathbf{A}\mathbf{z} = \mathbf{b}_1 - \mathbf{b}_1 - \mathbf{0} = \mathbf{0}$ , and  $\mathbf{x}^*$  is a solution.

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.

Suppose  $\mathbf{x}^* = \mathbf{A}^\dagger\mathbf{b} + \mathbf{z}$ , with  $\mathbf{z} \in \ker(\mathbf{A})$  is a solution. To show the minimum norm property  $\|\mathbf{x}^*\|_2 \geq \|\mathbf{A}^\dagger\mathbf{b}\|_2$  we 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,  $\mathbf{V}_2^H \mathbf{V}_1 = \mathbf{0}$  since  $\mathbf{V}$  has orthonormal columns. If  $\mathbf{A}^\dagger = \mathbf{V}_1 \Sigma^{-1} \mathbf{U}_1^H$  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

$$\mathbf{z}^H \mathbf{A}^\dagger \mathbf{b} = \mathbf{y}^H \mathbf{V}_2^H \mathbf{V}_1 \Sigma^{-1} \mathbf{U}_1^H \mathbf{b} = \mathbf{0}.$$

But then  $\|\mathbf{x}^*\|_2^2 = \|\mathbf{A}^\dagger\mathbf{b} + \mathbf{z}\|_2^2 = \|\mathbf{A}^\dagger\mathbf{b}\|_2^2 + \|\mathbf{z}\|_2^2 \geq \|\mathbf{A}^\dagger\mathbf{b}\|_2^2$ .  $\square$

The linear system

$$\mathbf{A}^H \mathbf{A} \mathbf{x} = \mathbf{A}^H \mathbf{b}$$

is called the **normal equations**. It is a linear system of  $n$  equations in  $n$  unknowns. If  $\mathbf{A}$  is real then the coefficient matrix  $\mathbf{A}^T \mathbf{A}$  is nonsingular and hence symmetric positive definite if and only if  $\mathbf{A}$  has linearly independent columns.

Before discussing numerical methods for solving the least squares problem we consider some examples.

## 17.1 Examples

**Example 17.2** We choose  $n$  functions  $\phi_1, \phi_2, \dots, \phi_n$  defined for  $t \in \{t_1, t_2, \dots, t_m\}$ . Typical examples of functions might be polynomials, trigonometric functions, exponential functions, or splines. We want to find  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  such that

$$E(\mathbf{x}) = \sum_{i=1}^m w_i \left[ \sum_{j=1}^n x_j \phi_j(t_i) - y_i \right]^2$$

is as small as possible. Let  $p(t) := \sum_{j=1}^n x_j \phi_j(t)$ . The weights  $w_i$  are positive numbers. If  $y_i$  is an accurate observation, we can choose a large weight  $w_i$ . This will force  $p(t_i) - y_i$  to be small. Similarly, a small  $w_i$  will allow  $p(t_i) - y_i$  to be large. If an estimate for the standard deviation  $\delta y_i$  in  $y_i$  is known for each  $i$ , we can choose  $w_i = 1/(\delta y_i)^2$ ,  $i = 1, 2, \dots, m$ . Let  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  have elements  $a_{i,j} = \sqrt{w_i} \phi_j(t_i)$  and  $b_i = \sqrt{w_i} y_i$ . Then

$$(\mathbf{A}\mathbf{x})_i = \sqrt{w_i} \sum_{j=1}^n x_j \phi_j(t_i),$$

$$E(\mathbf{x}) = \sum_{i=1}^m [(\mathbf{A}\mathbf{x})_i - b_i]^2 = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

and we have a least squares problem.

The  $i, j$  element  $b_{i,j}$  in  $\mathbf{B} = \mathbf{A}^T \mathbf{A}$  and the  $i$ th component  $c_i$  in  $\mathbf{c} = \mathbf{A}^T \mathbf{b}$  take the form

$$\begin{aligned} b_{i,j} &= \sum_{k=1}^m a_{k,i} a_{k,j} = \sum_{k=1}^m w_k \phi_i(t_k) \phi_j(t_k), \\ c_i &= \sum_{k=1}^m w_k^{1/2} y_k \phi_i(t_k). \end{aligned} \quad (17.1)$$

In particular, if  $n = 2$ ,  $w_i = 1$ ,  $i = 1, \dots, m$ ,  $\phi_1(t) = 1$ , and  $\phi_2(t) = t$ , the normal equations can be written

$$\begin{bmatrix} m & \sum t_k \\ \sum t_k & \sum t_k^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \sum y_k \\ \sum t_k y_k \end{bmatrix}. \quad (17.2)$$

Here  $k$  ranges from 1 to  $m$  in the sums. This  $2 \times 2$  system is symmetric positive definite and is easily solved for  $x_1$  and  $x_2$ .

**Example 17.3** With the data

$x$	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 (17.2). In this case we obtain

$$\begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 10.1 \end{bmatrix}. \quad (17.3)$$

The solution is  $x_1 = 3.95$  and  $x_2 = -0.98$ . The data and the polynomial  $p(t)$  are shown in Figure 17.1.

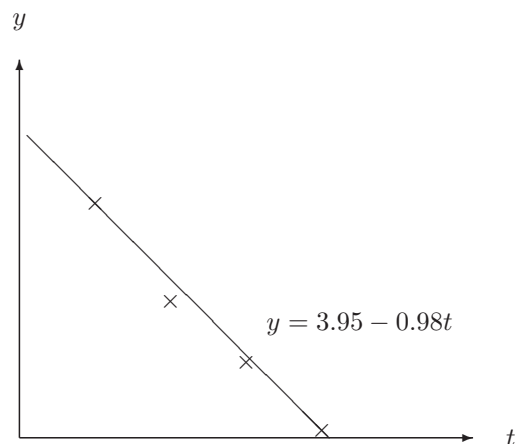
**Example 17.4** Suppose we have a simple input/output model. To every input  $\mathbf{u} \in \mathbb{R}^n$  we obtain an output  $y \in \mathbb{R}$ . Assuming we have a linear relation

$$y = \mathbf{u}^T \mathbf{x} = \sum_{i=1}^n u_i x_i,$$

between  $\mathbf{u}$  and  $y$ , how can we determine  $\mathbf{x}$ ?

Performing  $m \geq n$  experiments we obtain a table of values

$\mathbf{u}$	$\mathbf{u}_1$	$\mathbf{u}_2$	$\cdots$	$\mathbf{u}_m$
$y$	$y_1$	$y_2$	$\cdots$	$y_m$



**Figure 17.1.** A least squares fit to data.

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  $\min \|\mathbf{Ax} - \mathbf{b}\|_2^2$ .

**Exercise 17.5** Suppose  $(t_i, y_i)_{i=1}^m$  are  $m$  points in the plane. We consider the overdetermined systems

$$\begin{array}{ll} \text{(i)} & \begin{array}{l} x_1 = y_1 \\ x_1 = y_2 \\ \vdots \\ x_1 = y_m \end{array} & \text{(ii)} & \begin{array}{l} x_1 + t_1 x_2 = y_1 \\ x_1 + t_2 x_2 = y_2 \\ \vdots \\ x_1 + t_m x_2 = y_m \end{array} \end{array}$$

- 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 17.6** Related to (ii) in Exercise 17.5 we have the overdetermined system

$$\text{(iii)} \quad x_1 + (t_i - \hat{t})x_2 = y_i, \quad i = 1, 2, \dots, m,$$

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 17.7** 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

$$(t_i - c_1)^2 + (y_i - c_2)^2 = r^2, \quad i = 1, \dots, m, \quad (17.4)$$

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

$$t_i x_1 + y_i x_2 + x_3 = t_i^2 + y_i^2, \quad i = 1, \dots, m, \quad (17.5)$$

and then setting  $c_1 = x_1/2$ ,  $c_2 = x_2/2$  and  $x_3 = r^2 - x_1^2 - x_2^2$ .

- a) Derive (17.5) from (17.4). Explain how we can find  $c_1, c_2, r$  once  $[x_1, x_2, x_3]$  is determined.
- b) Formulate (17.5) 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 (17.5) to find the circle passing through the three points  $(1, 4)$ ,  $(3, 2)$ ,  $(1, 0)$ .

## 17.2 Numerical Solution using the Normal Equations

We assume that  $\mathbf{A}$  and  $\mathbf{b}$  are real and that  $\mathbf{A}$  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  $\mathbf{R}^T \mathbf{R}$ -factorization of  $\mathbf{B}$ .

Consider forming the normal equations. We can use either a column oriented- or a row oriented approach. To derive these we partition  $\mathbf{A}$  in terms of columns or rows as

$$\mathbf{A} = [\mathbf{a}_{.1}, \dots, \mathbf{a}_{.n}] = \begin{bmatrix} \mathbf{a}_{1.}^T \\ \vdots \\ \mathbf{a}_{m.}^T \end{bmatrix}.$$

We then find

1.  $(\mathbf{A}^T \mathbf{A})_{ij} = \mathbf{a}_{.i}^T \mathbf{a}_{.j}, \quad (\mathbf{A}^T \mathbf{b})_i = \mathbf{a}_{.i}^T \mathbf{b},$  (inner product form),
2.  $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^m \mathbf{a}_{i.} \mathbf{a}_{i.}^T, \quad \mathbf{A}^T \mathbf{b} = \sum_{i=1}^m b_i \mathbf{a}_{i.},$  (outer product form).

The outer product form is suitable for large problems since it uses only one pass through the data.

Consider the number of operations to compute  $\mathbf{B} := \mathbf{A}^T \mathbf{A}$ . We need  $2m$  flops to find each  $\mathbf{a}_{i.}^T \mathbf{a}_{.j}$ . 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  $O(mn^2)$  flops. The computation of  $\mathbf{B}$  using outer products can also be done in  $O(mn^2)$  flops by computing only one half of  $\mathbf{A}$ . In conclusion the number of operations are  $O(mn^2)$  to find  $\mathbf{B}$ ,  $2mn$  to find  $\mathbf{A}^T \mathbf{b}$ ,  $O(n^3/3)$  to find  $\mathbf{R}$ ,  $O(n^2)$  to solve  $\mathbf{R}^T \mathbf{y} = \mathbf{c}$  and  $O(n^2)$  to solve  $\mathbf{R}\mathbf{x} = \mathbf{y}$ . Since  $m \geq n$ , the bulk of the work is to find  $\mathbf{B}$ .

A problem with the normal equation 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 singular values of  $\mathbf{B}$  are the square of the singular values of  $\mathbf{A}$ . If  $\mathbf{A}$  is ill-conditioned, this could make the normal equation approach problematic. One difficulty which can be encountered is that the computed  $\mathbf{A}^T \mathbf{A}$  might not be positive definite. See Problem 17.18 for an example.

### 17.3 Numerical Solution using the QR Factorization

Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$  has rank  $n$  and let  $\mathbf{b} \in \mathbb{R}^m$ . The QR factorization can be used to solve the least squares problem  $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ . 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

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b} \implies \mathbf{R}_1 \mathbf{x} = \mathbf{Q}_1^T \mathbf{b}.$$

We have the following method to solve the least squares problem.

1. Find a QR factorization  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  of  $\mathbf{A}$ .
2. Solve  $\mathbf{R}_1 \mathbf{x} = \mathbf{Q}_1^T \mathbf{b}$  for the least squares solution  $\mathbf{x}$ .

**Example 17.8** Consider the least squares problem with

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

This is the matrix in Example 16.2. 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  $\mathbf{x} = [1, 0, 0]^T$ .



In the following algorithm we use Householder transformations to solve the least squares problem. We first find a QR decomposition  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  of  $\mathbf{A}$  and compute  $\mathbf{Q}^T\mathbf{b}$ . We then solve the system  $\mathbf{R}_1\mathbf{x} = \mathbf{Q}_1^T\mathbf{b}$ , where  $\mathbf{A} = \mathbf{Q}_1\mathbf{R}_1$  is the QR factorization of  $\mathbf{A}$ . For each Householder transformation we need to compute the update  $\hat{\mathbf{H}}_k\mathbf{D}_k$  (cf. Section 16.3). We do not form the matrix  $\hat{\mathbf{H}}_k$ , but compute the product as

$$(\mathbf{I} - \hat{\mathbf{u}}_k\hat{\mathbf{u}}_k^T)\mathbf{D}_k = \mathbf{D}_k - \hat{\mathbf{u}}_k(\hat{\mathbf{u}}_k^T\mathbf{D}_k). \quad (17.6)$$

**Algorithm 17.9 (Solving least squares by Householder QR)** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$  has rank  $n$  and that  $\mathbf{b} \in \mathbb{R}^m$ . This algorithm uses Householder transformations to solve the least squares problem  $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  if  $m > n$  and the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  if  $m = n$ . It uses Algorithm 16.12 and the built in MATLAB function `linsolve`.

```
function x = houselsq(A,b)
[m,n]=size(A); A=[A b];
for k=1:min(n,m-1)
    [v,A(k,k)]=housegen(A(k:m,k));
    A(k:m,k+1:n+1)=A(k:m,k+1:n+1)-v*(v'*A(k:m,k+1:n+1));
end
opts.UT = true;
x=linsolve(A(1:n,1:n),A(1:n,n+1),opts);
```

The function `housegen(x)` returns a Householder transformation for any  $\mathbf{x} \in \mathbb{R}^n$ . Thus in Algorithm 17.9 we obtain a QR decomposition  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q} = \mathbf{H}_1 \dots \mathbf{H}_r$ , is orthogonal and  $r = \min\{n, m-1\}$ . Thus a QR factorization always exists and we have proved

**Theorem 17.10** Suppose  $m \geq n \geq 1$  and  $\mathbf{A} \in \mathbb{R}^{m,n}$ . Then  $\mathbf{A}$  has a QR decomposition and a QR factorization.

Algorithm 17.9 is a useful alternative to normal equations for solving full rank least squares problems. The condition number for the system  $\mathbf{R}\mathbf{x} = \mathbf{c}$  is  $K_2(\mathbf{R}) = K_2(\mathbf{Q}\mathbf{R}) = K_2(\mathbf{A})$ , 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 17.9 is quite stable. But using Householder transformations requires more work. The leading term in the number of flops in Algorithm 17.9 can be estimated from (17.6). We use the following lemma.

**Lemma 17.11** Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{u} \in \mathbb{R}^m$  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  $O(4mn)$  flops.

**Proof.** It costs  $O(2mn)$  flops to compute  $\mathbf{w}^T := \mathbf{u}^T\mathbf{A}$ ,  $O(mn)$  flops to compute  $\mathbf{W} = \mathbf{u}\mathbf{w}^T$  and  $O(mn)$  flops for the final subtraction  $\mathbf{A} - \mathbf{W}$ , a total of  $O(4mn)$  flops. Taking transpose we obtain the same count for  $\mathbf{A} - (\mathbf{A}\mathbf{v})\mathbf{v}^T$ .  $\square$

Since  $\mathbf{D}_k \in \mathbb{R}^{m-k+1, n-k+1}$  and  $m \geq n$  the cost of computing (17.6) is  $4(m-k)(n-k)$  flops. This implies that the work in Algorithm 17.9 can be estimated as

$$\int_0^n 4(m-k)(n-k)dk = 2mn^2 - \frac{2}{3}n^3. \quad (17.7)$$

When  $m$  is large compared to  $n$  the term  $2mn^2$  dominates. Now forming the normal equations and taking advantage of the symmetry requires  $O(mn^2)$  flops. Thus Algorithm 17.9 requires approximately twice as many flops as using the normal equations when  $m$  is large compared to  $n$ .

### 17.3.1 QR and Linear Systems

Algorithm 17.9 can be used to solve linear systems. If  $\mathbf{A}$  is nonsingular and  $m = n$  then the output  $\mathbf{x}$  will be the solution of  $\mathbf{Ax} = \mathbf{b}$ . This follows since the QR decomposition and QR factorization are the same when  $\mathbf{A}$  is square. Therefore, if  $\mathbf{Rx} = \mathbf{Q}^T \mathbf{b}$  then

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \|\mathbf{QRx} - \mathbf{b}\|_2 = \|\mathbf{Rx} - \mathbf{Q}^T \mathbf{b}\|_2 = 0.$$

So Algorithm 17.9 can be used as 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? Linear systems can be constructed where Gaussian elimination with partial pivoting will fail numerically. On the other hand the transformations used in Householder triangulation are orthonormal so the method is quite stable. So why is Gaussian elimination more popular than Householder triangulation? One reason is that the number of flops in (17.7) when  $m = n$  is given by  $4n^3/3$ , while the count for Gaussian elimination is half of that. 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. Tradition might also play a role.

## 17.4 Numerical Solution using the Singular Value Factorization

This method can be used even if  $\mathbf{A}$  does not have full rank. It requires knowledge of the pseudo-inverse of  $\mathbf{A}$ . By Theorem 17.1

$$\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$$

is a least squares solution for any  $\mathbf{z} \in \ker(\mathbf{A})$ .

**Example 17.12** Consider the singular value factorization derived from the singular value decomposition of Example 11.8.

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \end{bmatrix} [2] [1/\sqrt{2} \ 1/\sqrt{2}].$$

Then

$$\mathbf{A}^\dagger = \mathbf{V}_1 \boldsymbol{\Sigma}_1^{-1} \mathbf{U}_1^T = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix} [1/2] [1/\sqrt{2} \ 1/\sqrt{2} \ 0] = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}.$$

Moreover,  $[-1, 1]^T$  is a basis for  $\ker(\mathbf{A})$ . If  $\mathbf{b} = [b_1, b_2, b_3]^T$ , then for any  $z \in \mathbb{R}$  the vector

$$\mathbf{x} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + z \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

is a solution of  $\min \|\mathbf{Ax} - \mathbf{b}\|_2$  and this gives all solutions.

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 minimal norm solution  $\mathbf{A}^\dagger \mathbf{b}$ . (Cf. Theorem 17.1.)

## 17.5 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  $\mathbf{A}^H \mathbf{Ax} = \mathbf{A}^H \mathbf{b}$ ) as

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}, \quad \mathbf{A}^\dagger := (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H.$$

### 17.5.1 Perturbing the Right Hand Side

Let us now consider the effect of a perturbation in  $\mathbf{b}$  on  $\mathbf{x}$ .

**Theorem 17.13** Suppose  $\mathbf{A} \in \mathbb{C}^{m,n}$  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  $\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2$  and  $\min_{\mathbf{y}} \|\mathbf{Ay} - \mathbf{b} - \mathbf{e}\|_2$ . Finally, let  $\mathbf{b}_1, \mathbf{e}_1$  be the projections of  $\mathbf{b}$  and  $\mathbf{e}$  on  $\text{span}(\mathbf{A})$ . If  $\mathbf{b}_1 \neq \mathbf{0}$ , we have for any operator norm  $\|\cdot\|$  subordinate to a vector norm that is also denoted by  $\|\cdot\|$

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|. \quad (17.8)$$

**Proof.** Subtracting  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$  from  $\mathbf{y} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{A}^\dagger \mathbf{e}$  we have  $\mathbf{y} - \mathbf{x} = \mathbf{A}^\dagger \mathbf{e}$ . Since  $\mathbf{A}^H \mathbf{e} = \mathbf{A}^H \mathbf{e}_1$ , we find  $\mathbf{A}^\dagger \mathbf{e} = \mathbf{A}^\dagger \mathbf{e}_1$ . 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.  $\square$

(17.8) is analogous to the bound (cf. (12.12))

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|.$$

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 is that  $\|\mathbf{e}\|/\|\mathbf{b}\|$  has been replaced by the projections  $\|\mathbf{e}_1\|/\|\mathbf{b}_1\|$ . If  $\mathbf{b}$  lies almost entirely in  $\ker(\mathbf{A}^H)$ , 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 17.2. 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.

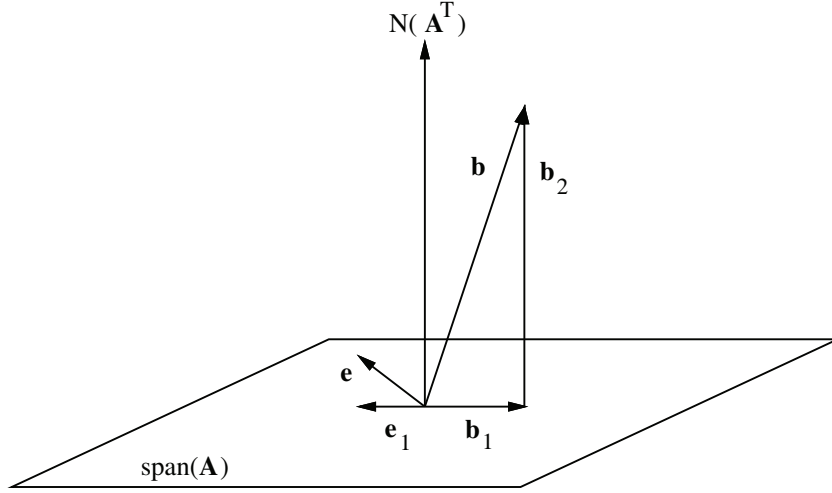


Figure 17.2. Graphical interpretation of the bounds in Theorem 17.13.

**Example 17.14** 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,

$$\begin{aligned} \mathbf{A}^T \mathbf{A} &= \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad (\mathbf{A}^T \mathbf{A})^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} \\ \mathbf{A}^\dagger &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix}. \end{aligned}$$

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  $\mathbf{b}_2$ . We have  $\text{span}(\mathbf{A}) = \text{span}(\mathbf{e}_1, \mathbf{e}_2)$ , and  $\ker(\mathbf{A}^T) = \text{span}(\mathbf{e}_3)$ . The projection  $\mathbf{b}_1$  on  $\text{span}(\mathbf{A})$  is  $\mathbf{b}_1 = [10^{-4}, 0, 0]^T$ . Since  $\|\mathbf{b}\|_\infty / \|\mathbf{b}_1\|_\infty = 10^4$  is large, we expect that the solutions  $\mathbf{x}$  and  $\mathbf{y}$  of  $\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  and  $\min \|\mathbf{A}\mathbf{y} - \mathbf{b} - \mathbf{e}\|_2$  will differ by much more than  $\|\mathbf{e}\|_\infty / \|\mathbf{b}\|_\infty = 10^{-6}$ . To check this we compute  $\mathbf{x}$  and  $\mathbf{y}$ . These can be found by either solving the normal

equations or by solving  $\mathbf{Ax} = \mathbf{b}_1$ ,  $\mathbf{Ay} = \mathbf{b}_1 + \mathbf{e}_1$ . This gives  $\mathbf{x} = [10^{-4}, 0]^T$  and  $\mathbf{y} = [10^{-4} + 10^{-6}, 0]^T$ . We find

$$\frac{\|\mathbf{x} - \mathbf{y}\|_\infty}{\|\mathbf{x}\|_\infty} = \frac{10^{-6}}{10^{-4}} = 10^{-2},$$

and this is indeed much larger than  $\|\mathbf{e}\|_\infty / \|\mathbf{b}\|_\infty = 10^{-6}$ , but equals  $\|\mathbf{e}_1\|_\infty / \|\mathbf{b}_1\|_\infty$ . (17.8) takes the form

$$\frac{1}{4}10^{-2} \leq \frac{\|\mathbf{y} - \mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \leq 4 \cdot 10^{-2}$$

and this captures pretty well the true states of affairs.

**Exercise 17.15** 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 (17.8). The lower bound is best possible in a similar way.

**Exercise 17.16** a) Let  $\mathbf{A} \in \mathbb{C}^{m,n}$ . Show that we have equality to the right in (17.8) if  $\mathbf{b} = \mathbf{Ay}_A$ ,  $\mathbf{e}_1 = \mathbf{y}_{A^\dagger}$  where  $\|\mathbf{Ay}_A\| = \|\mathbf{A}\|$ ,  $\|\mathbf{A}^\dagger \mathbf{y}_{A^\dagger}\| = \|\mathbf{A}^\dagger\|$ .

- 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 17.14. Find extremal  $\mathbf{b}$  and  $\mathbf{e}$  when the  $l_\infty$  norm is used.

## 17.5.2 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 [12, p. 51]. For other estimates see [2] and [19].

**Theorem 17.17** Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{m,n}$ ,  $m > n$ , 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  $\ker(\mathbf{A}^H)$  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

$$\rho = \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{\alpha} K(1 + \beta K) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2}, \quad \beta = \frac{\|\mathbf{b}_2\|_2}{\|\mathbf{b}_1\|_2}, \quad K = \|\mathbf{A}\|_2 \|\mathbf{A}^\dagger\|_2. \quad (17.9)$$

(17.9) 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  $\mathbf{A}$ . Suppose  $\|\mathbf{E}\|_2$  is so small that  $\alpha \approx 1$ . If  $\beta K$  is small compared to  $K$  then the condition number is approximately  $K$  as for linear systems. But in the worst case the bounding term is dominated by  $\beta K^2$ . Thus  $K$  is squared in addition to the effect of a small  $\mathbf{b}_1$  compared to  $\mathbf{b}_2$ .

**Exercise 17.18** Consider the least squares problems where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1+\epsilon \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 2 \\ 3 \\ 2 \end{bmatrix}, \quad \epsilon \in \mathbb{R}.$$

- Find the normal equations and the exact least squares solution.
- Suppose  $\epsilon$  is small and we replace the  $(2, 2)$  element  $3+2\epsilon+\epsilon^2$  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). 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).

## 17.6 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 11.36. If  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m,n}$  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 17.19** Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m,n}$  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

$$|\alpha_j - \beta_j| \leq \|\mathbf{A} - \mathbf{B}\|_2, \quad \text{for } j = 1, 2, \dots, n. \quad (17.10)$$

**Proof.** Fix  $j$  and let  $\mathcal{S}$  be the  $n-j+1$  dimensional subspace for which the minimum in Theorem 11.35 is obtained for  $\mathbf{A}$ . Then

$$\alpha_j = \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{B} + (\mathbf{A} - \mathbf{B}))\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{B}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \beta_j + \|\mathbf{A} - \mathbf{B}\|_2.$$

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 12.33.

**Theorem 17.20** *Let  $\mathbf{A}, \mathbf{E} \in \mathbb{R}^{m,n}$  have singular values  $\alpha_1 \geq \cdots \geq \alpha_n$  and  $\epsilon_1 \geq \cdots \geq \epsilon_n$ . If  $\|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2 < 1$  then*

1.  $\text{rank}(\mathbf{A} + \mathbf{E}) \geq \text{rank}(\mathbf{A})$ ,
2.  $\|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2} = \frac{1}{\alpha_r - \epsilon_1}$ ,

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  $\beta_1 \geq \cdots \geq \beta_n$ . 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 17.19 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

$$\|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{1}{\beta_r} \leq \frac{1}{\alpha_r - \epsilon_1} = \frac{1/\alpha_r}{1 - \epsilon_1/\alpha_r} = \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2}.$$

$\square$





## **Part VI**

# **Eigenvalues and Eigenvectors**



## Chapter 18

# Numerical Eigenvalue Problems

In this and the next chapter we consider numerical methods for finding one or more of the eigenvalues and eigenvectors of a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$ . 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  $\lambda_j = 10^{-1}e^{2\pi ij/16}$ ,  $j = 1, \dots, 16$ . 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.

## 18.1 Perturbation of Eigenvalues

In this section we study the following problem. Given matrices  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n,n}$ , 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 12.42, 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  $\mathbf{A} + \mathbf{E}$  is  $\pi(\lambda) := \lambda^n - (-1)^n \epsilon$ , 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 18.1 (Elsner's Theorem)** *Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n,n}$ . To every  $\mu \in \sigma(\mathbf{A} + \mathbf{E})$  there is a  $\lambda \in \sigma(\mathbf{A})$  such that*

$$|\mu - \lambda| \leq (\|\mathbf{A}\|_2 + \|\mathbf{A} + \mathbf{E}\|_2)^{1-1/n} \|\mathbf{E}\|_2^{1/n}. \quad (18.1)$$

**Proof.** Suppose  $\mathbf{A}$  has eigenvalues  $\lambda_1, \dots, \lambda_n$  and let  $\lambda_1$  be one which is closest to  $\mu$ . Let  $\mathbf{u}_1$  with  $\|\mathbf{u}_1\|_2 = 1$  be an eigenvector corresponding to  $\mu$ , 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 (16.1) 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 18.4. For a more general discussion see [19].

**Theorem 18.2** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\mathbf{E} \in \mathbb{C}^{n,n}$  and let  $\mu$  be an eigenvalue of  $\mathbf{A} + \mathbf{E}$ . Then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{E}\|_p, \quad 1 \leq p \leq \infty, \text{ where } K_p(\mathbf{X}) := \|\mathbf{X}\|_p \|\mathbf{X}^{-1}\|_p. \quad (18.2)$$

The equation (18.2) 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 (18.2) 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 (18.2) implies the following result.

**Theorem 18.3** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is normal and let  $\mu$  be an eigenvalue of  $\mathbf{A} + \mathbf{E}$  for some  $\mathbf{E} \in \mathbb{C}^{n,n}$ . Then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that  $|\lambda - \mu| \leq \|\mathbf{E}\|_2$ .

For an even stronger result for Hermitian matrices see Corollary 10.16. 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.

### 18.1.1 Gerschgorin's Theorem

The following theorem is useful for locating eigenvalues of an arbitrary square matrix.

**Theorem 18.4 (Gerschgorin's Circle Theorem)** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$ . 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  $(\lambda, \mathbf{x})$  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

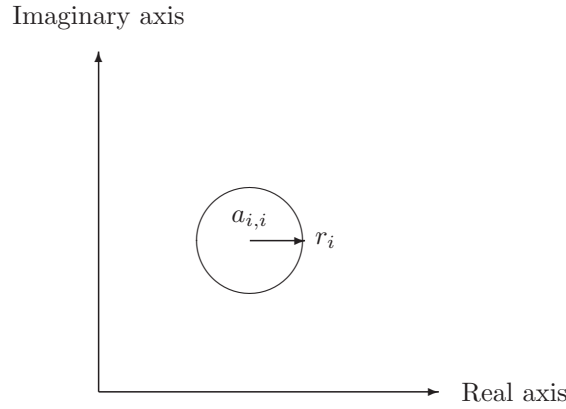
$(\lambda - a_{ii})x_i = \sum_{j \neq i} a_{ij}x_j$ . Dividing by  $x_i$  and taking absolute values we find

$$|\lambda - a_{ii}| = \left| \sum_{j \neq i} a_{ij}x_j/x_i \right| \leq \sum_{j \neq i} |a_{ij}| |x_j/x_i| \leq r_i$$

since  $|x_j/x_i| \leq 1$  for all  $j$ . Thus  $\lambda \in R_i$ .

Since  $\lambda$  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 18.1.  $R_i$  is called a (Gerschgorin) row disk.



**Figure 18.1.** The Gerschgorin disk  $R_i$ .

An eigenvalue  $\lambda$  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 18.5** Let  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m,m}$  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

$$R_1 = R_m = \{z \in \mathbb{R} : |z-2| \leq 1\}, \text{ and } R_i = \{z \in \mathbb{R} : |z-2| \leq 2\}, \quad i = 2, 3, \dots, m-1.$$

We conclude that  $\lambda \in [0, 4]$  for any eigenvalue  $\lambda$  of  $\mathbf{T}$ . To check this, we recall that by Lemma 8.11 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 Gerschgorin's theorem gives a remarkably good estimate for large  $m$ .

Sometimes some of the Gerschgorin disks are distinct and we have

**Corollary 18.6** *If  $p$  of the Gerschgorin 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 18.1, see Exercise 18.7. The row disks  $R_i(t)$  of  $\mathbf{A}(t)$  have radius proportional to  $t$ , indeed

$$R_i(t) = \{z \in \mathbb{C} : |z - a_{ii}| \leq tr_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|.$$

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 lose or gain eigenvalues. It follows that  $R^p(1)$  must contain  $p$  eigenvalues of  $\mathbf{A}$ .  $\square$

**Exercise 18.7** *Suppose  $t_1, t_2 \in [0, 1]$  and that  $\mu$  is an eigenvalue of  $\mathbf{A}(t_2)$ . Show, using Theorem 18.1 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  $\lambda$  such that*

$$|\lambda - \mu| \leq C(t_2 - t_1)^{1/n}, \quad \text{where } C \leq 2(\|\mathbf{D}\|_2 + \|\mathbf{A} - \mathbf{D}\|_2).$$

*Thus, as a function of  $t$ , every eigenvalue of  $\mathbf{A}(t)$  is a continuous function of  $t$ .*

**Example 18.8** *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 18.6 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 18.9** 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 Gerschgorin's theorem that  $\mathbf{A}$  is nonsingular.

**Exercise 18.10** Show using Gerschgorin'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.

## 18.2 Unitary Similarity Transformation of a Matrix into Upper Hessenberg or Tridiagonal 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 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 done.

Recall that a matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  is upper Hessenberg if  $a_{i,j} = 0$  for  $j = 1, 2, \dots, i-2$ ,  $i = 3, 4, \dots, n$ . We will reduce  $\mathbf{A} \in \mathbb{R}^{n,n}$  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

$$\mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \in \mathbb{R}^{n,n}.$$



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 16.12 to find  $\mathbf{v}_k$  and  $\alpha_k$ . 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 18.11 (Householder reduction to Hessenberg form)** This algorithm uses Householder similarity transformations to reduce a matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  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 16.12 is used in each step of the reduction.

```
function [L,B] = hesshousegen(A)
n=length(A); L=zeros(n,n); B=A;
for k=1:n-2
    [v,B(k+1:k,n)]=housegen(B(k+1:n,k));
    L(k+1:n,k)=v; B(k+2:n,k)=zeros(n-k-1,1);
    C=B(k+1:n,k+1:n); B(k+1:n,k+1:n)=C-v*(v'*C);
    C=B(1:n,k+1:n); B(1:n,k+1:n)=C-(C*v)*v';
end
```

**Exercise 18.12** Show that the number of flops for Algorithm 18.11 is  $O(\frac{10}{3}n^3)$ .

We can use the output of Algorithm 18.11 to assemble the matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$  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  $\begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_k \end{bmatrix}$ , where  $\mathbf{U}_k \in \mathbb{R}^{n-k, n-k}$ . Then

$$\mathbf{Q}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{v}_k \mathbf{v}_k^T \end{bmatrix} * \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_k \end{bmatrix} = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_k - \mathbf{v}_k (\mathbf{v}_k^T \mathbf{U}_k) \end{bmatrix}.$$

This leads to the following algorithm.

**Algorithm 18.13 (Assemble Householder transformations)** Suppose  $[L, B] = \text{hesshousegen}(A)$  is the output of Algorithm 18.11. This algorithm assembles an orthonormal matrix  $Q$  from the columns of  $L$  such that  $B = Q^T A Q$  is upper Hessenberg.

```
function Q = accumulateQ(L)
n=length(L); Q=eye(n);
for k=n-2:-1:1
    v=L(k+1:n,k); C=Q(k+1:n,k+1:n);
    Q(k+1:n,k+1:n)=C-v*(v'*C);
end
```

**Exercise 18.14** Show that the number of flops required by Algorithm 18.13 is  $O(\frac{4}{3}n^3)$ .

**Exercise 18.15** If  $A$  is symmetric we can modify Algorithm 18.11 as follows. To find  $A_{k+1}$  from  $A_k$  we have to compute  $V_k E_k V_k$  where  $E_k$  is symmetric. Dropping subscripts we have to compute a product of the form  $G = (I - vv^T)E(I - vv^T)$ . Let  $w := Ev$ ,  $\beta := \frac{1}{2}v^T w$  and  $z := w - \beta v$ . Show that  $G = E - vz^T - zv^T$ . Since  $G$  is symmetric, only the sub- or superdiagonal elements of  $G$  need to be computed. Computing  $G$  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 [18] for a detailed algorithm.

### 18.3 Computing a Selected Eigenvalue of a Symmetric Matrix

Let  $A \in \mathbb{R}^{n,n}$  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  $\lambda_m$  for some  $1 \leq m \leq n$ . Using Householder similarity transformations as outlined in the previous section we can assume that  $A$  is symmetric and tridiagonal.

$$A = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}. \quad (18.3)$$

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

$$\mathbf{A}_1 = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{i-2} & d_{i-1} & c_{i-1} \\ & & & c_{i-1} & d_i \end{bmatrix} \quad \text{and} \quad \mathbf{A}_2 = \begin{bmatrix} d_{i+1} & c_{i+1} & & & \\ c_{i+1} & d_{i+2} & c_{i+2} & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}.$$

Thus  $\mathbf{A}$  is block diagonal and each diagonal block is tridiagonal. By 6. of Theorem 5.3 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 18.16** *An irreducible, tridiagonal and symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  has  $n$  real and distinct eigenvalues.*

**Proof.** Let  $\mathbf{A}$  be given by (18.3). By Theorem 10.7 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

$$p_{k+1}(x) = (x - d_{k+1})p_k(x) - c_k^2 p_{k-1}(x). \quad (18.4)$$

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

$$p_{k+1}(y_j) = (-1)^{k+1-j} |p_{k+1}(y_j)| \neq 0, \quad \text{for } j = 0, 1, \dots, k+1.$$

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.  $\square$

### 18.3.1 The Inertia Theorem

We say that two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$  are **congruent** if  $\mathbf{A} = \mathbf{E}^H \mathbf{B} \mathbf{E}$  for some nonsingular matrix  $\mathbf{E} \in \mathbb{C}^{n,n}$ . By Theorem 10.5 a Hermitian matrix  $\mathbf{A}$  is both congruent and similar to a diagonal matrix  $\mathbf{D}$ ,  $\mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{D}$  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 18.17 (Sylvester's Inertia Theorem)** *If  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$  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  $\mathbf{A} = \mathbf{E}^H \mathbf{B} \mathbf{E}$ , 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 3.5). 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  $\mu_1, \dots, \mu_n$ , where  $\mu_i \leq 0$  for  $i \geq m + 1$  then

$$\mathbf{y}^H \mathbf{A} \mathbf{y} = \sum_{i=1}^n \lambda_i |y_i|^2 = \sum_{i=1}^k \lambda_i |x_i|^2 > 0.$$

But

$$\mathbf{y}^H \mathbf{A} \mathbf{y} = \mathbf{y}^H \mathbf{E}^H \mathbf{B} \mathbf{E} \mathbf{y} = \mathbf{z}^H \mathbf{B} \mathbf{z} = \sum_{i=m+1}^n \mu_i |z_i|^2 \leq 0,$$

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  $\mathbf{U}_1^H \mathbf{A} \mathbf{U}_1 = \mathbf{D}_1$  and  $\mathbf{U}_2^H \mathbf{B} \mathbf{U}_2 = \mathbf{D}_2$  where  $\mathbf{D}_1$  and  $\mathbf{D}_2$  are diagonal matrices. Since  $\mathbf{A} = \mathbf{E}^H \mathbf{B} \mathbf{E}$ , we find  $\mathbf{D}_1 = \mathbf{F}^H \mathbf{D}_2 \mathbf{F}$  where  $\mathbf{F} = \mathbf{U}_2^H \mathbf{E} \mathbf{U}_1$  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  $\zeta$  and  $v$ .  $\square$

**Corollary 18.18** Suppose  $\mathbf{A} = \text{tridiag}(c_i, d_i, c_i) \in \mathbb{R}^{n,n}$  is symmetric and that  $\alpha \in \mathbb{R}$  is such that  $\mathbf{A} - \alpha \mathbf{I}$  has an  $\mathbf{LDL}^T$  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  $\alpha$  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

$$d_1(\alpha) = d_1 - \alpha, \quad d_k(\alpha) = d_k - \alpha - c_{k-1}^2/d_{k-1}(\alpha), \quad k = 2, 3, \dots, n. \quad (18.5)$$

**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 (18.5) follows immediately from (6.25). 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  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  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  $\alpha$ .  $\square$

**Exercise 18.19** Consider the matrix in Exercise 18.9. Determine the number of eigenvalues greater than 4.5.

**Exercise 18.20** Let for  $n \in \mathbb{N}$

$$\mathbf{A}_n = \begin{bmatrix} 10 & 1 & 0 & \cdots & 0 \\ 1 & 10 & 1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & 1 & 10 & 1 \\ 0 & \cdots & 0 & 1 & 10 \end{bmatrix}.$$

- Let  $d_k$  be the diagonal elements of  $\mathbf{D}$  in an  $\mathbf{LDL}^T$  factorization of  $\mathbf{A}_n$ . Show that  $5 + \sqrt{24} < d_k \leq 10$ ,  $k = 1, 2, \dots, n$ .
- 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 18.21** (Simultaneous diagonalization of two symmetric matrices by a congruence transformation). Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,n}$  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  $\mathbf{D}^{-1/2} = \text{diag}(d_1^{-1/2}, \dots, d_n^{-1/2})$ .

- 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  $\mathbf{E} = \mathbf{U}^T \mathbf{D}^{-1/2} \hat{\mathbf{U}}^T$ .
- Show that  $\mathbf{E}$  is nonsingular and that  $\mathbf{E}^T \mathbf{A} \mathbf{E} = \hat{\mathbf{D}}$ ,  $\mathbf{E}^T \mathbf{B} \mathbf{E} = \mathbf{I}$ .

### 18.3.2 Approximating $\lambda_m$

Corollary 18.18 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 Gerschgorin'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  $\lambda_m$  we generate a sequence  $\{[a_j, b_j]\}$  of intervals, each containing  $\lambda_m$  and  $b_j - a_j = 2^{-j}(b - a)$ .

As it stands this method will fail if in (18.5) 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  $\epsilon_M$  is the Machine epsilon, typically  $2 \times 10^{-16}$  for Matlab. This replacement is done if  $|d_k(\alpha)| < |\delta_k|$ .

**Exercise 18.22** 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  $\lambda_m$  for a given  $m$  using bisection and the method outlined in Section 18.3.2.

- 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.
- 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  $\lambda_m$ . Iterate until  $b_k - a_k \leq (b - a)\epsilon_M$ , where  $\epsilon_M$  is Matlab's machine epsilon `eps`. Typically  $\epsilon_M \approx 2.22 \times 10^{-16}$ .
- Test the program on  $\mathbf{T} := \text{tridiag}(-1, 2, -1)$  of size 100. Compare the exact value of  $\lambda_5$  with your result and the result obtained by using Matlab's built-in function `eig`.

**Exercise 18.23** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is upper Hessenberg and  $x \in \mathbb{C}$ . We will study two algorithms to compute  $f(x) = \det(\mathbf{A} - x\mathbf{I})$ .

- Show that Gaussian elimination without pivoting requires  $O(1 * n^2)$  flops.
- Show that the number of flops is the same if partial pivoting is used.
- Estimate the number of flops if Givens's rotations are used.
- Compare the two methods discussing advantages and disadvantages.

## 18.4 Perturbation Proofs

We first show that the  $p$ -norm of a diagonal matrix is equal to its spectral radius.

**Lemma 18.24** *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 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  $\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p} \leq \rho(\mathbf{A})$ . But from Theorem 12.42 we have  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|_p$  and the proof is complete for  $p < \infty$ .  $\square$

**Exercise 18.25** *Give a direct proof that  $\|\mathbf{A}\|_\infty = \rho(\mathbf{A})$  if  $\mathbf{A}$  is diagonal.*

Suppose now  $(\mu, \mathbf{x})$  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 18.26 (Absolute errors)** *Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\lambda$  of  $\mathbf{A}$  such that*

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{r}\|_p, \quad 1 \leq p \leq \infty, \quad (18.6)$$

where  $\mathbf{r} := \mathbf{Ax} - \mu\mathbf{x}$  and  $K_p(\mathbf{X}) := \|\mathbf{X}\|_p \|\mathbf{X}^{-1}\|_p$ . If for some  $\mathbf{E} \in \mathbb{C}^{n,n}$  it holds that  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} + \mathbf{E}$ , then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{E}\|_p, \quad 1 \leq p \leq \infty, \quad (18.7)$$

**Proof.** If  $\mu \in \sigma(\mathbf{A})$  then we can take  $\lambda = \mu$  and (18.6), (18.7) hold trivially. So assume  $\mu \notin \sigma(\mathbf{A})$ . Since  $\mathbf{A}$  is nondefective it can be diagonalized, we have  $\mathbf{A} = \mathbf{XD}\mathbf{X}^{-1}$ , 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 18.24 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 (18.6) 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 (18.6) proves (18.7).  $\square$

For the accuracy of an eigenvalue of small magnitude we are interested in the size of the relative error.

**Theorem 18.27 (Relative errors)** *Suppose in Theorem 18.26 that  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\lambda$  of  $\mathbf{A}$  such that*

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{r}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (18.8)$$

where  $\mathbf{r} := \mathbf{A}\mathbf{x} - \mu\mathbf{x}$ . If for some  $\mathbf{E} \in \mathbb{C}^{n,n}$  it holds that  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} + \mathbf{E}$ , then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})\|\mathbf{A}^{-1}\mathbf{E}\|_p \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{E}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (18.9)$$

**Proof.** Applying Theorem 12.42 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 (18.8) follows from (18.6). To prove (18.9) we define the matrices  $\mathbf{B} := \mu\mathbf{A}^{-1}$  and  $\mathbf{F} := -\mathbf{A}^{-1}\mathbf{E}$ . If  $(\lambda_j, \mathbf{x})$  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  $(\mu, \mathbf{x})$  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 18.26 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 (18.9). The second inequality in (18.9) follows from the submultiplicativity of the  $p$ -norm.  $\square$



## Chapter 19

# The Power and QR Methods

### 19.1 The Power Method

Let  $\mathbf{A} \in \mathbb{C}^{n,n}$  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.
- (19.1)

The first assumption means that  $\mathbf{A}$  has a dominant eigenvalue  $\lambda_1$  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

$$\mathbf{z}_k := \mathbf{A}^k \mathbf{z}_0 = \mathbf{A} \mathbf{z}_{k-1}, \quad k = 1, 2, \dots \quad (19.2)$$

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 (19.1) 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

$$\mathbf{z}_k = c_1 \lambda_1^k \mathbf{v}_1 + c_2 \lambda_2^k \mathbf{v}_2 + \dots + c_n \lambda_n^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (19.3)$$

Dividing by  $\lambda_1^k$  we find

$$\frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1 + c_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k \mathbf{v}_2 + \dots + c_n \left(\frac{\lambda_n}{\lambda_1}\right)^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (19.4)$$

Assumption (i) of (19.1) implies that  $(\lambda_j/\lambda_1)^k \rightarrow 0$  as  $k \rightarrow \infty$  for all  $j \geq 2$  and we obtain

$$\lim_{k \rightarrow \infty} \frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1, \quad (19.5)$$

the dominant eigenvector of  $\mathbf{A}$ . It can be shown that this also holds for defective matrices as long as (i) and (ii) of (19.1) hold, see for example page 58 of [18].

In practice we need to scale the iterates  $\mathbf{z}_k$  somehow and we normally do not know  $\lambda_1$ . 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:

$$\begin{aligned} (i) \quad & \mathbf{y}_k = \mathbf{A}\mathbf{x}_{k-1} \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k/\|\mathbf{y}_k\|. \end{aligned} \tag{19.6}$$

**Lemma 19.1** *Suppose (19.1) 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

$$\mathbf{x}_k = \frac{\mathbf{z}_k}{\|\mathbf{z}_k\|} = \frac{c_1 \lambda_1^k}{|c_1 \lambda_1^k|} \frac{\mathbf{v}_1 + \frac{c_2}{c_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \dots + \frac{c_n}{c_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n}{\|\mathbf{v}_1 + \frac{c_2}{c_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \dots + \frac{c_n}{c_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n\|}, \quad k = 0, 1, 2, \dots,$$

and this implies the lemma.  $\square$

Suppose we know an approximate eigenvector  $\mathbf{u}$  of  $\mathbf{A}$ , but not the corresponding eigenvalue  $\mu$ . One way of estimating  $\mu$  is to minimize the Euclidian norm of the residual  $r(\lambda) := \mathbf{A}\mathbf{u} - \lambda\mathbf{u}$ .

**Theorem 19.2** *Let  $\mathbf{A} \in \mathbb{C}^{n,n}$ ,  $\mathbf{u} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ , and let  $\rho : \mathbb{C} \rightarrow \mathbb{R}$  be given by  $\rho(\lambda) = \|\mathbf{A}\mathbf{u} - \lambda\mathbf{u}\|_2$ . Then  $\rho$  is minimized when  $\lambda := \frac{\mathbf{u}^H \mathbf{A} \mathbf{u}}{\mathbf{u}^H \mathbf{u}}$ , the Rayleigh quotient for  $\mathbf{A}$ .*

**Proof.** It is equivalent to minimize  $E(\lambda) := \rho^2(\lambda)$ . Now

$$E(\lambda) = \mathbf{u}^T \mathbf{u} \lambda^2 - 2\mathbf{u}^T \mathbf{A} \mathbf{u} \lambda + \mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u}.$$

We see that  $E$  is a quadratic polynomial and since  $\mathbf{u}^T \mathbf{u} > 0$ ,  $E$  has a unique minimum  $\lambda$ , where  $E'(\lambda) = 0$ . The solution of  $E'(\lambda) = 0$  is given by  $\lambda = \frac{\mathbf{u}^H \mathbf{A} \mathbf{u}}{\mathbf{u}^H \mathbf{u}}$ .  $\square$

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. The estimate (18.8) can give us some insight. Recall that if  $\mathbf{A}$  is nonsingular and nondefective with eigenvector matrix  $\mathbf{X}$  and  $(\mu, \mathbf{u})$  is an approximate eigenpair with  $\|\mathbf{u}\|_2 = 1$ , then we can find an eigenvalue  $\lambda$  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.

**Algorithm 19.3 (The Power Method)** Given  $\mathbf{A} \in \mathbb{C}^{n,n}$ , 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  $\|\mathbf{A}\mathbf{x} - l\mathbf{x}\|_2 / \|\mathbf{A}\|_F < tol$ . If no such eigenpair is found in  $K$  iterations the value  $it = K + 1$  is returned.

```
function [l, x, it] = powerit(A, z, K, tol)
af = norm(A, 'fro'); x = z / norm(z);
for k = 1:K
    y = A*x; l = x' * y;
    if norm(y - l*x) / af < tol
        it = k; x = y / norm(y); return
    end
    x = y / norm(y);
end
it = K + 1;
```

**Example 19.4** We try `powerit` on the three matrices

$$\mathbf{A}_1 := \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 1.7 & -0.4 \\ 0.15 & 2.2 \end{bmatrix}, \quad \text{and} \quad \mathbf{A}_3 = \begin{bmatrix} 1 & 2 \\ -3 & 4 \end{bmatrix}.$$

In each case we start with the random vector  $\mathbf{z} = [0.6602, 0.3420]$  and  $tol = 10^{-6}$ . 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 (19.1) 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 (19.4). 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  $\lambda$  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 19.4 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.

### 19.1.1 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  $(\lambda, \mathbf{x})$  of  $\mathbf{A}$  as long as  $\lambda$  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  $\lambda_1$  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  $\lambda_1$  the inverse power method will converge to that eigenvalue.

For the inverse power method (19.6) is replaced by

$$\begin{aligned} (i) \quad & (\mathbf{A} - s\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1} \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|. \end{aligned} \tag{19.7}$$

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 iteration.

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

$$\begin{aligned} (i) \quad & (\mathbf{A} - s_{k-1}\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1}, \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|, \\ (iii) \quad & s_k = \mathbf{x}_k^H \mathbf{A} \mathbf{x}_k, \\ (iv) \quad & \mathbf{r}_k = \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k. \end{aligned}$$

We can avoid the calculation of  $\mathbf{A} \mathbf{x}_k$  in (iii) and (iv). Let

$$\rho_k := \frac{\mathbf{y}_k^H \mathbf{x}_{k-1}}{\mathbf{y}_k^H \mathbf{y}_k}, \quad \mathbf{w}_k := \frac{\mathbf{x}_{k-1}}{\|\mathbf{y}_k\|_2}.$$

Then

$$s_k = \frac{\mathbf{y}_k^H \mathbf{A} \mathbf{y}_k}{\mathbf{y}_k^H \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^H (\mathbf{A} - s_{k-1} \mathbf{I}) \mathbf{y}_k}{\mathbf{y}_k^H \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^H \mathbf{x}_{k-1}}{\mathbf{y}_k^H \mathbf{y}_k} = s_{k-1} + \rho_k,$$

$$\mathbf{r}_k = \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k = \frac{\mathbf{A} \mathbf{y}_k - (s_{k-1} + \rho_k) \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \frac{\mathbf{x}_{k-1} - \rho_k \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \mathbf{w}_k - \rho_k \mathbf{x}_k.$$

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.

We obtain the following algorithm.

**Algorithm 19.5 (Rayleigh quotient iteration)** Given an approximation  $(s, \mathbf{x})$  to an eigenpair  $(\lambda, \mathbf{v})$  of a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$ . This algorithm computes a hopefully better approximation to  $(\lambda, \mathbf{v})$  by doing one Rayleigh quotient iteration. The length  $nr$  of the new residual is also returned

```
function [x, s, nr] = rayleight(A, x, s)
n=length(x);
y=(A-s*eye(n,n))\x;
yn=norm(y);
w=x/yn;
x=y/yn;
rho=x'*w;
s=s+rho;
nr=norm(w-rho*x);
```

Since the shift changes from iteration to iteration the computation of  $\mathbf{y}$  in **rayleight** will require  $O(n^3)$  flops for a full matrix. For such a matrix it might pay to reduce it to a 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 19.6** The matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$  has an eigenvalue  $\lambda_1 = (5 - \sqrt{33})/2 \approx -0.37$ . We test the rate of convergence by calling **rayleight** 5 times starting with  $\mathbf{x} = [1, 1]^T$  and  $s = 0$ . The results are shown in Table 19.7. The errors are approximately squared in each iteration indicating quadratic convergence.

## 19.2 The QR Algorithm

The QR algorithm is an iterative method to compute all eigenvalues and eigenvectors of a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$ . The matrix is reduced to triangular form by a sequence

$k$	1	2	3	4	5
$\ \mathbf{r}\ _2$	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 19.7.** Quadratic convergence of Rayleigh quotient iteration.

of unitary similarity transformations computed from the QR factorization of  $\mathbf{A}$ . Recall that for a square matrix the QR factorization and the QR decomposition are the same. If  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  is a QR factorization then  $\mathbf{Q} \in \mathbb{C}^{n,n}$  is unitary,  $\mathbf{Q}^H \mathbf{Q} = \mathbf{I}$  and  $\mathbf{R} \in \mathbb{C}^{n,n}$  is upper triangular.

The basic QR algorithm takes the following form:

$$\begin{array}{l}
 \mathbf{A}_1 = \mathbf{A} \\
 \text{for } k = 1, 2, \dots \\
 \quad \mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k \quad (\text{QR factorization of } \mathbf{A}_k) \\
 \quad \mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k. \\
 \text{end}
 \end{array} \tag{19.8}$$

Here are two examples to illustrate the convergence.

**Example 19.8** 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

$$\mathbf{A}_2 = \mathbf{R}_1 \mathbf{Q}_1 = \frac{1}{5} \begin{bmatrix} -5 & -4 \\ 0 & 3 \end{bmatrix} * \begin{bmatrix} -2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 2.8 & -0.6 \\ -0.6 & 1.2 \end{bmatrix}.$$

Continuing we find

$$\mathbf{A}_3 = \begin{bmatrix} 2.997 & -0.074 \\ -0.074 & 1.0027 \end{bmatrix}, \quad \mathbf{A}_9 = \begin{bmatrix} 3.0000 & -0.0001 \\ -0.0001 & 1.0000 \end{bmatrix}$$

$\mathbf{A}_9$  is almost diagonal and contains the eigenvalues  $\lambda_1 = 3$  and  $\lambda_2 = 1$  on the diagonal.

**Example 19.9** Applying the QR iteration (19.8) 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

$$A_{14} = \left[ \begin{array}{c|c|c|c} 2.323 & 0.047223 & -0.39232 & -0.65056 \\ \hline -2.1e-10 & 0.13029 & 0.36125 & 0.15946 \\ -4.1e-10 & -0.58622 & 0.052576 & -0.25774 \\ \hline 1.2e-14 & 3.3e-05 & -1.1e-05 & 0.22746 \end{array} \right].$$

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  $1 \times 1$  blocks give us two real eigenvalues  $\lambda_1 \approx 2.323$  and  $\lambda_4 \approx 0.2275$ . The middle  $2 \times 2$  block has complex eigenvalues resulting in  $\lambda_2 \approx 0.0914 + 0.4586i$  and  $\lambda_3 \approx 0.0914 - 0.4586i$ . From Gerschgorin's circle theorem 18.4 and Corollary 18.6 it follows that the approximations to the real eigenvalues are quite accurate. We would also expect the complex eigenvalues to have small absolute errors.

### 19.2.1 The Relation to the Power Method

In the basic QR algorithm we obtain the QR factorization of the powers  $\mathbf{A}^k$  as follows:

**Theorem 19.10** Let  $\mathbf{Q}_1, \dots, \mathbf{Q}_k$  and  $\mathbf{R}_1, \dots, \mathbf{R}_k$  be the matrices generated by the basic QR algorithm (19.8). Then the products

$$\tilde{\mathbf{Q}}_k := \mathbf{Q}_1 \cdots \mathbf{Q}_k \text{ and } \tilde{\mathbf{R}}_k := \mathbf{R}_k \cdots \mathbf{R}_1 \text{ for } k \geq 1 \quad (19.9)$$

are the matrices in a QR factorization  $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$  of  $\mathbf{A}^k$ .

**Proof.** 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  $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k = \tilde{\mathbf{Q}}_{k-1}^H \mathbf{A} \tilde{\mathbf{Q}}_{k-1}$  we find

$$\tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k = \tilde{\mathbf{Q}}_{k-1} (\mathbf{Q}_k \mathbf{R}_k) \tilde{\mathbf{R}}_{k-1} = \tilde{\mathbf{Q}}_{k-1} \mathbf{A}_k \tilde{\mathbf{R}}_{k-1} = (\tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{Q}}_{k-1}^H) \mathbf{A} \tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{R}}_{k-1} = \mathbf{A}^k.$$

□

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 (19.6) 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  $\lambda_1$  is a dominant eigenvalue of  $\mathbf{A}$ . This is clearly what happens in Examples 19.8 and 19.9.

### 19.2.2 A convergence theorem

There is no theorem which proves convergence of the QR algorithm in general. The following theorem shows convergence under somewhat restrictive assumptions.

**Theorem 19.11** *Suppose in the basic QR algorithm (19.8) that*

1.  $\mathbf{A} \in \mathbb{R}^{n,n}$  can be diagonalized,  $\mathbf{X}^{-1}\mathbf{A}\mathbf{X} = \mathbf{\Lambda} := \text{diag}(\lambda_1, \dots, \lambda_n)$ .
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  $\mathbf{X}^{-1} = \mathbf{L}\mathbf{R}$ .

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  $\pm 1$  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  $\mathbf{Q}$ -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  $\mathbf{X} = \mathbf{Q}\mathbf{R}$  be the QR factorization of  $\mathbf{X}$ . We observe that  $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$  is upper triangular. For since  $\mathbf{X}^{-1}\mathbf{A}\mathbf{X} = \mathbf{\Lambda}$  we have  $\mathbf{R}^{-1}\mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{R} = \mathbf{\Lambda}$  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  $\pm 1$ .

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 19.10  $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$  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  $\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k}$  is lower triangular with elements  $(\frac{\lambda_i}{\lambda_j})^k l_{ij}$  on and under the diagonal. Thus for  $n = 3$

$$\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} = \begin{bmatrix} 1 & 0 & 0 \\ (\frac{\lambda_2}{\lambda_1})^k l_{21} & 1 & 0 \\ (\frac{\lambda_3}{\lambda_1})^k l_{31} & (\frac{\lambda_3}{\lambda_2})^k l_{32} & 1 \end{bmatrix}.$$

By Assumption 2, it follows that  $\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \rightarrow \mathbf{I}$ , and hence  $\mathbf{F}_k \rightarrow \mathbf{I}$ . Since  $\hat{\mathbf{R}}_k^T \hat{\mathbf{R}}_k$  is the Cholesky factorization of  $\mathbf{F}_k^T \mathbf{F}_k$  it follows that  $\hat{\mathbf{R}}_k^T \hat{\mathbf{R}}_k \rightarrow \mathbf{I}$ . By the continuity



of the Cholesky factorization it holds  $\hat{\mathbf{R}}_k \rightarrow \mathbf{I}$  and hence  $\hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ . But then  $\hat{\mathbf{Q}}_k = \mathbf{F}_k \hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ .  $\square$

**Exercise 19.12** Use Theorem 12.33 to show that  $\hat{\mathbf{R}}_k \rightarrow \mathbf{I}$  implies  $\hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ .

### 19.2.3 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:

$$\begin{aligned}
 &\mathbf{A}_1 = \mathbf{A} \\
 &\text{for } k = 1, 2, \dots \\
 &\quad \text{Choose a shift } s_k \\
 &\quad \mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k - s_k \mathbf{I} \quad (\text{QR factorization of } \mathbf{A}_k - s_k \mathbf{I}) \\
 &\quad \mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I}. \\
 &\text{end}
 \end{aligned} \tag{19.10}$$

We will not develop the practical details of an implementation of this algorithm. We content ourselves with the following remarks. See [18] for a detailed discussion and algorithms.

1.  $\mathbf{A}_{k+1}$  is unitary similar to  $\mathbf{A}_k$ . For since  $\mathbf{R}_k = \mathbf{Q}_k^H (\mathbf{A}_k - s_k \mathbf{I})$  we find  $\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I} = \mathbf{Q}_k^H (\mathbf{A}_k - s_k \mathbf{I}) \mathbf{Q}_k + s_k \mathbf{I} = \mathbf{Q}_k^H \mathbf{A}_k \mathbf{Q}_k$ .
2. Before applying this algorithm we reduce  $\mathbf{A}$  to upper Hessenberg form using Algorithm 18.11.
3. If  $\mathbf{A}$  is upper Hessenberg then all matrices  $\{\mathbf{A}_k\}_{k \geq 1}$  will be upper Hessenberg. This follows since  $\mathbf{Q}_k = (\mathbf{A}_k - s_k \mathbf{I}) \mathbf{R}_k^{-1}$  implies  $\mathbf{A}_{k+1} = \mathbf{R}_k (\mathbf{A}_k - s_k \mathbf{I}) \mathbf{R}_k^{-1} + s_k \mathbf{I} = \mathbf{R}_k \mathbf{A}_k \mathbf{R}_k^{-1}$ . This product of two upper triangular matrices and an upper Hessenberg matrix is upper Hessenberg.
4. Givens rotations is used to compute the QR factorization of  $\mathbf{A}_k - s_k \mathbf{I}$ .
5. To compute  $\mathbf{A}_{k+1}$  from  $\mathbf{A}_k$  requires  $O(n^2)$  flops if  $\mathbf{A}_k$  is upper Hessenberg and  $O(n)$  flops if  $\mathbf{A}_k$  is tridiagonal.
6. The shifted QR algorithm is related to the power method, cf. Theorem 19.10.
7. The equation  $\mathbf{A} - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k$  implies that  $(\mathbf{A} - s_k \mathbf{I})^T \mathbf{q}_k = r_{nn}^k \mathbf{e}_n$ , where  $\mathbf{q}_k$  is the last column of  $\mathbf{Q}_k$  and  $r_{nn}^k$  is the  $(n, n)$  element in  $\mathbf{R}_k$ . Thus  $\mathbf{q}_k$  is the result of one iteration of the inverse power method to  $\mathbf{A}^T$  with shift  $s_k$ .
8. 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. This splitting occurs often in practice and can with a proper implementation reduce the computation time considerably.

9. The shift  $s_k := \mathbf{e}_n^T \mathbf{A}_k \mathbf{e}_n$  is called the **Rayleigh quotient shift**.
10. The eigenvalue of the lower right  $2 \times 2$  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.
11. The convergence is very fast and at least quadratic both for the Rayleigh quotient shift and the Wilkinson shift.
12. 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**.
13. 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  $\mathbf{A}$ .
14. 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.

# **Part VII**

## **Appendix**



## Appendix A

# Gaussian Elimination

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

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

and in matrix form

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

We recall (see Definition 3.6 and Theorem 3.7) 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 3.9) that a square matrix is invertible 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  $\mathbf{A}$ . However, for large systems it is inefficient to compute  $\mathbf{x}$  in this way. For an example see (6.30) and the discussion about the matrix  $\mathbf{T}$  there. We also note (Lemma 3.8) that if  $\mathbf{A} = \mathbf{BC}$ , where  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  are square matrices, then  $\mathbf{A}$  is nonsingular if and only if both  $\mathbf{B}$  and  $\mathbf{C}$  are nonsingular and in that case  $\mathbf{A}^{-1} = \mathbf{C}^{-1}\mathbf{B}^{-1}$ .

The elements of  $\mathbf{A}$  and  $\mathbf{b}$  can be either real or complex numbers. For simplicity and ease of exposition we assume real elements.

In Gaussian elimination with no row interchanges we compute a triangular factorization of the coefficient matrix  $\mathbf{A}$ . This factorization is known as an  $LU$ -factorization<sup>3</sup> of  $\mathbf{A}$ . In this chapter we discuss some theoretical and algorithmic

<sup>3</sup>We normally denote an upper triangular matrix by  $\mathbf{R}$ , but we respect common practice and most often we refer to the factorization  $\mathbf{A} = \mathbf{LR}$  as an LU factorization of  $\mathbf{A}$

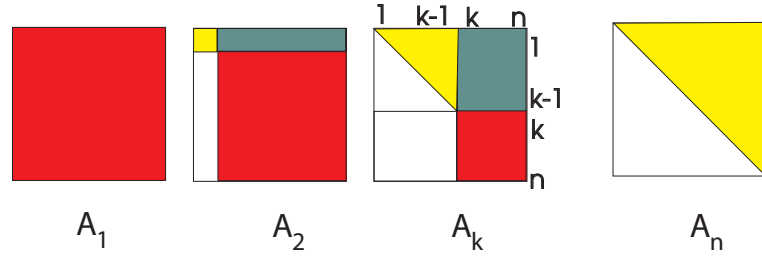


Figure A.1. Gaussian elimination

aspects of Gaussian elimination. We consider also Gaussian elimination with row interchanges.

## A.1 Gaussian Elimination and LU factorization

In **Gaussian elimination** without row interchanges we 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  $\mathbf{A}^{(n)}\mathbf{x} = \mathbf{b}^{(n)}$  is easy to solve. The process is illustrated in Figure A.1.

The matrix  $\mathbf{A}^{(k)}$  takes the form

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

The process transforming  $\mathbf{A}^{(k)}$  into  $\mathbf{A}^{(k+1)}$  for  $k = 1, \dots, n-1$  can be described as follows.

for  $i = k+1 : n$   
 $l_{ik}^k = a_{ik}^k / a_{kk}^k$   
 for  $j = k : n$   
 $a_{ij}^{k+1} = a_{ij}^k - l_{ik}^k a_{kj}^k$

(A.2)

For  $j = k$  it follows from (A.2) 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 (A.2) are called **multipliers**.

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 A.1** 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  $\mathbf{M}_k \in \mathbb{R}^{n,n}$  by

$$\mathbf{M}_k := \mathbf{I} - \begin{bmatrix} \mathbf{0} \\ \mathbf{l}_k \end{bmatrix} \mathbf{e}_k^T = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & \ddots & \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & -l_{k+1,k} & 1 & \cdots & 0 \\ \vdots & & & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -l_{n,k} & 0 & \cdots & 1 \end{bmatrix}, \quad (\text{A.3})$$

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

We have

$$\mathbf{A}^{(k+1)} = \mathbf{M}_k \mathbf{A}^{(k)}, \text{ for } k = 1, \dots, n-1, \quad (\text{A.4})$$

where  $\mathbf{M}_k \in \mathbb{R}^{n,n}$  is an elementary lower triangular matrix of the form (A.3) with  $l_{ik} = l_{ik}^k$  given by (A.2) for  $i = k+1, \dots, n$ .

**Exercise A.2** Show (A.4).

Gaussian elimination with no row interchanges is valid if and only if the **pivots**  $a_{kk}^k$  are nonzero for  $k = 1, \dots, n-1$ .

**Theorem A.3** We have  $a_{k,k}^k \neq 0$  for  $k = 1, \dots, n-1$  if and only if the leading principal submatrices

$$\mathbf{A}_k = \begin{bmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{bmatrix}$$

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

**Proof.** Let  $\mathbf{B}_k = \mathbf{A}_{k-1}^{(k)}$  be the upper left  $k-1$  corner of  $\mathbf{A}^{(k)}$  given by (A.1). Observe that the elements of the matrix  $\mathbf{B}_k$  is computed from  $\mathbf{A}$  by using only elements from  $\mathbf{A}_{k-1}$  and that only row-operations preserving nonsingularity are used. It follows that  $\mathbf{A}_{k-1}$  is nonsingular if and only if  $\mathbf{B}_k$  is nonsingular. By Lemma 6.32  $\mathbf{B}_k$  is nonsingular if and only if  $a_{ii}^{(i)} \neq 0$ ,  $i = 1, \dots, k-1$ . We conclude that  $\mathbf{A}_1, \dots, \mathbf{A}_{n-1}$  are nonsingular if and only if  $\mathbf{B}_2, \dots, \mathbf{B}_n$  are nonsingular which is equivalent to  $a_{kk}^{(k)} \neq 0$  for  $k = 1, \dots, n-1$ .  $\square$

Gaussian elimination is a way to compute the LU factorization of the coefficient matrix.

**Theorem A.4** Suppose  $\mathbf{A} \in \mathbb{R}^{n,n}$  and that  $\mathbf{A}_k$  is nonsingular for  $k = 1, \dots, n-1$ . Then Gaussian elimination with no row interchanges results in an LU factorization of  $\mathbf{A} \in \mathbb{R}^{n,n}$ . In particular  $\mathbf{A} = \mathbf{L}\mathbf{R}$ , where

$$\mathbf{L} = \begin{bmatrix} 1 & & & \\ l_{21}^1 & 1 & & \\ \vdots & & \ddots & \\ l_{n1}^1 & l_{n2}^2 & \cdots & 1 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} a_{11}^1 & \cdots & a_{1n}^1 \\ & \ddots & \vdots \\ & & a_{nn}^n \end{bmatrix}, \quad (\text{A.5})$$

where the  $l_{ij}^j$  and  $a_{ij}^i$  are given by (A.2).

**Proof.** From (A.2) we have for all  $i, j$

$$l_{ik}a_{kj}^k = a_{ij}^k - a_{ij}^{k+1} \text{ for } k < \min(i, j), \text{ and } l_{ij}a_{jj}^j = a_{ij}^j \text{ for } i > j.$$

Thus for  $i \leq j$  we find

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

while for  $i > j$

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

□

Note that this Theorem holds even if  $\mathbf{A}$  is singular. Since  $\mathbf{L}$  is nonsingular the matrix  $\mathbf{R}$  is singular, and we must have  $a_{nn}^n = 0$  when  $\mathbf{A}$  is singular.

### A.1.1 Algorithms

Consider next an algorithm to find the LU factorization of  $\mathbf{A}$  using Gaussian elimination with no row interchanges. Storing both the elements  $l_{ij}^j$  and  $a_{ij}^i$  in  $\mathbf{A}$  we can write (A.2) as follows for  $k = 1, \dots, n-1$ .

for  $i = k+1 : n$   
 $a_{ik} = a_{ik}/a_{kk}$   
 for  $j = k+1 : n$   
 $a_{ij} = a_{ij} - a_{ik}a_{kj}$

(A.6)

We can write (A.6) using outer product notation. We have

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



The result is a matrix of order  $n - k$ .

This leads to the following algorithm.

**Algorithm A.5 (lufactor)** Given  $\mathbf{A} \in \mathbb{R}^{n,n}$  with  $\mathbf{A}_k \in \mathbb{R}^{k,k}$  nonsingular for  $k = 1, \dots, n-1$ . This algorithm computes an LU factorization of  $\mathbf{A}$  using Gaussian elimination without row interchanges.

```
function [L,R]=lufactor(A)
n=length(A); for k=1:n-1
    kn=k+1:n;
    A(kn,k)=A(kn,k)/A(k,k);
    A(kn,kn)=A(kn,kn)-A(kn,k)*A(k,kn);
end
L=eye(n,n)+tril(A,-1);
R=triu(A);
```

Once we have an LU factorization of  $\mathbf{A}$  the system  $\mathbf{Ax} = \mathbf{b}$  is solved easily in two steps. Since  $\mathbf{LRx} = \mathbf{b}$  we have  $\mathbf{Ly} = \mathbf{b}$ , where  $\mathbf{y} := \mathbf{Rx}$ . We first solve  $\mathbf{Ly} = \mathbf{b}$  for  $\mathbf{y}$  and then  $\mathbf{Rx} = \mathbf{y}$  for  $\mathbf{x}$ . Consider solving a system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is lower triangular with nonzero diagonal elements. For  $n = 3$  we have

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

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 and we arrive at the following algorithm.

**Algorithm A.6 (forwardsolve)** Given a nonsingular lower triangular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  and  $\mathbf{b} \in \mathbb{R}^n$ . An  $\mathbf{x} \in \mathbb{R}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```
function x=forwardsolve(A,b)
n=length(b); x=b(:);
for k=1:n
    x(k)=(x(k)-A(k,1:k-1)*x(1:k-1))/A(k,k);
end
```

A system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is upper triangular must be solved 'bottom-up'. We first find  $x_n$  from the last equation and then move upwards for the remaining unknowns. We have the following algorithm.

**Algorithm A.7 (backsolve)** Given a nonsingular upper triangular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  and  $\mathbf{b} \in \mathbb{R}^n$ . An  $\mathbf{x} \in \mathbb{R}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```
function x=backsolve(A,b)
n=length(b); x=b(:);
for k=n:-1:1
    x(k)=(x(k)-A(k,k+1:n)*x(k+1:n))/A(k,k);
end
```

### A.1.2 Operation Count

We define a **flop** (floating point operation) as one of the floating point arithmetic operations, ie. multiplication, division, addition and subtraction. We denote by **nflops** the total number of flops in an algorithm, i.e. the the sum of all multiplications, divisions, additions and subtractions. For a problem of size  $n$  the number nflops will often be a polynomial in  $n$ . For example, we will show below that an LU factorization requires  $\frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{n}{6}$  flops. For large values of  $n$  the highest term  $\frac{2}{3}n^3$  dominates and we usually say that  $\text{nflops} = O(\frac{2}{3}n^3)$  ignoring lower order terms. We sometimes say that  $\text{nflops} = O(n^3)$  if we do not bother with the constant (in this case  $2/3$ ) in front of the  $n^3$  term.

In many implementations the computing time  $T_A$  for an algorithm  $A$  applied to a large problem is proportional to  $N_A := \text{nflops}$ . If this is true then we typically have  $T_A = \alpha N_A$ , where  $\alpha$  is in the range  $10^{-12}$  to  $10^{-9}$  on a modern computer.

Consider now  $N_{LU} := \text{nflops}$  for LU factorization. Let  $M, D, A, S$  be the number of multiplications, divisions, additions, and subtractions. We first do an exact count. From (A.6) we find

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

Thus

$$N_{LU} = M + D + A + S = \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n = O(\frac{2}{3}n^3)$$

There is a quick way to arrive at the leading term  $2n^3/3$ . We only consider the arithmetic operations contributing to the leading term. 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  $N_{LU} = O(2n^3/3)$  which is reasonably correct for large values of  $n$ .

Consider next forward and backward substitution. Counting flops and letting  $N_S := N_F + N_B$  we find

$$N_S \approx \int_1^n 2(k-1)dk + \int_1^n 2(n-k)dk \approx \int_0^n 2kdk + \int_0^n 2(n-k)dk = 2n^2.$$

Comparing  $N_{LU}$  and  $N_S$  we see that LU factorization is an  $O(n^3)$  process while the solution stage only require  $O(n^2)$  flops. This leads to dramatic differences in computing time as illustrated in the following table:

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

Here we have assumed that the computing time for the LU factorization is  $T_{LU} = 10^{-9}n^3$  and the computing time for the forward and backwards substitution is  $T_S = 3 \times 10^{-9}n^2$  corresponding to  $\alpha = 3 \times 10^{-9}/2$ .

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  $\mathbf{A}_j \in \mathbb{R}^{n,n}$  and  $\mathbf{b}_j \in \mathbb{R}^n$ . We need  $m(\frac{2}{3}n^3 + 2n^2)$  flops 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  $\mathbf{A}\mathbf{X} = \mathbf{B}$ , where  $\mathbf{A} \in \mathbb{R}^{n,n}$ ,  $\mathbf{B} \in \mathbb{R}^{n,m}$  and the matrix  $\mathbf{X} \in \mathbb{R}^{n,m}$  is the unknown. To solve  $\mathbf{A}\mathbf{X} = \mathbf{B}$  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.

## A.2 Pivoting

We have seen that Gaussian elimination without row interchanges is only well defined if the leading principal minors  $\mathbf{A}_k \in \mathbb{R}^{k,k}$  are nonsingular for  $k = 1, \dots, n-1$ .

Suppose now  $\mathbf{A} \in \mathbb{R}^{n,n}$  is nonsingular. We can still solve a linear system with  $\mathbf{A}$  if we incorporate row interchanges. 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. Gaussian elimination with row interchanges 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 (A.2).

We have seen that the elimination step can be described as multiplying the current matrix by an elementary transformation matrix  $\mathbf{M}_k$  given by (A.3). But before we can multiply by  $\mathbf{M}_k$  we have to interchange rows. This can be described in terms of permutation matrices.

### A.2.1 Permutation Matrices

**Definition A.8** A **permutation matrix** is a matrix of the form

$$P = I(:, \mathbf{p}) = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_n}] \in \mathbb{R}^{n,n},$$

where  $\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n}$  is a permutation of the unit vectors  $\mathbf{e}_1, \dots, \mathbf{e}_n \in \mathbb{R}^n$ .

Every permutation  $\mathbf{p} = [i_1, \dots, i_n]^T$  of the integers  $1, 2, \dots, n$  gives rise to a permutation matrix and vice versa. Post-multiplying a matrix  $\mathbf{A}$  by a permutation matrix results in a permutation of the columns, while pre-multiplying by a permutation matrix gives a permutation of the rows. In symbols

$$\mathbf{A}\mathbf{P} = \mathbf{A}(:, \mathbf{p}), \quad \mathbf{P}^T \mathbf{A} = \mathbf{A}(\mathbf{p}, :). \quad (\text{A.7})$$

Indeed,  $\mathbf{A}\mathbf{P} = (\mathbf{A}\mathbf{e}_{i_1}, \dots, \mathbf{A}\mathbf{e}_{i_n}) = \mathbf{A}(:, \mathbf{p})$  and  $\mathbf{P}^T \mathbf{A} = (\mathbf{A}^T \mathbf{P})^T = (\mathbf{A}^T(:, \mathbf{p}))^T = \mathbf{A}(\mathbf{p}, :)$ .

Since  $\mathbf{P}^T \mathbf{P} = \mathbf{I}$  the inverse of  $\mathbf{P}$  is equal to its transpose,  $\mathbf{P}^{-1} = \mathbf{P}^T$  and  $\mathbf{P}\mathbf{P}^T = \mathbf{I}$  as well. Thus a permutation matrix is an orthonormal matrix.

We will use a particularly simple permutation matrix.

**Definition A.9** 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.

### A.2.2 Gaussian Elimination Works Mathematically

The process going from  $\mathbf{A}^{(k)}$  to  $\mathbf{A}^{(k+1)}$  can be written

$$\mathbf{A}^{(k+1)} = \mathbf{M}_k \mathbf{P}_k \mathbf{A}^{(k)}, \text{ for } k = 1, \dots, n-1, \quad (\text{A.8})$$

where  $\mathbf{P}_k = \mathbf{I}_{r_k, k} \in \mathbb{R}^{n,n}$  is a permutation matrix interchanging rows  $k$  and  $r_k$  of  $\mathbf{A}^{(k)}$  and  $\mathbf{M}_k \in \mathbb{R}^{n,n}$  is an elementary lower triangular matrix of the form (A.3) with  $l_{ik} = l_{ik}^k$  given by (A.2) for  $i = k+1, \dots, n$ .

If  $\mathbf{A}$  is nonsingular then Gaussian elimination can always be carried to completion by using suitable row interchanges. To show this, suppose by induction on  $k$  that  $\mathbf{A}^{(k)}$  is nonsingular. Since  $\mathbf{A}^{(1)} = \mathbf{A}$  this holds for  $k = 1$ . By Lemma 6.31 the lower right diagonal block in  $\mathbf{A}^{(k)}$  is nonsingular. But then at least one element in the first column of that block must be nonzero and it follows that  $r_k$  exists so that  $a_{r_k, k}^k \neq 0$ . But then the matrices  $\mathbf{P}_k$  and  $\mathbf{M}_k$  in (A.8) are well defined. By Lemma 6.32 the matrix  $\mathbf{M}_k$  is nonsingular and since a permutation matrix is nonsingular it follows from Lemma 3.8 that  $\mathbf{A}^{(k+1)}$  is nonsingular. We conclude that  $\mathbf{A}^{(k)}$  is nonsingular for  $k = 1, \dots, n$ .

### A.2.3 Pivot Strategies

Up to now we have said nothing about what rows in  $\mathbf{A}$  to interchange during the elimination. We start with an example illustrating that small pivots should be avoided.

**Example A.10** *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.

We briefly describe the two most common pivoting strategies. The choice

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

with  $r_k$  the smallest such index in case of a tie, is known as **partial pivoting**. It 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 stable, but requires a lot of search and is seldom used in practice.

### A.3 The PLU Factorization

Consider now Gaussian elimination with row pivoting. We can keep track of the row interchanges using **pivot vectors**  $\mathbf{p}_k$ . We define

$$\mathbf{p} := \mathbf{p}_n, \text{ where } \mathbf{p}_1 := [1, 2, \dots, n]^T, \text{ and } \mathbf{p}_{k+1} := \mathbf{I}_{r_k, k} \mathbf{p}_k \text{ for } k = 1, \dots, n-1. \quad (\text{A.9})$$

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

There is a close relation between the pivot vectors  $\mathbf{p}_k$  and the corresponding interchange matrices  $\mathbf{P}_k := \mathbf{I}_{r_k, k}$ . Since  $\mathbf{P}_k \mathbf{I}(\mathbf{p}_k, :) = \mathbf{I}(\mathbf{P}_k \mathbf{p}_k, :) = \mathbf{I}(\mathbf{p}_{k+1}, :)$  we obtain

$$\mathbf{P}^T := \mathbf{P}_{n-1} \cdots \mathbf{P}_1 = \mathbf{I}(\mathbf{p}, :), \quad \mathbf{P} := \mathbf{P}_1 \mathbf{P}_2 \cdots \mathbf{P}_{n-1} = \mathbf{I}(:, \mathbf{p}). \quad (\text{A.10})$$

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 Gaussian elimination in Section A.1 with elements  $a_{ij}^1$  can be described as follows:

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

(A.11)

This leads to the following factorization:

**Theorem A.11** *Gaussian elimination with row pivoting on a nonsingular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  leads to a factorization  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{R}$ , where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is lower triangular with ones on the diagonal, and  $\mathbf{R}$  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

$$\mathbf{L} = \begin{bmatrix} 1 & & & \\ a_{p_2, 1}^1 & 1 & & \\ \vdots & & \ddots & \\ a_{p_n, 1}^1 & a_{p_n, 2}^2 & \cdots & 1 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} a_{p_1, 1}^1 & \cdots & a_{p_1, n}^1 \\ & \ddots & \vdots \\ & & a_{p_n, n}^n \end{bmatrix}, \quad (\text{A.12})$$

**Proof.** The proof is analogous to the proof for LU factorization without pivoting. From (A.11) we have for all  $i, j$

$$a_{p_i, k}^k a_{p_k, j}^k = a_{p_i, j}^k - a_{p_i, j}^{k+1} \text{ for } k < \min(i, j), \text{ and } a_{p_i, j}^k a_{p_j, j}^j = a_{p_i, j}^j \text{ for } i > j.$$

Thus for  $i \leq j$  we find

$$\begin{aligned} (\mathbf{LR})_{ij} &= \sum_{k=1}^n l_{i,k} u_{kj} = \sum_{k=1}^{i-1} a_{p_i, k}^k a_{p_k, j}^k + a_{p_i, j}^i \\ &= \sum_{k=1}^{i-1} (a_{p_i, j}^k - a_{p_i, j}^{k+1}) + a_{p_i, j}^i = a_{p_i, j}^1 = a_{p_i, j} = (\mathbf{P}^T \mathbf{A})_{ij}, \end{aligned}$$

while for  $i > j$

$$\begin{aligned} (\mathbf{LR})_{ij} &= \sum_{k=1}^n l_{i,k} u_{kj} = \sum_{k=1}^{j-1} a_{p_i, k}^k a_{p_k, j}^k + a_{p_i, j}^j a_{p_j, j}^j \\ &= \sum_{k=1}^{j-1} (a_{p_i, j}^k - a_{p_i, j}^{k+1}) + a_{p_i, j}^j = a_{p_i, j}^1 = a_{p_i, j} = (\mathbf{P}^T \mathbf{A})_{ij}. \end{aligned}$$

□

## A.4 An Algorithm for Finding the PLU Factorization

Using pivot vectors we can compute the PLU factorization of  $\mathbf{A}$  without physically interchanging the elements  $a_{ij}^k$ . As is clear from (A.12) we can store the elements of  $\mathbf{L}$  and  $\mathbf{R}$  in  $\mathbf{A}$  and work with  $\mathbf{A}(\mathbf{p}_k, :)$ . At the end the elements of  $\mathbf{L}$  and  $\mathbf{R}$  will be located under and above the diagonal.

In the following algorithm we use partial pivoting.

**Algorithm A.12 (PLU factorization)** Given a nonsingular  $\mathbf{A} \in \mathbb{R}^{n,n}$ . 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  $\mathbf{P} = \mathbf{I}(:, \mathbf{p})$ .

```
function [p,L,R] = plufactor(A)
n = length(A);
p = 1:n;
for k=1:n-1
    [maxv, r] = max(abs(A(p(k:n), k)));
    p([k r+k-1]) = p([r+k-1 k]);
    ps=p(k+1:n);
    A(ps, k) = A(ps, k)/A(p(k), k);
    A(ps, k+1:n) = A(ps, k+1:n) - A(ps, k)*A(p(k), k+1:n);
end
L = eye(n,n) + tril(A(p,:), -1);
R = triu(A(p,:));
```

Once we have a PLU factorization of  $\mathbf{A}$  the system  $\mathbf{Ax} = \mathbf{b}$  is solved easily in three steps. Since  $\mathbf{PLRx} = \mathbf{b}$  we have  $\mathbf{Pz} = \mathbf{b}$ ,  $\mathbf{Ly} = \mathbf{z}$ , and  $\mathbf{Rx} = \mathbf{y}$ . Using the output  $[p, L, R]$  of Algorithm A.12 the solution can be found from Algorithms A.6 and A.7 in two steps.

1.  $\mathbf{y} = \text{forwardsolve}(\mathbf{L}, \mathbf{b}(\mathbf{p}))$ ;
2.  $\mathbf{x} = \text{backsolve}(\mathbf{R}, \mathbf{y})$ ;

**Exercise A.13** In this exercise we develop column oriented vectorized versions of forward and backward substitution. Suppose  $\mathbf{L} \in \mathbb{R}^{n,n}$  is lower triangular and  $\mathbf{R} \in \mathbb{R}^{n,n}$  is upper triangular. Consider the system  $\mathbf{Lx} = \mathbf{b}$ . Suppose after  $k-1$  steps of the algorithm we have a reduced system in the form

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

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

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

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



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

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

**b)** Suppose now  $L \in \mathbb{R}^{n,n}$  is lower triangular,  $R \in \mathbb{R}^{n,n}$  is upper triangular and  $\mathbf{b} \in \mathbb{R}^n$ . Justify the following column oriented vectorized algorithms for solving  $Lx = b$  and  $Rx = b$ .

**Algorithm A.14 (Forward Substitution (column oriented))**

```

for  $k = 1 : n$ 
     $x(k) = b(k)/L(k, k);$ 
     $b(k+1:n) = b(k+1:n) - L(k+1:n, k) * x(k);$ 
end

```

**Algorithm A.15 (Backward Substitution (column oriented))**

```

for  $k = n : -1 : 1$ 
     $x(k) = b(k)/R(k, k);$ 
     $b(1:k-1) = b(1:k-1) - R(1:k-1, k) * x(k);$ 
end

```

Each algorithm requires  $n^2$  flops.



## 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  $\epsilon$  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** *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  $\rho_b$  is small then  $\rho_a$  is small and it does not matter whether we choose  $\rho_a$  or  $\rho_b$  to discuss relative error.

**Exercise B.4** Compare  $\rho_a$  and  $\rho_b$  when  $a := 3.1415.9265$  and  $b := 31415.8951$ .

## 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  $\beta$ . 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

$$(0.1)_{10} = (1.100110011001\dots)_2 \cdot 2^{-4}.$$

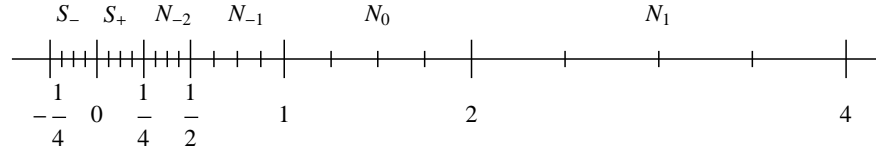
We turn now to a description of the floating-point numbers. We will only describe a **standard system**, namely the binary IEEE floating-point standard. Although it is not used by all systems, it has been widely adopted and is used in Matlab. For a more complete introduction to the subject see [8],[17].

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

$$\epsilon_M := 2^{-t}, \quad \text{machine epsilon,} \tag{B.1}$$

and

$$\begin{aligned}\mathcal{F} &:= \{0\} \cup \mathcal{S} \cup \mathcal{N}, \text{ where} \\ \mathcal{N} &:= \mathcal{N}_+ \cup \mathcal{N}_-, \quad \mathcal{N}_+ := \bigcup_{e=\underline{e}}^{\bar{e}} \mathcal{N}_e, \quad \mathcal{N}_- := -\mathcal{N}_+, \\ \mathcal{N}_e &:= \{(1.d_1 d_2 \dots d_t)_2\} * 2^e = \{1, 1 + \epsilon_M, 1 + 2\epsilon_M, \dots, 2 - \epsilon_M\} * 2^e, \\ \mathcal{S} &:= \mathcal{S}_+ \cup \mathcal{S}_-, \quad \mathcal{S}_+ := \{\epsilon_M, 2\epsilon_M, 3\epsilon_M, \dots, 1 - \epsilon_M\} * 2^{\underline{e}}, \quad \mathcal{S}_- := -\mathcal{S}_+.\end{aligned} \tag{B.2}$$



**Figure B.1.** *Distribution of some positive floating-point numbers*

**Example B.5** Suppose  $t := 2$ ,  $\bar{e} = 3$  and  $\underline{e} := -2$ . Then  $\epsilon_M = 1/4$  and we find

$$\begin{aligned}\mathcal{N}_{-2} &= \left\{\frac{1}{4}, \frac{5}{16}, \frac{3}{8}, \frac{7}{16}\right\}, & \mathcal{N}_{-1} &= \left\{\frac{1}{2}, \frac{5}{8}, \frac{3}{4}, \frac{7}{8}\right\}, & \mathcal{N}_0 &= \left\{1, \frac{5}{4}, \frac{3}{2}, \frac{7}{4}\right\}, \\ \mathcal{N}_1 &= \left\{2, \frac{5}{2}, 3, \frac{7}{2}\right\}, & \mathcal{N}_2 &= \{4, 5, 6, 7\}, & \mathcal{N}_3 &= \{8, 10, 12, 14\}, \\ \mathcal{S}_+ &= \left\{\frac{1}{16}, \frac{1}{8}, \frac{3}{16}\right\}, & \mathcal{S}_- &= \left\{-\frac{3}{16}, -\frac{1}{8}, -\frac{1}{16}\right\}.\end{aligned}$$

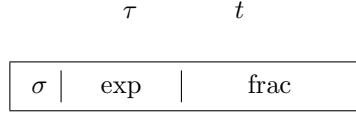
The position of some of these sets on the real line is shown in Figure B.1

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,  $+\text{Inf}$  and  $-\text{Inf}$ .
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,  $\tau$  bits for the exponent, and  $t$  bits for the fractional part of the mantissa.



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  $\text{frac}$  is the fractional part  $d_1 d_2 \dots d_t$  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  $\text{exp}$ , 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  $\tau$  and  $t$  are shown in the following table

precision	$\tau$	$t$	$b$	$\epsilon_M = 2^{-t}$	$r_m = 2^{1-b}$	$r_M$
half	5	10	15	$9.8 \times 10^{-4}$	$6.1 \times 10^{-5}$	$6.6 \times 10^4$
single	8	23	127	$1.2 \times 10^{-7}$	$1.2 \times 10^{-38}$	$3.4 \times 10^{38}$
double	11	52	1023	$2.2 \times 10^{-16}$	$2.2 \times 10^{-308}$	$1.8 \times 10^{308}$
quad	15	112	16383	$1.9 \times 10^{-34}$	$3.4 \times 10^{-4932}$	$1.2 \times 10^{4932}$

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**).

**Exercise B.6** Check the results of the following operations on your computer.  $1^{\text{Inf}}$ ,  $2^{\text{Inf}}$ ,  $e^{-\text{Inf}}$ ,  $\text{Inf}^0$ ,  $\log 0$ ,  $\sin(\text{Inf})$ ,  $\arctan(-\text{Inf})$ .

## 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  $\text{frac} = 0$ .

The following lemma gives a bound for the relative error in rounding.

**Theorem B.7** *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**.

**Exercise B.8** *Show that the upper bound for  $\delta$  is attained for  $x = (1 + \epsilon_M/2) * 2^e$  when round to zero is used. Compute  $\delta$  when  $x = (2 - \epsilon_M/2) * 2^e$ .*

### 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  $\alpha_1, \alpha_2$ , and in this context we call  $\delta$  the **backward error**. A similar interpretation is valid for the other arithmetic operations  $-, *, /, \sqrt{\phantom{x}}$ , 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  $\beta$  as an approximation to  $\alpha$  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  $\eta$  factors.



**Lemma B.9** Suppose for integers  $k, m$  with  $0 \leq m \leq k$  and  $k \geq 1$  that

$$1 + \eta_k := \frac{(1 + \delta_1) \cdots (1 + \delta_m)}{(1 + \delta_{m+1}) \cdots (1 + \delta_k)}, \quad |\delta_j| \leq u_M, \quad j = 1, \dots, k.$$

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  $\alpha$  can be found. [17] uses  $\alpha = 1/10$  giving  $u'_M = 1.12u_M$ , while in [8] the value  $\alpha = 0.01$  can be found. In the classical work [25] 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

$$|\eta_1| \leq (n - 1)u'_M, \quad |\eta_i| \leq (n - i + 1)u'_M, \quad \text{for } i = 2, \dots, n.$$

or more simply

$$|\eta_i| \leq (n - 1)u'_M, \quad \text{for } i = 1, \dots, n. \quad (\text{B.11})$$

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  $s = \alpha_1 + \cdots + \alpha_n$  from the computed sum  $\tilde{s} = \alpha_1(1 + \eta_1) + \cdots + \alpha_n(1 + \eta_n)$ , to get

$$|\tilde{s} - s| = |\alpha_1\eta_1 + \cdots + \alpha_n\eta_n| \leq (|\alpha_1| + \cdots + |\alpha_n|)(n-1)u'_M.$$

Thus the relative error in the computed sum of  $n$  numbers is bounded as follows

$$\left| \frac{\tilde{s} - s}{s} \right| \leq \kappa(n-1)u'_M, \text{ where } \kappa := \frac{|\alpha_1| + \cdots + |\alpha_n|}{\alpha_1 + \cdots + \alpha_n}. \quad (\text{B.12})$$

This bound shows that the backward error can be magnified by at most  $\kappa$ . The number  $\kappa$  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  $|\tilde{c}_{ij} - c_{ij}| \leq \kappa_{ij} |c_{ij}| n u'_M$ . Using the infinity matrix norm we find

$$\sum_j |\tilde{c}_{ij} - c_{ij}| \leq n u'_M \sum_j \kappa_{ij} |c_{ij}| \leq \kappa n u'_M \sum_j |c_{ij}| \leq \kappa n u'_M \|\mathbf{C}\|_\infty, \quad \text{all } i,$$

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.

**Exercise C.1** For  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$  show 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** Suppose  $f \in C^2(\Omega)$ , where  $\Omega \in \mathbb{R}^n$  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

$$\begin{aligned}g(0) &= f(\mathbf{x}) \quad g(1) = f(\mathbf{x} + \mathbf{h}), \\ g'(t) &= \sum_{i=1}^n h_i \frac{\partial f(\mathbf{x} + t\mathbf{h})}{\partial x_i} = \mathbf{h}^T \nabla f(\mathbf{x} + t\mathbf{h}), \\ g''(t) &= \sum_{i=1}^n \sum_{j=1}^n h_i h_j \frac{\partial^2 f(\mathbf{x} + t\mathbf{h})}{\partial x_i \partial x_j} = \mathbf{h}^T \nabla \nabla^T f(\mathbf{x} + t\mathbf{h}) \mathbf{h}.\end{aligned}$$

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** 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.

□





## Appendix D

# Some Inequalities

In this appendix we derive an inequality for convex functions called Jensen's inequality and use it to show Hölder's and Minkowski's inequalities.

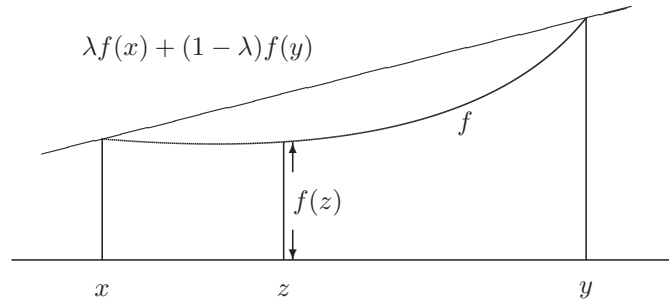
### D.1 Convexity

**Definition D.1 (Convex function)** Let  $I \subset \mathbb{R}$  be an interval. A function  $f : I \rightarrow \mathbb{R}$  is called *convex* if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for all  $x, y \in I$  and all  $\lambda \in [0, 1]$ . The sum  $\sum_{j=1}^n \lambda_j z_j$  is called a **convex combination** of  $z_1, \dots, z_n$  if  $\lambda_j \geq 0$  for  $j = 1, \dots, n$  and  $\sum_{j=1}^n \lambda_j = 1$ .

The condition is shown graphically in Figure D.1.



**Figure D.1.** A convex function.

We state without proof that if  $f \in C^2(I)$  and  $f''(x) \geq 0$  for  $x \in I$  then  $f$  is convex. It then follows that the function  $-\log x$  is convex on  $I = (0, \infty)$ .

## D.2 Inequalities

**Theorem D.2 (Jensen's Inequality)** Suppose  $I \subset \mathbb{R}$  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  $z_1, \dots, z_n \in I$  we have

$$f\left(\sum_{j=1}^n \lambda_j\right) \leq \sum_{j=1}^n \lambda_j f(z_j).$$

**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  $\lambda_j, z_j$  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 D.3 (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

$$a_1^{\lambda_1} a_2^{\lambda_2} \cdots a_n^{\lambda_n} \leq \sum_{j=1}^n \lambda_j a_j, \quad (\text{D.1})$$

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  $j$ . We use Jensen's inequality with the convex function  $f(x) = -\log x$  on  $I = (0, \infty)$ . Then

$$-\log\left(\sum_{j=1}^n \lambda_j a_j\right) \leq -\sum_{j=1}^n \lambda_j \log(a_j) = -\log(a_1^{\lambda_1} \cdots a_n^{\lambda_n})$$

and since the log function is monotone the inequality follows.  $\square$

Taking  $\lambda_j = \frac{1}{n}$  for all  $j$  in (D.1) we obtain the classical **geometric/arithmetic mean inequality**

$$(a_1 a_2 \cdots a_n)^{\frac{1}{n}} \leq \frac{1}{n} \sum_{j=1}^n a_j. \quad (\text{D.2})$$

**Corollary D.4 (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

$$a^{\frac{1}{p}} b^{\frac{1}{q}} \leq \frac{1}{p} a + \frac{1}{q} b, \text{ where } \frac{1}{p} + \frac{1}{q} = 1. \quad (\text{D.3})$$

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 D.3 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 D.5 (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.  $\square$



## Appendix E

# The Jordan Form

### E.1 The Jordan Form

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?

**Definition E.1** A **Jordan block**, denoted  $\mathbf{J}_m(\lambda)$  is an  $m \times m$  matrix of the form

$$\mathbf{J}_m(\lambda) := \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda & 1 & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \vdots & \\ 0 & 0 & 0 & \cdots & \lambda & 1 \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix}$$

A  $3 \times 3$  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  $\lambda$  is an eigenvalue of  $\mathbf{J}_m(\lambda)$  and any eigenvector must be a multiple of  $\mathbf{e}_1$ . 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.

**Theorem E.2** Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  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  $\mathbf{S} \in \mathbb{C}^{n,n}$  such that

$$\mathbf{J} := \mathbf{S}^{-1} \mathbf{A} \mathbf{S} = \text{diag}(\mathbf{U}_1, \dots, \mathbf{U}_k), \text{ with } \mathbf{U}_i \in \mathbb{C}^{a_i, a_i}, \quad (\text{E.1})$$

where each  $\mathbf{U}_i$  is block diagonal having  $g_i$  Jordan blocks along the diagonal

$$\mathbf{U}_i = \text{diag}(\mathbf{J}_{m_{i,1}}(\lambda_i), \dots, \mathbf{J}_{m_{i,g_i}}(\lambda_i)). \quad (\text{E.2})$$

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 (E.1) is called the **Jordan form** of  $\mathbf{A}$ . As an example consider the Jordan form

$$\mathbf{J} := \text{diag}(\mathbf{U}_1, \mathbf{U}_2) = \begin{bmatrix} 2 & 1 & 0 & & & & & \\ 0 & 2 & 1 & & & & & \\ 0 & 0 & 2 & & & & & \\ & & & 2 & 1 & & & \\ & & & 0 & 2 & & & \\ & & & & & 2 & & \\ & & & & & & 3 & 1 \\ & & & & & & 0 & 3 \end{bmatrix} \in \mathbb{R}^{8,8}. \quad (\text{E.3})$$

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  $\lambda_i$  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

$$\mathbf{J} := \begin{bmatrix} 3 & 1 & & & & & & \\ 0 & 3 & & & & & & \\ & & 2 & 1 & & & & \\ & & 0 & 2 & & & & \\ & & & & 2 & & & \\ & & & & & 2 & 1 & 0 \\ & & & & & 0 & 2 & 1 \\ & & & & & 0 & 0 & 2 \end{bmatrix}$$

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 (E.3) we have  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_8]$  and we find

$$\begin{aligned} \mathbf{As}_1 &= 2\mathbf{s}_1, & \mathbf{As}_2 &= 2\mathbf{s}_2 + \mathbf{s}_1, \\ \mathbf{As}_3 &= 2\mathbf{s}_3, \\ \mathbf{As}_4 &= 2\mathbf{s}_4, & \mathbf{As}_5 &= 2\mathbf{s}_5 + \mathbf{s}_4, & \mathbf{As}_6 &= 2\mathbf{s}_6 + \mathbf{s}_5, \\ \mathbf{As}_7 &= 3\mathbf{s}_7, & \mathbf{As}_8 &= 3\mathbf{s}_8 + \mathbf{s}_7, \end{aligned}$$

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 E.3** For the Jordan form of the matrix  $\mathbf{A} = \begin{bmatrix} 3 & 0 & 1 \\ -4 & 1 & -2 \\ -4 & 0 & -1 \end{bmatrix}$  we have  $\mathbf{J} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ . Find  $\mathbf{S}$ .

**Exercise E.4** Find the Jordan form of the matrix

$$\mathbf{A} = \frac{1}{9} \begin{bmatrix} 10 & 16 & -8 & -5 & 6 & 1 & -3 & 4 \\ -7 & 32 & -7 & -10 & 12 & 2 & -6 & 8 \\ -6 & 12 & 12 & -15 & 18 & 3 & -9 & 12 \\ -5 & 10 & -5 & -2 & 24 & 4 & -12 & 16 \\ -4 & 8 & -4 & -16 & 30 & 14 & -15 & 20 \\ -3 & 6 & -3 & -12 & 9 & 24 & -9 & 24 \\ -2 & 4 & -2 & -8 & 6 & -2 & 15 & 28 \\ -1 & 2 & -1 & -4 & 3 & -1 & -6 & 41 \end{bmatrix}. \quad (\text{E.4})$$

The following lemma is useful when studying powers of matrices.

**Lemma E.5** *Let  $\mathbf{J}$  be the Jordan form of a matrix  $\mathbf{A} \in \mathbb{C}^{n,n}$  as given in Theorem E.2. Then for  $r = 0, 1, 2, \dots$ ,  $m = 2, 3, \dots$ , and any  $\lambda \in \mathbb{C}$*

1.  $\mathbf{A}^r = \mathbf{S}\mathbf{J}^r\mathbf{S}^{-1}$ ,
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  $\mathbf{A}^2 = \mathbf{S}\mathbf{J}\mathbf{S}^{-1}\mathbf{S}\mathbf{J}\mathbf{S}^{-1} = \mathbf{S}\mathbf{J}^2\mathbf{S}^{-1}$  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

$$\mathbf{E}_m = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots & \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{m-1} \\ \mathbf{0} & \mathbf{0}^T \end{bmatrix}. \quad (\text{E.5})$$

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  $\mathbf{E}_m^m = \mathbf{0}$ .

□

**Exercise E.6** Determine  $\mathbf{J}_3^r$  for  $r \geq 1$ .

**Exercise E.7** Find  $\mathbf{J}^{100}$  and  $\mathbf{A}^{100}$  for the matrix in Exercise E.3.

### E.1.1 The Minimal Polynomial

Let  $\mathbf{J}$  be the Jordan form of  $\mathbf{A}$  given in Theorem E.2. 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 E.8** 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 E.9** What is the characteristic polynomial and the minimal polynomial of the matrix  $\mathbf{J}$  in (E.3)?

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

$$\pi_{\mathbf{A}}(\mathbf{A}) := \prod_{i=1}^k \prod_{j=1}^{g_i} (\lambda_i \mathbf{I} - \mathbf{A})^{m_{ij}}, \quad \mu(\mathbf{A}) := \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{A})^{m_i}. \quad (\text{E.6})$$

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,  $\mu(\mathbf{A}) = \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{SJS}^{-1})^{m_i} = \mathbf{S}\mu(\mathbf{J})\mathbf{S}^{-1}$ , 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

$$\begin{aligned} (\lambda_r \mathbf{I} - \mathbf{U}_r)^{m_r} &= \text{diag}((\lambda_r \mathbf{I} - \mathbf{J}_{m_{r1}})^{m_r}, \dots, (\lambda_r \mathbf{I} - \mathbf{J}_{m_{r,g_r}})^{m_r}) \\ &= \text{diag}(\mathbf{E}_{m_{r1}}^{m_r}, \dots, \mathbf{E}_{m_{r,g_r}}^{m_r}) = \mathbf{0}, \end{aligned}$$



by Lemma E.5 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  $\mu$ . 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 E.10** 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 E.11** What is the minimal polynomial of the unit matrix and more generally of a diagonalizable matrix?



# 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] Golub, G. H., and C. F. Van Loan, *Matrix Computations*, John Hopkins University Press, Baltimore, MD, third edition, 1996.
- [4] Greenbaum, Anne, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [5] Hackbusch, Wolfgang, *Iterative Solution of Large Sparse Systems of Equations*, Springer-Verlag, Berlin, 1994.
- [6] Hestenes, Magnus, *Conjugate Direction Methods in Optimization*, Springer-Verlag, Berlin, 1980.
- [7] 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.
- [8] Higham, Nicloas J., *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.
- [9] Horn, Roger A. and Charles R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [10] Horn, Roger A. and Charles R. Johnson, *Topics in Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1991.
- [11] Kato, *Perturbation Theory for Linear Operators*, Pringer.
- [12] Lawson, C.L. and R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, Englewood Cliffs, N.J, 1974.
- [13] Lax, Peter D., *Linear Algebra*, John Wiley & Sons, New York, 1997.
- [14] Leon, Steven J., *Linear Algebra with Applications*, Prentice Hall, NJ, Seventh Edition, 2006.

- 
- [15] Meyer, Carl D., *Matrix Analysis and Applied Linear Algebra*, Siam Philadelphia, 2000.
  - [16] Steel, J. Michael, *The Cauchy-Schwarz Master Class*, Cambridge University Press, Cambridge, UK, 2004.
  - [17] Stewart, G. G., *Matrix Algorithms Volume I: Basic Decompositions*, Siam Philadelphia, 1998.
  - [18] Stewart, G. G., *Matrix Algorithms Volume II: Eigensystems*, Siam Philadelphia, 2001.
  - [19] Stewart, G. G. and Ji-guang Sun, *Matrix Perturbation Theory*, Academic Press, San Diego, 1990.
  - [20] Stewart, G. G., *Introduction to Matrix Computations*, Academic press, New York, 1973.
  - [21] Trefethen, Lloyd N., and David Bau III, *Numerical Linear Algebra*, Siam Philadelphia, 1997.
  - [22] Tveito, A., and R. Winther, *Partial Differential Equations*, Springer, Berlin.
  - [23] Van Loan, Charles, *Computational Frameworks for the Fast Fourier Transform*, Siam Philadelphia, 1992.
  - [24] Varga, R. S., *Matrix Iterative Analysis/ 2nd Edn.*, Springer Verlag, New York, 2000.
  - [25] Wilkinson, J. H., *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.
  - [26] Young, D. M., *Iterative Solution of Large Linear Systems*, Academic Press, New York, 1971.

# Index

- (leading) principal minor, 82
- 1D test matrix, 100
- 2D test matrix, 100
- A-norm, 189
- A-orthogonal, 189
- abelian group, 8
- absolute error, 156, 277
- adjoint matrix, 43
- adjusted rounding unit, 283
- algebraic multiplicity, 55
- averaging matrix, 100
- backward error, 282
- backward stable, 282
- banded matrix, 25
  - symmetric LU factorization, 94
- banded symmetric LU factorization, 94
- biharmonic equation, 107
  - fast solution method, 118
  - nine point rule, 118
- Cauchy sequence, 17, 160
- Cauchy-Binet formula, 48
- Cauchy-Schwarz inequality, 20
- Cayley Hamilton Theorem, 299
- characteristic equation, 49
- characteristic polynomial, 49
- Chebyshev polynomial, 195
- Cholesky factorization, 89
- cofactor, 43
- column operations, 45
- companion matrix, 53
- complete pivoting, 272
- computer arithmetic, 277
- condition number, 284
  - ill-conditioned, 155
- congruent matrices, 246
- conjugate gradient method, 185
  - A-norm, 189
  - convergence, 193
  - derivation, 189
  - energy norm, 189
  - Krylov subspace, 189
  - least squares problem, 192
  - preconditioning, 201
  - preconditioning algorithm, 203
  - preconditioning convergence, 203
- convergence
  - absolute, 18
- convex combination, 291
- convex function, 291
- Courant-Fischer theorem, 126
- Cramers rule, 43
- cubic spline, 62
- defective eigenvalue, 56
- defective matrix, 56
- deflation, 122
- determinant, 39
  - additivity, 40
  - block triangular, 40
  - Cauchy-Binet, 48
  - cofactor, 43
  - cofactor expansion, 45
  - homogeneity, 40
  - permutation of columns, 40
  - product rule, 40
  - singular matrix, 40
  - transpose, 40
  - triangular matrix, 39

- Van der Monde, 46
- direct sum, 13
- Discrete Fourier Transform, 113
  - Fourier matrix, 113
- Discrete Sine Transform, 112
- double precision, 280
- eigenpair, 49
  - left eigenpair, 57
  - orthonormal eigenpairs, 121
  - right eigenpair, 57
- eigenvalue, 49
  - algebraic multiplicity, 55
  - characteristic equation, 49
  - characteristic polynomial, 49
  - Courant-Fischer theorem, 126
  - defective, 56
  - geometric multiplicity, 55
  - Hoffman-Wielandt theorem, 128
  - Kronecker sum, 103
  - location, 239
  - Rayleigh quotient, 125
  - Schur form, real, 128
  - spectral theorem, 124
  - spectrum, 49
- eigenvector, 49
  - Kronecker sum, 103
  - left eigenvector, 57
  - right eigenvector, 57
- elementary divisors, 298
- elementary lower triangular matrix, 265
- elementary reflector, 214
- Elsner's theorem, 238
- energy norm, 189
- exchange matrix, 95
- Fast Fourier Transform, 114
  - recursive FFT, 116
- field, 7
- fill-inn, 110
- finite difference method, 69
- fixed-point, 173
- fixed-point iteration, 173
- floating-point number
  - bias, 280
  - denormalized, 279
  - double precision, 280
  - exponent part, 279
  - guard digits, 281
  - half precision, 280
  - Inf, 279
  - mantissa, 279
  - NaN, 279
  - normalized, 279
  - overflow, 279
  - quadruple precision, 280
  - round away from zero, 281
  - round to zero, 281
  - rounding, 281
  - rounding unit, 281
  - single precision, 280
  - subnormal, 279
- flops, 268
- Fourier matrix, 113
- Fredholm's alternative, 144
- fundamental subspaces, 31
- Gaussian elimination, 264
  - complete pivoting, 272
  - elementary lower triangular matrix, 265
  - flops, 268
  - interchange matrix, 270
  - partial pivoting, 272
  - pivot, 269
  - pivot vector, 272
  - pivoting, 269
- geometric multiplicity, 55
- Gerschgorin's theorem, 239
- Given's rotation, 218
- gradient, 287
- group, 7
- guard digits, 281
- Hölder's inequality, 15, 293
- Hadamard's inequality, 212
- half precision, 280
- hessian, 287
- Hilbert matrix, 48
- Hoffman-Wielandt theorem, 128

- Householder transformation, 214
- identity matrix, 4
- ill-conditioned, 284
- ill-conditioned problem, 155
- inequality, 291
  - geometric/arithmetic mean, 292
  - Hölder, 293
  - Jensen, 291
  - Minkowski, 293
- Inf, 279
- inner product, 19
  - inner product norm, 19
  - standard inner product in  $\mathbb{C}^n$ , 19
  - standard inner product in  $\mathbb{R}^n$ , 19
- inner product space
  - linear projection operator, 24
  - orthogonal basis, 21
  - orthogonal complement, 24
  - orthogonal decomposition, 24
  - orthonormal basis, 21
- interchange matrix, 270
- inverse power method, 254
- iterative method
  - convergence, 173
  - Gauss-Seidel, 168
  - Jacobi, 168
  - SOR, 168
  - SOR, convergence, 179
  - SSOR, 168
- iterative methods, 167
- Jacobian, 287
- Jensen's inequality, 291
- Jordan form, 296
  - elementary divisors, 298
  - Jordan block, 295
  - Jordan canonical form, 295
  - principal vectors, 296
- Kronecker product, 101
  - eigenvalues, 102
  - eigenvectors, 102
  - inverse, 103
  - left product, 101
  - mixed product rule, 102
  - nonsingular, 103
  - positive definite, 103
  - right product, 101
  - symmetry, 103
  - transpose, 102
- Kronecker sum, 101
  - eigenvalues, 103
  - eigenvectors, 103
  - nonsingular, 103
  - positive definite, 103
  - symmetry, 103
- Krylov subspace, 189
- Laplacian, 287
- leading principal block submatrices, 84
- leading principal submatrices, 82
- least squares
  - error analysis, 229
  - normal equations, 222
- left eigenpair, 57
- left eigenvector, 57
- left triangular, 81
- linear combination, 9
- linear mapping, 33
- linear system
  - homogenous, 28
  - overdetermined, 28
  - residual vector, 158
  - square, 28
  - underdetermined, 28
- linear transformation
  - kernel, 34
  - span, 34
- LU factorization, 81
  - symmetric, 85
- LU factorization, see also LR-factorization, 263
- mantissa, 279
- matrix
  - addition, 25
  - adjoint, 43

- block lower triangular, 25
- block matrix, 77
- block upper triangular, 25
- blocks, 77
- companion matrix, 53
- conjugate transpose, 27
- defective, 56
- deflation, 122
- diagonal, 25
- diagonalizable, 54
- diagonally dominant, 71
- element-by-element operations, 26
- equivalent, 32
- fundamental subspaces, 31
- Hadamard product, 26
- Hermitian transpose, 27
- Hilbert, 48
- ill-conditioned, 156
- inverse, 29
- invertible, 29
- leading principal submatrices, 82
- left inverse, 29
- left triangular, 25
- lower banded, 25
- lower Hessenberg, 25
- lower triangular, 25
- multiplication, 26
- nilpotent, 52
- nonsingular, 28
- normal, 121, 123
- nullity, 31
- outer product expansion, 78
- permutation, 270
- pseudo-inverse, 140
- quasi-triangular, 123
- rank, 31
- right inverse, 29
- right triangular, 25
- scalar multiplication, 25
- Schur factorization, 122
- Schur product, 26
- second derivative, 69
- similar matrices, 53
- similarity transformation, 53
- singular, 28
- spectral radius, 160, 161
- strictly diagonally dominant, 67
- test matrix, 1D , 100
- test matrix, 2D , 100
- trace, 51
- transpose, 26
- tridiagonal, 25
- unitary similar, 121
- upper banded, 25
- upper Hessenberg, 25
- upper triangular, 25
- well-conditioned, 156
- matrix norm
  - consistent norm, 149
  - Frobenius norm, 148
  - max norm, 148
  - operator norm, 150
  - spectral norm, 152
  - subordinate norm, 149
  - sum norm, 148
  - two-norm, 152
- Minkowski's inequality, 15, 293
- mixed product rule, 102
- NaN, 279
- natural ordering, 98
- negative (semi)definite, 86
- Neumann Series, 163
- nflops, 268
- nilpotent matrix, 52
- norm
  - $l_1$ -norm, 15
  - $l_2$ -norm, 15
  - $l_\infty$ -norm, 15
  - absolute norm, 155
  - Euclidian norm, 15
  - infinity-norm, 15
  - max norm, 15
  - monotone norm, 155
  - one-norm, 15
  - triangle inequality, 15
  - two-norm, 15
- normal equations, 222
- normal matrix, 121, 123



- not-a-knot, 67
- nullity, 31
- operation count, 268
- optimal relaxation parameter, 178
- orthogonal matrix, see orthonormal matrix, 34
- orthogonal projection, 23
- orthonormal eigenpairs, 121
- orthonormal matrix, 34
- overflow, 279
- paraboloid, 200
- partial pivoting, 272
- permutation, 37
  - identity, 37
  - inversion, 38
  - sign, 38
  - symmetric group, 39
- permutation matrix, 95, 270
- perpendicular vectors, 21
- pivot vector, 272
- pivots, 265
- plane rotation, 218
- PLU factorization, 83, 95
- Poisson matrix, 99
- Poisson problem, 97
  - five point stencil, 98
  - nine point scheme, 107
  - Poisson matrix, 99
  - variable coefficients, 204
- positive definite, 86
- positive semidefinite, 86
- power method, 251
  - inverse, 254
  - Rayleigh quotient iteration, 254
  - shifted, 254
- pp representation, 63
- preconditioning, 201
- principal submatrix, 82
- principal vectors, 296
- pseudo-inverse, 140
- QR algorithm
  - implicit shift, 260
  - Rayleigh quotient shift, 260
  - shifted, 259
  - Wilkinson shift, 260
- QR decomposition, 211
- QR factorization, 211
- quadratic form, 86
- quadruple precision, 280
- quotient space, 14
- rank, 31
- rate of convergence, 174
- Rayleigh quotient, 125
- Rayleigh quotient iteration, 254
- relative error, 156, 277
- residual vector, 158
- right eigenpair, 57
- right eigenvector, 57
- right triangular, 81
- rotation in the  $i, j$ -plane, 219
- rounding unit, 281
- rounding-error analysis
  - adjusted rounding unit, 283
  - backward error, 282
  - backward stable, 282
  - condition number, 284
  - ill-conditioned, 284
- row operations, 45
- RTR factorization, 89
- scalar product, 19
- Schur factorization, 122
- Schur form, real, 128
- second derivative matrix, 69
- semi-Cholesky factorization, 89
- sequence
  - bounded sequence of vectors, 17
  - subsequence, 18
- Sherman-Morrison formula, 30
- shifted power method, 254
- similar matrices, 53
- similarity transformation, 53
- single precision, 280
- singular value
  - Courant-Fischer theorem, 144
  - error analysis, 232

- Hoffman-Wielandt theorem, 145
- singular values, 132
- singular vector
  - left singular vectors, 138
  - right singular vectors, 138
- spectral radius, 160, 161
- spectral theorem, 124
- spectrum, 49
- steepest descent, 201
- stencil, 98
- Sylvester's inertia theorem, 246
- symmetric positive semidefinite, 86
- trace, 51
- triangle inequality, 15
- triangular matrix
  - left triangular, 81
  - right triangular, 81
- unit vectors, 4
- unitary similar, 121
- vector
  - addition, 7
  - angle, 21
  - linearly dependent, 10
  - linearly independent, 10
  - orthogonal, 21
  - orthonormal, 21
  - scalar multiplication, 7
- vector space, 8
  - basis, 10
  - complementary, 13
  - complete, 17
  - complex inner product space, 19
  - dimension, 11
  - direct sum, 13
  - finite dimensional, 9
  - intersection, 12
  - normed, 15
  - quotient space, 14
  - real inner product space, 19
  - subspace, 8
  - sum, 12
  - trivial, 8
  - union, 12
- vectorization, 98
- vectornorm, 14
- Wilkinson diagram, 218