

COMPUTATIONAL LINEAR ALGEBRA

Les S. Jennings
les@maths.uwa.edu.au
<http://maths.uwa.edu.au/~les/>

Notes for third year and honours courses in
Numerical Analysis
Printed 3 May 2001

The University of Western Australia
DEPARTMENT of MATHEMATICS AND STATISTICS

CONTENTS.

| | |
|---|-----------|
| 1. Introduction | 1 |
| 2. Errors and Computer Arithmetic | 2 |
| 2.1 Accuracy..... | 2 |
| 2.2 Precision..... | 2 |
| 2.3 Arithmetic unit errors..... | 3 |
| 2.4 Backward error analysis..... | 3 |
| 2.5 Loss of precision..... | 3 |
| 2.6 Exercises..... | 3 |
| 3. Euclidean Vector Spaces | 4 |
| 3.1 Revision of the vector spaces \mathbb{R}^n | 4 |
| 3.2 Orthogonality..... | 5 |
| 3.3 Exercises..... | 6 |
| 4. Linear Mappings, Matrices | 7 |
| 4.1 Linear mapping (linear transformation)..... | 7 |
| 4.2 Some alternate definitions..... | 7 |
| 4.3 Inverse mapping..... | 8 |
| 4.4 Orthogonal matrices..... | 8 |
| 4.5 Matrix norms..... | 8 |
| 4.6 Change of basis..... | 8 |
| 4.7 Determinant of a matrix..... | 9 |
| 4.8 Eigenvalues and eigenvectors..... | 9 |
| 4.9 Exercises..... | 10 |
| 5. Areas of Application, Types of Matrices | 12 |
| 5.1 Classification of problems..... | 12 |
| 5.2 Types of Matrices..... | 12 |
| 5.3 Exercises..... | 13 |
| 6. Elementary Operations and Elementary Matrices | 14 |
| 6.1 Elementary operations..... | 14 |
| 6.2 Multiplying a row by a non-zero scalar..... | 14 |
| 6.3 Interchanging rows..... | 14 |
| 6.4 Subtracting a multiple of one row from other rows..... | 14 |
| 6.5 Householder transformations..... | 15 |
| 6.6 Plane rotations..... | 16 |
| 6.7 Exercises..... | 16 |
| 7. Decompositions (Factorizations) | 18 |
| 7.1 Singular value decomposition..... | 18 |
| 7.2 Q–U factorization (Q–R)..... | 18 |
| 7.3 L–U factorization..... | 19 |
| 7.4 Choleski factorization and variants..... | 22 |
| 7.5 Banded matrices..... | 23 |
| 7.6 Eigenvalue–eigenvector decomposition..... | 23 |
| 7.7 Exercises..... | 24 |

| | |
|---|-----------|
| 8. The Inverse Mapping Problem | 26 |
| 8.1 Sensitivity analysis | 26 |
| 8.2 Direct methods | 29 |
| 8.3 Least squares methods | 32 |
| 8.4 Iterative methods | 33 |
| 8.5 Exercises | 35 |
| 9. The Generalized Inverse Problem | 38 |
| 9.1 The geometry of a linear mapping | 38 |
| 9.2 Natural orthogonal bases for a mapping | 38 |
| 9.3 Canonical form of a linear mapping | 39 |
| 9.4 Generalized inverse mapping | 40 |
| 9.5 Best rank s approximation to a matrix | 42 |
| 9.6 Generalized inverse and rounding error | 42 |
| 9.7 Scaling considerations | 43 |
| 9.8 Generalized inverses under scalings | 44 |
| 9.9 Approximations to the generalized inverse | 44 |
| 9.10 Projection operators | 46 |
| 9.11 Exercises | 47 |
| 10. Eigenvalue Computations | 48 |
| 10.1 Introduction | 48 |
| 10.2 Classes of algorithm | 48 |
| 10.3 Classes of matrices | 49 |
| 10.4 Real symmetric matrices (Hermitian matrices) | 49 |
| 10.5 Tridiagonal matrices | 51 |
| 10.6 Inverse iteration | 52 |
| 10.7 Reduction to upper triangular form | 52 |
| 10.8 Balancing a matrix | 53 |
| 10.9 EISPACK | 54 |
| 10.10 Other eigenvalue problems | 54 |
| 10.12 Exercises | 56 |
| References | 57 |

APPENDICES.

| | |
|--|----|
| A. Software Development — a Brief Personal History | 59 |
| B.1 LAPACK | ? |
| B.1 LINALG | ? |
| B.3 LINPACK, IMSL, NAG, Numerical Recipes | ? |
| C. Description of the Q–U package | ? |
| D. Eigenvalue/Eigenvector software | ? |
| E. NA-NET libraries | ? |

Copyright: These notes may be reproduced provided the copies are not sold for profit. Contact the author at les@maths.uwa.edu.au for an up to date plain \TeX source.

1. Introduction.

This set of notes on Computational Linear Algebra contains a mathematical description of operations with matrices and vectors on a modern computer. They have been used for courses in numerical analysis at the University of Western Australia for some time. In particular, the solution of sets of linear equations is studied and an overview of eigenvalue–eigenvector computations is given. Some of the material, for example, singular value decomposition and generalized inverses is not relevant to third year students, being more advanced, while some early chapters revise second year work. Chapters 6–8 contain the mathematical description of Gauss Elimination and the Q-U factorization. Students are expected to consult other more complete references where applicable.

References to the computing library packages are also included and it is expected that these will be consulted when necessary. Note that the top quality computing libraries are written in FORTRAN but that on some computers the Pascal and C compilers have means to communicate with FORTRAN subroutines and functions. There is a slow trickle of conversions of the FORTRAN libraries into C language libraries. One problem when using both FORTRAN and Pascal or C is that arrays are stored differently; in FORTRAN by columns, in Pascal by rows, and in C by user choice but usually rows. The easiest way to overcome this problem for Pascal is to transpose the matrix in Pascal. Care has to be taken with the FORTRAN column length variable which equates with the Pascal row length. Similar interface problems exist with the C language. Users of C libraries must inform themselves of how a matrix is stored into an appropriate C data structure.

Many of the computational exercises contained herein could be computed with MATLAB, a commercially available high level system for numerical computing. In particular, MATLAB is a good interface to the NSF funded software libraries LINPACK (linear equations package) and EISPACK (eigen systems package), and now LAPACK (linear algebra package) which supercedes LINPACK and EISPACK. There is a companion library of basic linear algebra subroutines (BLAS) which can be tailored to particular machine architectures for efficiency. The NSF funded libraries are available over the network, by email,

from
`netlib@ornl.gov`

and by the WEB, at

`http://www.netlib.org/index.html`

MATLAB information is available on the WEB at

`http://www.mathworks.com`

There are some freely available MATLAB look-alikes.

SCINET is available from

`http://www-rocq.inria.fr/scilab`

and OCTAVE from

`http://www.che.wisc.edu/octave/`

Appendices contain references to software and descriptions of some local software useful for teaching purposes.

2. Errors and Computer Arithmetic.

Suggested extra reading in Atkinson (1989), Dahlquist *et al.* (1974), Stewart (1973). The beginnings of this work can be found in Wilkinson (1963, 1965).

2.1 Accuracy.

Suppose x represents an exact quantity and \bar{x} an approximation to it. Define the following types of accuracy with respect to some small quantity ϵ .

- (i) Absolute accuracy: \bar{x} is within ϵ of x absolutely if

$$|x - \bar{x}| < \epsilon.$$

This gives a certain number of decimal places accuracy (the same in both \bar{x} and x) determined by ϵ .

- (ii) Relative accuracy: \bar{x} is within ϵ of x in a relative sense if

$$|x - \bar{x}| < \epsilon|x|, \quad x \neq 0.$$

This gives a certain number of significant digits accuracy, determined by ϵ .

- (iii) Mixed test:

$$|x - \bar{x}| < \epsilon(a + |x|)$$

This test is useful if it is known that x is of magnitude a . If $|x| \gg a$ it acts as a relative test while if $|x| \ll a$ it acts as a relative test against a . For most purposes the relative or mixed test should be used.

2.2 Precision.

Whenever a number expressed in decimal notation is stored on a computer (in binary or hexadecimal form) truncation (chopping) or rounding occurs in the last binary bit. Hence one computer number will represent a range of numbers on the real line. The relative (compared to the number) length of this range is the *precision* of the computer. The error between the number and its computer representation is known as *rounding error* and is a function of computer hardware. A suitable definition of precision of a computer (language/compiler) is the smallest positive computer number ϵ such that

$$1 + \epsilon \neq 1,$$

when executed in that computer (language/compiler).

Precision of FORTRAN variables on computers.

| computer | single precision | double precision | long precision |
|-----------------|-------------------------|-----------------------------|------------------------------|
| CDC 6000 series | $2^{-47} \sim 10^{-14}$ | $2^{-95} \sim 10^{-29}$ | — |
| INTEL 80x87 | $2^{-23} \sim 10^{-7}$ | $2^{-52} \sim 10^{-16}$ | — |
| SUN SPARC | $2^{-23} \sim 10^{-7}$ | $2^{-52} \sim 10^{-16}$ | — |
| HP 9000 | $2^{-23} \sim 10^{-7}$ | $2^{-52} \sim 10^{-16}$ | — |
| DEC MIPS | $2^{-23} \sim 10^{-7}$ | $2^{-52} \sim 10^{-16}$ | — |
| DEC VAX | $2^{-26} \sim 10^{-8}$ | $2^{-58} \sim 10^{-18}$ (D) | — |
| | — | $2^{-52} \sim 10^{-16}$ (G) | $2^{-102} \sim 10^{-31}$ (G) |
| DEC ALPHA | $2^{-26} \sim 10^{-8}$ | $2^{-58} \sim 10^{-18}$ (D) | — |
| | — | $2^{-52} \sim 10^{-16}$ (G) | $2^{-102} \sim 10^{-31}$ (G) |

Note that rounding error and precision are expressed as a relative error (significant binary digits). The precision of a computer should not be confused with the smallest positive number of a computer. There is now an IEEE standard for storage of floating point numbers and for their arithmetic. Not all manufacturers use this standard. Most numerical computation should be done at a precision smaller than 10^{-10} , which means that except in special circumstances, single precision should not be used on most machines.

2.3 Arithmetic unit errors.

Actual machine arithmetic generates further errors. Let x and y be two computer numbers and θ be one of the operations $+$, $-$, $*$ or $/$. Then the computer number stored which represents $x\theta y$ is denoted $fl(x\theta y)$ and,

$$fl(x\theta y) = (x\theta y)(1 + \epsilon(x, y)),$$

where $\epsilon(x, y)$ represents the relative error. The object of computer manufacturers is to have $|\epsilon(x, y)| \leq$ machine precision for all choices of computer numbers x and y . (See Kahan (1966)) Most computers use double length accumulators for arithmetic so effectively it can be assumed that no errors are created in the arithmetic. But when the double length result is stored back to single length in memory an error is created. This is not necessarily true for all computers, especially at double precision.

2.4 Backward error analysis.

This relates the error in a computed result back to contrived errors in the operands, that is,

$$fl(x\theta y) = (x + \delta x)\theta(y + \delta y).$$

The actual result is considered to be the exact result using perturbed data. Sometimes it is easier to put bounds on δx and δy than on $\epsilon(x, y)$. Here the operands may be matrices and θ represents a matrix operation.

2.5 Loss of precision.

- (i) *Cancellation* (of significant digits)— occurs when two nearly equal numbers are subtracted, for example,

$$f'(x) = (f(x+h) - f(x))/h.$$

This is the most serious mechanism of loosing accuracy when executing millions of instructions. It could be hidden deeply within an algorithm.

- (ii) *Underflow* — occurs when summing numbers which differ vastly in magnitude. For example on 32 bit hexadecimal machines

$$fl(1 + 10^{-7}) = 1.$$

This can effect statistical computation and the like in that the result of summing thousands of numbers depends on the order in which they are summed.

$$e.g. \quad fl\left(1.0 + \sum_{i=1}^{10^7} 10^{-7}\right) \neq 2 \neq fl\left(\sum_{i=1}^{10^7} 10^{-7} + 1.0\right).$$

- (iii) *Scaling* problems are usually critical. Basically it means trying to bring significant figure accuracy to be the same as decimal place accuracy. If possible all numbers should be scaled to be of magnitude 1.0. For example the matrix $\begin{pmatrix} 1.111 & 1.234 \times 10^{-7} \\ 2.222 & 1.234 \times 10^{-7} \end{pmatrix}$ could be considered singular, but if the numbers are accurate to four significant figures the second column should be scaled to give the matrix $\begin{pmatrix} 1.111 & 1.234 \\ 2.222 & 1.234 \end{pmatrix}$ which is far from singular. Care must be taken to use the full significant accuracy of any given data, especially in matrix computations.

2.6 Exercises.

1. Find the error bounds both absolute and relative for storing (from a decimal number), computing and storing the result of $a + b$, $a - b$, $a * b$, and a/b in terms of a computer with precision ϵ .
2. Given two vectors \mathbf{a} and \mathbf{b} and a computer with precision ϵ , what is an upper bound for the error in storing, computing and storing the result of $\mathbf{a} + \mathbf{b}$ and $\mathbf{a}^t \mathbf{b}$ in terms of the norms of \mathbf{a} and \mathbf{b} ? You need to follow the detailed algorithm for these two operations on the vectors.
3. For a new computer chip using 52 bit mantissa, how long would it take to check every mantissa bit combination for one of the four basic arithmetic operations? Use a realistic time to test a pair of numbers, and, assume there are 10^7 seconds in a year.

3. Euclidean Vector Spaces.

Consult any linear algebra texts, such as, Anton and Rorres (1987), Lipschutz (1968), Strang (1988), Noble (1969), Noble and Daniel (1977) for a fuller discussion of the material in chapters 3-5.

3.1 Revision of the vector spaces \mathbb{R}^n .

3.1.1 Notation.

The finite dimensional vector spaces considered are the Euclidean vector spaces of dimension n , denoted \mathbb{R}^n or \mathbb{E}^n . They have elements \mathbf{x} written as column vectors so that

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad x_i \in \mathbb{R} = (-\infty, \infty).$$

For convenience, write $\mathbf{x} = (x_1, x_2, \dots, x_n)^t$, where \cdot^t represents the operation ‘transpose’. Similarly \mathbb{C}^n denotes the n dimensional vector space where $x_i \in \mathbb{C}$, the field of complex numbers.

- (i) *The Unit Vectors:* \mathbf{e}_i is the i -th unit vector of *the standard basis* for \mathbb{R}^n .

$$\mathbf{e}_i = (0, 0, \dots, 1, \dots, 0)^t$$

where the ‘1’ is in the i -th position. Note that x_i and \mathbf{x}_i are different, one being the i -th component of a vector, the other being the i -th vector of a sequence of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots\}$.

- (ii) The vector $\mathbf{e} = (1, 1, 1, \dots, 1)^t$ will occur occasionally.
 (iii) Lower case greek letters are always scalars.

3.1.2 Inner product.

An inner product on \mathbb{R}^n , denoted $\langle \mathbf{x}, \mathbf{y} \rangle$ is a binary operation on two vectors, producing a scalar, with the properties, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\alpha, \beta \in \mathbb{R}$;

- (i) $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ (commutative),
 (ii) $\langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle = \alpha \langle \mathbf{x}, \mathbf{z} \rangle + \beta \langle \mathbf{y}, \mathbf{z} \rangle$,
 (iii) $\langle \mathbf{x}, \mathbf{x} \rangle > 0$, $\mathbf{x} \neq \mathbf{0}$.

The standard inner product for \mathbb{R}^n is

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i = \mathbf{x}^t \mathbf{y} = \mathbf{y}^t \mathbf{x}.$$

The definition is changed for \mathbb{C}^n , the vector spaces over the field of complex numbers. Item (i) is changed to $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$ where \bar{z} is the complex conjugate of the complex number z . Combined with (ii) above this gives $\langle \mathbf{x}, \beta \mathbf{y} \rangle = \bar{\beta} \langle \mathbf{x}, \mathbf{y} \rangle$. The standard inner product for \mathbb{C}^n is $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^t \bar{\mathbf{y}} = \overline{\mathbf{y}^t \mathbf{x}}$. The operation of conjugate transpose is sometimes denoted $\mathbf{x}^* = \overline{\mathbf{x}^t}$.

3.1.3 Vector norms.

The norm of a vector \mathbf{x} , is a function from $\mathbb{R}^n \rightarrow \mathbb{R}$ (or $\mathbb{C}^n \rightarrow \mathbb{R}$), denoted $\|\mathbf{x}\|$, if the following properties hold:

- (i) $\|\mathbf{x}\| > 0$ for $\mathbf{x} \neq \mathbf{0}$.
 (ii) $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$.
 (iii) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality).

Examples:

$$\begin{aligned}\|\mathbf{x}\|_2 &= \left\{ \sum_{i=1}^n x_i^2 \right\}^{\frac{1}{2}} = \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}} & : \text{Euclidean norm or 2-norm} \\ \|\mathbf{x}\|_\infty &= \max_i |x_i| & : \infty\text{-norm} \\ \|\mathbf{x}\|_1 &= \sum_{i=1}^n |x_i| & : 1\text{-norm} \\ \|\mathbf{x}\|_p &= \left\{ \sum_{i=1}^n |x_i|^p \right\}^{\frac{1}{p}}, 1 < p < \infty & : p\text{-norm}\end{aligned}$$

3.1.4 Linear independence.

Vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r$, ($\{\mathbf{x}_i\}_{i=1}^r$) are linearly independent if

$$\sum_{i=1}^r a_i \mathbf{x}_i = \mathbf{0} \Rightarrow a_i = 0, \quad i = 1, 2, \dots, r.$$

This means no non-zero linear combination of the $\{\mathbf{x}_i, i = 1, \dots, r\}$ sum to zero.

3.1.5 Linear dependence.

At least one of the \mathbf{x}_i can be expressed as a linear combination of the others,

$$\mathbf{x}_i = \sum_{\substack{j=1 \\ j \neq i}}^r a_j \mathbf{x}_j.$$

This means there exists a non zero linear combination of the \mathbf{x}_i which sum to zero.

3.1.6 Span.

Span $\{\mathbf{x}_i\}_{i=1}^r$ is the subspace formed by all linear combinations of the vectors \mathbf{x}_i .

$$\text{span}\{\mathbf{x}_i\}_{i=1}^r = \{\mathbf{y} \mid \mathbf{y} = \sum_{i=1}^r a_i \mathbf{x}_i, a_i \in \mathbb{R}\}$$

3.1.7 Basis.

A *basis* of a vector space is a set of linearly independent vectors which have a span equal to the vector space.

3.1.8 Dimension.

The *dimension* of a subspace is the number of vectors in a basis for the subspace.

3.2 Orthogonality

3.2.1 Orthogonal vectors.

Two non-zero vectors \mathbf{x}, \mathbf{y} are orthogonal if their inner product $\langle \mathbf{x}, \mathbf{y} \rangle$ is zero, for example \mathbf{e}_i and \mathbf{e}_j are orthogonal if $i \neq j$, for the standard inner product.

3.2.2 Orthogonal subspaces.

Two subspaces X and Y are orthogonal if for all $\mathbf{x} \in X$ and for all $\mathbf{y} \in Y$, $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

3.2.3 Orthogonal complement.

The *orthogonal complement* of a subspace X , denoted X^\perp is defined as

$$X^\perp = \{\mathbf{y} \mid \langle \mathbf{x}, \mathbf{y} \rangle = 0, \text{ for all } \mathbf{x} \in X\}.$$

3.2.4 Orthogonal basis.

An *orthogonal basis* of a subspace is a basis $\{\mathbf{x}_i\}_{i=1}^r$ in which the basis elements are mutually orthogonal, that is

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0, \quad i \neq j.$$

3.2.5 Orthonormal basis.

An *orthonormal basis* of a subspace is an orthogonal basis in which the basis elements are normalized to unit norm,

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \delta_{ij} = \begin{cases} 1, & i = j, \\ 0, & i \neq j. \end{cases}$$

3.3 Exercises.

1. Find out how to find an orthonormal basis using Gram-Schmidt orthogonalization.
2. Show $(X^\perp)^\perp = X$ and $X^\perp \cap X = \{\mathbf{0}\}$.
3. If the subspace $X \subset \mathbb{R}^n$ show $X \oplus X^\perp = \mathbb{R}^n$ where $Y \oplus Z$ is the set $\{\mathbf{x} \mid \mathbf{x} = \mathbf{y} + \mathbf{z}, \mathbf{y} \in Y, \mathbf{z} \in Z\}$.
4. Show that $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1$, for all $\mathbf{x} \in \mathbb{R}^n$.
5. Draw the sets $\{\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_k = 1\}$, for $k = 1, 2$ and ∞ , and $n = 2$ and 3 .
6. Show that any inner product (not just the standard one) on \mathbb{R}^n defines a natural norm, called the inner product norm, by the definition

$$\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}}.$$

4. Linear Mappings, Matrices.

4.1 Linear mapping (linear transformation).

Let \mathcal{A} be a mapping of \mathbb{R}^n to \mathbb{R}^m . \mathcal{A} is linear if $\mathcal{A} : \mathbf{x} \rightarrow \mathbf{w}$, $\mathcal{A} : \mathbf{y} \rightarrow \mathbf{z}$ then

$$\mathcal{A} : \alpha\mathbf{x} + \beta\mathbf{y} \rightarrow \alpha\mathbf{w} + \beta\mathbf{z}.$$

A linear mapping of $\mathbb{R}^n \rightarrow \mathbb{R}^m$ is represented in the usual basis by an $m \times n$ matrix \mathbf{A} (say). Write

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

to describe $\mathcal{A} : \mathbf{x} \rightarrow \mathbf{y}$.

4.1.1 Notation.

- (i) A_{ij} (or $(A)_{ij}$) represents the element in the i -th row and j -th column of the matrix \mathbf{A} .
- (ii) $\rho_i(\mathbf{A})$ represents the i -th row of \mathbf{A} as a row vector.
- (iii) $\kappa_j(\mathbf{A})$ represents the j -th column of \mathbf{A} as a column vector.
- (iv) The transpose of the matrix \mathbf{A} is the matrix \mathbf{A}^t (or \mathbf{A}') defined,

$$(\mathbf{A}^t)_{ij} = A_{ji}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

Note that \mathbf{A}^t represents a mapping from \mathbb{R}^m to \mathbb{R}^n .

- (v) The matrix-vector product $\mathbf{A}\mathbf{x} = \sum_{i=1}^n x_i \kappa_i(\mathbf{A})$, that is, a linear combination of columns of the matrix \mathbf{A} .

4.1.2 Adjoint.

The *adjoint* mapping \mathbf{A}^* is defined by the relation

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{A}^*\mathbf{y} \rangle, \quad \forall \mathbf{x} \in \mathbb{R}^n \quad \text{and} \quad \mathbf{y} \in \mathbb{R}^m.$$

4.1.3 Null space (kernel).

The *null space* (or *kernel*) of a mapping denoted $\mathcal{N}(\mathbf{A})$ (or $\ker \mathbf{A}$) is

$$\mathcal{N}(\mathbf{A}) = \{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{0}\} \subset \mathbb{R}^n.$$

4.1.4 Range space.

The *range* (or *image*) of a mapping, denoted $\mathcal{R}(\mathbf{A})$ (or $\mathcal{I}m\mathbf{A}$) is

$$\mathcal{R}(\mathbf{A}) = \{\mathbf{y} \mid \mathbf{y} = \mathbf{A}\mathbf{x} \quad \text{for some} \quad \mathbf{x} \in \mathbb{R}^n\} \subset \mathbb{R}^m.$$

4.1.5 Rank.

The *rank* of a matrix \mathbf{A} is the dimension of $\mathcal{R}(\mathbf{A})$.

4.2 Some alternative definitions.

Let the set of vectors $\{\mathbf{x}_i\}_{i=1}^r$ be represented by the $n \times r$ matrix \mathbf{X} where $\kappa_i(\mathbf{X}) = \mathbf{x}_i$, for $i = 1, 2, \dots, r$. Then the following are alternative definitions.

Note :

$$\sum_{i=1}^r a_i \mathbf{x}_i = \mathbf{X}\mathbf{a}, \quad \mathbf{a} = (a_1, a_2, \dots, a_r)^t.$$

- (i) Linear independence; \nexists non zero \mathbf{a} such that $\mathbf{X}\mathbf{a} = \mathbf{0}$, or $\mathcal{N}(\mathbf{X}) = \{\mathbf{0}\}$.

- (ii) Linear dependence; \exists a non zero \mathbf{a} such that $\mathbf{X}\mathbf{a} = \mathbf{0}$, or $\mathcal{N}(\mathbf{X}) \neq \{\mathbf{0}\}$.
- (iii) Span $(\mathbf{X}) = \{\mathbf{y} \mid \mathbf{y} = \mathbf{X}\mathbf{a}, \mathbf{a} \in \mathbb{R}^r\}$, ($= \mathbf{X}(\mathbb{R}^r)$).
- (iv) Orthogonal basis; $\mathbf{X}^t\mathbf{X} = \mathbf{D} = \text{diag}[d_1, d_2, \dots, d_r]$, $d_i \neq 0$, $i = 1, 2, \dots, r$.
- (v) Orthonormal basis; $\mathbf{X}^t\mathbf{X} = \mathbf{I}_r$, the identity matrix of order r .

4.3 Inverse mapping. ($\mathcal{A}^{-1} : \mathbb{R}^m \rightarrow \mathbb{R}^n$)

$$\begin{aligned} \mathcal{A}^{-1}\mathcal{A} &= \mathcal{A}\mathcal{A}^{-1} = \text{Identity mapping} \\ \mathbf{A}^{-1}\mathbf{A} &= \mathbf{A}\mathbf{A}^{-1} = \mathbf{I} \end{aligned}$$

The inverse mapping \mathcal{A}^{-1} exists and is represented by the inverse matrix \mathbf{A}^{-1} , in the usual basis, provided \mathcal{A} is 1:1 and onto (and hence $m = n$).

$$\begin{aligned} \text{onto} &\Leftrightarrow \mathcal{R}(\mathbf{A}) = \mathbb{R}^m \\ 1 : 1 &\Leftrightarrow \mathcal{N}(\mathbf{A}) = \{\mathbf{0}\} \end{aligned}$$

4.4 Orthogonal matrices.

\mathbf{Q} is an orthogonal matrix if

$$\mathbf{Q}\mathbf{Q}^t = \mathbf{Q}^t\mathbf{Q} = \mathbf{I},$$

that is $\mathbf{Q}^{-1} = \mathbf{Q}^t$. This means the rows and columns of \mathbf{Q} are orthonormal.

4.5 Matrix norms.

A matrix norm measures the ‘magnitude’ of a matrix .

- (a) The function $\|\mathbf{A}\|$ is a matrix norm if the four properties below are true;
 - (i) $\|\mathbf{A}\| > 0$, if $\mathbf{A} \neq \mathbf{0}$,
 - (ii) $\|\alpha\mathbf{A}\| = |\alpha| \|\mathbf{A}\|$,
 - (iii) $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$,
 - (iv) $\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$.
- (b) Operator norms (subordinate norms). Given norms for \mathbb{R}^n and \mathbb{R}^m ,

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$$

is an operator norm or subordinate norm.

- (i) $\|\mathbf{A}\|_1 = \max_{\|\mathbf{x}\|_1=1} \|\mathbf{A}\mathbf{x}\|_1 = \max_j \sum_{i=1}^m |A_{ij}|$. (Maximum column sum.)
- (ii) $\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \max_i \{\lambda_i(\mathbf{A}^t\mathbf{A})\}^{\frac{1}{2}}$. The 2-norm of a matrix is the square root of the largest eigenvalue of $\mathbf{A}^t\mathbf{A}$.
- (iii) $\|\mathbf{A}\|_\infty = \max_{\|\mathbf{x}\|_\infty=1} \|\mathbf{A}\mathbf{x}\|_\infty = \max_i \sum_{j=1}^n |A_{ij}|$. (Maximum row sum.)
- (c) Euclidean norm of a matrix (Frobenius norm)

$$\|\mathbf{A}\|_E = \left\{ \sum_{i,j} A_{ij}^2 \right\}^{\frac{1}{2}}.$$

- (d) If $\mathbf{y} = \mathbf{A}\mathbf{x}$, then $\|\mathbf{y}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$, for appropriate ‘same’ norms, as in (b).

4.6 Change of basis.

Many of the factorizations of Chapter 6 have an explanation in terms of a change of basis, either of the domain or range of a mapping, or of both. Chapter 9 requires a good understanding of this section.

4.6.1 Coordinates of a vector.

All vectors are usually expressed in terms of *the unit vectors* of that space, for example

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} a + \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} b + \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} c = \mathbf{I} \begin{pmatrix} a \\ b \\ c \end{pmatrix}.$$

Here a , b and c are the coordinates of the vector. More generally, in the standard basis

$$\mathbf{x} = \sum_{i=1}^n \mathbf{e}_i x_i \quad \left(= \sum_{i=1}^n \langle \mathbf{e}_i, \mathbf{x} \rangle \mathbf{e}_i \right) = \mathbf{I} \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n.$$

4.6.2 Change of basis.

Choose another basis $\{\mathbf{b}_i\}_{i=1}^r$ which can be represented (in ‘old’ coordinates) by the $n \times r$ matrix \mathbf{B} , where $\kappa_i(\mathbf{B}) = \mathbf{b}_i$, $i = 1, 2, \dots, r$. The equation

$$\mathbf{x} = \mathbf{B} \mathbf{a}, \quad \mathbf{a} \in \mathbb{R}^r,$$

means that \mathbf{x} can be represented in the basis \mathbf{B} by the coordinate r -tuple \mathbf{a} . Note that \mathbf{x} must be an element of $\mathcal{R}(\mathbf{B})$. If \mathbf{B} has an inverse ($r = n$) then \mathbf{a} can be found as $\mathbf{B}^{-1} \mathbf{x}$, where \mathbf{x} is the n -tuple for the vector \mathbf{x} in terms of the standard basis. For vectors in \mathbb{R}^n no distinction between the vector and its coordinate n -tuple is usually made, only because the standard basis has been chosen to represent the vector.

4.6.3 Rotation of axes.

This can be considered as a change of basis in which another orthonormal basis has been chosen. Let \mathbf{Q} be a matrix such that $\mathbf{q}_i = \kappa_i(\mathbf{Q})$ and $\{\mathbf{q}_i\}_{i=1}^r$ be an orthonormal basis ($\mathbf{Q}^t \mathbf{Q} = \mathbf{I}_r$). Then

$$\mathbf{x} = \sum_{i=1}^r \mathbf{q}_i a_i = \mathbf{Q} \mathbf{a}, \quad \mathbf{a} \in \mathbb{R}^r.$$

The term *rotation of axes* usually refers to the case $r = n$, so that any vector $\mathbf{x} \in \mathbb{R}^n$ can be represented in the new basis. In this case $\mathbf{x} = \mathbf{Q} \mathbf{a}$, so that $\mathbf{Q}^t \mathbf{x} = \mathbf{a}$ as $\mathbf{Q}^t \mathbf{Q} = \mathbf{Q} \mathbf{Q}^t = \mathbf{I}_n$, that is, $a_i = \langle \mathbf{q}_i, \mathbf{x} \rangle$ and

$$\mathbf{x} = \sum_{i=1}^n \langle \mathbf{q}_i, \mathbf{x} \rangle \mathbf{q}_i.$$

Hence the use of orthogonal matrices can be thought of as a change of orthonormal basis or rotation of axes. As such, the lengths of co-ordinate vectors ($\|\mathbf{x}\|_2$) remains unchanged under an orthogonal change of basis. In particular any vector of errors has its Euclidean size unchanged when multiplied by an orthogonal matrix.

4.7 Determinant of a matrix

See one of the texts Anton and Rorres (1987) or Lipschutz (1968) for the definition of the determinant of a matrix. The determinant of a matrix is only defined for a square matrix and the following properties are relevant to this set of notes. Let \mathbf{A} and \mathbf{B} be $n \times n$ matrices.

- (i) $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$.
- (ii) $\det(\mathbf{A}) \neq 0$ iff \mathbf{A} is invertible.
- (iii) $\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i(\mathbf{A})$, the product of all the eigenvalues of \mathbf{A} , counting repeated eigenvalues.
- (iv) $\det(c\mathbf{A}) = c^n \det(\mathbf{A})$, c a scalar.

4.8 Eigenvalues and eigenvectors

See one of the texts Anton and Rorres (1987) or Lipschutz (1968) for a fuller discussion. The eigenvalue–eigenvector pair of a matrix \mathbf{A} is defined to be the solution of the equation

$$\mathbf{A} \mathbf{v} = \lambda \mathbf{v}, \quad \mathbf{v} \neq \mathbf{0},$$

or equivalently, the eigenvalues are those values of λ which make the matrix $\mathbf{A} - \lambda\mathbf{I}$ singular, and the eigenvector \mathbf{v} is any nonzero vector in $\mathcal{N}(\mathbf{A} - \lambda\mathbf{I})$.

4.9 Exercises.

- Check $A_{ij} = \mathbf{e}_i^t \mathbf{A} \mathbf{e}_j$, $\rho_i(\mathbf{A}) = \mathbf{e}_i^t \mathbf{A}$, $\kappa_j(\mathbf{A}) = \mathbf{A} \mathbf{e}_j$.
- Prove $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^t)$, and $\mathcal{N}(\mathbf{A}^t)^\perp = \mathcal{R}(\mathbf{A})$.
- (i) Show that $\mathbf{A}^* = \mathbf{A}^t$ if the mapping, represented by the $m \times n$ matrix \mathbf{A} is real, and the usual inner products are used in \mathbb{R}^m and \mathbb{R}^n .
(ii) Show that $(\mathbf{A}\mathbf{B})^t = \mathbf{B}^t \mathbf{A}^t$.
- (i) Show that \mathcal{A}^{-1} exists iff $m = n = r$ (the rank of \mathcal{A}).
(ii) Show that $(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$.
- Show that for \mathbf{Q} orthogonal, $\|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$, $\|\mathbf{Q}\|_2 = 1$, $\|\mathbf{Q}\mathbf{A}\|_2 = \|\mathbf{A}\|_2$ and $\|\mathbf{Q}\mathbf{A}\|_E = \|\mathbf{A}\|_E$.
- Show that $\|\mathbf{A}\|_2 = \|\mathbf{A}^t\|_2$, and hence, from 5 above, that $\|\mathbf{A}\mathbf{Q}\|_2 = \|\mathbf{A}\|_2$, for \mathbf{Q} orthogonal.
- Show that $\|\mathbf{A}\|_E$ is not a subordinate norm for $m, n > 1$. (Hint: use $\mathbf{A} = \mathbf{I}$).
- Show that the 1, 2 and ∞ norms for $n \times 1$ matrices agree with the 1, 2 and ∞ norms for column vectors. What about row vectors?
- Show that any finite product of orthogonal matrices is also orthogonal. (What about an infinite product?)
- A symmetric matrix \mathbf{S} is positive definite if $\mathbf{x}^t \mathbf{S} \mathbf{x} > 0$, $\forall \mathbf{x} \neq \mathbf{0}$.
(i) Show that $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{S}} = \mathbf{x}^t \mathbf{S} \mathbf{y}$ is an inner product if \mathbf{S} is positive definite.
(ii) Let \mathbf{A} be an $m \times n$ real matrix and let the inner product on \mathbb{R}^m be defined as $\langle \mathbf{y}_1, \mathbf{y}_2 \rangle = \mathbf{y}_1^t \mathbf{S} \mathbf{y}_2$, \mathbf{S} positive definite, and the innerproduct on \mathbb{R}^n be defined as $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \mathbf{x}_1^t \mathbf{T} \mathbf{x}_2$, \mathbf{T} positive definite. What is the adjoint mapping of \mathbf{A} ?
- From the definition of a subordinate norm show that the expressions in §4.5(b)(i), (ii) and (iii) do satisfy the four properties of a matrix norm.
- (a) Consider the set $S_1 = \{\mathbf{x} | \mathbf{x} \in \mathbb{R}^2, \|\mathbf{x}\|_1 = 1\}$ (draw it) and the set of points $\{\mathbf{A}\mathbf{x} | \mathbf{x} \in S_1\}$ for a 2×2 matrix \mathbf{A} . Convince yourself that the vectors $\mathbf{x} = \pm \mathbf{e}_i \in S$, $i = 1, 2$ are the only vectors to consider in the definition of $\|\mathbf{A}\|_1$, and hence that the 1-norm of a matrix is the maximum column sum as in §4.5(b)(i).
(b) Repeat part (a) for the 2-norm. The relevant vectors are different, and the eigenvalues and eigenvectors of $\mathbf{A}^t \mathbf{A}$ need to be considered. A linear programming exercise is part of the proof.
(c) Repeat part (a) for the ∞ -norm. The relevant vectors are different and of course the ∞ -norm is defined in terms of the maximum row sum.
- Consider the solution set S of the set of equations $\mathbf{A}\mathbf{x} = \mathbf{y}$.
(a) Show that if \mathbf{E} is a compatible matrix then the solution set of $\mathbf{E}\mathbf{A}\mathbf{x} = \mathbf{E}\mathbf{y}$ contains S .
(b) Show that if \mathbf{E} is invertible then the solution set of $\mathbf{E}\mathbf{A}\mathbf{x} = \mathbf{E}\mathbf{y}$ equals S .
- (a) Show that if $\|\mathbf{A}\| < 1$, then $\|\mathbf{A}^n\| \rightarrow 0$ as $n \rightarrow \infty$ for any matrix norm.
(b) Show that if $\mathbf{A} = \begin{pmatrix} \lambda & a \\ 0 & \lambda \end{pmatrix}$, then $\mathbf{A}^n = \begin{pmatrix} \lambda^n & n\lambda^{n-1}a \\ 0 & \lambda^n \end{pmatrix}$. Hence show that for all $K > 0$, there exists \mathbf{A} such that $\|\mathbf{A}\| > K$ and $\|\mathbf{A}^n\| \rightarrow 0$ as $n \rightarrow \infty$. That is $\|\mathbf{A}\| > 1$ does not imply $\|\mathbf{A}^n\| \rightarrow \infty$.
(c) Does there exist a non-zero matrix \mathbf{A} such that $\mathbf{A}^k = \mathbf{0}$ for $k < n$?
- Show that while $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1$, (exercise 3.3.4), no such relation holds for matrix norms, even though the corresponding matrix norms are defined in terms of vector norms. Experiment with the matrices, $\begin{pmatrix} 3 & 1 \\ 1 & 3 \end{pmatrix}$, $\begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}$, $\begin{pmatrix} 1 & a \\ 0 & 2 \end{pmatrix}$, $\begin{pmatrix} 2 & a \\ 0 & 1 \end{pmatrix}$. Maple or MATLAB would be useful to find the 2-norms.
- Find the co-ordinates of the vector $[2, 3]^t$, in terms of the orthogonal basis

$$\left\{ \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\}.$$

- The two diagonal matrices of order 20,

$$\mathbf{D}_1 = 10^{-1} \mathbf{I}, \quad \text{and} \quad \mathbf{D}_2 = \text{diag}(1, 1, \dots, 1, 10^{-20}),$$

both have the same determinant $10^{-20} \approx 0$. Which of these matrices is closer to being exactly singular? Think about this in terms of the size of perturbation required to make the matrix singular. This shows that the determinant of a matrix is not a good numerical indicator of closeness to singularity.

5. Areas of Application, Types of Matrices.

5.1 Classification of problems.

There are four major areas of application or interest of computational linear algebra.

- (i) Solving equations: Given \mathbf{A} and \mathbf{y} , find \mathbf{x} such that $\mathbf{Ax} = \mathbf{y}$. The methods employed depend on the type of matrix (see §5.2) and the order of the matrix. See also Forsythe and Moler (1967), Golub and van Loan (1989), Stewart (1973), Dongarra *et al.* (1979) and Anderson *et al.* (1995).
- (ii) Solving eigenvalue-eigenvector problems: Find (λ, \mathbf{x}) such that

$$(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}.$$

Once again, the methods used depend on the type of matrix, its order and how many eigenvalues and or eigenvectors are required. More generally, given $\mathbf{M}(\lambda)$ a matrix being a function of λ , find λ, \mathbf{x} such that

$$\mathbf{M}(\lambda)\mathbf{x} = \mathbf{0}.$$

See also Wilkinson (1965), Golub and van Loan (1989), Garbow *et al.* (1977), Smith *et al.* (1976) and Anderson *et al.* (1995).

- (iii) Computation of matrix expressions like projection operators (e.g. $\mathbf{A}(\mathbf{A}^t\mathbf{A})^{-1}\mathbf{A}^t$) and updating formulae (e.g. $(\mathbf{A} + \mathbf{u}\mathbf{v}^t)^{-1}$ in terms of \mathbf{A}^{-1}) in optimization and linear programming. Very often the algebraic expressions used to represent these operations do not provide an efficient and numerically stable method of computation. Factorizations employed in (i) above are used to simplify these expressions. See also Golub and van Loan (1989), or Trefethen and Bau (1997), or Demmel (1997).
- (iv) Error Analysis: Finding expressions for error bounds for the various procedures used above, with respect to error in data and the rounding error of computers. See also Wilkinson (1965).

5.2 Types of Matrices.

5.2.1 Non-square matrices (rectangular matrices).

These have application mainly in statistics, econometrics and statistical modelling. The inverse problem (§4.3) has to be re-stated to give a meaningful result and is re-defined in chapter 9. If \mathbf{A} is an $m \times n$ matrix then the system of equations, given \mathbf{A} and \mathbf{y} , find \mathbf{x} ,

$$\mathbf{Ax} : \mathbf{y}$$

is called over-determined if $m > n$, and under-determined if $m < n$.

5.2.2 Square matrices.

The following divisions of square matrices do not divide the set of square matrices into disjoint sets. The classifications are useful for the different uses as described in §5.1, and is not exhaustive.

- (i) Orthogonal matrices: $\mathbf{Q}^t\mathbf{Q} = \mathbf{Q}\mathbf{Q}^t = \mathbf{I}$, if \mathbf{Q} is real.
Unitary matrices: $\overline{\mathbf{Q}}^t\mathbf{Q} = \mathbf{Q}\overline{\mathbf{Q}}^t = \mathbf{I}$, if \mathbf{Q} is complex.
- (ii) Real symmetric matrices: Defined by $\mathbf{A}^t = \mathbf{A}$, or $A_{ij} = A_{ji}$, while in the complex case they are called Hermitian matrices defined by $\overline{\mathbf{A}}^t = \mathbf{A}$, or $\overline{A}_{ij} = A_{ji}$.
 - (a) Positive definite (symmetric) matrices: $\mathbf{x}^t\mathbf{Ax} > 0$, $\mathbf{x} \neq \mathbf{0}$ and real, if \mathbf{A} is real. (And $\overline{\mathbf{x}}^t\mathbf{Ax} > 0$, $\mathbf{x} \neq \mathbf{0}$, if \mathbf{A} is complex Hermitian.)
 - (b) Positive semi-definite matrices: $\mathbf{x}^t\mathbf{Ax} \geq 0$, $\mathbf{x} \neq \mathbf{0}$.
 - (c) Negative definite matrices: $\mathbf{x}^t\mathbf{Ax} < 0$, $\mathbf{x} \neq \mathbf{0}$.
- (iii) Anti-symmetric matrices: defined by the property $\mathbf{A}^t = -\mathbf{A}$ or $A_{ij} = -A_{ji}$.
- (iv) Full matrices: order n^2 non zero elements.
- (v) Sparse matrices: order n non zero elements.
 - (a) Banded matrices: $A_{ij} = 0, |i - j| > k$ give a banded matrix with band width $2k + 1$. Special cases are the bidiagonal and tridiagonal matrices. For non symmetric matrices the upper and lower band widths are usually specified as generally they would be different.

- (b) Structured sparse matrices (e.g. from finite element models).
- (c) Random sparse matrices.
- (vi) Upper or lower triangular matrices: $U_{ij} = 0, i > j; L_{ij} = 0, i < j$.
Strictly upper or lower triangular: $U_{ij} = 0, i \geq j; L_{ij} = 0, i \leq j$.
- (vii) Upper or lower Hessenberg matrices $U_{ij} = 0, i > j + 1; L_{ij} = 0, i < j - 1$.
- (viii) Stable matrix: for linear constant coefficient differential equations, the real part of all eigenvalues is negative, while for linear constant coefficient difference equations the real part of all eigenvalues must lie within $(-1, 1)$.

5.3 Exercises.

1. Describe (a) a positive definite upper triangular matrix, (b) a symmetric upper Hessenberg matrix.
2. (i) Show that the product of lower triangular matrices is lower triangular. (Similarly upper triangular matrices.)
(ii) Is there a similar property for upper or lower Hessenberg matrices? What about the product of an upper triangular matrix and an upper Hessenberg matrix (two cases, pre and post multiplication)?
(iii) Show that the inverse of an upper triangular matrix, if it exists, is also upper triangular.
(iv) Let U_1 and U_2 be $n \times n$ upper triangular matrices. Show that if either U_1 or U_2 is invertible and $U_1 A = U_2$, then A is upper triangular. Show by example (2×2 is enough) that if either U_1 or U_2 is not invertible then A may not be upper triangular.
3. Show that the product of two tridiagonal matrices is a pentadiagonal matrix. Is there a general rule for the product of two band matrices with band width $2k + 1$ and $2\ell + 1$ say?
4. Is the product of two symmetric matrices also symmetric? If not, what other additional properties are needed?
5. Under what conditions are the eigenvalues of the product of two matrices equal to the product of the eigenvalues of the matrices?
6. Show that the rank of an $n \times n$ upper triangular matrix is less than n if at least one diagonal u_{ii} is zero. Is the rank of U equal to the number of non-zero diagonal elements? Is the rank of U equal to the number of nonzero eigenvalues?
7. Show that every matrix can be written (uniquely) as the sum of a symmetric matrix and an anti-symmetric matrix.
8. Show that every square matrix can be written as the product of two symmetric matrices.

6. Elementary Operations and Elementary Matrices.

6.1 Elementary operations.

The reader is no doubt familiar with the following invertible elementary operations on a given matrix used to solve a set of equations, by the procedure called Gauss elimination.

- (i) Multiplying a row (column) by a non-zero scalar.
- (ii) Interchanging rows (columns).
- (iii) Adding a multiple of a row (column) to another row (column).
- (iv) An extension of (iii) is to replace a set of rows (columns) by linear combinations of that set of rows (columns).

In this chapter, the above processes are described in terms of pre (rows) or post (columns) multiplications of a given matrix by a matrix which has the same effect as an elementary operation. The matrices used to do this are called elementary matrices. The operations will be described in terms of row operations which are equivalent to pre-multiplication by an elementary matrix. The extension to column operations is straight forward, as a post multiplication by an elementary matrix. If \mathbf{A} is the $m \times n$ matrix to be operated on then the process

$$\mathbf{B} = \mathbf{EA},$$

where \mathbf{E} is an elementary matrix, can be written

$$\kappa_i(\mathbf{B}) = \mathbf{E}\kappa_i(\mathbf{A}), \quad i = 1, 2, \dots, n.$$

That is, for pre-multiplication, the matrix \mathbf{A} is considered as a sequence of independent columns. Hence only operations of the type,

$$\mathbf{b} = \mathbf{Ea},$$

are considered, where \mathbf{a} and \mathbf{b} are m -vectors. Most of the elementary operations can be expressed as the matrix $\mathbf{I} + \mathbf{uv}^t$ for some non-zero vectors \mathbf{u} and \mathbf{v} , and the matrix must be invertible. See exercise 4.9.13.

6.2 Multiplying a row by a non-zero scalar.

Find \mathbf{D}_i such that

$$\mathbf{b} = \mathbf{D}_i\mathbf{a},$$

where $b_j = a_j$, $j \neq i$, $b_i = \alpha a_i$. Obviously $\mathbf{D}_i = \text{diag}(1, 1, \dots, 1, \alpha, 1, \dots, 1)$, where α is in the i -th position. Note that \mathbf{D}_i can also be written as $\mathbf{D}_i = \mathbf{I} + (\alpha - 1)\mathbf{e}_i\mathbf{e}_i^t$, and its inverse as $\mathbf{D}_i^{-1} = \mathbf{I} + (\frac{1}{\alpha} - 1)\mathbf{e}_i\mathbf{e}_i^t$. This operation is usually used to produce a 'nice' number like 1 in some appropriate position when working by hand. It is usually inefficient to do this in a computer program.

6.3 Interchanging rows (elementary permutation matrices).

If $\mathbf{P}_{ij} = \mathbf{I} - (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^t$, and $\mathbf{b} = \mathbf{P}_{ij}\mathbf{a}$, then $b_k = a_k$, $k \neq i$ or j , and $b_i = a_j$, $b_j = a_i$. \mathbf{P}_{ij} is an identity matrix with rows i and j interchanged. Note that \mathbf{P}_{ij} is symmetric and orthogonal so that $\mathbf{P}_{ij}\mathbf{P}_{ij} = \mathbf{I}$. A product of elementary permutation matrices is called a permutation matrix. Note that interchanging two rows of a matrix takes a sizeable amount of time on a computer due to the memory accesses. Many computer programs do not interchange, rather they keep track of which rows were done in order.

6.4 Subtracting a multiple of one row from other rows.

Usually this is done to introduce zeros in various positions of a column of a matrix. That is, find a matrix \mathbf{E}_i of the form

$$\mathbf{E}_i = \mathbf{I} + \mathbf{uv}^t,$$

such that if $\mathbf{b} = \mathbf{E}_i\mathbf{a}$, then $b_j = a_j$, $j = 1, 2, \dots, i$, and $b_{i+1} = b_{i+2} = \dots = b_m = 0$. (Zeroing elements below the i -th diagonal element of the i -th column as in Gauss elimination.)

$$\begin{aligned} \mathbf{E}_i\mathbf{a} &= (\mathbf{I} + \mathbf{uv}^t)\mathbf{a} \\ \mathbf{b} &= \mathbf{a} + (\mathbf{v}^t\mathbf{a})\mathbf{u}. \end{aligned}$$

Now choose the sign of α so that subtractive cancellation does not take place in (6.5.3). Hence α has opposite sign to a_i ,

$$\alpha = -\text{sign}(a_i) \left(\sum_{j=i}^m a_j^2 \right)^{\frac{1}{2}}, \tag{6.5.4}$$

and the sequence of calculation is (i) (6.5.1) $|\alpha|$, (ii) (6.5.4) α , (iii) (6.5.3) γ , (iv) (6.5.2) w_j , $j = i, i + 1, \dots, m$.

6.6 Plane rotations (Given's rotations).

These elementary operations replace two rows by linear combinations of the two rows. Let \mathbf{P}^{ij} denote the rotation of the i -th and j -th rows of a vector \mathbf{a} . \mathbf{P}^{ij} is orthogonal and is usually used to zeroise one element of a vector (say a_j). Hence

$$\mathbf{P}^{ij} \mathbf{a} = \mathbf{b} \quad \text{say,}$$

is defined by $a_k = b_k$, $k \neq i$ or j , $b_j = 0$. \mathbf{P}^{ij} is orthogonal so $b_i^2 = a_i^2 + a_j^2$ as these are the only elements of \mathbf{a} to change. \mathbf{P}^{ij} is an identity matrix with the i -th and j -th row and column intersections replaced with one of

$$\begin{matrix} & \begin{matrix} i & j \end{matrix} \\ \begin{matrix} i \\ j \end{matrix} & \begin{pmatrix} c & s \\ s & -c \end{pmatrix} \end{matrix} \quad \text{or} \quad \begin{matrix} & \begin{matrix} i & j \end{matrix} \\ \begin{matrix} i \\ j \end{matrix} & \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \end{matrix}$$

rotation + reflection pure rotation.

It is best to use the symmetric one as the sign in the anti-symmetric one is bothersome when considering inverses. Only consider that part of $\mathbf{b} = \mathbf{P}^{ij} \mathbf{a}$ which is involved in the computation, that is,

$$\begin{pmatrix} c & s \\ s & -c \end{pmatrix} \begin{pmatrix} a_i \\ a_j \end{pmatrix} = \begin{pmatrix} b_i \\ 0 \end{pmatrix}$$

$$b_i = \sqrt{a_i^2 + a_j^2}, \Rightarrow c = a_i/b_i \text{ and } s = a_j/b_i.$$

Note that $c^2 + s^2 = 1$, hence the notation, $c = \cos \theta$, $s = \sin \theta$, where θ is the angle of rotation of the i -th and j -th plane.

It can be shown that the plane rotation plus reflection matrix is also a Householder transformation (see exercise 6.7.4). Generally if many elements are to be zeroed the Householder transformation is more efficient than using many plane rotations by a factor of 2. The reader might like to compare the plane rotations with the elementary matrix of Gauss elimination type for zeroing one element of a column.

6.7 Exercises.

1. (i) If $\mathbf{u} \neq \mathbf{0}$ and $\alpha \mathbf{u} + \beta \mathbf{x} = \mathbf{0}$, show that either α and β both equal zero or $\mathbf{x} = \gamma \mathbf{u}$, $\gamma \in \mathbb{R}$.
 (ii) If $\mathbf{u} \neq \mathbf{0}$, $\mathbf{v} \neq \mathbf{0}$ and $\mathbf{u}\mathbf{v}^t + \mathbf{a}\mathbf{b}^t = \mathbf{0}$, using (i) show that there exists nonzero scalars α and β such that $\mathbf{a} = \alpha \mathbf{u}$ and $\mathbf{b} = \beta \mathbf{v}$.
 (iii) What is the inverse of $\mathbf{E} = \mathbf{I} + \mathbf{u}\mathbf{v}^t$ and when does it exist? Assume $\mathbf{E}^{-1} = \mathbf{I} + \mathbf{a}\mathbf{b}^t$ and use (ii).
 (iv) What are the eigenvalues and eigenvectors of $\mathbf{E} = \mathbf{I} + \mathbf{u}\mathbf{v}^t$? Use (i) on $(\mathbf{I} + \mathbf{u}\mathbf{v}^t)\mathbf{x} = \lambda \mathbf{x}$.
2. (i) Repeat 1(iii) and 1(iv) for $\mathbf{E} = \mathbf{E}_i = \mathbf{I} - \mathbf{m}_i \mathbf{e}_i^t$.
 (ii) Examine the i -th major step of Gauss elimination (which is to zero elements below the i -th diagonal) as you have done it in the past and compare it to using \mathbf{E}_i , that is, expand and interpret $\rho_j(\mathbf{E}_i \mathbf{A})$ in terms of rows. See 4.9.1.
 (iii) What effect does \mathbf{E}_i have on \mathbf{e}_j , $j \leq i$ on both pre-multiplication ($\mathbf{E}_i \mathbf{e}_j$) and post-multiplication ($\mathbf{e}_j^t \mathbf{E}_i$)?
 (iv) Find an expression for $\|\mathbf{E}_i\|_\infty$.
 (v) Find an expression for $\mathbf{E}_i \mathbf{E}_j$ and show the pattern of zero and non-zero elements in the product matrix. (Cases $i < j$, $i = j$, $i > j$.)
 (vi) Find an expression for $\det(\mathbf{E}_i)$.

- (vii) What form does $\mathbf{E}_i = \mathbf{I} + \mathbf{w}\mathbf{v}^t$ take in Gauss-Jordan elimination (which zeros above and below the diagonal)?
3. (i) Show that if $\mathbf{I} + \mathbf{w}\mathbf{v}^t$ is orthogonal then $\mathbf{u} = -\mathbf{v}$ and $\|\mathbf{u}\| = \sqrt{2}$.
- (ii) Where is $\|\mathbf{w}\|_2 = 1$ used in computing \mathbf{H}_i ?
- (iii) Repeat the calculation of \mathbf{H}_i using $\mathbf{H}_i = \mathbf{I} - \mathbf{w}\mathbf{w}^t$, $\|\mathbf{w}\|_2 = \sqrt{2}$.
- (iv) Can \mathbf{H}_i ever be undefined? If so, how?
- (v) Write down \mathbf{H}_i^{-1} in terms of \mathbf{w} .
- (vi) What are the eigenvalues and eigenvectors of \mathbf{H}_i .
- (vii) Show $\mathbf{H}_i\mathbf{e}_j = \mathbf{e}_j$ if $j < i$.
- (viii) If α is not chosen by (6.5.4), show that it is possible for $\gamma = 0$, and hence for \mathbf{H}_i to be undefined. But does it matter in this case?
- (ix) Find an expression for $\det(\mathbf{H}_i)$.
- (x) The complex case has $\mathbf{H}_i = \mathbf{I} - 2\mathbf{w}\bar{\mathbf{w}}^t$. Repeat the computation for a complex vector \mathbf{a} . Note that \mathbf{H}_i is now unitary.
4. (i) Show that $\begin{pmatrix} c & s \\ s & -c \end{pmatrix}$ (actually a rotation and reflection) can be written in the form $\mathbf{I} - 2\mathbf{w}\mathbf{w}^t$.
- (ii) Show that $\begin{pmatrix} c & s \\ -s & c \end{pmatrix}$ (a rotation) can be written in the form $\mathbf{J} - 2\mathbf{w}\mathbf{w}^t\mathbf{J}$ or $\mathbf{J} - 2\mathbf{J}\mathbf{w}\mathbf{w}^t$ where $\mathbf{J} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$, a reflection.
- (iii) If $\theta = \pi/2$ what is the effect of the plane rotation, and, plane rotation and reflection.
5. (i) Find the plane rotations which will zero the second and third elements of the vector $[1, 2, 2]^t$.
- (ii) Find the Householder matrix to cause the same effect.
6. Let \mathbf{H} be a Householder matrix defined by $\mathbf{H}\mathbf{a} = \|\mathbf{a}\|\mathbf{e}_1$. Let $\mathbf{H}\mathbf{b} = \mathbf{c}$. Show that $c_1 = \langle \mathbf{a}, \mathbf{b} \rangle / \|\mathbf{a}\|$, that is, the projection of \mathbf{b} onto a unit vector in the direction \mathbf{a} , and hence that $\mathbf{H}[c_1, 0, 0, \dots, 0]^t$ is parallel to \mathbf{a} . Show that $\mathbf{H}[0, c_2, c_3, \dots, c_n]^t$ is perpendicular to \mathbf{a} .
7. Show that each elementary matrix can be obtained by doing the elementary operation on an identity matrix.
8. (i) Compare the number of arithmetic operations to compute \mathbf{E}_i and \mathbf{H}_i .
- (ii) Compare the number of arithmetic operations to apply \mathbf{E}_i and \mathbf{H}_i to an arbitrary vector.
- (iii) Compare the number of arithmetic operations to compute and apply a Householder transformation to zero elements below the i -th with the corresponding sequence of plane rotations which zero the same elements.

7. Decompositions, Factorizations of Matrices.

7.1 The singular value decomposition (SVD).

This factorization is very useful as a tool to gain some geometric insight into a linear mapping from \mathbb{R}^n to \mathbb{R}^m . It was not used in computation to a great extent, but with faster workstations it is now used more often.

Theorem 7.1.1 An arbitrary $m \times n$ real matrix \mathbf{A} can be factored to the matrix product

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^t$$

where \mathbf{V} is $m \times m$ orthogonal, \mathbf{U} is $n \times n$ orthogonal, and \mathbf{D} is $m \times n$, diagonal with one of the forms,

$$\left(\begin{array}{cccc|c} d_1 & & & & \\ & d_2 & & & \\ & & \ddots & & \\ & & & d_m & \\ \hline & & & & \circ \end{array} \right) \quad \text{or} \quad \left(\begin{array}{cccc} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \\ \hline - & - & - & - \\ & & \circ & \end{array} \right).$$

Further, $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_\ell)$, $\ell = \min(m, n)$, and d_i are ordered so that $d_1 \geq d_2 \geq \dots \geq d_r > 0$, $d_{r+1} = \dots = d_\ell = 0$, where $r = \text{rank}(\mathbf{A})$. The d_i are known as the singular values of \mathbf{A} . The SVD can also be written

$$\mathbf{A} = \sum_{i=1}^r d_i \mathbf{v}_i \mathbf{u}_i^t,$$

that is, the sum of r rank 1 matrices $\mathbf{v}_i \mathbf{u}_i^t$, where $\mathbf{v}_i = \kappa_i(\mathbf{V})$ and $\mathbf{u}_i = \kappa_i(\mathbf{U}) = \rho_i(\mathbf{U}^t)$. ■

See Golub and Kahan (1965) for the original modern algorithm, which involves an implicit eigenvalue/eigenvector computation, and Eckart and Young (1936) for an original use in statistics.

7.2 The Q–U factorization (or Q–R factorization).

Theorem 7.2.1

(a) Every $m \times n$ matrix \mathbf{A} , $m \geq n$, with rank n (full rank) can be factored to the form

$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{U} \\ \mathbf{0} \end{pmatrix},$$

where \mathbf{Q} is orthogonal and \mathbf{U} is $n \times n$ upper triangular and non-singular.

(b) Every $m \times n$ matrix \mathbf{A} , $m \leq n$, with rank m (full rank) can be factored to the form

$$\mathbf{A} = (\mathbf{L} \ \mathbf{0}) \mathbf{Q},$$

where \mathbf{Q} is $n \times n$ orthogonal and \mathbf{L} is $m \times m$ lower triangular with rank m .

(c) Every $m \times n$ matrix \mathbf{A} of rank r can be factored to the form

$$\mathbf{A} = \mathbf{Q}_1 \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}_2,$$

where \mathbf{Q}_1 is $m \times m$ orthogonal, \mathbf{Q}_2 is $n \times n$ orthogonal and \mathbf{L} is $r \times r$ lower triangular of rank r .

Proof:

(a) Algorithm:

$$\mathbf{A}_0 = \mathbf{A}.$$

For $i = 1, 2, \dots, n$

$$\text{do } \mathbf{A}_i = \mathbf{H}_i \mathbf{A}_{i-1};$$

where H_i is a Householder transformation which zeros column i of A_{i-1} below the diagonal element. Then

$$\begin{aligned} A_n &= H_n A_{n-1} \\ &= H_n H_{n-1} \cdots H_1 A \\ &= Q^t A, \end{aligned}$$

where $Q = H_1 H_2 \cdots H_n$ is an orthogonal matrix. The matrix A_n has the form $\begin{pmatrix} U \\ 0 \end{pmatrix}$, hence

$$A = Q \begin{pmatrix} U \\ 0 \end{pmatrix}.$$

Note:

- (i) Verify that the application of H_i to A_{i-1} does not destroy the zeros already produced in columns $1, 2, \dots, i-1$ of A_{i-1} .
 - (ii) Is it possible that one of the H_i are undefined? (See §6.7 exercise 3(v).)
- (b) Apply (a) to A^t .
- (c) Try this as an exercise being careful with the rank question as in (ii) above. Some column interchanges might be necessary. ■

Note: Part (a) is a statement of Gram-Schmidt orthogonalization (see any introductory linear algebra text, for example, Anton and Rorres (1987)) but the Gram-Schmidt algorithm is numerically unstable, while the Householder algorithm is stable. Björck (1968) has an alternative stable algorithm.

7.3 The L–U factorization.

Theorem 7.3.1 Under certain conditions (stated in proof) an $n \times n$ matrix A of rank n can be factored to the form

$$A = LU,$$

where L is lower triangular with diagonal elements unity and U is upper triangular of rank n . Note that this is an algebraic statement of the Gauss elimination algorithm with no row interchanges.

Proof:

- (i) Algorithm:

$$A_0 = A$$

For $i = 1, 2, \dots, n-1$

$$\text{do } A_i = E_i A_{i-1};$$

where E_i is the elementary matrix of type discussed in §6.4. E_i zeros the elements in the i -th column of A_{i-1} below the diagonal. Hence

$$\begin{aligned} A_{n-1} &= E_{n-1} A_{n-2} \\ &= E_{n-1} E_{n-2} \cdots E_1 A, \end{aligned}$$

and

$$A = E_1^{-1} E_2^{-1} \cdots E_{n-1}^{-1} A_{n-1}.$$

- (ii) It remains to show that A_{n-1} is upper triangular of full rank, and the product $E_1^{-1} \cdots E_{n-1}^{-1}$ is lower triangular. At each step, E_i does not affect the zeros already introduced in columns $1, 2, \dots, i-1$ of A_{i-1} (exercise 2(iv) in §6.7). This shows that A_{n-1} is upper triangular. A_{n-1} has full rank because it can be shown that the product $E_1^{-1} \cdots E_{n-1}^{-1}$ is nonsingular.
- (iii) When is E_i defined? Consider the i -th step of the algorithm. A_{i-1} has the form (let the (k, j) -th element of A_{i-1} be denoted $a_{k,j}$)

$$A_{i-1} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1i} & \cdots & a_{1n} \\ & a_{22} & \cdots & a_{2i} & \cdots & a_{2n} \\ & & \ddots & \vdots & & \vdots \\ & & & \bigcirc & a_{ii} & \cdots & a_{in} \\ & & & & a_{i+1,i} & & \\ & & & & \vdots & & \vdots \\ & & & & a_{n,i} & \cdots & a_{nn} \end{pmatrix}$$

Proof:

(i) Algorithm:

$$\mathbf{A}_0 = \mathbf{A}$$

For $i = 1, 2, \dots, n - 1$

$$\text{do } \mathbf{A}_i = \mathbf{E}_i \mathbf{P}_i \mathbf{A}_{i-1}.$$

(a) \mathbf{P}_i is an elementary permutation matrix which interchanges rows of \mathbf{A}_{i-1} so that the (i, i) -th element of $\mathbf{P}_i \mathbf{A}_{i-1}$ is the largest (in modulus) of the elements of the i -th column below and on the diagonal, i.e.

$$|(\mathbf{P}_i \mathbf{A}_{i-1})_{ii}| \geq |(\mathbf{P}_i \mathbf{A}_{i-1})_{ji}|, \quad j = i, i + 1, \dots, n.$$

(b) \mathbf{E}_i is defined to be \mathbf{I} if $(\mathbf{P}_i \mathbf{A}_{i-1})_{i,i}$ is zero, otherwise \mathbf{E}_i is the usual elementary matrix for zeroing the i -th column below the diagonal

$$\mathbf{E}_i = \mathbf{I} - \mathbf{m}_i \mathbf{e}_i^t,$$

where now \mathbf{m}_i is based on the i -th column of $\mathbf{P}_i \mathbf{A}_{i-1}$. Note that the elements of \mathbf{m}_i are all less than or equal to 1 in magnitude. This pivoting procedure takes account of any rank deficiency in \mathbf{A} (why?).

(ii) The algorithm produces

$$\mathbf{A}_{n-1} = \mathbf{E}_{n-1} \mathbf{P}_{n-1} \mathbf{E}_{n-2} \mathbf{P}_{n-2} \cdots \mathbf{E}_2 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 \mathbf{A},$$

and

$$\mathbf{A}_{n-1} (= \mathbf{U}) = \mathbf{L}^{-1} \mathbf{P}^{-1} \mathbf{A},$$

is required in order to prove the theorem. This can be done provided all \mathbf{P}_i can be shifted to the inside and leave the \mathbf{E}_i essentially untouched in the positions of non-zero elements.

(a) Consider $\mathbf{P}_i \mathbf{E}_j \mathbf{P}_i$, $i > j$. \mathbf{P}_i interchanges rows and columns of \mathbf{E}_j using rows and columns bigger than the i -th.

$$\begin{aligned} \mathbf{P}_i (\mathbf{I} - \mathbf{m}_j \mathbf{e}_j^t) \mathbf{P}_i &= \mathbf{I} - \mathbf{P}_i \mathbf{m}_j \mathbf{e}_j^t, \text{ as } \mathbf{e}_j^t \mathbf{P}_i = \mathbf{e}_j^t \\ &= \mathbf{I} - \overline{\mathbf{m}}_j \mathbf{e}_j^t, \text{ say,} \end{aligned}$$

where $\overline{\mathbf{m}}_j$ differs from \mathbf{m}_j in that two non-zero elements have been interchanged. Hence

$$\mathbf{P}_i \mathbf{E}_j \mathbf{P}_i = \overline{\mathbf{E}}_j,$$

where $\overline{\mathbf{E}}_j$ has the same form as \mathbf{E}_j in terms of zeros and non-zeros.

(b) Now write (using (a) above)

$$\begin{aligned} \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 &= \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_2 \mathbf{P}_2 \mathbf{P}_1 \\ &= \overline{\mathbf{E}}_1 \mathbf{P}_2 \mathbf{P}_1, \end{aligned}$$

and further

$$\mathbf{U} = \mathbf{E}_{n-1} \overline{\mathbf{E}}_{n-2} \overline{\overline{\mathbf{E}}}_{n-3} \cdots \overline{\overline{\overline{\mathbf{E}}}}_1 \mathbf{P}_{n-1} \mathbf{P}_{n-2} \cdots \mathbf{P}_1 \mathbf{A}.$$

Hence $\mathbf{A} = \mathbf{P} \mathbf{L} \mathbf{U}$, where

$$\begin{aligned} \mathbf{P} &= \mathbf{P}_1 \mathbf{P}_2 \cdots \mathbf{P}_{n-1}, \\ \mathbf{L} &= \overline{\overline{\overline{\overline{\mathbf{E}}}}}^{-1}_1 \cdots \overline{\overline{\mathbf{E}}}_{n-2}^{-1} \mathbf{E}_{n-1}^{-1}, \end{aligned}$$

and is still lower triangular. ■

The above is known as partial pivoting or row pivoting. The result (ii)(b) above,

$$\mathbf{P}^{-1} \mathbf{A} = \mathbf{L} \mathbf{U},$$

is saying that if the row interchanges which need to be done were done first and then the algorithm of Theorem 7.3.1 applied to the new matrix $\mathbf{P}^{-1} \mathbf{A}$, there would be no difficulty with zero divisors.

The next theorem gives the result for full pivoting where at each stage, rows and columns are interchanged in order that $(\mathbf{A}_{i-1})_{ii}$ is the largest element of the lower right partition of \mathbf{A}_{i-1} .

Theorem 7.3.3 (Full pivoting) Every $m \times n$ matrix A , of rank r , can be factored to a form,

$$A = PLUQ$$

where P and Q are permutation matrices and L is $m \times r$ lower triangular and U is $r \times n$ upper triangular.

$$A = P \begin{pmatrix} 1 & & \circ \\ \vdots & \ddots & \\ L_{r1} & & 1 \\ \vdots & & \vdots \\ L_{m1} & \cdots & L_{mr} \end{pmatrix} \begin{pmatrix} U_{11} & \cdots & U_{1r} & \cdots & U_{1n} \\ & \ddots & & & \vdots \\ \circ & & U_{rr} & \cdots & U_{rn} \end{pmatrix} Q$$

Proof: — Exercise! ■

From $P^{-1}AQ^{-1} = LU$, if the appropriate permutations of rows and columns were done before Gauss elimination, then the biggest element of the remaining lower sub-matrix would be in the top left position, the pivot position.

7.4 Choleski factorization and variants.

Theorem 7.4.1 A positive definite matrix (symmetric and square) has a factorization of the form

$$A = LL^t,$$

where L is lower triangular and invertible.

Proof: See for example the book by Stewart and exercise 7.7.7.

Algorithm: By solving the $n(n + 1)/2$ defining equations in a particular order all elements of L are easily found.

| Defining Equations | Component of L |
|---|---|
| (i) $A_{11} = L_{11}^2$ | $L_{11} = \sqrt{A_{11}}$ |
| (ii) For $j = 2, 3, \dots, n$ do $A_{j1} = L_{j1}L_{11}$ | $L_{j1} = A_{j1}/L_{11}$ |
| (iii) For $i = 2, \dots, n$ do $A_{ii} = \sum_{k=1}^i L_{ik}^2$ For $j = i + 1, i + 2, \dots, n$ do $A_{ji} = \sum_{k=1}^i L_{jk}L_{ik}$ | $L_{ii} = \sqrt{(A_{ii} - \sum_{k=1}^{i-1} L_{ik}^2)}$ $L_{ji} = (A_{ji} - \sum_{k=1}^{i-1} L_{jk}L_{ik})/L_{ii}$ ■ |

Corollary 7.4.1 A positive definite matrix has a factorization

$$A = LDL^t$$

where L is lower triangular with $L_{ii} = 1$, $i = 1, 2, \dots, n$ and D is a positive diagonal matrix.

Proof: Exercise. ■

There is no need to consider pivoting in this algorithm because the matrix is non-singular, and the error analysis shows that pivoting is not needed to reduce error magnification. The upper triangular matrix DL^t is equal to U of the L-U factorization, and the L matrices of both factorizations are equal.

Corollary 7.4.2 Some symmetric $n \times n$ matrices A with rank n , can be factored to

$$A = LDL^t,$$

where now D is diagonal but not necessarily positive definite, and $L_{ii} = 1$.

Proof: Exercise. ■

If D has both positive and negative diagonal elements the factorization is unstable with respect to roundoff error. A negative definite matrix will have all $D_{ii} < 0$. There is no need to have alternative software for negative definite matrices A , as $-A$ is positive definite.

Corollary 7.4.3 If A is symmetric of rank n then A can be factored to

$$A = LBL^t,$$

where \mathbf{L} is as in Corollary 7.4.1 and \mathbf{B} is block diagonal, with either a 1×1 block or a 2×2 block with negative determinant.

Proof: See Dongarra *et al.* (1979). Exercise. ■

This factorization is stable, and is used in LAPACK.

7.5 Banded matrices.

7.5.1 The L–U factorization.

Let \mathbf{B} be a banded matrix with ℓ bands below the diagonal and p bands above the diagonal ($\ell + p + 1$ bands altogether). Provided no pivoting is required then \mathbf{B} has an L–U factorization

$$\mathbf{B} = \mathbf{L}\mathbf{U},$$

where \mathbf{L} is lower triangular with a diagonal of ones and ℓ bands below it, and, \mathbf{U} is upper triangular with p bands above the diagonal. That is, the structure of \mathbf{B} is preserved in \mathbf{L} and \mathbf{U} . Proof of this result follows as an application of Theorem 7.3.1.

If pivoting is used the structure of \mathbf{L} and \mathbf{U} is destroyed. However if only row pivoting (partial pivoting) is allowed then

$$\mathbf{B} = \mathbf{P}\mathbf{L}\mathbf{U},$$

where \mathbf{P} is a permutation matrix, \mathbf{L} as above and \mathbf{U} is upper triangular with $p + \ell$ bands above the diagonal. The proof follows by considering Theorem 7.3.2 and the worst possible pivoting strategy in terms of “fill-in”.

7.5.2 Choleski.

If \mathbf{B} is positive definite banded (symmetric) with $2n + 1$ bands altogether, then \mathbf{B} has a Choleski factorization,

$$\mathbf{B} = \mathbf{L}\mathbf{L}^t,$$

where \mathbf{L} is lower triangular with $n + 1$ bands including the diagonal.

It is possible to re-order the rows and columns of some sparse matrices to make the resulting set of equations have as small a band structure as possible. This helps reduce the number of computations involved in solving the set of equations. This is important in applications where the same matrix or same structured matrix is inverted many times, for example, iterative eigenvalue/eigenvector computations. See Duff (1977) and George and Liu (1981).

7.6 Eigenvalue — eigenvector decomposition.

Let \mathbf{A} be an $n \times n$ matrix. The eigenvalues (characteristic values) of \mathbf{A} are defined as the roots of the polynomial in λ

$$\det|\mathbf{A} - \lambda\mathbf{I}| = 0.$$

There exists n linear factors $(\lambda - \lambda_i)$, $i = 1, 2, \dots, n$ of this polynomial. As $\mathbf{A} - \lambda_i\mathbf{I}$ is singular it must have a null space. Let $\{\mathbf{s}_i\}$ be a set of vectors spanning these spaces. Let \mathbf{S} be the matrix of columns \mathbf{s}_i .

$$(\mathbf{A} - \lambda_i\mathbf{I})\mathbf{s}_i = \mathbf{0},$$

$$\begin{aligned} \mathbf{A}(\mathbf{s}_1 \quad \mathbf{s}_2 \quad \dots \quad \mathbf{s}_n) &= (\lambda_1\mathbf{s}_1 \quad \lambda_2\mathbf{s}_2 \quad \dots \quad \lambda_n\mathbf{s}_n) \\ \mathbf{A}\mathbf{S} &= \mathbf{S}\mathbf{\Lambda}, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \end{aligned}$$

If \mathbf{S}^{-1} exists, this is now written as

$$\mathbf{A} = \mathbf{S}\mathbf{\Lambda}\mathbf{S}^{-1}.$$

It is only under certain conditions that \mathbf{S}^{-1} will exist.

The following examples show how roundoff error might affect decisions about existence of eigenvectors.

(a) $\mathbf{A} = \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}$, $\lambda_1 = \lambda_2 = 1$, $\mathbf{s}_1 = \mathbf{s}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ if $a \neq 0$.

$$(b) \mathbf{A} = \begin{pmatrix} 1 & 1 \\ 0 & 1 + \epsilon \end{pmatrix}, \quad \lambda_1 = 1, \quad \lambda_2 = 1 + \epsilon, \quad \mathbf{s}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{s}_2 = \begin{pmatrix} 1 \\ \epsilon \end{pmatrix},$$

$$\mathbf{S}^{-1} = \begin{pmatrix} 1 & -\epsilon^{-1} \\ 0 & \epsilon^{-1} \end{pmatrix}.$$

Of course a further computational problem is that if \mathbf{A} is real, λ_i may be complex.

7.7 Exercises.

1. (i) Write down all singular value decompositions of a 1×1 matrix. Hence, how many different SVD are there for an $n \times n$ matrix where the d_i are distinct.
 - (ii) What is the SVD of an $m \times 1$ matrix?
 - (iii) The SVD is not uniquely defined. Why and in what way?
 Use the notation of Theorem 7.1.1 in the next eleven questions.
 - (iv) Let \mathbf{A} represent the mapping \mathcal{A} in the usual basis. What is the matrix representation of \mathcal{A} if the basis of \mathbb{R}^n is $\{\mathbf{u}_i\}_{i=1}^n$ and of \mathbb{R}^m is $\{\mathbf{v}_i\}_{i=1}^m$.
 - (v) Show that the number of non-zero singular values is the rank of \mathbf{A} .
 - (vi) Show that $\mathcal{R}(\mathbf{A}) = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$.
 - (vii) Show that $\mathcal{R}(\mathbf{A}^t) = \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$.
 - (viii) Show that $\mathcal{N}(\mathbf{A}) = \text{span}\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_n\}$.
 - (ix) Show that $\mathcal{N}(\mathbf{A}^t) = \text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_m\}$.
 - (x) Show that $\mathbf{A}\mathbf{u}_i = d_i\mathbf{v}_i$.
 - (xi) Show that $\mathbf{A}^t\mathbf{v}_i = d_i\mathbf{u}_i$.
 - (xii) Show that $\|\mathbf{A}\|_2 = \{\lambda_{\max}(\mathbf{A}^t\mathbf{A})\}^{\frac{1}{2}} = d_1$. What is $\|\mathbf{A}^{-1}\|_2$, if \mathbf{A}^{-1} exists?
 - (xiii) Show that $\|\mathbf{A}\|_E = \{\sum_{i=1}^r d_i^2\}^{\frac{1}{2}}$.
 - (xiv) Show that $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_E \leq \sqrt{n}\|\mathbf{A}\|_2$.
2. Show that the result of Theorem 7.2.1(a) can be interpreted as Gram-Schmidt orthogonalization by writing it in matrix form.
3. (i) (Crout factorization) Suppose $\mathbf{A}(n \times n)$ can be factored to the product $\mathbf{L}\mathbf{U}$ with no problems of zero divisors. Show that the elements of \mathbf{L} and \mathbf{U} can be found by solving the n^2 equations

$$A_{ij} = \sum_k L_{ik}U_{kj},$$

in an appropriate order. (Let $L_{ii} = 1$, $i = 1, \dots, n$.) Count the number of multiplications and divisions.

- (ii) Carry out the steps of Gauss elimination, the L-U factorization via elementary matrices and the Crout method on the matrix.

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 3 & 1 & 3 & 4 \\ 2 & -3 & 4 & 8 \end{pmatrix}.$$

(Do not use pivoting in Gauss elimination, and do not divide the pivot row by the diagonal element before doing the row operation.) Can you see that these three methods compute the same numbers, but in different orders?

4. Using the notation of §7.6, prove the following:
 - (i) All λ_i distinct $\Rightarrow \mathbf{S}^{-1}$ exists.
 - (ii) \mathbf{A} symmetric $\Rightarrow \mathbf{S}^{-1}$ exists.
 - (iii) If \mathbf{A} is real and λ_i is not real then $\bar{\lambda}_i$ is also an eigenvalue.
 - (iv) If \mathbf{A} is real symmetric (Hermitian) then \mathbf{S} may be chosen and scaled so that it is orthogonal (unitary).
 - (v) \mathbf{A} is real symmetric (Hermitian) $\Rightarrow \lambda_i$ are real. \mathbf{A} is positive definite $\Rightarrow \lambda_i$ are positive.
 - (vi) What is the relation between the SVD and the eigen decomposition of, (a) a symmetric positive definite matrix, (b) a symmetric matrix.
 - (vii) Show that $\mathbf{A}\mathbf{A}^t$ and $\mathbf{A}^t\mathbf{A}$ have the same non-zero eigenvalues.

- (viii) Examine the relation between the SVD of \mathbf{A} and the eigen decomposition of $\mathbf{A}\mathbf{A}^t$ and $\mathbf{A}^t\mathbf{A}$.
 - (ix) If \mathbf{A} is real symmetric and all eigenvalues are positive, show that \mathbf{A} is positive definite.
 - (x) If \mathbf{A} is anti-symmetric, show that $i\mathbf{A}$ is Hermitian, and hence show that the eigenvalues of \mathbf{A} are purely imaginary.
 - (xi) Show that if \mathbf{A} is anti-symmetric and has odd order, it is singular.
5. Define recursive relations for the elements of \mathbf{L} and \mathbf{U} for tridiagonal band matrices, (without pivoting) both positive definite and otherwise.
 6. Define recursive relations for the elements of \mathbf{L} and \mathbf{U} for pentadiagonal band matrices, (without pivoting) both positive definite and otherwise.
 7. Show that the Choleski factorization $\mathbf{A} = \mathbf{L}\mathbf{L}^t$ exists if and only if \mathbf{A} is positive definite. One part of this proof requires an induction argument on the order of the matrix, and in the process gives an alternative definition of the Choleski algorithm. See question 7.7.12.
 8. Describe Gauss–Jordan elimination in terms of elementary matrices (zeros above and below the diagonal), and count the number of multiplications and divisions required.
 9. Expand $\mathbf{e}_k^t \mathbf{E}_1 \mathbf{A}$, for $k = 1, \dots, n$, where \mathbf{A} is an $n \times n$ matrix and \mathbf{E}_1 is the first elementary matrix of the L-U factorization. Describe the resulting rows of $\mathbf{E}_1 \mathbf{A}$ in terms of rows of \mathbf{A} .
 10. Examine the Simplex tableau algorithm for linear programming and explain it in terms of elementary matrices operating on the tableau. Note that as each major step of Gaussian elimination is done there is no account taken of pivoting requirements. Hence the basis columns of the Simplex algorithm have to be re-inverted, with pivoting, every so many basis change steps, so as to keep some accuracy.
 11. Consider the following partitioned matrix equation, which can be used to describe the L-U factorization without pivoting. Show how \mathbf{x} , \mathbf{y} and z can be computed from \mathbf{A} , \mathbf{a} , \mathbf{b} , c and the already computed factors of \mathbf{A} , namely \mathbf{L} and \mathbf{U} .

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{a}^t & c \end{bmatrix} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{x}^t & 1 \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{y} \\ \mathbf{0}^t & z \end{bmatrix}$$

12. Consider the following partitioned matrix equation, which can be used to describe the Choleski factorization of a symmetric positive definite matrix. Show how \mathbf{x} and z can be computed from \mathbf{A} , \mathbf{a} and c and the already computed factors of \mathbf{A} , namely \mathbf{L} .

$$\begin{bmatrix} \mathbf{A} & \mathbf{a} \\ \mathbf{a}^t & c \end{bmatrix} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{x}^t & z \end{bmatrix} \begin{bmatrix} \mathbf{L}^t & \mathbf{x} \\ \mathbf{0}^t & z \end{bmatrix}$$

13. Use questions 11 and 12 to count the number of operations in L-U and Choleski factorizations, given that it takes $n(n+1)/2$ operations to compute the solution of a set of equations involving an $n \times n$ lower or upper triangular matrix.

8. The Inverse Problem.

The procedures used to solve the set of linear equations,

$$\mathbf{A}\mathbf{x} = \mathbf{y}, \tag{8.0.1}$$

where \mathbf{A} is an $n \times n$ nonsingular matrix, $\mathbf{y} \in \mathbb{R}^n$, a vector termed the right hand side and $\mathbf{x} \in \mathbb{R}^n$, the solution vector are the subjects of this chapter. Both the algorithms for finding \mathbf{x} , and the error analysis of each method are of interest. This chapter builds on the work of chapters 6 and 7. However before considering the methods, some theory is necessary to understand how close to the exact solution a computed solution can be expected to be. This analysis is called a sensitivity analysis or condition number analysis. Part of this analysis gives a clue as to how to make a numerical decision about when a matrix is close to being singular.

8.1 Sensitivity analysis.

Primarily because of error in the data (\mathbf{A} and \mathbf{y}), and to a lesser extent roundoff error in the computer and error generated by the arithmetic unit of the computer, the solution to the perturbed equation

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}, \tag{8.1.1}$$

needs to be considered. The perturbation in \mathbf{x} , namely $\delta\mathbf{x}$ is expressed in terms of the perturbations in the data ($\delta\mathbf{A}$ and $\delta\mathbf{y}$). Under what conditions on $\delta\mathbf{A}$ does $(\mathbf{A} + \delta\mathbf{A})^{-1}$ exist?

Lemma 8.1.1

$(\mathbf{I} + \mathbf{X})^{-1}$ exists for all \mathbf{X} such that $\|\mathbf{X}\| < 1$ where $\|\cdot\|$ is a subordinate norm. ($\|\mathbf{I}\| = 1$) Think about this statement for a 1×1 matrix.

Proof: Note: $(\mathbf{I} + \mathbf{X})\mathbf{x} = \mathbf{0}$ has only one solution, $\mathbf{x} = \mathbf{0}$, if and only if $(\mathbf{I} + \mathbf{X})^{-1}$ exists. Suppose $\mathbf{x} \neq \mathbf{0}$ satisfies $(\mathbf{I} + \mathbf{X})\mathbf{x} = \mathbf{0}$. Then $\mathbf{x} = -\mathbf{X}\mathbf{x}$ and

$$\|\mathbf{x}\| \leq \|\mathbf{X}\| \|\mathbf{x}\|.$$

As $\|\mathbf{X}\| < 1$, and $\|\mathbf{x}\| \neq 0$

$$\|\mathbf{x}\| < \|\mathbf{x}\|.$$

This is a contradiction, so $\mathbf{x} = \mathbf{0}$ is the only solution. Hence $(\mathbf{I} + \mathbf{X})^{-1}$ exists if $\|\mathbf{X}\| < 1$. ■

Lemma 8.1.2

$(\mathbf{A} + \delta\mathbf{A})^{-1}$ exists for all $\delta\mathbf{A}$ such that $\|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1$.

Proof:

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

By lemma 8.1.1, $(\mathbf{I} + \delta\mathbf{A}\mathbf{A}^{-1})^{-1}$ exists if

$$\|\delta\mathbf{A}\mathbf{A}^{-1}\| < 1,$$

which is true if $\|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1$. Hence $(\mathbf{A} + \delta\mathbf{A})^{-1}$ exists if $\|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1$. ■

Lemma 8.1.3

$$\frac{1}{1 + \|\mathbf{X}\|} \leq \|(\mathbf{I} + \mathbf{X})^{-1}\| \leq \frac{1}{1 - \|\mathbf{X}\|}, \text{ if } \|\mathbf{X}\| < 1,$$

where $\|\cdot\|$ is a subordinate norm.

Proof:

$$\begin{aligned} \mathbf{I} &= (\mathbf{I} + \mathbf{X})^{-1}(\mathbf{I} + \mathbf{X}) & (8.1.2) \\ \|\mathbf{I}\| &\leq \|(\mathbf{I} + \mathbf{X})^{-1}\| \|\mathbf{I} + \mathbf{X}\| \\ &\leq \|(\mathbf{I} + \mathbf{X})^{-1}\| (\|\mathbf{I}\| + \|\mathbf{X}\|) \end{aligned}$$

Rearranging,

$$\frac{1}{1 + \|\mathbf{X}\|} \leq \|(\mathbf{I} + \mathbf{X})^{-1}\|.$$

From (8.1.2)

$$\begin{aligned} (\mathbf{I} + \mathbf{X})^{-1} &= \mathbf{I} - (\mathbf{I} + \mathbf{X})^{-1} \mathbf{X} \\ \therefore \|(\mathbf{I} + \mathbf{X})^{-1}\| &\leq 1 + \|(\mathbf{I} + \mathbf{X})^{-1}\| \|\mathbf{X}\|, \end{aligned}$$

and rearranging,

$$\|(\mathbf{I} + \mathbf{X})^{-1}\| \leq \frac{1}{1 - \|\mathbf{X}\|}. \quad \blacksquare$$

Lemma 8.1.4

$$\|(\mathbf{A} + \delta\mathbf{A})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\|} \quad \text{if } \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1.$$

Proof:

$$\begin{aligned} \|(\mathbf{A} + \delta\mathbf{A})^{-1}\| &= \|\mathbf{A}^{-1}(\mathbf{I} + \delta\mathbf{A}\mathbf{A}^{-1})^{-1}\| \\ &\leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\delta\mathbf{A}\mathbf{A}^{-1}\|}, \quad \text{by Lemma 8.1.3} \\ &\leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\|}, \quad \text{as } \|\delta\mathbf{A}\mathbf{A}^{-1}\| \leq \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1 \end{aligned} \quad \blacksquare$$

Theorem 8.1.1

The relative error in the solution of (8.1.1) expressed as $\|\delta\mathbf{x}\|/\|\mathbf{x}\|$ can be related to the relative errors in \mathbf{A} and \mathbf{y} , namely $\|\delta\mathbf{A}\|/\|\mathbf{A}\|$ and $\|\delta\mathbf{y}\|/\|\mathbf{y}\|$ by the expression

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\chi(\mathbf{A})}{1 - \chi(\mathbf{A}) \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left\{ \frac{\|\delta\mathbf{y}\|}{\|\mathbf{y}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right\} \quad (8.1.3)$$

if $\|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| < 1$, where $\chi(\mathbf{A})$ is the condition number of the matrix \mathbf{A} , defined by

$$\chi(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|.$$

Note the norm for matrices is a subordinate norm to the vector norm.

Proof: Subtract (8.0.1) from (8.1.1) to give after rearranging

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})\delta\mathbf{x} &= \delta\mathbf{y} - \delta\mathbf{A}\mathbf{x} \\ \delta\mathbf{x} &= (\mathbf{A} + \delta\mathbf{A})^{-1} \{\delta\mathbf{y} - \delta\mathbf{A}\mathbf{x}\} \\ \|\delta\mathbf{x}\| &\leq \|(\mathbf{A} + \delta\mathbf{A})^{-1}\| \{ \|\delta\mathbf{y}\| + \|\delta\mathbf{A}\| \|\mathbf{x}\| \} \end{aligned}$$

From Lemma 8.1.4

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\|} \left\{ \frac{\|\delta\mathbf{y}\|}{\|\mathbf{x}\|} + \|\delta\mathbf{A}\| \right\}.$$

Inserting $\|\mathbf{A}\|$ in appropriate places,

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\| \|\mathbf{A}\|}{1 - \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{A}\|}{\|\mathbf{A}\|}} \left\{ \frac{\|\delta\mathbf{y}\|}{\|\mathbf{A}\| \|\mathbf{x}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right\}.$$

As $\|\mathbf{y}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$,

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\chi(\mathbf{A})}{1 - \chi(\mathbf{A}) \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left\{ \frac{\|\delta\mathbf{y}\|}{\|\mathbf{y}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right\}. \quad \blacksquare$$

In practice the quantity $1 - \chi(\mathbf{A})\|\delta\mathbf{A}\|/\|\mathbf{A}\|$ will be of magnitude 1. The only time when this can become small is when the bound on $\|\delta\mathbf{A}\|$ is close to being violated. Assume

$$\frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \ll \frac{1}{\chi(\mathbf{A})},$$

that is, the relative error in \mathbf{A} is much smaller than the inverse of the condition number, then the quantity $\chi(\mathbf{A})$ gives the magnification factor between the error in the data and the error in the solution. As

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I} \Rightarrow \|\mathbf{A}\| \|\mathbf{A}^{-1}\| = \chi(\mathbf{A}) \geq \|\mathbf{I}\| = 1,$$

the relative error bound on the solution is always greater than the relative error bounds on the data.

A common mis-understanding is that various other quantities associated with a matrix give some idea of error magnification. The two most common mis-understandings involve the determinant of \mathbf{A} (see question 4.9.17) and the ratio of largest to smallest eigenvalue of \mathbf{A} . These do not determine condition numbers except for special types of matrices. For example the bidiagonal matrix, $A_{ii} = 1$, $A_{i,i+1} = -1$, $i = 1, 2, \dots, n$, and all other $A_{ij} = 0$, has determinant 1, eigenvalue ratio 1 but condition number $\chi(\mathbf{A}) = O(n^2)$. The above lemmas and theorem are not correct for the Euclidean norm of a matrix as $\|\mathbf{I}\|_E = \sqrt{n}$. Subroutines may however calculate $\chi(\mathbf{A})$ in terms of the Euclidean norm as the 2-norm is too costly to evaluate.

8.1.1 Scaling.

It is possible to effectively reduce the condition number of a matrix by diagonal scaling. Suppose the set of equations $\mathbf{A}\mathbf{x} = \mathbf{y}$, is replaced by

$$\underbrace{\mathbf{D}_1 \mathbf{A} \mathbf{D}_2}_{\overline{\mathbf{A}}} \underbrace{\mathbf{D}_2^{-1} \mathbf{x}}_{\overline{\mathbf{x}}} = \underbrace{\mathbf{D}_1 \mathbf{y}}_{\overline{\mathbf{y}}}$$

$$\overline{\mathbf{A}} \overline{\mathbf{x}} = \overline{\mathbf{y}},$$

where \mathbf{D}_1 and \mathbf{D}_2 are diagonal matrices with elements which are powers of 2. (A computer is capable of multiplying by powers of 2 with no further rounding error). \mathbf{D}_1 and \mathbf{D}_2 are chosen to make $\chi(\overline{\mathbf{A}}) < \chi(\mathbf{A})$. It is not obvious what the optimal choice of \mathbf{D}_1 and \mathbf{D}_2 are, however a good strategy appears to be to make all elements of $\overline{\mathbf{A}}$ as close as possible to the same order. For example $\begin{pmatrix} 2.345 & 1234 \\ 1234 & 5.678 \times 10^6 \end{pmatrix}$ should be replaced with $\begin{pmatrix} 2.345 & 1.234 \\ 1.234 & 5.678 \end{pmatrix}$ where $\mathbf{D}_1 = \mathbf{D}_2 = \text{diag}(1, 10^{-3})$. Perhaps a better way of expressing the requirement is that the absolute error in matrix elements should be uniform over the elements. Note that $\mathbf{D}_1 = \mathbf{D}_2$ if \mathbf{A} is symmetric, so as to make $\overline{\mathbf{A}}$ symmetric.

8.1.2 Backward error analysis.

In order to find the error bound due to the rounding error of a computer the computed solution $\mathbf{x} + \delta\mathbf{x}$ is related to contrived errors $\delta\mathbf{A}$, $\delta\mathbf{y}$, that is, the computed solution should be the exact solution of

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = (\mathbf{y} + \delta\mathbf{y}).$$

This procedure of finding perturbations in the data which match the actual computed solution is known as backward error analysis. (The forward error analysis is somewhat difficult.) Perturbation bounds for the L-U, Choleski and Q-U factorizations are known. Results are quoted later when discussing the solution of equations using the factorizations. Remember that usually error in the original data \mathbf{A} and \mathbf{y} will be larger than the roundoff error backward perturbations.

8.2 Direct methods for the solution of linear equations.

In this section the factorizations discussed in §7 are summarized and used to solve the set of linear equations

$$\mathbf{Ax} = \mathbf{y} \quad (8.2.0)$$

where \mathbf{A} is a square $n \times n$ matrix. As well, the problem of finding the solution \mathbf{x} which minimizes

$$\|\mathbf{Ax} - \mathbf{y}\|_2,$$

where \mathbf{A} is $m \times n$, $m > n$, $\text{rank}(\mathbf{A}) = n$, is discussed. This is the linear least squares problem.

8.2.1 Inversion costing n^2 operations.

There is of course a reason for using factorizations involving upper and lower triangular matrices. Equations in which the matrix is upper or lower triangular can be solved directly using the process known as forward or back substitution. As well the inversion of equations involving an orthogonal matrix is readily accomplished as the inverse is known with no further computation.

(a) **Forward substitution.** The set of equations

$$\mathbf{Lx} = \mathbf{y}$$

where \mathbf{L} is lower triangular and invertible, has a solution \mathbf{x} where formally

$$\mathbf{x} = \mathbf{L}^{-1}\mathbf{y}.$$

Components of \mathbf{x} can be computed in the order x_1, x_2, \dots, x_n . In detail

$$\begin{pmatrix} L_{11} & & & & \\ L_{21} & L_{22} & & & \\ L_{31} & L_{32} & L_{33} & & \\ \vdots & & & \ddots & \\ L_{n1} & L_{n2} & L_{n3} & \dots & L_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix},$$

and

$$\begin{aligned} x_1 &= y_1/L_{11} \\ x_2 &= (y_2 - L_{21}x_1)/L_{22} \\ x_3 &= (y_3 - L_{31}x_1 - L_{32}x_2)/L_{33} \\ &\vdots \\ x_n &= (y_n - \sum_{i=1}^{n-1} L_{ni}x_i)/L_{nn}. \end{aligned}$$

Note that \mathbf{L}^{-1} does not get computed. If \mathbf{L} is written as $\mathbf{D} + \mathbf{S}$ where \mathbf{D} is the diagonal of \mathbf{L} and \mathbf{S} is the strictly lower triangular part then the above computation can be written formally as

$$\mathbf{x} = \mathbf{D}^{-1}(\mathbf{y} - \mathbf{Sx}).$$

(b) **Back substitution.** The set of equations of the form

$$\mathbf{Ux} = \mathbf{y}$$

has a solution \mathbf{x} , where

$$\mathbf{x} = \mathbf{U}^{-1}\mathbf{y}.$$

Here the components of \mathbf{x} can be computed in the order $x_n, x_{n-1}, \dots, x_2, x_1$, where,

$$\begin{aligned} x_n &= y_n/U_{nn}, \\ x_{n-1} &= (y_{n-1} - U_{n-1,n}x_n)/U_{n-1,n-1}, \\ &\vdots \\ x_1 &= (y_1 - \sum_{i=2}^n U_{1i}x_i)/U_{11}. \end{aligned}$$

Note that U^{-1} does not get computed. If U is written as $D + S$ where D is the diagonal of U and S is the strictly upper triangular part then the above computation can be written formally as

$$x = D^{-1}(y - Sx).$$

Both of these “inversions” cost $\frac{1}{2}n(n + 1)$ multiplications and divisions and this cost compares favourably with the cost of a factorization which is $\mathcal{O}(n^3)$.

(c) Multiplying by orthogonal matrices. The orthogonal matrices in the factorizations are either stored in full or as a finite sequence of Householder transformations or plane rotations. Hence it is easy to “invert” such systems. Note if $Q = H_1 H_2 \cdots H_p$, then $Q^{-1} = Q^t = H_p \cdots H_2 H_1$.

8.2.2 Solution of general equations.

The detail of solving (8.2.0) has already been given in §7 and §8.2.1. As such it only remains to describe the major steps and results of the solution procedure. The factorization step is completely independent of any right hand side y and hence can be done independently of y , as was done in §7. Methods used in elementary texts, for example, Gauss elimination, are usually described using augmented matrices with row operations done on both the matrix and right hand side at the same time. Factorization descriptions simply give a picture of this process with the independent parts dissected. Included in the descriptions below are the cost (number of multiplications and divisions), the error analysis results, convenient approximations to $\chi(A)$, and the singularity test. Only the results for real matrices are given as the complex case is similar. In the error analysis, ϵ represents the machine precision and $f(n)$ represents polynomials in n , usually of order 3 (cubic) or less, corresponding to the number of arithmetic operations.

(a) Using the SVD.

$$Ax = y, \quad VDU^t x = y \quad \Rightarrow \quad x = UD^{-1}V^t y.$$

(i) Algorithm — compute in turn:

- 1. $A \rightarrow VDU^t$ $\mathcal{O}(n^3)$
- 2. $V^t y$ n^2
- 3. $D^{-1}(V^t y)$ n
- 4. $U(D^{-1}V^t y)$ n^2

(ii) Storage requirements for A, V, D, U , are $3n^2 + n$, which can be reduced to $2n^2$ by not storing V , and computing $V^t y$ as part of the algorithm.

(iii) Error analysis: If $\bar{V}, \bar{D}, \bar{U}$ are the computed V, D, U then

$$\|A - \bar{V}\bar{D}\bar{U}^t\| \leq \epsilon f(n)\|A\|,$$

even though \bar{V}, \bar{D} and \bar{U} may not be close to V, D and U respectively, even assuming the same ordering and signs of columns of V and U .

(iv) $\chi(A) = d_1/d_n$ ($= \|A\|_2 \|A^{-1}\|_2$).

(v) The smallest singular value close to zero is the test for close to singularity. This test would usually be done as a relative test, against the size of the largest singular value.

(b) Using the Q–U factorization.

$$Ax = y, \quad QUx = y \quad \Rightarrow \quad x = U^{-1}Q^t y.$$

(i) Algorithm — compute in turn:

- 1. $A \rightarrow QU$ $\frac{2}{3}n^3$
- 2. $Q^t y$ n^2
- 3. $U^{-1}(Q^t y)$ $\frac{1}{2}n^2$

(ii) Storage — requires one extra n vector to store Q and U over the matrix A .

(iii) Error analysis

$$\begin{aligned} \|\delta A\| &\leq \epsilon f_1(n)\|A\| \\ \|\delta y\| &\leq \epsilon f_2(n)\|y\| \end{aligned}$$

- (iv) If column pivoting is used then $|U_{11}/U_{nn}| \leq \chi(\mathbf{A})$ and is a good estimate, otherwise $\chi(\mathbf{A}) = \chi(\mathbf{U}) = \|\mathbf{U}\| \|\mathbf{U}^{-1}\|$ which costs $\frac{1}{6}n^3$ to compute. $\chi(\mathbf{A})$ can be estimated by solving with trial right hand sides to give maximum amplification as in LINPACK and LAPACK, that is, attempt to find that \mathbf{x} which maximizes $\|\mathbf{Ax}\|/\|\mathbf{x}\|$. This procedure works well if \mathbf{A} is ill-conditioned, that is, in the case that matters.
- (v) Any U_{ii} close to zero is the test for close to singularity. This is more accurate if column pivoting is used. The relative test against the largest U_{ii} is usually used.

(c) Using the L–U factorization with pivoting.

$$\mathbf{Ax} = \mathbf{y}, \quad (\mathbf{PL})\mathbf{Ux} = \mathbf{y} \quad \Rightarrow \quad \mathbf{x} = \mathbf{U}^{-1}(\mathbf{PL})^{-1}\mathbf{y}.$$

- (i) Algorithm — compute in turn:
1. $\mathbf{A} \rightarrow (\mathbf{PL})\mathbf{U}$ $\frac{1}{3}n^3$
 2. $(\mathbf{PL})^{-1}\mathbf{y}$ $\frac{1}{2}n^2$
 3. $\mathbf{U}^{-1}((\mathbf{PL})^{-1}\mathbf{y})$ $\frac{1}{2}n^2$
- (ii) No extra storage is needed, other than an integer array of length n to store the pivot information. \mathbf{U} is stored over the upper triangle of \mathbf{A} , while \mathbf{L} is stored in the strictly lower triangle of \mathbf{A} , the diagonal of 1's needing no storage.
- (iii) Error analysis:

$$\|\delta\mathbf{A}\| \leq \epsilon f(n)(\|\mathbf{A}\| + \|\mathbf{U}\|).$$

Unfortunately $\|\mathbf{U}\|$ can be much larger than $\|\mathbf{A}\|$,

$$\begin{aligned} \mathbf{U} &= (\mathbf{PL})^{-1}\mathbf{A}, \\ \therefore \|\mathbf{U}\| &\leq \|(\mathbf{PL})^{-1}\| \|\mathbf{A}\|, \end{aligned}$$

but $\|\mathbf{L}^{-1}\| \leq 2^{n-1}$ if pivoting is used (prove it using $\|\cdot\|_\infty$), and this upper bound can be reached.

- (iv) $\chi(\mathbf{A})$ could be indicated by $\max_i |U_{ii}| / \min_i |U_{ii}|$, but a better estimate is

$$\chi(\mathbf{A}) \leq \chi(\mathbf{L})\chi(\mathbf{U}) = \|\mathbf{L}\| \|\mathbf{L}^{-1}\| \|\mathbf{U}\| \|\mathbf{U}^{-1}\|$$

at a cost of $\frac{1}{3}n^3$ operations for the inverses. $\chi(\mathbf{A})$ can be estimated by solving with trial right hand sides to give maximum amplification as in LINPACK and LAPACK, that is, attempt to find that \mathbf{x} which maximizes $\|\mathbf{Ax}\|/\|\mathbf{x}\|$. This procedure works well if \mathbf{A} is ill-conditioned, that is, in the case that matters.

- (v) Any U_{ii} close to zero, using row pivoting, is the test for closeness to singularity. A relative test against the norm of the matrix is usually required.

(d) Using the Choleski factorization.

$$\mathbf{Ax} = \mathbf{y}, \quad \mathbf{A} \text{ pos.def.} \quad \mathbf{LL}^t\mathbf{x} = \mathbf{y} \quad \Rightarrow \quad \mathbf{x} = \mathbf{L}^{-t}\mathbf{L}^{-1}\mathbf{y}.$$

- (i) Algorithm — compute in turn:
1. $\mathbf{A} \rightarrow \mathbf{LL}^t$ $\frac{1}{6}n^3$
 2. $\mathbf{L}^{-1}\mathbf{y}$ $\frac{1}{2}n^2$
 3. $\mathbf{L}^{-t}(\mathbf{L}^{-1}\mathbf{y})$ $\frac{1}{2}n^2$
- (ii) No extra storage is needed, as \mathbf{L} is stored over \mathbf{A} .
- (iii) Error analysis: very stable.

$$\begin{aligned} \|\delta\mathbf{A}\| &\leq \epsilon \|\mathbf{A}\| \\ \|\delta\mathbf{y}\| &\leq \epsilon f(n) \|\mathbf{y}\|. \end{aligned}$$

- (iv) $\chi(\mathbf{A})$ could be indicated by $(\max_i L_{ii} / \min_i L_{ii})^2$, or using $\frac{1}{6}n^3$ operations

$$\chi(\mathbf{A}) \leq (\|\mathbf{L}\| \|\mathbf{L}^{-1}\|)^2.$$

- (v) Any L_{ii} close to zero, compared to the norm of the matrix, is the test for closeness to singularity.

8.3. Least squares methods.

See Lawson and Hanson (1974) for a more complete discussion and Jennings (1980) for applications in econometrics. Björk (1968) has an alternative stable algorithm for least squares. Note that Gram Schmidt orthogonalization is numerically unstable and should never be used unless orthogonalizing two or three vectors.

Write an overdetermined set of linear equations as

$$\mathbf{Ax} = \mathbf{y} - \mathbf{r},$$

where \mathbf{A} is $m \times n$, $m > n$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y}, \mathbf{r} \in \mathbb{R}^m$. The residual vector \mathbf{r} is defined for every \mathbf{x} as $\mathbf{r} = \mathbf{y} - \mathbf{Ax}$. The least squares solution is that \mathbf{x} which minimizes $\mathbf{r}^t \mathbf{r} = \|\mathbf{r}\|_2^2$. Let the function $f(\mathbf{x}) = \mathbf{r}^t \mathbf{r}$.

$$\begin{aligned} f(\mathbf{x}) &= (\mathbf{Ax} - \mathbf{y})^t (\mathbf{Ax} - \mathbf{y}) \\ &= \mathbf{x}^t \mathbf{A}^t \mathbf{Ax} - 2\mathbf{x}^t \mathbf{A}^t \mathbf{y} + \mathbf{y}^t \mathbf{y} \end{aligned}$$

At a turning point $\nabla f = \mathbf{0}$, so that $2\mathbf{A}^t \mathbf{Ax} - 2\mathbf{A}^t \mathbf{y} = \mathbf{0}$, that is,

$$(\mathbf{A}^t \mathbf{A})\mathbf{x} = \mathbf{A}^t \mathbf{y}.$$

These are known as the normal equations. The condition for a unique minimum is that $\mathbf{A}^t \mathbf{A}$ be positive definite which is true if \mathbf{A} has rank n .

8.3.1. Choleski algorithm for least squares.

- (i) Compute $\mathbf{A}^t \mathbf{A}$ $\frac{1}{2}n^2m$
- (ii) Factor $\mathbf{A}^t \mathbf{A} \rightarrow \mathbf{LL}^t$ $\frac{1}{6}n^3$
- (iii) Compute $\mathbf{L}^{-t}(\mathbf{L}^{-1}(\mathbf{A}^t \mathbf{y}))$ $nm + n^2$

The error analysis for the Choleski algorithm shows that

$$\|\delta \mathbf{x}\| \leq \epsilon f(n) \chi^2(\mathbf{A}) \|\mathbf{y}\| \quad \text{as} \quad \chi(\mathbf{A}^t \mathbf{A}) = \chi^2(\mathbf{A}).$$

Most of the accumulated error for this algorithm is in the computation of $\mathbf{A}^t \mathbf{A}$.

8.3.2. Golub algorithm for least squares.

See Golub (1966). If \mathbf{A} is factored using the Q-U factorization, then the normal equations become

$$\mathbf{U}^t \mathbf{U} \mathbf{x} = (\mathbf{U}^t \quad \mathbf{0}) \mathbf{Q}^t \mathbf{y},$$

and writing $\mathbf{Q}^t \mathbf{y}$ as $\begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix}$, where $\mathbf{y}_1 \in \mathbb{R}^n$,

$$\mathbf{U}^t \mathbf{U} \mathbf{x} = \mathbf{U}^t \mathbf{y}_1 \quad \Rightarrow \quad \mathbf{x} = \mathbf{U}^{-1} \mathbf{y}_1.$$

Algorithm:

- (i) Factorize $\mathbf{A} \rightarrow \mathbf{Q} \begin{pmatrix} \mathbf{U} \\ \mathbf{0} \end{pmatrix}$ $\frac{2}{3}n^2m$
- (ii) Compute $\mathbf{Q}^t \mathbf{y}$ $2nm - \frac{1}{2}n^2$
- (iii) Compute $\mathbf{U}^{-1} \mathbf{y}_1$ $\frac{1}{2}n^2$

Note: $\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|_2 \equiv \min_{\mathbf{x}} \left\| \begin{pmatrix} \mathbf{U} \\ \mathbf{0} \end{pmatrix} \mathbf{x} - \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} \right\|_2 = \|\mathbf{y}_2\|_2$ if $\mathbf{U} \mathbf{x} = \mathbf{y}_1$. The residual vector is

$$\mathbf{r} = \mathbf{Q} \begin{pmatrix} \mathbf{0} \\ \mathbf{y}_2 \end{pmatrix}.$$

The error analysis for the Golub algorithm shows that

$$\|\delta \mathbf{x}\| \leq \epsilon f_1(n) \chi(\mathbf{A}) \|\mathbf{y}\| + \epsilon f_2(n) \chi^2(\mathbf{A}) \|\mathbf{r}\|.$$

See Jennings and Osborne (1974) for details. The Golub algorithm is better in terms of roundoff error but involves about twice the computation than the Choleski method.

Estimates of $\chi(\mathbf{A})$ follow from earlier discussion for square matrices for the Choleski method, while for the Golub algorithm the condition number of \mathbf{A} is equal to that of \mathbf{U} .

8.4 Iterative methods.

The methods for solving a set of linear equations to be discussed in this section apply mainly to systems which are large ($n > 500$) but have a sparse coefficient matrix. Usually in this case the cost of finding and storing a factorization with its resultant “fill in” is prohibitive. Further details of iterative methods can be found in Young (1971) and Varga (1962), while sparse matrix research was reviewed in Duff (1977). Blackford *et al.* (1997) discusses the LAPACK library modifications to parallel processing for all types of matrices. Dongarra *et al.* (1990) discusses the impact of parallelization on matrix computations.

8.4.1 Successive iteration (linear)

Suppose \mathbf{A} and \mathbf{y} are given and a guess at the solution \mathbf{x} , say \mathbf{x}_1 is known. Then if

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{y},$$

the iteration

$$\mathbf{x}_{i+1} = \mathbf{G}_i\mathbf{x}_i + \mathbf{r}_i, \quad i = 1, 2, \dots,$$

is proposed where \mathbf{G}_i is a matrix and \mathbf{r}_i a vector. Convergence requires that

$$\lim_{i \rightarrow \infty} \|\mathbf{x}_i - \mathbf{x}\| = 0.$$

Consistency: The solution must satisfy the iteration

$$\mathbf{A}^{-1}\mathbf{y} = \mathbf{G}_i\mathbf{A}^{-1}\mathbf{y} + \mathbf{r}_i,$$

so that \mathbf{G}_i and \mathbf{r}_i are related.

$$\mathbf{r}_i = (\mathbf{I} - \mathbf{G}_i)\mathbf{A}^{-1}\mathbf{y}.$$

Convergence: Let $\mathbf{z}_i = \mathbf{x}_i - \mathbf{x} = \mathbf{x}_i - \mathbf{A}^{-1}\mathbf{y}$.

$$\begin{aligned} \mathbf{z}_{i+1} &= \mathbf{x}_{i+1} - \mathbf{x} \\ &= \mathbf{G}_i\mathbf{x}_i + (\mathbf{I} - \mathbf{G}_i)\mathbf{x} - \mathbf{x} \\ &= \mathbf{G}_i(\mathbf{x}_i - \mathbf{x}) \\ &= \mathbf{G}_i\mathbf{z}_i \\ &= \mathbf{G}_i\mathbf{G}_{i-1} \cdots \mathbf{G}_1\mathbf{z}_1 \\ &= \mathbf{H}_i\mathbf{z}_1 \quad \text{say} \end{aligned}$$

Now if $\mathbf{H}_i \rightarrow 0$ as $i \rightarrow \infty$ then $\mathbf{x}_i \rightarrow \mathbf{x}$. If $\mathbf{G}_i = \mathbf{G}$, (constant) then $\mathbf{H}_i = \mathbf{G}^i$ and convergence is assured if any of; $\|\mathbf{G}\| < 1$, for any matrix norm, or, $\max_j |\lambda_j| < 1$, where λ_j is an eigenvalue of \mathbf{G} .

Jacobi iteration.

Write the matrix \mathbf{A} as

$$\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U},$$

where \mathbf{L} is strictly lower triangular, \mathbf{D} the diagonal elements of \mathbf{A} , and \mathbf{U} is strictly upper triangular. The Jacobi method is

$$\begin{aligned} \mathbf{D}\mathbf{x}_{i+1} &= \mathbf{y} - (\mathbf{L} + \mathbf{U})\mathbf{x}_i \\ \therefore \mathbf{G} &= -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}), \quad \mathbf{r} = \mathbf{D}^{-1}\mathbf{y}. \end{aligned}$$

It is easily shown that if \mathbf{A} is diagonally dominant, that is

$$\sum_{\substack{i \neq j \\ i=1}}^n |A_{ij}| \leq |A_{jj}| \quad \text{for all } j \quad \text{or} \quad \sum_{\substack{i \neq j \\ j=1}}^n |A_{ij}| \leq |A_{ii}| \quad \text{for all } i$$

then this method is convergent. (Prove it!!)

Gauss-Siedel iteration.

Here, essentially a forward substitution is computed.

$$\begin{aligned}(\mathbf{L} + \mathbf{D})\mathbf{x}_{i+1} &= \mathbf{y} - \mathbf{U}\mathbf{x}_i \\ \mathbf{G} &= -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}, \quad \mathbf{r} = (\mathbf{L} + \mathbf{D})^{-1}\mathbf{y} \\ \mathbf{x}_{i+1} &= \mathbf{D}^{-1}(\mathbf{y} - \mathbf{U}\mathbf{x}_i - \mathbf{L}\mathbf{x}_{i+1})\end{aligned}$$

This can be thought of as “if better values for \mathbf{x} are available in \mathbf{x}_{i+1} , why not use them”. This has better convergence properties than Jacobi in that it is faster. Jacobi and Gauss Siedel iterations converge or diverge together, depending on the properties of \mathbf{A} . There are variants of this algorithm, for example, doing a back substitution instead, which corresponds to computing the new vector in reverse order. Alternating directions methods combine forward and backward methods alternately for faster convergence, as well as re-orderings of equations and variables in some problems.

Relaxation of Gauss-Siedel.

Consider $\mathbf{x}_{i+1} - \mathbf{x}_i$ as a direction of search. Then define a next iterate as

$$\mathbf{x}_i + w(\mathbf{x}_{i+1} - \mathbf{x}_i).$$

If $w > 1$ this is called over-relaxation while if $w < 1$, under relaxation. Formally

$$\mathbf{x}_{i+1} = w\mathbf{D}^{-1}(\mathbf{y} - \mathbf{U}\mathbf{x}_i - \mathbf{L}\mathbf{x}_{i+1}) + (1 - w)\mathbf{x}_i$$

but this is usually best explained in terms of elements of vectors. For $j = 1, 2, \dots, n$

$$\mathbf{x}_{i+1}(j) = wD_{jj}^{-1}\{y_j - \sum_{k>j} U_{jk}\mathbf{x}_i(k) - \sum_{k<j} L_{jk}\mathbf{x}_{i+1}(k)\} + (1 - w)\mathbf{x}_i(j).$$

Some problems, especially in the numerical solution of partial differential equations are very efficiently solved using variants of relaxation and alternating directions.

8.4.2 The conjugate gradient method.

See Hestines and Steifel (1952) for the original detail. Let the matrix \mathbf{A} be positive definite. Then the quadratic form

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^t \mathbf{A}\mathbf{x} - \mathbf{y}^t \mathbf{x},$$

has a minimum when

$$(\nabla f)^t = \mathbf{A}\mathbf{x} - \mathbf{y} = \mathbf{0}.$$

The conjugate gradient method is used to find \mathbf{x} which gives the minimum of $f(\mathbf{x})$ by creating search directions from some starting point \mathbf{x}_0 say. At the i -th step form a search direction \mathbf{p}_i and move along it a distance α_i

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{p}_i,$$

where α_i is chosen so that \mathbf{x}_{i+1} is a minimum of $f(\mathbf{x})$ in this direction. In particular the gradient of $f(\mathbf{x})$ at \mathbf{x}_{i+1} will be perpendicular to \mathbf{p}_i ,

$$\nabla f_{i+1} \mathbf{p}_i = 0, \quad \text{where} \quad \nabla f_{i+1} = \nabla f(\mathbf{x}_{i+1}).$$

This gives

$$\alpha_i = \frac{\nabla f_i \mathbf{p}_i}{\mathbf{p}_i^t \mathbf{A} \mathbf{p}_i}.$$

Now

$$\mathbf{x}_n = \mathbf{x}_i + \sum_{j=i}^{n-1} \alpha_j \mathbf{p}_j$$

and

$$\mathbf{A}\mathbf{x}_n - \mathbf{y} = \mathbf{A}\mathbf{x}_i - \mathbf{y} + \sum_{j=i}^{n-1} \alpha_j \mathbf{A}\mathbf{p}_j$$

or,

$$\nabla f_n^t = \nabla f_i^t + \sum_{j=i}^{n-1} \alpha_j \mathbf{A} \mathbf{p}_j.$$

In particular,

$$\mathbf{p}_{i-1}^t \nabla f_n^t = \mathbf{p}_{i-1}^t \nabla f_i^t + \sum_{j=i}^{n-1} \alpha_j \mathbf{p}_{i-1}^t \mathbf{A} \mathbf{p}_j, \tag{8.4.1}$$

and as $\mathbf{p}_{i-1}^t \nabla f_i^t = 0$, then $\mathbf{p}_{i-1}^t \nabla f_n^t = 0$, if $\mathbf{p}_{i-1}^t \mathbf{A} \mathbf{p}_j = 0$, $j = i, i + 1, \dots, n - 1$. Hence \mathbf{p}_j is chosen such that

$$\mathbf{p}_i^t \mathbf{A} \mathbf{p}_j = 0 \quad i \neq j.$$

This is known as *A-conjugacy*, and the set of vectors $\{\mathbf{p}_j\}_{j=0}^{n-1}$ is said to be *A-conjugate*. It is easy to show the set of non-zero vectors $\{\mathbf{p}_j\}_{j=0}^{n-1}$ form a basis for \mathbb{R}^n and hence from (8.4.1)

$$\nabla f_n = \mathbf{0}$$

that is, \mathbf{x}_n is the point giving a minimum of $f(\mathbf{x})$. To find the vectors \mathbf{p}_j , suppose $\mathbf{p}_0 = -\nabla f_0$, and let

$$\mathbf{p}_{i+1} = -\nabla f_{i+1}^t + \beta_i \mathbf{p}_i, \quad i = 0, 1, 2, \dots, n - 1.$$

Now as

$$\begin{aligned} \mathbf{p}_i^t \mathbf{A} \mathbf{p}_{i+1} &= 0 \\ -\mathbf{p}_i^t \mathbf{A} \nabla f_{i+1}^t + \beta_i \mathbf{p}_i^t \mathbf{A} \mathbf{p}_i &= 0 \\ \Rightarrow \beta_i &= \frac{-\mathbf{p}_i^t \mathbf{A} \nabla f_{i+1}^t}{\mathbf{p}_i^t \mathbf{A} \mathbf{p}_i} \\ &= \frac{-\mathbf{p}_i^t \mathbf{A} (\mathbf{A} \mathbf{x}_{i+1} - \mathbf{y})}{\mathbf{p}_i^t \mathbf{A} \mathbf{p}_i} \end{aligned}$$

8.5 Exercises.

- (i) Solve these simultaneous equations exactly. This is $\mathbf{A} \mathbf{x} = \mathbf{y}$.

$$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 2 & 1 \\ 2 & 1 & 0 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 3 \end{pmatrix}.$$

- (ii) Solve these simultaneous equations exactly. This is $(\mathbf{A} + \delta \mathbf{A})(\mathbf{x} + \delta \mathbf{x}) = \mathbf{y}$.

$$\begin{pmatrix} \epsilon & 1 & 1 \\ 1 & 2 & 1 \\ 2 & 1 & 0 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 3 \end{pmatrix}.$$

Is there much change in solution for small ϵ ? Expand your answer as a Taylor series to first order.

- (iii) Solve with no pivoting using three significant figure arithmetic (chop not round) the system in (ii) with $\epsilon = 10^{-4}$. Follow a strict Gaussian elimination algorithm where the row operations are $row_j \leftarrow row_j - m_j row_i$, $j > i$, that is, subtract a multiple of the pivot row from row_j . Does the solution change very much from the answer of (ii)?
- (iv) Is the system ill-conditioned compared to the small change to the (1,1) element? Why doesn't your solution to (iii) satisfy Theorem 8.1.1?
- (i) Indicate how you would solve the matrix system (a rank one correction or update)

$$(\mathbf{A} + \mathbf{u} \mathbf{v}^t) \mathbf{x} = \mathbf{y},$$

given the factors of $\mathbf{A} = (\mathbf{P} \mathbf{L}) \mathbf{U}$, in order n^2 operations. You know the inverse of $\mathbf{I} + \mathbf{w} \mathbf{v}^t$ from exercise 6.7.1. Illustrate the algorithm by solving

$$\left(\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} + \begin{pmatrix} 0 & 10^{-4} \\ 0 & 0 \end{pmatrix} \right) \mathbf{x} = \begin{pmatrix} 2 \\ 3 \end{pmatrix}.$$

- (ii) Generalize to $(\mathbf{A} + \mathbf{ZV}^t)\mathbf{x} = \mathbf{y}$ where \mathbf{Z} and \mathbf{V} are $n \times r$ matrices ($r < n$). This is the Sherman–Morrison–Woodbury formula for r rank one updates simultaneously. Hint: show that the inverse of $\mathbf{I} + \mathbf{WV}^t$ is $\mathbf{I} - \mathbf{WRV}^t$, where $\mathbf{R} = (\mathbf{I} + \mathbf{V}^t\mathbf{W})^{-1}$, an $r \times r$ matrix.
- 3. Suggest an efficient method for solving a matrix equation of order 1000 which is pentadiagonal except that $A_{1,1000}$ and $A_{1000,1}$ are non-zero. Show that this can be done in about 19–25 thousand operations depending on whether row pivoting needs to be used.
- 4. Consider the equation

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & 0.01 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 1 \\ 0.1 \\ 0.01 \end{pmatrix}.$$

- Show that the relative error bounds for $\|\delta\mathbf{x}\|/\|\mathbf{x}\|$ in (8.1.3) can be reached by choosing appropriate $\delta\mathbf{A}$ with $\delta\mathbf{y} = 0$ and choosing appropriate $\delta\mathbf{y}$ with $\delta\mathbf{A} = 0$, that is, first simplify the theorem by putting in turn one of $\delta\mathbf{y}$ and $\delta\mathbf{A}$ to zero. What is the biggest perturbation $\delta\mathbf{A}$ so that all smaller perturbations ensure $(\mathbf{A} + \delta\mathbf{A})^{-1}$ exists? Hint: use the infinity norm and choose $\delta\mathbf{A} = -\alpha\mathbf{e}_3\mathbf{e}_3^t$ for some suitable range of α values. This shows that Theorem 8.1.1 cannot be improved.
- 5. Let \mathbf{A} be an $m \times n$ matrix $m > n$ and let

$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{U} \\ \mathbf{0} \end{pmatrix}.$$

Consider the set of least squares equations, $\mathbf{Ax} = \mathbf{y} - \mathbf{r}$. If \mathbf{x} is chosen to minimize $\|\mathbf{r}\|_2$ show that

$$\begin{aligned} \mathbf{r} &= (\mathbf{I} - \mathbf{A}(\mathbf{A}^t\mathbf{A})^{-1}\mathbf{A}^t)\mathbf{y} \\ &= \mathbf{Q} \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-n} \end{pmatrix} \mathbf{Q}^t\mathbf{y}. \end{aligned}$$

- 6. The linear system of equations

$$\begin{pmatrix} 1 & -a \\ -a & 1 \end{pmatrix} \mathbf{x} = \mathbf{y}$$

where a is real, can, under certain conditions be solved by the iterative method

$$\begin{pmatrix} 1 & 0 \\ -wa & 1 \end{pmatrix} \mathbf{x}_{i+1} = \begin{pmatrix} 1-w & wa \\ 0 & 1-w \end{pmatrix} \mathbf{x}_i + w\mathbf{y}.$$

- for which values of a is the method convergent for $w = 1$? Use Maple to look for the best value of w in terms of a , by considering the matrix \mathbf{G} defined in the notes. Can you generalise this to the corresponding $n \times n$ symmetric tridiagonal matrix?
- 7. Compute the first few iterations of the Jacobi and Gauss-Siedel iterations on the system,

$$\begin{pmatrix} 4 & 1 & 0 & 0 & 0 \\ 1 & 4 & 1 & 0 & 0 \\ 0 & 1 & 4 & 1 & 0 \\ 0 & 0 & 1 & 4 & 1 \\ 0 & 0 & 0 & 1 & 4 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 3 \\ -2 \\ 2 \\ -2 \\ 3 \end{pmatrix},$$

- starting at $\mathbf{x} = \mathbf{0}$. (Solution is $[1, -1, 1, -1, 1]^t$.) Compute the two over-relaxation iterations using $w = 1.5$. (If you write a simple 10 minute program to do this, compare the convergence rates.)
- 8. Write an algorithm for the conjugate gradient method. Because of the poor numerical performance of this algorithm it is usual to restart it every n iterations.
 - 9a. Let \mathbf{A} be an $n \times n$ matrix which is invertible.

- (i) Show the inverse of the partitioned matrix $\begin{pmatrix} \mathbf{A} & \mathbf{a} \\ \mathbf{b}^t & z \end{pmatrix}$ is

$$\begin{pmatrix} \mathbf{A}^{-1} + w\mathbf{A}^{-1}\mathbf{a}\mathbf{b}^t\mathbf{A}^{-1} & -w\mathbf{A}^{-1}\mathbf{a} \\ -w\mathbf{b}^t\mathbf{A}^{-1} & w \end{pmatrix}$$

where $w = (z - \mathbf{b}^t\mathbf{A}^{-1}\mathbf{a})^{-1}$.

- (ii) What is the inverse (if it exists) if \mathbf{A}^{-1} does not exist and $z \neq 0$.

- (iii) Does the inverse exist if $z = 0$ and \mathbf{A}^{-1} does not exist?
 (iv) In (i) if \mathbf{A} has been factored to $(\mathbf{P}\mathbf{L})\mathbf{U}$ how could you solve

$$\begin{pmatrix} \mathbf{A} & \mathbf{a} \\ \mathbf{b}^t & z \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ u \end{pmatrix} = \begin{pmatrix} \mathbf{y} \\ v \end{pmatrix}$$

in $\mathcal{O}(n^2)$ operations?

- 9b. Continuing this question let \mathbf{A} be a symmetric positive definite matrix and $\mathbf{b} = \mathbf{a}$.
 (i) What are the conditions on \mathbf{a} and z to keep the $(n+1) \times (n+1)$ matrix positive definite?
 (ii) If \mathbf{A} has been factored to $\mathbf{L}\mathbf{L}^t$ how could you solve the order $n+1$ system in $\mathcal{O}(n^2)$ operations?
 10. Find an expression for the inverse of the (not symmetric) partitioned matrix

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}$$

where \mathbf{A} and \mathbf{D} are square and not necessarily of the same order. You should only have two inverses in your expression. Note the conditions for existence of the computations. Compare the operation cost for solving the $2n \times 2n$ system

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix}$$

using L–U factorization, with solving the system in partitioned form using L–U factorization on $n \times n$ systems. Note the common factors in the expressions for \mathbf{x}_1 and \mathbf{x}_2 , and that $\mathbf{A}\mathbf{B}\mathbf{x}$ is more efficiently computed as $\mathbf{A}(\mathbf{B}\mathbf{x})$ and not $(\mathbf{A}\mathbf{B})\mathbf{x}$.

11. Given the two systems of equations $\mathbf{A}\mathbf{x} = \mathbf{y}$ and $\mathbf{A}\bar{\mathbf{x}} = \mathbf{y} + \mathbf{z}$, show that the minimum value of $\|\mathbf{x} - \bar{\mathbf{x}}\|_2$ is $\|\mathbf{z}\|/d_1$ and the maximum value is $\|\mathbf{z}\|/d_n$. Here d_1 and d_n are the largest and smallest singular values of \mathbf{A} .
 12. Show that if partial pivoting is used in the L–U factorization

$$\|\mathbf{U}\|_\infty \leq 2^{n-1} \|\mathbf{A}\|_\infty.$$

(It is easy to verify that the 1-norm gives a larger bound while it is somewhat harder to show that the 2-norm gives a larger bound.)

13. Show that the eigenvectors of the system

$$(\mathbf{A} - \lambda\mathbf{B})\mathbf{v} = \mathbf{0},$$

where \mathbf{A} and \mathbf{B} are symmetric, are both \mathbf{A} -conjugate and \mathbf{B} -conjugate, if all the eigenvalues are distinct.

14. Consider the system

$$\begin{pmatrix} 1 & 0 & -1/4 & -1/4 \\ 0 & 1 & -1/4 & -1/4 \\ -1/4 & -1/4 & 1 & 0 \\ -1/4 & -1/4 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \end{pmatrix}$$

and use MATLAB to do the following computations.

- (i) Starting with $\mathbf{x} = \mathbf{0}$, iterate the Jacobi method until iterates are within 10^{-10} of each other. Note the number of iterations to achieve this accuracy.
 (ii) Repeat (i) using the Gauss-Seidel method.
 (iii) For both methods find the matrix \mathbf{G} of the notes and find the norm and the spectral radius (largest in magnitude eigenvalue) and hence explain the different convergence behaviour of (i) and (ii).
 (iv) Is there a relationship between the norm of the distance of an iterate from the true solution and the norm of the difference of two successive iterates?

15. Show that it costs $n^3/6$ multiplications and divisions to compute the inverse matrix of an upper triangular matrix. Remember that the inverse is also upper triangular so there is no need to compute the zeros in the lower triangle.
16. What are the operation counts for computing the following most efficiently?
 - (a) The product of two upper triangular matrices.
 - (b) The product of an upper triangular and lower triangular matrix.
 - (c) The product of the inverse of a lower triangular matrix and an upper triangular matrix, given the lower triangular matrix. Is this best done via Question 15 or by Gauss elimination in a tableau method?
 - (d) The product of the inverse of a lower triangular matrix and a lower triangular matrix, given the first lower triangular matrix. Is this best done via Question 15 or by Gauss elimination in a tableau method?

9. The Generalized Inverse Problem.

The concept of a generalized inverse is examined with the view of providing an ‘inverse’ for any linear mapping from \mathbb{R}^n to \mathbb{R}^m . The geometric aspects are emphasized and used to explain the computation of solutions to any system of simultaneous linear equations in both the analytical and practical sense. Further material can be found in Ben-Israel and Greville (1974). Original work in this area goes back to Moore (1920) and Penrose (1956).

9.1 The geometry of a linear mapping.

Let a linear map from \mathbb{R}^n to \mathbb{R}^m be represented by the $m \times n$ matrix \mathbf{A} in the usual orthogonal bases. The abstract distinction between the linear map \mathcal{A} , and the matrix \mathbf{A} is not made, hence the linear map will be called \mathbf{A} . Let $\langle \cdot, \cdot \rangle$ be the usual inner product and $\|\cdot\|$ be the Euclidean norm for the space \mathbb{R}^n or \mathbb{R}^m . Recall,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^t \mathbf{y} = \sum_{i=1}^n x_i y_i \quad (= \mathbf{y}^t \mathbf{x}) = \langle \mathbf{y}, \mathbf{x} \rangle,$$

and

$$\langle \mathbf{x}, \mathbf{x} \rangle^{1/2} = \|\mathbf{x}\|_2 = \|\mathbf{x}\|.$$

This inner product and norm are invariant to multiplication by an orthogonal matrix.

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{y} \rangle$$

$$\|\mathbf{x}\| = \|\mathbf{Q}\mathbf{x}\|$$

From §5 this means that this inner product and norm are invariant to change of orthonormal basis

$$\mathbf{I}\mathbf{I}^t \mathbf{x} = \sum_{i=1}^n \langle \mathbf{e}_i, \mathbf{x} \rangle \mathbf{e}_i = \sum_{i=1}^n \langle \mathbf{q}_i, \mathbf{x} \rangle \mathbf{q}_i = \mathbf{Q}\mathbf{Q}^t \mathbf{x}$$

Note that $(\mathbf{I}\mathbf{x})_i = \langle \mathbf{e}_i, \mathbf{x} \rangle$ is the i -th coordinate of \mathbf{x} with respect to the standard basis and $(\mathbf{Q}^t \mathbf{x})_i = \langle \mathbf{q}_i, \mathbf{x} \rangle$ is the i -th coordinate of \mathbf{x} with respect to the orthogonal basis $[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$. From §4 the adjoint of a real mapping is the transpose matrix as for all \mathbf{x} and \mathbf{y}

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{A}^t \mathbf{y} \rangle.$$

The mapping \mathbf{A} divides \mathbb{R}^n and \mathbb{R}^m into ‘natural’ orthogonal subspaces, namely

$$\begin{aligned} \mathbb{R}^n &= \mathcal{N}(\mathbf{A}) \oplus \mathcal{R}(\mathbf{A}^t), \\ \mathbb{R}^m &= \mathcal{N}(\mathbf{A}^t) \oplus \mathcal{R}(\mathbf{A}). \end{aligned}$$

It is easy to show that

$$\begin{aligned} \mathcal{N}(\mathbf{A})^\perp &= \mathcal{R}(\mathbf{A}^t), \\ \mathcal{N}(\mathbf{A}^t)^\perp &= \mathcal{R}(\mathbf{A}). \end{aligned}$$

Note the dimensions of $\mathcal{R}(\mathbf{A})$ and $\mathcal{R}(\mathbf{A}^t)$ are the same and is the rank of \mathbf{A} denoted r . Hence the dimensions of $\mathcal{N}(\mathbf{A})$ is $n - r$ and that of $\mathcal{N}(\mathbf{A}^t)$ is $m - r$. It is easy to show that \mathbf{A} maps all elements of $\mathcal{R}(\mathbf{A}^t)$ to $\mathcal{R}(\mathbf{A})$, while \mathbf{A}^t maps all elements of $\mathcal{R}(\mathbf{A})$ to $\mathcal{R}(\mathbf{A}^t)$.

9.2 Natural orthogonal bases for a mapping.

9.2.1 The Q–U factorization.

Recall that an $m \times n$ matrix of rank r can be factored to

$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{P}^t$$

where Q is $m \times m$ orthogonal, P is $n \times n$ orthogonal and L is $r \times r$ lower triangular (of rank r). If Q and P are partitioned appropriately then

$$A = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \begin{pmatrix} L & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} P_1^t \\ P_2^t \end{pmatrix} = Q_1 L P_1^t$$

where Q_1 is $m \times r$ and P_1 is $n \times r$. If q_i are columns of Q and p_i are columns of P it is easy to show that

$$\begin{aligned} \mathcal{R}(A) &= \text{span}\{q_1, q_2, \dots, q_r\}, \\ \mathcal{N}(A^t) &= \text{span}\{q_{r+1}, q_{r+2}, \dots, q_m\}, \\ \mathcal{R}(A^t) &= \text{span}\{p_1, p_2, \dots, p_r\}, \\ \mathcal{N}(A) &= \text{span}\{p_{r+1}, p_{r+2}, \dots, p_n\}. \end{aligned}$$

So the two orthogonal matrices Q and P of the Q-U factorization give a natural orthonormal basis for the linear mapping A . The matrix of the mapping in the new bases is of course

$$\begin{pmatrix} L & 0 \\ 0 & 0 \end{pmatrix}.$$

since

$$AP = Q \begin{pmatrix} L & 0 \\ 0 & 0 \end{pmatrix}.$$

Q_1 and P_1 and L are unique up to multiplication by diagonal matrices of ± 1 .

9.2.2 The singular value decomposition.

The $m \times n$ matrix A of rank r can be factored to

$$A = VDU^t$$

where V is $m \times m$ orthogonal, U is $n \times n$ orthogonal and D is $m \times n$ diagonal with $D_{ii} = d_i$, $i = 1, 2, \dots, \min(m, n)$ and $d_1 \geq d_2 \geq \dots \geq d_r > 0$, $d_{r+1} = d_{r+2} = \dots = d_{\min(m, n)} = 0$. If $v_i = k_i(V)$ and $u_i = k_i(U)$ then

$$A = \sum_{i=1}^r d_i v_i u_i^t,$$

that is, A is the sum of r rank one matrices which in a sense are orthogonal to each other. Partitioning according to the zeros of D produces

$$A = \begin{pmatrix} V_1 & V_2 \end{pmatrix} \begin{pmatrix} \hat{D} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} U_1^t \\ U_2^t \end{pmatrix} = V_1 \hat{D} U_1^t,$$

where V_1 is $m \times r$, U_1 is $n \times r$ and \hat{D} is $r \times r$ diagonal. It is easy to show that

$$\begin{aligned} \mathcal{R}(A) &= \text{span}\{v_1, v_2, \dots, v_r\}, \\ \mathcal{N}(A^t) &= \text{span}\{v_{r+1}, \dots, v_m\}, \\ \mathcal{R}(A^t) &= \text{span}\{u_1, u_2, \dots, u_r\}, \\ \mathcal{N}(A) &= \text{span}\{u_{r+1}, \dots, u_n\}. \end{aligned}$$

9.3 Canonical form of a linear map.

By choosing appropriate bases in \mathbb{R}^n and \mathbb{R}^m every linear mapping from \mathbb{R}^n to \mathbb{R}^m can be represented by a diagonal matrix. The basis for \mathbb{R}^n is the columns of U and the basis for \mathbb{R}^m is the columns of V . Hence the $m \times n$ matrix representing the mapping is the diagonal matrix of the SVD,

$$AU = VD,$$

$$A(u_1 \ u_2 \ \dots \ u_n) = (v_1 \ v_2 \ \dots \ v_m) D,$$

$$\text{or } Au_i = \begin{cases} d_i v_i & i = 1, 2, \dots, r, \\ 0 & i = r + 1, \dots, n. \end{cases}$$

Likewise the mapping $\mathbf{A}^t : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is diagonal and represented by \mathbf{D}^t if the above bases are chosen

$$\mathbf{A}^t = \mathbf{U}\mathbf{D}^t\mathbf{V}^t$$

or

$$\mathbf{A}^t\mathbf{V} = \mathbf{U}\mathbf{D}^t$$

or

$$\mathbf{A}^t\mathbf{v}_i = \begin{cases} d_i\mathbf{u}_i & i = 1, 2, \dots, r, \\ \mathbf{0} & i = r + 1, \dots, m. \end{cases}$$

The projection of vectors in the usual bases to the new bases represented by \mathbf{U} and \mathbf{V} are,

$$\mathbf{x} = \mathbf{U}\mathbf{U}^t\mathbf{x} = \sum_{i=1}^n \langle \mathbf{u}_i, \mathbf{x} \rangle \mathbf{u}_i,$$

$$\mathbf{y} = \mathbf{V}\mathbf{V}^t\mathbf{y} = \sum \langle \mathbf{v}_i, \mathbf{y} \rangle \mathbf{v}_i,$$

and using $\mathbf{A} = \sum_{i=1}^r d_i\mathbf{v}_i\mathbf{u}_i^t$,

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \sum_{i=1}^r d_i\mathbf{v}_i\mathbf{u}_i^t \sum_{j=1}^n \langle \mathbf{u}_j, \mathbf{x} \rangle \mathbf{u}_j, \\ &= \sum_{i=1}^r d_i \langle \mathbf{u}_i, \mathbf{x} \rangle \mathbf{v}_i. \end{aligned}$$

Hence if $\mathbf{y} = \mathbf{A}\mathbf{x}$, equating coefficients of \mathbf{v}_i gives

$$\begin{aligned} \langle \mathbf{v}_i, \mathbf{y} \rangle &= d_i \langle \mathbf{u}_i, \mathbf{x} \rangle, \quad i = 1, 2, \dots, r, \\ \langle \mathbf{v}_i, \mathbf{y} \rangle &= 0, \quad i = r + 1, \dots, m. \end{aligned}$$

So the component of $\mathbf{A}\mathbf{x}$ in the direction \mathbf{v}_i is related only to the component of \mathbf{x} in the direction \mathbf{u}_i . Similarly the mapping \mathbf{A}^t ,

$$\begin{aligned} \mathbf{A}^t\mathbf{y} \quad (= \mathbf{x}) &= \sum_{i=1}^r d_i \langle \mathbf{v}_i, \mathbf{y} \rangle \mathbf{u}_i, \\ \langle \mathbf{u}_i, \mathbf{x} \rangle &= d_i \langle \mathbf{v}_i, \mathbf{y} \rangle, \quad i = 1, 2, \dots, r \\ \langle \mathbf{u}_i, \mathbf{x} \rangle &= 0, \quad i = r + 1, \dots, n. \end{aligned}$$

So the component of $\mathbf{A}^t\mathbf{y}$ in the direction \mathbf{u}_i is related only to the component of \mathbf{y} in the direction \mathbf{v}_i .

9.4 Generalized Inverse Mapping

9.4.1 An inverse map by restricting the domain and range.

Suppose the mapping \mathbf{A} is restricted to $\mathcal{R}(\mathbf{A}^t)$ as its domain and $\mathcal{R}(\mathbf{A})$ as its range. As these both have the same dimension the mapping is now 1 : 1 and onto, that is,

$$\mathbf{A}\mathbf{u}_i \rightarrow d_i\mathbf{v}_i, \quad i = 1, 2, \dots, r,$$

and $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ is a basis of $\mathcal{R}(\mathbf{A}^t)$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ is a basis for $\mathcal{R}(\mathbf{A})$. As such an inverse map exists, say \mathbf{A}^r , which is described by

$$\mathbf{A}^r\mathbf{v}_i \rightarrow \frac{1}{d_i}\mathbf{u}_i, \quad i = 1, 2, \dots, r.$$

It is easy to check that $\mathbf{A}^r\mathbf{A} = \mathbf{A}\mathbf{A}^r = \mathbf{I}_r$ for $\mathbf{A} : \mathcal{R}(\mathbf{A}^t) \rightarrow \mathcal{R}(\mathbf{A})$.

As the inverse map is linear it can be represented by an $r \times r$ matrix in bases \mathbf{U}_1 and \mathbf{V}_1 as $\text{diag} [d_1^{-1}, d_2^{-1}, \dots, d_r^{-1}]$,

$$\mathbf{A}^r\mathbf{V}_1 = \mathbf{U}_1 \begin{pmatrix} d_1^{-1} & & & \\ & d_2^{-1} & \circ & \\ & \circ & \ddots & \\ & & & d_r^{-1} \end{pmatrix}.$$

The definition of the inverse map can be extended so that

$$\mathbf{A}^r \mathbf{v}_i \rightarrow \mathbf{0} \quad , \quad i = r + 1, r + 2, \dots, m,$$

and extending to the basis of \mathbb{R}^n

$$\begin{aligned} \mathbf{A}^r (\mathbf{V}_1 \quad \mathbf{V}_2) &= (\mathbf{U}_1 \quad \mathbf{U}_2) \begin{pmatrix} d_1^{-1} & & & & \\ & d_2^{-1} & & & \\ & & \ddots & & \\ & & & d_r^{-1} & \\ & & & & \circ \\ & & & & \circ \end{pmatrix} \\ &= (\mathbf{U}_1 \mathbf{U}_2) \mathbf{D}^r \quad (\text{say}). \end{aligned}$$

Hence

$$\mathbf{A}^r = \mathbf{U} \mathbf{D}^r \mathbf{V}^t,$$

and \mathbf{A}^r maps $\mathcal{N}(\mathbf{A}^t)$ to the zero vector, that is

$$\mathcal{N}(\mathbf{A}^t) = \mathcal{N}(\mathbf{A}^r),$$

and

$$\mathcal{R}(\mathbf{A}^t) = \mathcal{R}(\mathbf{A}^r).$$

The mappings \mathbf{A}^r and \mathbf{A}^t are similar except

$$\mathbf{A}^t \mathbf{v}_i \rightarrow d_i \mathbf{u}_i \quad i = 1, 2, \dots, r,$$

and

$$\mathbf{A}^r \mathbf{v}_i \rightarrow \frac{1}{d_i} \mathbf{u}_i \quad i = 1, 2, \dots, r.$$

9.4.2 Moore-Penrose generalized inverse.

The mapping \mathbf{A}^+ from \mathbb{R}^m to \mathbb{R}^n is defined to be the unique solution to the equations,

$$\begin{aligned} \mathbf{A} \mathbf{A}^+ \mathbf{A} &= \mathbf{A} \\ \mathbf{A}^+ \mathbf{A} \mathbf{A}^+ &= \mathbf{A}^+ \\ \mathbf{A} \mathbf{A}^+ &= (\mathbf{A} \mathbf{A}^+)^t \\ \mathbf{A}^+ \mathbf{A} &= (\mathbf{A}^+ \mathbf{A})^t. \end{aligned} \tag{9.4.1}$$

It is easy to show by substitution that if

$$\mathbf{A} = \mathbf{Q} \mathbf{B} \mathbf{P}^t,$$

where \mathbf{Q} and \mathbf{P} are orthogonal then

$$\mathbf{A}^+ = \mathbf{P} \mathbf{B}^+ \mathbf{Q}^t,$$

that is, this inverse is invariant to change of orthogonal basis. Likewise for the partitioned $m \times n$ matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

where \mathbf{B} is nonsingular ($r \times r$),

$$\mathbf{A}^+ = \begin{pmatrix} \mathbf{B}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

where \mathbf{A}^+ is $n \times m$. Hence for the two factorizations,

SVD :
$$\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{U}^t, \quad \mathbf{A}^+ = \mathbf{U} \mathbf{D}^+ \mathbf{V}^t,$$

where $\mathbf{D}^+ = \mathbf{D}^r$, and

Q-U :
$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{P}^t, \quad \mathbf{A}^+ = \mathbf{P} \begin{pmatrix} \mathbf{L}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^t.$$

Note that $\mathbf{A}^r = \mathbf{A}^+$, that is, the Moore-Penrose generalized inverse is simply the inverse mapping restricted to the appropriate subspaces and mapping $\mathcal{N}(\mathbf{A}^t)$ to zero. The name pseudo-inverse is also given to this inverse. Other generalized inverses of \mathbf{A} satisfy a subset of the four defining equations (9.4.1).

9.4.3 Minimum norm least squares solution.

Define a mapping $\mathbf{A}^\dagger : \mathbb{R}^m \rightarrow \mathbb{R}^n$ by the rule that $\mathbf{y} \in \mathbb{R}^m$ maps to \mathbf{x}^\dagger which is that \mathbf{x} of minimum Euclidean norm of all \mathbf{x} which minimize $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$. It is easy to show \mathbf{x}^\dagger exists and is unique and that the mapping is linear. Use the SVD and change bases to \mathbf{U} and \mathbf{V} . Let $\mathbf{x} = \mathbf{U}\bar{\mathbf{x}}$, $\mathbf{y} = \mathbf{V}\bar{\mathbf{y}}$. The $\bar{\mathbf{x}}$ which has minimum norm ($\|\bar{\mathbf{x}}\| = \|\mathbf{x}\|$) of all $\bar{\mathbf{x}}$ which minimize

$$\begin{aligned} \|\mathbf{V}\mathbf{D}\mathbf{U}^t\mathbf{U}\bar{\mathbf{x}} - \mathbf{V}\bar{\mathbf{y}}\|^2 &= \|\mathbf{D}\bar{\mathbf{x}} - \bar{\mathbf{y}}\|^2 \\ &= \sum_{i=1}^r (d_i \bar{x}_i - \bar{y}_i)^2 + \sum_{i=r+1}^m \bar{y}_i^2 \end{aligned}$$

is chosen by $\bar{x}_i = \bar{y}_i/d_i$, $i = 1, 2, \dots, r$, and $\bar{x}_{r+1}, \bar{x}_{r+2}, \dots, \bar{x}_n$ any value. This gives the coordinates of vectors $\bar{\mathbf{x}}$ minimizing the sum of squares of residuals. Now minimizing $\|\mathbf{x}\|$ gives that

$$\bar{x}_{r+1} = \bar{x}_{r+2} = \dots = \bar{x}_n = 0.$$

Hence

$$\bar{\mathbf{x}}^\dagger = \mathbf{D}^r \bar{\mathbf{y}} = \mathbf{D}^\dagger \bar{\mathbf{y}},$$

and

$$\mathbf{x}^\dagger = \mathbf{A}^\dagger \mathbf{y} = \mathbf{A}^r \mathbf{y} = \mathbf{A}^+ \mathbf{y}.$$

So these three different approaches produce the same unique generalized inverse.

9.5 Best rank s approximation to a matrix.

Let \mathbf{A} be an $m \times n$ matrix with SVD, $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^t$. Then of all matrices \mathbf{B} of rank $s \leq r$, the matrix which minimizes

$$\|\mathbf{A} - \mathbf{B}\|_E = \left\{ \sum_{i,j} (a_{ij} - b_{ij})^2 \right\}^{1/2}$$

is

$$\mathbf{A}_T = \mathbf{V}\mathbf{D}_T\mathbf{U}^t,$$

where \mathbf{D}_T is a diagonal matrix equal to

$$\text{diag}[d_1, d_2, \dots, d_s, 0, \dots, 0],$$

and assuming $d_1 \geq d_2 \geq \dots \geq d_r > 0$. See Hoffman and Wielandt (1953) for a proof. \mathbf{A}_T is called the truncated estimate of rank s of \mathbf{A} . Note that as the Euclidean norm is invariant to multiplication by orthogonal matrices

$$\|\mathbf{A} - \mathbf{A}_T\|_E = \|\mathbf{D} - \mathbf{D}_T\|_E = \left(\sum_{i=s+1}^r d_i^2 \right)^{1/2}.$$

Also $\|\mathbf{A} - \mathbf{A}_T\|_2 = d_{s+1}$. This replaces the smallest singular values by zero, and this has relevance in numerical work as seldom can singular values $d_i = 0$, $i = r + 1, r + 2, \dots, \min(m, n)$ be calculated accurately.

9.6 Generalized inverse and rounding error.

In computing a generalized inverse a decision has been made of when a singular value is zero, given that rounding error effects all the factors of the SVD or Q-U factorization. In linear algebra use is made of the ‘backward error analysis’ approach. The actual computed factors (denoted by $\bar{\cdot}$ in this section) are assigned to be the exact factors of some perturbed original matrix. So in terms of the SVD

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^t,$$

$$\bar{\mathbf{A}} = \bar{\mathbf{V}} \bar{\mathbf{D}} \bar{\mathbf{U}}^t,$$

and a stable algorithm for calculating the SVD will mean that

$$\|\mathbf{A} - \bar{\mathbf{A}}\| \leq p(m, n)\epsilon,$$

where ϵ is the machine precision and $p(m, n)$ is a low degree polynomial in m and n . Note this does not mean that the factors $\bar{\mathbf{V}}, \bar{\mathbf{D}}, \bar{\mathbf{U}}$ are necessarily close, respectively to \mathbf{V}, \mathbf{D} and \mathbf{U} .

In particular if $\bar{\mathbf{D}} = \text{diag}(\bar{d}_1, \bar{d}_2, \dots, \bar{d}_\ell)$, where $\ell = \min(m, n)$, some of the singular values of \mathbf{A} which might have been zero are probably of order the machine precision ϵ in $\bar{\mathbf{D}}$. In order to determine the rank of \mathbf{A} a decision about small singular values of $\bar{\mathbf{D}}$ needs to be made. The neatest procedure is to replace $\bar{\mathbf{D}}$ with its truncated estimate $\bar{\mathbf{D}}_T$ so that the effective matrix is

$$\bar{\mathbf{A}}_T = \bar{\mathbf{V}} \bar{\mathbf{D}}_T \bar{\mathbf{U}}^t.$$

If $\bar{\mathbf{A}}_T$ has rank s then (approximately as $\bar{\mathbf{V}}$ and $\bar{\mathbf{U}}$ are not orthogonal)

$$\|\mathbf{A} - \bar{\mathbf{A}}_T\| \leq p(m, n)\epsilon + \bar{d}_{s+1}$$

and $\bar{\mathbf{A}}_T$ will still be close to \mathbf{A} . If the non zero singular values of \mathbf{A} are well separated from zero, relative to the machine precision, so that $d_i \gg \epsilon$, $i = 1, 2, \dots, r$, then as $|d_i - \bar{d}_i| \leq \kappa\epsilon$ it is easy to choose an s equal to the rank of \mathbf{A} , that is,

$$\begin{aligned} \bar{d}_i &\gg \epsilon & i = 1, 2, \dots, r \\ \bar{d}_i &\approx \epsilon & i = r + 1, \dots, \ell. \end{aligned}$$

The above ideas also apply to the Q-U factorization in that the rank of the matrix has to be decided at some stage during the factorization. Column pivoting is needed during the factorization;

$$\mathbf{A} \longrightarrow \mathbf{Q} \begin{pmatrix} \mathbf{U} \\ \mathbf{0} \end{pmatrix} \mathbf{P},$$

where \mathbf{P} represents column interchanges. In practice before the i -th elementary Hermitian stage

$$\mathbf{A} \longrightarrow \bar{\mathbf{Q}} \begin{pmatrix} \bar{\mathbf{U}} \\ \mathbf{0} \end{pmatrix} \bar{\mathbf{P}}.$$

If $\|\delta\bar{\mathbf{U}}\| < \kappa\epsilon$ for some constant κ , $\delta\bar{\mathbf{U}}$ may be considered to be entirely round off error and an effective $\bar{\mathbf{A}}$ is

$$\bar{\mathbf{A}} = \bar{\mathbf{Q}} \begin{pmatrix} \bar{\mathbf{U}} \\ \mathbf{0} \end{pmatrix} \bar{\mathbf{P}}.$$

Column pivoting ensures that as \mathbf{U} gains its i -th row, $|U_{ii}|$ is larger than or equal to any succeeding $|U_{jj}|$. So $|U_{11}| \geq |U_{22}| \geq \dots \geq |U_{\ell\ell}|$. Also $|U_{ii}| \geq |U_{ij}|$, $j = i + 1, \dots, n$. Approximately (because $\bar{\mathbf{Q}}$ is not exactly orthogonal)

$$\|\mathbf{A} - \bar{\mathbf{A}}\| \leq \|\delta\bar{\mathbf{U}}\|,$$

so this procedure can be looked upon as an approximation to the truncated estimate. The Q-U factorization is faster to compute than the SVD, but this has become less of a problem with modern workstations.

9.7 Scaling considerations.

It is possible in practical situations to have a matrix made from data which differs vastly in magnitude but which is accurate to only four significant digits (not decimal places). For example in econometrics, an $m \times n$ data matrix \mathbf{A} has n columns each of which represent an economic variable measured over m time periods. The norm of a column depends on the units of the column, for example, population measured as a single person, 10^3 people, 10^6 people or 10^9 people, will give vastly different numbers to represent the population of China. Similarly for the rows, because of inflation, if money variables are expressed in dollars, the last row is expected to be larger than the first. These two scaling problems need to be handled differently but if a meaningful rank decision is to be made about \mathbf{A} they must be taken into account.

Taking the row problem first. Usually a least squares solution is found to a set of equations $\mathbf{Ax} = \mathbf{y} + \mathbf{r}$ ($m > n$), that is, the sum of squares of residuals of each time period ($\mathbf{r}^t \mathbf{r}$) is minimized. This is a valid statistical procedure if the residuals are independent and homoscedastic, that is $\mathcal{E}(\mathbf{r} \mathbf{r}^t) = \sigma^2 \mathbf{I}_m$, where $\mathbf{r} = \mathbf{Ax} - \mathbf{y}$. If $\mathcal{E}(\mathbf{r} \mathbf{r}^t) = \mathbf{V}$ (variance-covariance of residuals) then $\mathbf{r}^t \mathbf{V}^{-1} \mathbf{r}$ should be minimized instead. If the residuals are independent but not homoscedastic (heteroscedastic) then \mathbf{V} is a diagonal matrix and the rows of the equation are weighted inversely with the expected standard error of the i -th residual. Hence any row scaling is usually governed by statistical considerations in this example.

The column scalings are obviously free to be chosen to make the numerical computation as stable as possible. Consider the 3×2 matrix

$$\begin{pmatrix} 1.123 & 2.345 \times 10^9 \\ 1.123 & 2.567 \times 10^9 \\ 1.123 & 2.789 \times 10^9 \end{pmatrix}$$

representing a constant column and a variable in units of dollars say. This matrix is essentially singular as one column is less than 10^{-9} times the norm of the matrix. In fact one singular value will be of order 10^9 , the other of order 1. Yet both columns have just as much significant data in them. If another column equal to $\kappa_1(\mathbf{A}) + 10^{-9} \kappa_2(\mathbf{A})$ is added then for the 3×3 matrix the third singular value is zero. But on a 32 bit machine one singular value of order 10^9 is computed, the other two computed are of magnitude $10-100$, due to roundoff error. As the choice of units is arbitrary, units to make each column unity can be chosen. This essentially means that significant figure accuracy lines up with decimal place accuracy. For our scaled 3×3 matrix two singular values of magnitude 1 and one of magnitude the machine precision are computed.

In terms of the system of equations

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\| \tag{9.7.1}$$

$$\begin{aligned} &\equiv \min_{\mathbf{x}} \|\mathbf{ADD}^{-1} \mathbf{x} - \mathbf{y}\| \\ &\equiv \min_{\hat{\mathbf{x}}} \|\hat{\mathbf{A}} \hat{\mathbf{x}} - \mathbf{y}\| \quad , \quad \hat{\mathbf{x}} = \mathbf{D}^{-1} \mathbf{x} \quad , \quad \hat{\mathbf{A}} = \mathbf{AD}. \end{aligned} \tag{9.7.2}$$

If \mathbf{A} has full rank, so that $\mathbf{A}^+ = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t = \mathbf{D} \hat{\mathbf{A}}^+$ then the same vector \mathbf{x} is computed for both problems (9.7.1) and (9.7.2). However if \mathbf{A} does not have rank n and the minimum norm solution of the least squares solution is calculated then in (9.7.1) $\|\mathbf{x}\|$ is minimized and in (9.7.2) $\|\hat{\mathbf{x}}\| = \|\mathbf{D}^{-1} \mathbf{x}\|$ is minimized. In terms of the previous 3×2 example (9.7.2) is likely to be preferred anyway as it appears to give equal relative weighting to all \hat{x}_i values.

9.8 Generalized inverses under scalings.

Let \mathbf{R} and \mathbf{S} be symmetric positive definite matrices of sizes m and n respectively. The inner products $\langle \mathbf{a}, \mathbf{b} \rangle_{\mathbf{S}} = \mathbf{a}^t \mathbf{S} \mathbf{b}$ and hence norms $\|\mathbf{a}\|_{\mathbf{S}} = \langle \mathbf{a}, \mathbf{a} \rangle_{\mathbf{S}}^{1/2} = (\mathbf{a}^t \mathbf{S} \mathbf{a})^{1/2}$ can be used to define generalized inverses under scalings. Define a generalized inverse \mathbf{A}^+ as the mapping of $\mathbf{y} \in \mathbb{R}^m$ to the unique vector \mathbf{x}_G which is the smallest vector in terms of $\|\cdot\|_{\mathbf{S}}$ of all vectors \mathbf{x} which minimize $\|\mathbf{Ax} - \mathbf{y}\|_{\mathbf{R}}$. If $\mathbf{r} = \mathbf{y} - \mathbf{Ax}$, this is finding the smallest vector $\hat{\mathbf{x}}_G$ in terms of $\|\cdot\|_2$ of all vectors $\hat{\mathbf{x}}$ which minimize $\mathbf{r}^t \mathbf{R} \mathbf{r}$ where $\hat{\mathbf{x}} = \mathbf{S}^{1/2} \mathbf{x}$ and $(\mathbf{S}^{1/2})^t \mathbf{S}^{1/2} = \mathbf{S}$.

In short, scaling is equivalent to putting different metrics on spaces for generalised inverse computations. These different metrics can be more appropriate than a user's arbitrary choice of units. Computation within statistical or econometric packages should be done in terms of the scaled variables.

9.9 Approximations to the generalized inverse.

9.9.1 Limiting approximations.

If \mathbf{A} is $m \times n$, $m > n$ of rank n (full rank) then

$$\mathbf{A}^+ = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t,$$

and if $m < n$, rank m

$$\mathbf{A}^+ = \mathbf{A}^t (\mathbf{A} \mathbf{A}^t)^{-1}.$$

In the case of not full rank, $r < \min(m, n)$, we have

$$\mathbf{A}^+ = \lim_{\epsilon \rightarrow 0} (\mathbf{A}^t \mathbf{A} + \epsilon^2 \mathbf{I})^{-1} \mathbf{A}^t = \lim_{\epsilon \rightarrow 0} \mathbf{A}^t (\mathbf{A} \mathbf{A}^t + \epsilon^2 \mathbf{I})^{-1},$$

which is easily proved using the SVD of \mathbf{A} and noting that

$$\lim_{\epsilon \rightarrow 0} \frac{d_i}{d_i^2 + \epsilon^2} = d_i^+ = \begin{cases} d_i^{-1}, & d_i \neq 0, \\ 0, & d_i = 0. \end{cases}$$

Of course in numerical computation $\epsilon \rightarrow 0$ is not possible as an increasingly ill-conditioned inversion needs to be computed. A similar effect to the truncated estimate generalized inverse is preferred. Note that the singular values of $(\mathbf{A}^t \mathbf{A} + \epsilon^2 \mathbf{I}) \mathbf{A}^t$ are

$$\frac{d_i}{d_i^2 + \epsilon^2} \approx \begin{cases} d_i^{-1} & \text{if } d_i \gg \epsilon, \\ 0 & \text{if } d_i \ll \epsilon. \end{cases}$$

Hence by keeping ϵ to be of magnitude of the cut off for the truncated estimate (perhaps a bit smaller) a similar effect is achieved.

9.9.2 Damped least squares.

Consider the problem

$$\begin{aligned} & \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \epsilon^2 \|\mathbf{x}\|_2^2, \\ & \equiv \min_{\mathbf{x}} \left\| \begin{pmatrix} \mathbf{A} \\ \epsilon \mathbf{I} \end{pmatrix} \mathbf{x} - \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \right\|_2^2. \end{aligned}$$

As $\begin{pmatrix} \mathbf{A} \\ \epsilon \mathbf{I} \end{pmatrix}$ has full rank so the solution is given by,

$$\left((\mathbf{A}^t \quad \epsilon \mathbf{I}) \begin{pmatrix} \mathbf{A} \\ \epsilon \mathbf{I} \end{pmatrix} \right)^{-1} (\mathbf{A}^t \quad \epsilon \mathbf{I}) \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} = (\mathbf{A}^t \mathbf{A} + \epsilon^2 \mathbf{I})^{-1} \mathbf{A}^t \mathbf{y}.$$

A way to compute this solution is to use the Q-U factorization on $\begin{pmatrix} \mathbf{A} \\ \epsilon \mathbf{I} \end{pmatrix}$.

9.9.3 Underdetermined case.

Consider the system $\mathbf{A}\mathbf{X} = \mathbf{y}$, \mathbf{A} is $m \times n$, $m < n$, and the augmented system

$$(\mathbf{A} \quad \epsilon \mathbf{I}) \begin{pmatrix} \mathbf{x} \\ \mathbf{x}_r \end{pmatrix} = \mathbf{A}\mathbf{x} + \epsilon \mathbf{x}_r = \mathbf{y},$$

which has full rank, so that

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{x}_r \end{pmatrix} = \begin{pmatrix} \mathbf{A}^t \\ \epsilon \mathbf{I} \end{pmatrix} (\mathbf{A} \mathbf{A}^t + \epsilon^2 \mathbf{I})^{-1} \mathbf{y},$$

or

$$\mathbf{x} = \mathbf{A}^t (\mathbf{A} \mathbf{A}^t + \epsilon^2 \mathbf{I})^{-1} \mathbf{y},$$

and the residual is given by $\epsilon \mathbf{x}_r = \epsilon^2 (\mathbf{A} \mathbf{A}^t + \epsilon^2 \mathbf{I})^{-1} \mathbf{y}$.

9.9.4 Rutishauser's doubly relaxed least squares.

See Rutishauser (1968) for the details of this approximation.

$$\mathbf{A}^+ = \lim_{\epsilon \rightarrow 0} \{ (\mathbf{A}^t \mathbf{A} + \epsilon^2 \mathbf{I}) + \epsilon^2 (\mathbf{A}^t \mathbf{A} + \epsilon^2 \mathbf{I})^{-1} \}^{-1} \mathbf{A}^t.$$

9.10 Projection operators.

9.10.1 Projection onto a subspace.

Let \mathcal{X} be a subspace of \mathbb{R}^m . $\mathbf{P}_{\mathcal{X}}$ is a projection operator from \mathbb{R}^m to \mathcal{X} if

$$\mathcal{R}(\mathbf{P}_{\mathcal{X}}) = \mathcal{X},$$

and, if $\mathbf{x} \in \mathcal{X}$, then $\mathbf{P}_{\mathcal{X}}\mathbf{x} = \mathbf{x}$. So if $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{y} = \mathbf{x} + \mathbf{z}$, $\mathbf{x} \in \mathcal{X}$, $\mathbf{z} \in \mathcal{X}^\perp$ then

$$\mathbf{P}_{\mathcal{X}}\mathbf{y} = \mathbf{x}.$$

9.10.2 Spaces represented by matrices.

Let \mathcal{X} be the space of the span of the columns of an $m \times n$ matrix \mathbf{A} , that is, $\mathcal{X} = \mathcal{R}(\mathbf{A}) \subset \mathbb{R}^m$. Then $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^t)$, and any vector in \mathbb{R}^m can be expressed as the sum of a component in $\mathcal{R}(\mathbf{A})$ and a component in $\mathcal{N}(\mathbf{A}^t)$. $\mathbf{P}_{\mathcal{X}} = \mathbf{P}_{\mathcal{R}(\mathbf{A})}$ has to choose the component in $\mathcal{R}(\mathbf{A})$. So

$$\mathbf{P}_{\mathcal{X}}\kappa_j(\mathbf{A}) = \kappa_j(\mathbf{A}) \quad , \quad j = 1, 2, \dots, n,$$

or

$$\mathbf{P}_{\mathcal{X}}\mathbf{A} = \mathbf{A}. \tag{9.10.1}$$

This means that $\mathbf{P}_{\mathcal{X}}$ must be represented by an $m \times m$ matrix, with each column a linear combination of columns of \mathbf{A} . Hence $\mathbf{P}_{\mathcal{X}} = \mathbf{A}\mathbf{B}$ for some $n \times m$ matrix \mathbf{B} . So from (9.10.1)

$$\mathbf{A}\mathbf{B}\mathbf{A} = \mathbf{A},$$

and a choice for \mathbf{B} is \mathbf{A}^+ , that is

$$\mathbf{P}_{\mathcal{X}} = \mathbf{A}\mathbf{A}^+.$$

In the full rank cases, if \mathbf{A} has rank n then $\mathbf{P}_{\mathcal{X}} = \mathbf{A}(\mathbf{A}^t\mathbf{A})^{-1}\mathbf{A}^t$, while if the rank is m , then $\mathbf{P}_{\mathcal{X}} = \mathbf{A}\mathbf{A}^t(\mathbf{A}\mathbf{A}^t)^{-1} = \mathbf{I}$.

It is easy to show that

$$\mathbf{P}_{\mathcal{X}^\perp} = \mathbf{I} - \mathbf{A}\mathbf{A}^+$$

is the projection onto $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^t)$. There is no need to find \mathbf{A}^+ in order to compute $\mathbf{P}_{\mathcal{X}}$ or $\mathbf{P}_{\mathcal{X}^\perp}$, as the orthogonal factorizations give relevant expressions in terms of parts of orthogonal matrices. The actual rank of \mathbf{A} needs to be determined.

9.10.3 Projections in terms of the SVD.

Let $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^t$, be of rank r . Then $\mathbf{v}_1\mathbf{v}_2, \dots, \mathbf{v}_r$ is an orthogonal basis for $\mathcal{R}(\mathbf{A})$. So if $\mathbf{y} \in \mathbb{R}^m$,

$$\begin{aligned} \mathbf{y} &= \sum_{i=1}^m \langle \mathbf{v}_i, \mathbf{y} \rangle \mathbf{v}_i, \\ \mathbf{P}_{\mathcal{X}}\mathbf{y} &= \sum_{i=1}^r \langle \mathbf{v}_i, \mathbf{y} \rangle \mathbf{v}_i, \end{aligned}$$

and

$$\mathbf{P}_{\mathcal{X}^\perp}\mathbf{y} = \sum_{i=r+1}^m \langle \mathbf{v}_i, \mathbf{y} \rangle \mathbf{v}_i.$$

Writing $\mathbf{V} = (\mathbf{V}_1 \quad \mathbf{V}_2)$, where \mathbf{V}_1 is $m \times r$,

$$\mathbf{y} = \mathbf{V}\bar{\mathbf{y}} \quad , \quad \bar{\mathbf{y}} = \mathbf{V}^t\mathbf{y}$$

and

$$\begin{aligned} \mathbf{P}_{\mathcal{X}}\mathbf{y} &= \mathbf{V}_1 \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \end{pmatrix} \bar{\mathbf{y}} \\ &= \mathbf{V}_1 \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \end{pmatrix} \mathbf{V}^t\mathbf{y} \\ &= \mathbf{V}_1\mathbf{V}_1^t\mathbf{y} \quad \text{or} \quad \mathbf{V} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^t\mathbf{y}, \end{aligned}$$

and similarly

$$\mathbf{P}_{\mathcal{X}^\perp}\mathbf{y} = \mathbf{V}_2\mathbf{V}_2^t\mathbf{y} = \mathbf{V} \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-r} \end{pmatrix} \mathbf{V}^t\mathbf{y}.$$

If bases are changed to the columns of \mathbf{V} then these two projection operators are now represented by

$$\begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-r} \end{pmatrix}.$$

Substitution of the SVD of \mathbf{A} and \mathbf{A}^+ into $\mathbf{A}\mathbf{A}^+$ and $\mathbf{I} - \mathbf{A}\mathbf{A}^+$ will also yield the above formulae. See Jennings (1980) for an application in econometrics.

9.11 Exercises

1. Find projection operators $\mathbf{A}\mathbf{A}^+$ and $\mathbf{I} - \mathbf{A}\mathbf{A}^+$ in terms of the Q–U factorization.
2. Prove the Rutishauser formula and compare its approximation to the truncated estimate generalized inverse.
3. Show $\mathbf{A}^+\mathbf{A}\mathbf{A}^t = \mathbf{A}^t$, and $(\mathbf{A}^+)^t = (\mathbf{A}^t)^+$.
4. Show that a projection matrix \mathbf{P} is idempotent, that is, $\mathbf{P}^2 = \mathbf{P}$.
5. Using the notation of §9.8, show that

$$\mathbf{x}_G = \lim_{\epsilon \rightarrow 0} (\mathbf{A}^t \mathbf{R} \mathbf{A} + \epsilon^2 \mathbf{S})^{-1} \mathbf{A}^t \mathbf{R} \mathbf{y},$$

which is the solution of the generalized damped least squares problem

$$\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_R^2 + \epsilon^2 \|\mathbf{x}\|_S.$$

10. Eigenvalue Computations – An Overview.

10.1 Introduction.

In general $M(\lambda)$ an $n \times n$ matrix being a function of one parameter λ . The objective is to find values of λ such that

$$\det |M(\lambda)| = 0.$$

At these points a representation of $\mathcal{N}(M(\lambda))$ is required, that is, non-zero \mathbf{v} such that

$$M(\lambda)\mathbf{v} = \mathbf{0}.$$

Classical linear case.

Here $M(\lambda) = \mathbf{A} - \lambda\mathbf{I}$, where \mathbf{A} is an $n \times n$ matrix. The problem is stated as: find λ such that $\det |\mathbf{A} - \lambda\mathbf{I}| = 0$, and non-zero vector \mathbf{v} such that $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$. The scalar λ is known as the eigenvalue and \mathbf{v} its corresponding eigenvector. As $\det |\mathbf{A} - \lambda\mathbf{I}|$ is a polynomial in λ of degree n there exists n linear factors over \mathbb{C} . For each of these λ_i a non-zero \mathbf{v}_i exists. A matrix \mathbf{B} is similar to \mathbf{A} if there exists a nonsingular matrix \mathbf{T} such that

$$\mathbf{B} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1}.$$

In this case \mathbf{A} and \mathbf{B} have the same eigenvalues as $\det |\mathbf{A} - \lambda\mathbf{I}| = \det |\mathbf{B} - \lambda\mathbf{I}|$. If \mathbf{v} is an eigenvector of \mathbf{A} then $\mathbf{T}\mathbf{v}$ is an eigenvector of \mathbf{B} .

Jordan Canonical Form.

Theorem: Every $n \times n$ matrix \mathbf{A} is similar to a matrix of the form

$$\mathbf{J} = \begin{pmatrix} \mathbf{J}_1 & & & \\ & \mathbf{J}_2 & & \\ & & \ddots & \\ & & & \mathbf{J}_r \end{pmatrix},$$

where each \mathbf{J}_i has one of the forms λ_i , $\begin{pmatrix} \lambda_i & 1 \\ 0 & \lambda_i \end{pmatrix}$, $\begin{pmatrix} \lambda_i & 1 & 0 \\ 0 & \lambda_i & 1 \\ 0 & 0 & \lambda_i \end{pmatrix}$, etc. There is only one null space vector of the sub-matrix $(\mathbf{J}_r - \lambda_r\mathbf{I})$, even though the eigenvalue is a repeated root if \mathbf{J}_r has order ≥ 2 . The invertible matrix \mathbf{T} where $\mathbf{A} = \mathbf{T}\mathbf{J}\mathbf{T}^{-1}$ still needs n linearly independent columns, even though only r of them are eigenvectors. See Lipschutz (1968) for the detail. ■

10.2 Classes of algorithm.

10.2.1 Eigenvalues: Classes of methods.

- (i) Methods which transform \mathbf{A} to \mathbf{B} in a finite number of elementary similarity transformations so that the polynomial

$$\det |\mathbf{B} - \lambda\mathbf{I}|$$

is readily evaluated and hence its zeros found using a zero finding technique, form the first class. These are good methods for finding one or two eigenvalues. Their efficiency is not good for finding all the eigenvalues.

- (ii) The second class of methods are those which transform \mathbf{A} to a form where the eigenvalues can be read off, for example diagonal or upper triangular.
- (iii) There are some methods distinct from the above used mainly in the nonlinear eigenvalue case, but which can be used in the linear case to improve the accuracy of estimates of λ and \mathbf{v} .

10.2.2 Eigenvectors.

In most cases the method of *inverse iteration* is used to find eigenvectors. Usually some accumulation of a finite sequence of similarity transformations is necessary for efficiency.

10.3 Classes of matrices.

In order of difficulty of computation

- (i) Complex general: The eigenvalues and eigenvectors are in general complex and n linearly independent eigenvectors may not exist if eigenvalues are repeated.
- (ii) Real general: The eigenvalues are either real or in complex conjugate pairs as are the eigenvectors and n linearly independent eigenvectors may not exist if eigenvalues are repeated.
- (iii) Hessenberg matrices: Matrices in the above two classes are similarity transformed to a real or complex Hessenberg matrix as part of the algorithm for finding eigenvalues and eigenvectors.
- (iv) Hermitian complex: The eigenvalues are real and the eigenvectors can be complex but can be chosen and scaled so that they form an orthonormal basis for \mathbb{C}^n .
- (v) Hermitian real (real symmetric): The eigenvalues are real and the eigenvectors can be chosen and scaled to form an orthonormal basis for \mathbb{R}^n .
- (vi) Symmetric banded: Special algorithms can be used to similarity transform these matrices to the next class.
- (vii) Symmetric tridiagonal: These are Hessenberg matrices and are formed as part of algorithms for computing eigenvalues and eigenvectors of classes (iv), (v) and (vi).

10.4 Real symmetric matrices (Hermitian matrices).

We will only discuss the real case. A recent history and summary can be found in Parlett (1981).

10.4.1 Jacobi method.

This is an old method not in use today but it is an example of the effectiveness of elementary plane rotations. Essentially the matrix A is transformed to diagonal form by an (infinite) sequence of similarity transformations.

$$\dots Q_k \dots Q_2 Q_1 A Q_1^t Q_2^t \dots Q_k \dots = \Lambda.$$

$$Q A Q^t = \Lambda.$$

(i) Basic algorithm.

1. Find the largest off-diagonal element, suppose it is (i, j) .
2. Apply a similarity transformation (plane rotation) to zero the element (i, j) (and (j, i) as well since A is symmetric).

$$\bar{A} = P A P^t,$$

$$\bar{A}_{ij} = 0 \Rightarrow (P A P^t)_{ij} = 0.$$

$$\rho_i(PA) = c\rho_i(A) + s\rho_j(A)$$

$$\rho_j(PA) = s\rho_i(A) - c\rho_j(A).$$

$$\kappa_i(PAP^t) = c\kappa_i(PA) + s\kappa_j(PA)$$

$$\kappa_j(PAP^t) = s\kappa_i(PA) - c\kappa_j(PA).$$

$$\therefore \begin{aligned} \bar{A}_{ii} &= c^2 A_{ii} + 2scA_{ij} + s^2 A_{jj} \\ \bar{A}_{jj} &= c^2 A_{jj} - 2scA_{ij} + s^2 A_{ii} \\ \bar{A}_{ij} &= (s^2 - c^2)A_{ij} + sc(A_{ii} - A_{jj}) = 0. \end{aligned}$$

$$\therefore \frac{2sc}{c^2 - s^2} = \frac{2A_{ij}}{A_{ii} - A_{jj}}.$$

If $c = \cos \theta$, $s = \sin \theta$,

$$\tan 2\theta = \frac{2A_{ij}}{A_{ii} - A_{jj}}.$$

Calculation of s and c has to be done carefully to avoid subtractive cancellation. See Boothroyd (1968).

3. If the sum of squares of off-diagonal elements $< \epsilon$ then stop, else goto 1.

(ii) Convergence:

Theorem: Let $\mathbf{A}^{(k)} = \mathbf{P}\mathbf{A}^{(k-1)}\mathbf{P}^t$ be the sequence of plane rotations ($\mathbf{A}^0 = \mathbf{A}$), then if E_k is the sum of squares of off-diagonal elements of $\mathbf{A}^{(k)}$ then $E_k \rightarrow 0$ as $k \rightarrow \infty$.

Proof:

$\|\mathbf{A}^{(k)}\|_E = \|\mathbf{A}^{(k-1)}\|_E = \|\mathbf{A}\|_E$. The sum of squares of the diagonal elements is shown to be monotone increasing and as it is bounded above by $\|\mathbf{A}\|_E$, converges to $\|\mathbf{A}\|_E$. This of course implies that $E_k \rightarrow 0$. Note

$$a^2 + b^2 = \frac{1}{2}\{(a+b)^2 + (a-b)^2\}. \quad (10.4.1)$$

Denote $\bar{A}_{ij} = (\mathbf{A}^{(k)})_{ij}$, $A_{ij} = (\mathbf{A}^{(k-1)})_{ij}$

$$\begin{aligned} \bar{A}_{ii} + \bar{A}_{jj} &= A_{ii} + A_{jj} \\ \bar{A}_{ii} - \bar{A}_{jj} &= (c^2 - s^2)(A_{ii} - A_{jj}) + 4csA_{ij} \\ &= (A_{ii} - A_{jj}) \sec 2\theta. \end{aligned}$$

$$\begin{aligned} \bar{A}_{ii}^2 + \bar{A}_{jj}^2 &= \frac{1}{2}\{(A_{ii} + A_{jj})^2 + (A_{ii} - A_{jj})^2 \sec^2 2\theta\} && \text{from(10.4.1)} \\ &= A_{ii}^2 + A_{jj}^2 + \frac{1}{2}(A_{ii} - A_{jj})^2(\sec^2 2\theta - 1) && \text{from(10.4.1)} \\ &= A_{ii}^2 + A_{jj}^2 + 2A_{ij}^2 \end{aligned}$$

Hence $E_k = E_{k-1} - 2(A_{ij}^{(k-1)})^2$, which implies $E_k \rightarrow 0$ as $k \rightarrow \infty$. ■

(iii) Rate of convergence:

$$\begin{aligned} E_k &= E_{k-1} - 2(A_{ij}^{(k-1)})^2 \\ &\leq \left(1 - \frac{2}{n(n-1)}\right)E_{k-1}, \text{ as } E_{k-1} \leq n(n-1)(A_{ij}^{(k-1)})^2. \end{aligned}$$

Let $N = n(n-1)/2$

$$\begin{aligned} E_{k+N} &\leq \left(1 - \frac{1}{N}\right)^N E_k \\ &\approx e^{-1}E_k. \end{aligned}$$

Stronger result: if all eigenvalues are strongly separated ($|\lambda_i - \lambda_j| > \delta$) then over every N iterations quadratic convergence is achieved, so that

$$E_{k+N} \leq KE_k^2.$$

(iv) Algorithm in practice.

Sequential algorithm — zero elements in turn with no searching for the maximum off-diagonal element. One pass through the off-diagonal elements is known as a sweep.

Threshold algorithm — sweeps are performed but miss all elements less than some tolerance. Once all off diagonals are less then the tolerance decrease it and begin again.

1. $tol = 0.1$
2. Sweep until all off-diagonals are $< tol$.
3. $tol = tol^2$
if $tol > \text{machine precision}$ goto 2
else stop.

(v) The eigenvectors are found by storing the plane rotations on an identity matrix.

10.4.2 Reduction to tridiagonal form.

(i) Givens method — using plane rotations.

For $i = 1$ to $n - 2$ do

For $j = i + 2$ to n do

plane rotate rows and columns $i + 1$ and j to make $A_{ji} = A_{ij} = 0$;

As the zeros are introduced past zeros are not changed. This takes $\frac{8}{3}n^3$ multiplications and $O(n^2)$ square roots.

(ii) Householder method — using Householder transformations.

For $i = 1$ to $n - 2$ do

transform row and column i into their $(i + 1)$ -th element, zeroing elements $i + 2$ to n .

This takes $\frac{4}{3}n^3$ operations and $n - 2$ square roots.

10.5 Tridiagonal matrices.

10.5.1 Non-symmetric matrices.

Under certain conditions a non-symmetric tridiagonal matrix can be made diagonally similar to a symmetric tridiagonal matrix.

$$\mathbf{A} = \begin{pmatrix} \alpha_1 & \beta_2 & & & & \\ \gamma_2 & \alpha_2 & \beta_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & \ddots & \ddots & \beta_n \\ & & & & \gamma_n & \alpha_n \end{pmatrix}$$

Choose $\mathbf{D} =$ diagonal such that

$$\mathbf{D}^{-1}\mathbf{A}\mathbf{D} = \mathbf{T},$$

where \mathbf{T} is symmetric tridiagonal. The condition is that $\gamma_i\beta_i > 0$, $i = 2, 3, \dots, n$.

10.5.2 Symmetric tridiagonal matrices.

It is an easy matter to evaluate the determinant of a symmetric tridiagonal matrix,

$$\mathbf{A} = \begin{pmatrix} \alpha_1 & \beta_2 & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & \ddots & \ddots & \beta_n \\ & & & & \beta_n & \alpha_n \end{pmatrix}.$$

Assume no $\beta_i = 0$, as the matrix could be partitioned if $\beta_i = 0$ for some i . Consider $P_n(\lambda) = \det |\mathbf{A} - \lambda\mathbf{I}|$. Expand $\det |\mathbf{A} - \lambda\mathbf{I}|$ about the last row. Let $P_i(\lambda)$ be the determinant of the i -th principle minor of $(\mathbf{A} - \lambda\mathbf{I})$,

$$P_n(\lambda) = (\alpha_n - \lambda)P_{n-1}(\lambda) - \beta_n^2 P_{n-2}(\lambda).$$

Similarly if $P_0(\lambda) = 1$, $P_{-1}(\lambda) = 0$, define

$$P_i(\lambda) = (\alpha_i - \lambda)P_{i-1}(\lambda) - \beta_i^2 P_{i-2}(\lambda), \quad i = 1, 2, \dots, n.$$

Using this recursion formulae to evaluate $P_n(\lambda)$ any root finding procedure could be used to find eigenvalues.

10.5.3 Sturm sequence property

This property is used to generate a bisection method of zero finding.

Lemma: Let the quantities $P_0(x), \dots, P_n(x)$ be evaluated at x . If $s(x)$ is the number of agreements in sign of consecutive members of the sequence, then $s(x)$ is also the number of eigenvalues of \mathbf{A} which are strictly greater than x .

Proof: by induction, see Wilkinson (1965) p300, or Froberg (1965) p115, or Ralston (1965) p494. ■

Numerical stability: Wilkinson (1965) p302 – zeros found using this technique are the zeros of a matrix $\mathbf{A} + \mathbf{E}$ where \mathbf{E} is of order the machine precision.

10.6 Inverse iteration.

This is a very effective method of finding an eigenvector given a good estimate of the eigenvalue.

10.6.1 Algorithm

Suppose λ approximates λ_i .

- (i) Solve $(\mathbf{A} - \lambda\mathbf{I})\mathbf{v}_1 = \mathbf{b}$, \mathbf{b} is chosen as \mathbf{e} say.
- (ii) Solve $(\mathbf{A} - \lambda\mathbf{I})\mathbf{v}_2 = \mathbf{v}_1/\|\mathbf{v}_1\|$.

The normalized vector $\mathbf{v}_2/\|\mathbf{v}_2\|$ is a good estimate of the eigenvector corresponding to λ_i . Note that a system which is ‘singular’ is solved, but of course with roundoff errors there is a low probability $(\mathbf{A} - \lambda\mathbf{I})$ is exactly singular numerically. If this happens, replace the zero element with the machine precision and continue the division.

10.6.2 Analysis of algorithm:

Suppose n linearly independent eigenvectors \mathbf{x}_i exist. Let $\mathbf{X} = [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_n]$, then $\mathbf{A} = \mathbf{X}\mathbf{\Lambda}\mathbf{X}^{-1}$. Write $\mathbf{b} = \sum_{i=1}^n c_i\mathbf{x}_i = \mathbf{X}\mathbf{c}$. Consider

$$\begin{aligned} (\mathbf{A} - \lambda\mathbf{I})\mathbf{v}_1 &= \mathbf{b} \\ \mathbf{v}_1 &= (\mathbf{X}\mathbf{\Lambda}\mathbf{X}^{-1} - \lambda\mathbf{I})^{-1}\mathbf{b} \\ &= \mathbf{X}(\mathbf{\Lambda} - \lambda\mathbf{I})^{-1}\mathbf{c} \\ &= \sum_{i=1}^n \frac{c_i}{\lambda_i - \lambda}\mathbf{x}_i. \end{aligned}$$

Now let λ be close to λ_1 say, and partition the sum

$$\mathbf{v}_1 = \frac{c_1}{\lambda_1 - \lambda}\mathbf{x}_1 + \sum_{i=2}^n \frac{c_i}{\lambda_i - \lambda}\mathbf{x}_i,$$

and if λ is distant (relatively) from the other λ_i and if $c_1 \neq 0$ or not small, the number $c_1/(\lambda_1 - \lambda)$ will be much larger than the rest of the $c_i/(\lambda_i - \lambda)$. (For example $\lambda_1 - \lambda \approx 10^{-13}$, $\lambda_i - \lambda \approx 1$ for $i = 2, \dots, n$, and all $c_i \approx 1.0$) On step (ii) we have (without normalization)

$$\mathbf{v}_2 \propto \frac{c_1}{(\lambda_1 - \lambda)^2}\mathbf{x}_1 + \sum_{i=2}^n \frac{c_i}{(\lambda_i - \lambda)^2}\mathbf{x}_i,$$

so that the major component of \mathbf{v}_2 is \mathbf{x}_1 . This algorithm is usually applied when $(\mathbf{A} - \lambda\mathbf{I})^{-1}\mathbf{b}$ is easily calculated (in $\leq \mathcal{O}(n^2)$ operations), say when \mathbf{A} is upper Hessenberg or upper triangular. This can also be used to improve the estimate of the eigenvalue, see §10.10.4.

10.7 Reduction to upper triangular form.

10.7.1 Theorem:

Every $n \times n$ matrix \mathbf{A} is unitary similar to an upper triangular matrix

$$\mathbf{A} = \mathbf{U}\mathbf{T}\mathbf{U}^*, \mathbf{U} \text{ unitary, } \mathbf{T} \text{ upper triangular.}$$

Proof: See Wilkinson (1965). ■

The eigenvalues of \mathbf{T} , namely T_{ii} are those of \mathbf{A} , while the eigenvectors of \mathbf{T} are related to those of \mathbf{A} via \mathbf{U} . The above computation is an infinite process.

10.7.2 Similarity reduction to upper Hessenberg form.

Algorithm:

For $i = 1$ to $n - 2$ do

Householder transform elements in column i below $(i + 1)$ -th row into element $(i + 1, i)$.

This is the same algorithm to reduce a symmetric matrix to symmetric tridiagonal, so it could be done with plane rotations as well. It is a finite procedure.

10.7.3 The L–R algorithm

- (i) $\mathbf{A}^{(0)} = \mathbf{A}$
- (ii) $\mathbf{A}^{(k)} \rightarrow \mathbf{L}^{(k)} \mathbf{R}^{(k)}$ (L–U factorization)
- (iii) $\mathbf{A}^{(k+1)} \leftarrow \mathbf{R}^{(k)} \mathbf{L}^{(k)}$

Note: $\mathbf{A}^{(k+1)} = \mathbf{L}^{(k)-1} \mathbf{A}^{(k)} \mathbf{L}^{(k)}$ so $\mathbf{A}^{(k+1)}$ is similar to $\mathbf{A}^{(k)}$. Under certain conditions on the separation of eigenvalues, $\mathbf{A}^{(k)}$ converges to an upper triangular matrix (or block upper triangular in the case of complex roots). The matrix $\mathbf{L}^{(k)}$ is the product of $n - 1$ elementary matrices so that if $\mathbf{A}^{(k)}$ is upper Hessenberg each factorization can be done in $O(n^2)$ operations, and $\mathbf{A}^{(k+1)}$ is upper Hessenberg. See Rutishauser (1958) for the original, and Rutishauser and Schwarz (1963) for the symmetric case.

10.7.4 Q–R algorithm

This is more stable as it uses the Q–U factorization. See Francis (1961) for the original.

- (i) $\mathbf{A}^{(0)} = \mathbf{A}$
- (ii) $\mathbf{A}^{(k)} \rightarrow \mathbf{Q}^{(k)} \mathbf{R}^{(k)}$
- (iii) $\mathbf{A}^{(k+1)} \rightarrow \mathbf{R}^{(k)} \mathbf{Q}^{(k)}$. ($\mathbf{A}^{(k+1)} = \mathbf{Q}^{(k)t} \mathbf{A}^{(k)} \mathbf{Q}^{(k)}$)

This gives slightly faster convergence to upper triangular form. See Parlett (1968) for the details. Check that if $\mathbf{A}^{(k)}$ is upper Hessenberg then so too is $\mathbf{A}^{(k+1)}$. If $\mathbf{A}^{(k)}$ is upper Hessenberg the matrix $\mathbf{Q}^{(k)}$ is the product of $n - 1$ plane rotations so the factorization can be done in $O(n^2)$ operations.

10.7.5 Single shifts— give faster convergence.

- (i) $\mathbf{A}^{(0)} = \mathbf{A}$
- (ii) $(\mathbf{A}^{(k)} - \lambda_k \mathbf{I}) \rightarrow \mathbf{Q}^{(k)} \mathbf{R}^{(k)}$
- (iii) $\mathbf{A}^{(k+1)} \leftarrow \mathbf{R}^{(k)} \mathbf{Q}^{(k)} + \lambda_k \mathbf{I}$, ($\mathbf{A}^{(k+1)} = \mathbf{Q}^{(k)t} \mathbf{A}^{(k)} \mathbf{Q}^{(k)}$).

At each iteration λ_k is chosen as an eigenvalue of the lower right 2×2 submatrix of $\mathbf{A}^{(k)}$. The iteration produces convergence of the (n, n) element first, which once converged, (the $(n, n - 1)$ element of the upper Hessenberg is close to zero), enables the algorithm to proceed on the remaining order $n - 1$ top left submatrix.

10.7.6 Double Shifts

This method gives even faster convergence and allows the implicit use of complex arithmetic for complex conjugate eigenvalues of real matrices.

- (i) $\mathbf{A}^{(0)} = \mathbf{A}$
- (ii) $(\mathbf{A}^{(k)} - \lambda_k \mathbf{I}) \rightarrow \mathbf{Q}^{(k)} \mathbf{R}^{(k)}$
- (iii) $\mathbf{A}^{(k+\frac{1}{2})} \leftarrow \mathbf{R}^{(k)} \mathbf{Q}^{(k)} + \lambda_k \mathbf{I}$
- (iv) $(\mathbf{A}^{(k+\frac{1}{2})} - \bar{\lambda}_k \mathbf{I}) \rightarrow \mathbf{Q}^{(k+\frac{1}{2})} \mathbf{R}^{(k+\frac{1}{2})}$
- (v) $\mathbf{A}^{(k+1)} \leftarrow \mathbf{R}^{(k+\frac{1}{2})} \mathbf{Q}^{(k+\frac{1}{2})} + \bar{\lambda}_k \mathbf{I}$

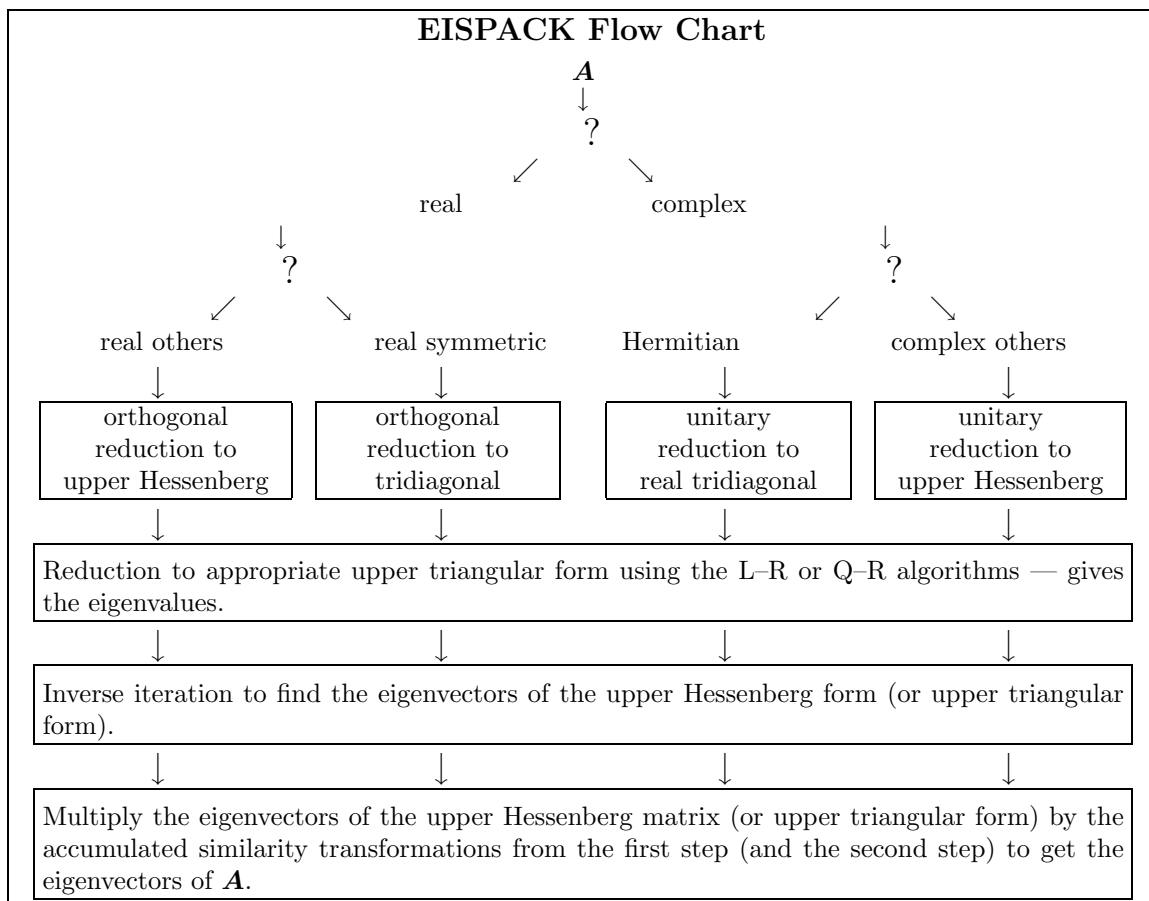
At each iteration λ_k and $\bar{\lambda}_k$ are the two eigenvalues of the lower right 2×2 submatrix of $\mathbf{A}^{(k)}$.

10.8 Balancing a matrix.

The idea here is to try and make the elements of the matrix as close to the same magnitude as possible using a diagonal similarity transformation.

$$\bar{\mathbf{A}} = \mathbf{D} \mathbf{A} \mathbf{D}^{-1}, \quad \mathbf{D} = \text{diag} [d_1, d_2, \dots, d_n],$$

where $d_i = 2^{k_i}$ for some k_i .



10.9 EISPACK.

EISPACK was written over a number of years around 1972 at a cost of some millions of dollars. As well as the published book, Garbow *et al.* (1977), Smith *et al.* (1976), there is machine readable documentation, a short description of the main driver subroutines and a large file with a description of all subroutines. Except for routines which supply specialized output most of the package can be described by the flow chart. If only eigenvalues are required then no accumulation of transformations is necessary. The algorithms are discussed in Wilkinson and Reinsch (1971).

EISPACK has now been superseded by LAPACK, Anderson *et al.* (1995), a combination of EISPACK and LINPACK.

10.10 Other eigenvalue problems.

10.10.1 General linear problem for Hermitian systems.

$$(A - \lambda B)v = 0, \quad B \text{ positive definite}$$

$$(A - \lambda LL^t)v = 0, \text{ where } B = LL^t \text{ (Choleski)}$$

$$(L^{-1}AL^{-t} - \lambda I)(L^t v) = 0.$$

That is, find eigenvalues of the symmetric matrix $L^{-1}AL^{-t}$. If A and B are banded then $L^{-1}AL^{-t}$ has complete fill-in so methods exist (L-Z and Q-Z algorithms) for the simultaneous reduction of A and B to upper triangular form. The L-Z and Q-Z algorithms also treat this problem more efficiently if B is ill-conditioned.

10.10.2 Singular value decomposition

The algorithm, Golub and Kahan (1965), uses an implicit Q-R algorithm, that is, it implicitly finds the square roots of the eigenvalues of, and the eigenvectors of AA^t and A^tA .

10.10.3 Band matrices

There is a method using plane rotations to reduce matrices of this class to tridiagonal form without causing massive “fill-in”. Basically an element on the edge of the band is zeroed, which introduces one non-zero outside the band. This is then “chased off” the matrix by further plane rotations before reducing another edge element to zero.

10.10.4 Nonlinear eigenvalue problems.

The ideas in this section can also be used to improve the estimates of eigenvalues. Suppose a nonlinear matrix function $M(\lambda)$ is given, where M is $n \times n$, and (λ, \mathbf{v}) is required such that

$$M(\lambda)\mathbf{v} = \mathbf{0}, \quad \mathbf{v} \neq \mathbf{0} \quad \text{or} \quad \det |M(\lambda)| = 0.$$

(i) Method (Osborne (1964)).

Consider fixed \mathbf{x} and \mathbf{s} , then the system

$$\begin{aligned} M(\lambda)\mathbf{v} &= \beta(\lambda)\mathbf{x}, \\ \mathbf{s}^t \mathbf{v} &= 1, \end{aligned}$$

defines

$$\beta(\lambda) = \frac{1}{\mathbf{s}^t M^{-1}(\lambda)\mathbf{v}}.$$

Lemma: If $\det |M(\lambda)| \rightarrow 0$ then $\beta(\lambda) \rightarrow 0$ for appropriate \mathbf{s} and \mathbf{x} .

Proof: Compare this with inverse iteration. Let $M(\lambda) = \mathbf{T}\mathbf{D}\mathbf{T}^{-1}$ and suppose that λ is close to λ^* where $M(\lambda^*)$ has its first eigenvalue $d_1(\lambda^*)$ equal to zero. Let $\mathbf{D}(\lambda) = \text{diag}[d_1(\lambda), d_2(\lambda), \dots, d_n(\lambda)]$.

$$M^{-1}(\lambda) = \mathbf{T}\mathbf{D}^{-1}\mathbf{T}^{-1},$$

$$\begin{aligned} \beta(\lambda) &= \frac{1}{(\mathbf{s}^t \mathbf{T})\mathbf{D}^{-1}(\mathbf{T}^{-1}\mathbf{x})} \\ &= \frac{1}{\sum_{i=1}^n d_i^{-1} \langle \mathbf{s}, \kappa_i(\mathbf{T}) \rangle \langle \rho_i(\mathbf{T}^{-1}), \mathbf{x} \rangle} \\ &= \frac{d_1}{\langle \mathbf{s}, \kappa_1(\mathbf{T}) \rangle \langle \rho_1(\mathbf{T}^{-1}), \mathbf{x} \rangle + \sum_{i=2}^n d_1/d_i \langle \cdot, \cdot \rangle}. \end{aligned}$$

Provided \mathbf{s} is not perpendicular to $\kappa_1(\mathbf{T})$ and \mathbf{x} is not perpendicular to $\rho_1(\mathbf{T}^{-1})$, $\beta(\lambda) \rightarrow 0$ as $\det |M(\lambda)| \rightarrow 0$. ■

Applying the Newton iteration to $\beta(\lambda)$ gives,

$$\begin{aligned} \lambda_{i+1} &= \lambda_i - \frac{\beta(\lambda_i)}{\beta'(\lambda_i)} \\ \beta'(\lambda) &= -(\mathbf{s}^t M^{-1}\mathbf{x})^{-2} \left(\mathbf{s}^t \frac{dM^{-1}(\lambda)}{d\lambda} \mathbf{x} \right) \\ &= +(\mathbf{s}^t M^{-1}\mathbf{x})^{-2} \left(\mathbf{s}^t M^{-1} \frac{dM}{d\lambda} M^{-1}\mathbf{x} \right) \\ \therefore \frac{\beta(\lambda)}{\beta'(\lambda)} &= \frac{\mathbf{s}^t M^{-1}(\lambda)\mathbf{x}}{\mathbf{s}^t M^{-1} \frac{dM}{d\lambda} M^{-1}(\lambda)\mathbf{x}} \end{aligned}$$

(ii) Algorithm.

1. $\lambda_0, \mathbf{x}_0, i = 1, \text{eps} = ?$
2. $\mathbf{y}_i = M^{-1}(\lambda_{i-1})\mathbf{x}_{i-1}$
3. $\mathbf{u}_i = M^{-1}(\lambda_{i-1}) \frac{dM}{d\lambda}(\lambda_{i-1})\mathbf{y}_i$
4. Choose $\mathbf{s} = \mathbf{e}_j$, where $|\mathbf{u}_i(j)|$ is the largest element of \mathbf{u}_i .
5. $\lambda_i = \lambda_{i-1} - \mathbf{s}^t \mathbf{y}_i / \mathbf{s}^t \mathbf{u}_i$
 $= \lambda_{i-1} - \mathbf{y}_i(j) / \mathbf{u}_i(j).$

6. $\mathbf{x}_i = \mathbf{u}_i / \mathbf{u}_i(j)$ (renormalize)
7. If $|\beta(\lambda)| = |1/\mathbf{y}_i(j)| < \text{eps}$, stop
else $i = i + 1$, goto 2.

Note: $\beta(\lambda)$ is redefined at each iteration, but in fact this improves convergence as $\mathbf{x}_i \rightarrow$ eigenvector. If $\mathbf{M}(\lambda)$ is symmetric then \mathbf{s} can be chosen as \mathbf{x}_i in step 4.

(iii) **An example:** A large $n \times n$ matrix with a banded structure and a few non-zeros in the top right and bottom left corners (periodic modelling of structures) can be written as $\mathbf{A} + \mathbf{W}\mathbf{W}^t$ where \mathbf{A} is a banded matrix and \mathbf{W} is $n \times m$, where $m \ll n$. The matrix \mathbf{B} is banded, so the large linear eigenvalue problem $(\mathbf{A} + \mathbf{W}\mathbf{W}^t - \lambda\mathbf{B})\mathbf{y} = \mathbf{0}$ can be changed to a small $m \times m$ non-linear problem as follows.

$$\begin{aligned}
 (\mathbf{A} + \mathbf{W}\mathbf{W}^t - \lambda\mathbf{B})\mathbf{y} &= \mathbf{0}. \\
 (\mathbf{A} - \lambda\mathbf{B})\mathbf{y} &= -\mathbf{W}\mathbf{W}^t\mathbf{y} \\
 \mathbf{y} &= -(\mathbf{A} - \lambda\mathbf{B})^{-1}\mathbf{W}\mathbf{W}^t\mathbf{y} \\
 \mathbf{W}^t\mathbf{y} &= -\mathbf{W}^t(\mathbf{A} - \lambda\mathbf{B})^{-1}\mathbf{W}\mathbf{W}^t\mathbf{y} \\
 \underbrace{(\mathbf{I}_m + \mathbf{W}^t(\mathbf{A} - \lambda\mathbf{B})^{-1}\mathbf{W})}_{\mathbf{M}(\lambda)}\mathbf{v} &= \mathbf{0}, \quad \mathbf{v} = \mathbf{W}^t\mathbf{y}.
 \end{aligned}$$

10.10.5 Matrix pencil problems

Here the matrices \mathbf{A} and \mathbf{B} are non-square and values of λ are required so that $\mathbf{A} = \lambda\mathbf{B}$ has less than full rank.

10.11 Exercises.

1. Show that if the matrix \mathbf{A} in the L-R or Q-R algorithm is upper Hessenberg then all iterates $\mathbf{A}^{(k)}$ are upper Hessenberg.
 - (i) Show that an L-U or Q-U factorization of an upper Hessenberg matrix cost $O(n^2)$ operations.
2. Show that the iterates $\mathbf{A}^{(k)}$ are similar in the case of double shifts.
3. If \mathbf{A} and \mathbf{B} are positive definite show that the eigenvectors of $(\mathbf{A} - \lambda\mathbf{B})\mathbf{v} = \mathbf{0}$ are both A-conjugate and B-conjugate.
4. Find $\frac{d\mathbf{M}}{d\lambda}$ if $\mathbf{M}(\lambda) = \mathbf{I} + \mathbf{W}^t(\mathbf{A} - \lambda\mathbf{B})^{-1}\mathbf{W}$. If \mathbf{A} and \mathbf{B} are $n \times n$ tridiagonal matrices and \mathbf{W} is a vector find the number of operations to calculate $\beta(\lambda)/\beta'(\lambda)$.
5. Let \mathbf{A} be positive definite with eigenvalues λ_i , $i = 1, 2, \dots, n$. Show that the eigenvalues of $\mathbf{A} + \mathbf{v}\mathbf{v}^t$ interlace those of \mathbf{A} , that is, between every consecutive pair of λ_i there is an eigenvalue of $\mathbf{A} + \mathbf{v}\mathbf{v}^t$.
6. Find the eigenvalues and eigenvectors of the symmetric tridiagonal matrix which has diagonal elements α and off-diagonal elements β . This has connections with the Discrete (Fast) Fourier Transform (see Briggs and Henson (1995)) because of the eigenvector structure.

Hint: Consider the associated second order difference equation with appropriate homogeneous boundary conditions.
7. Show that the non-symmetric tridiagonal matrix \mathbf{A} can be diagonally similarity transformed to a symmetric tridiagonal matrix if $\gamma_i\beta_i > 0$, $i = 2, \dots, n$ as in §10.5.1.

References

- Anderson E., Bai Z., Bischof C., Demmel J., Dongarra J., Du Croz J., Greenbaum A., Hammarling S., McKenney A., Ostrouchov S. and Sorenson D., *LAPACK Users' Guide, Second Edition*, SIAM, Philadelphia, 1995.
- Anton H. and Rorres C., *Elementary Linear Algebra with Applications*, Wiley, 1987.
- Atkinson K., *An Introduction to Numerical Analysis*, 2nd Ed., Wiley, 1989.
- Ben-Israel A. and Greville T.N.E., *Generalised Inverses: Theory and Applications*, Wiley, 1974.
- Björck A., Iterative refinement of linear least squares solutions II. *BIT*, **8**, pp8-30, 1968.
- Blackford L.S., Choi J., Cleary E., D'Azevedo E., Demmel J., Dhillon I., Dongarra J., Hammarling S., Henry G., Petitet A., Stanley K., Walker D. and Whaley R.C., *ScaLAPACK User's Guide*, SIAM, Philadelphia, 1997.
- J. Boothroyd., *Australian Computer Journal*, 1(1), 1968.
- Briggs W.L. and Van Emden Henson, *The DFT: An Owner's Manual for the Discrete Fourier Transform*, SIAM, Philadelphia, 1995.
- Dahlquist, G., Björck A. and Anderson N., *Numerical Methods*, Prentice-Hall, 1974.
- Demmel J.W., *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- Dongarra J.J., Bunch J.R., Moler C.B. and Stewart G.W., *LINPACK User's Guide*, SIAM, Philadelphia, 1979.
- Dongarra J.J., Duff I. and Sorenson D., *Solving Linear Systems on Vector and Shared Memory Computers*, SIAM, Philadelphia, 1990.
- Duff I., A survey of sparse matrix research, *Proceedings IEEE*, **65**, pp500-535, 1977.
- Forsythe G.E. and Moler C.B., *Computer Solution of Linear Algebraic Systems*, Prentice Hall Series in Automatic Computation, Englewood Cliffs, N.J., 1967.
- Francis J., The Q-R transformation. A unitary analogue to the LR transformation. *Computer Journal*, **4**, pp 265-271, 1961.
- Froberg C.E., *Introduction to Numerical Analysis*, Addison Wesley, 1965.
- Garbow B.S., Boyle J.M., Dongarra J.J. and Moler C.B., *Matrix Eigensystem Routines: EISPACK Guide Extension*, Lecture Notes in Computer Science, Springer-Verlag, 1977.
- George A. and Liu J., *Computer Solution of Large Sparse Positive Definite Systems*, Prentice-Hall, 1981.
- Golub, G.H., Iterative Refinement of least squares solutions, *Numerische Mathematik*, **9**, 139-148, 1966.
- Golub G.H. and Kahan W., Calculating the singular values and pseudo inverse of matrices. *SIAM Journal Numerical Analysis*, **2**, 205-224, 1965.
- Golub G.W. and van Loan C.F., *Matrix Computations*, John Hopkins University Press, Baltimore, Edition 2, 1989.
- Hestenes, M.R. and Steifel, E., Methods of Conjugate Gradients for Solving Linear Systems, *National Bureau Standards Report*, No. 1659, 1952.
- Hoffman A.J. and Wielandt H.W., The variation of the spectrum of a normal matrix, *Duke Mathematical Journal*, **20**, pp 37-39, 1953.
- Householder A., *The Theory of Matrices in Numerical Analysis*, Blaisdell, New York, 1964.
- Jennings L.S., Simultaneous Equation Estimation, Computational Aspects, *Journal of Econometrics*, **12**, pp 23-39, 1980.
- Jennings, L.S. and Osborne, M.R., A direct error analysis for least squares, *Numerische Mathematik*, **22**, 325-332, 1974.
- Kahan W., Numerical linear algebra, *Canadian Mathematics Bulletin*, **9**, pp 756-801, 1966.
- Lawson C.L. and Hanson R.J., *Solving Least Squares Problems*, Prentice-Hall, 1974.
- Lipschutz S., *Linear Algebra*, Schaum outline series, McGraw Hill, 1968.
- Moore E.H., On the reciprocal of the general algebraic matrix, *Bulletin American Mathematical Society*, **26**, pp394-395, 1920.
- Noble B., *Applied Linear Algebra*, Prentice-Hall, Englewood Cliffs, N.J., 1969.
- Noble B. and Daniel J., *Applied Linear Algebra*, Prentice-Hall, Englewood Cliffs, N.J., 1977.

- Osborne M.R., A new method for the solution of eigenvalue problems. *Computer Journal*, **7**, pp 358–362, 1964
- Parlett B.N., Global convergence of the basic QR algorithm on Hessenberg matrices, *Mathematics of Computation*, **22**, pp 803–817, 1968.
- Parlett B.N., *The Symmetric Eigenvalue Problem*, Prentice-Hall, 1981.
- Penrose R., On best approximate solutions of linear matrix equations, *Proceedings of the Cambridge Philosophical Society*, **52**, pp 17–19, 1956.
- Ralston A., *A First Course in Numerical Analysis*, McGraw-Hill, 1965.
- Rutishauser H., Solution of eigenvalue problems with the LR transformation. *Applied Mathematics Series, National Bureau Standards*, **49**, pp 47–81, 1958
- Rutishauser H., Once again: The least squares problem. *Linear Algebra and Its Applications*, **1**, pp 479–488, 1968.
- Rutishauser H. and Schwarz H.R., Handbook Series Linear Algebra. The LR transformation method for symmetric matrices. *Numerische Mathematik*, **5**, pp 273–289, 1963.
- Smith B.T., Boyle J., Garbow B., Ikebe Y., Klema V. and Moler C., Matrix Eigensystem Routines—EISPACK Guide, 2nd Ed., Lecture Notes in Computer Science, **6**, Springer Verlag, New York, 1976.
- Stewart G.W., *Introduction to Matrix Computations*, Academic Press, New York, 1973.
- Strang Gilbert, *Linear Algebra and its Applications*, 3rd ed., Harcourt Brace Jovanovich, 1988.
- Trefethen L.N. and Bau D., *Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- Varga R.S., *Matrix Iterative Analysis*, Prentice-Hall, 1962.
- Wilkinson J.H., Rounding Errors in Algebraic Processes, Prentice-Hall, Englewood Cliffs, N.J., 1963.
- Wilkinson J.H., *The Algebraic Eigenvalue Problem*, Clarendon Press, 1965.
- Wilkinson J.H. and Reinsch C., *Handbook for Automatic Computation*, Volumes I and II, Springer, 1971.
- Young D., *Iterative Solution of Large Linear Systems*, Academic Press, New York, 1971.