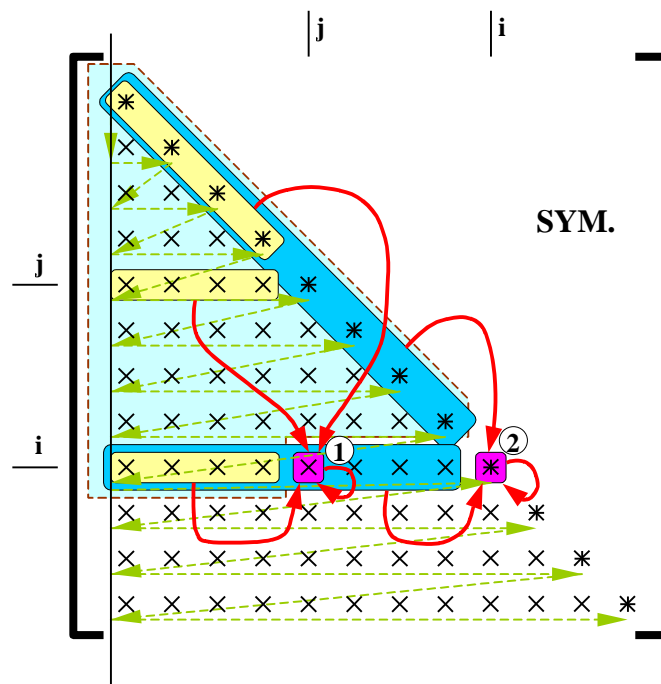


Matrix Decompositions With Implementation Remarks

For use in Nonlinear Optimization Software

Igor Grešovnik
May 2012

Electronic Book
Revision 4.2
(revision 1: August 2008)



This report is public.

Contents:

1	Introduction.....	1
2	Basics of Matrix Theory.....	2
2.1	Basic terms.....	2
2.2	Polar Decomposition	4
2.3	Eigenvalues and eigenvectors.....	5
2.4	Spectral theorem	7
2.5	Similarity transformations, Schur Decomposition & Spectral Decomposition	7
2.5.1	Schur Decomposition theorem.....	8
2.6	Singular Value Decomposition.....	10
2.7	Examples in physics	11
2.7.1	Finite deformation tensor	12
3	LU Decomposition.....	13
4	LDL^T Decomposition.....	15
4.1	Introduction: General on Matrix Products	15
4.1.1	Eigenvalues and Eigenvectors	17
4.1.2	General (asymmetric) matrices	18
4.1.3	Symmetric Real Matrices.....	18
4.2	LDLT Decomposition.....	19
4.2.1	LDLT Decomposition Using Lower Triangle.....	21
4.2.2	LDLT Decomposition Using Only Upper Triangle:	23
5	Cholesky Decomposition.....	25
6	Gram-Schmidt (GS) Orthogonalization	28
6.1.1	Basic Idea.....	28
6.2	Modified Gram-Schmidt Orthogonalization and QR Factorization	30
6.2.2	Remark: Extended QR decomposition	32
6.2.3	Solution of Systems of Equations with QR decomposition	33
6.3	Non-standard factorization by using the GS	35
6.3.1	Orthogonalization with normalization	38
6.3.2	Solution of Systems of Equations with Gram-Schmidt Orthogonalisation (non-standard form).....	39
6.3.3	Implementation of GS (non-standard form).....	40
7	Overdetermined Systems	41
7.1	Orthogonal methods.....	42
8	Special Systems of Equations	44
8.1	Orthogonal systems.....	44
8.2	Systems with Matrices that have Orthogonal Columns.....	45
8.2.1	Inverse of a Matrix with Orthogonal Columns	45
8.2.2	Component-wise verification.....	46

8.3	Systems of Equations with Matrix with Orthogonal Columns	48
8.4	Triangular Systems	49
8.4.1	Upper Triangular Systems	49
8.4.2	Lower Triangular Systems	50
9	<i>Number of Operatinos for Standard Matrix Operations</i>	<i>51</i>
9.1	Basic Operations	51
9.1.1	Matrix and vector multiplications:	51
9.2	Special systems of equations.....	51
9.3	Factorizations	52

1 INTRODUCTION

This document introduces the basics of matrix computations and especially matrix decompositions, which are also used in the Investigative Optimization Library (IOptLib).

The document was created by putting together a number of notes in MS Word created by the myself for memorizing a number of results from linear algebra that are relevant for optimization and other fields of numerical analysis.

I have created a large portion of these notes during my undergraduate study when I have practicing numerical software development for training, hobby and to earn some additional money that allowed me to cover more than just the basic needs, e.g. a bicycle or a one week trip to Paris. Another package of the notes was created when finishing the graduate degree on physics, when I have also put some hand written notes form the “Mathematics 1” subject into electronic form.

The collection was further supplemented during my Ph.D. study at the University of Swansea in the U.K. Then, one day when looking at all those scattered documents that contained short notes on individual topics, I realized that it would be much easier for myself to have all the material in a single document. At that time, the MS Word has also advanced enough to allow creation of longer documents with lots of formulas and images¹. So I have reserved a weekend and put all the documents in a single files, arranged titles and references and formatted the document in such a way that it acquired a relatively readable form.

After that, I was still adding notes from time to time when I needed something to memorize, and this was much easier to do in a settled document than in a group of scattered files. However, many portions of the script are not in the form I could recommend for studying. It must be understood that some of these notes emerged from yellowed old notes which represented just a very condensed hints on individual topics, and there are many portions that I didn’t touch (except for formatting) since they were put in electronic forms. But when I need to look back at something and I notice that this is very awkwardly written, I usually take some time to correct that. In this way I hope that the script will one day become a useful and easily digestible source of reference information for people like myself, who are not truly experts in linear algebra but use it in everyday life.

Author of this script, *Igor Grešounik*

¹ For a long time, at least over the decade, this was the major obstacle for writing scientific and engineering documents in MS Word, and this is probably the main reason that Tech remained the desktop publishing program of choice in many scientific and not so rare engineering environments. Because for people who did not have problems with learning programming languages and other computer related skills, it was far easier to work with a program where you couldn’t see the result of your work instantly than with a program that crashed for nothing every couple of minutes.

2 BASICS OF MATRIX THEORY

2.1 Basic terms

2.1.1.1 Inner product

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

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

2.1.1.2 Matrix basic terms

A *matrix* is a rectangular table of numbers or other quantities that can be added or multiplied. Most common is use of matrices that are defined over the real (real matrices) or complex (complex matrices) field. We will usually denote matrices by large bold letters, where dimensions (the number of rows and columns) will sometimes be specified in subscript, e.g. $\mathbf{A}_{m \times n}$ denotes a matrix with m rows and n columns. A matrix where one of the dimensions equals 1 is often called a vector, and is interpreted as an element of a coordinate space. A $1 \times n$ matrix is called a row vector and a $n \times 1$ matrix is called a column vector. Matrix components will usually be denoted by the corresponding letter, but not written in bold, with indices in the subscript, e.g. A_{ij} will mean the element in row i and column j of the matrix \mathbf{A} . Sometimes components will be denoted by small letters, i.e. a_{ij} . A *square matrix* is those that has the same number of rows and columns.

Matrix is *normal* if it commutes with its conjugate transpose:

$$\mathbf{A}^* \mathbf{A} = \mathbf{A} \mathbf{A}^* . \quad (2)$$

Adjoint matrix of a n -by- m matrix \mathbf{A} is its conjugate transpose:

$$\left(\mathbf{A}^* \right)_{i,j} = \overline{\mathbf{A}_{j,i}} \quad (3)$$

where subscripts denote the i,j -th component of the matrix, with $1 \leq i \leq n$ and $1 \leq j \leq m$. This is expressed as

$$\mathbf{A}^* = \overline{\mathbf{A}}^T. \quad (4)$$

The *conjugate transpose* of a matrix or its *adjoint matrix* is also denoted by \mathbf{A}^H (H coming from *Hermitian conjugate*)

2.1.1.3 Hermitian and unitary matrices

Hermitian or *self-adjoint matrix* is a square matrix that is equal to its conjugate transpose, i.e.

$$\mathbf{A} = \mathbf{A}^*. \quad (5)$$

Hermitian matrices are *normal*, and the finite dimensional spectral theorem applies, which means that every Hermitian matrix can be diagonalized by a unitary matrix. Eigenvalues of every Hermitian matrix are real and eigenvectors with distinct eigenvalues are orthogonal.

The sum of two Hermitian matrices is Hermitian, and inverse of an invertible Hermitian matrix is also Hermitian. The product of two Hermitian matrices is Hermitian only if they commute.

Anti-Hermitian or skew Hermitian matrix is a matrix for which $\mathbf{A}^* = -\mathbf{A}$. Entries in the main diagonal are pure imaginary. The same is true for eigenvalues.

Unitary matrix is a square matrix that satisfies the condition

$$\mathbf{U}^* \mathbf{U} = \mathbf{U} \mathbf{U}^* = \mathbf{I}, \quad (6)$$

where \mathbf{I} is the identity matrix. A matrix is unitary if it has an inverse which is equal to its conjugate transpose. Unitary matrix preserves the standard inner product on \mathbb{C}^n ,

$$\langle \mathbf{U}\mathbf{x}, \mathbf{U}\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle \quad (7)$$

If \mathbf{A} is a square matrix then the following conditions are equivalent:

1. \mathbf{A} is unitary
2. \mathbf{A}^* is unitary
3. The columns of \mathbf{A} form an *orthonormal basis* of \mathbb{C}^n with respect to the standard inner product on \mathbb{C}^n .

4. \mathbf{A} is an *isometry* (i.e. distance preserving isomorphism¹) with respect to the norm for this inner product.

A unitary matrices are called special if its determinant is 1.

Unitary matrices are *normal*, therefore the spectral theorem applies to them. Every unitary matrix \mathbf{U} has decomposition of the form

$$\mathbf{U} = \mathbf{V}\mathbf{\Sigma}\mathbf{V}^*, \quad (8)$$

where \mathbf{V} is unitary and $\mathbf{\Sigma}$ is diagonal and unitary.

Real unitary matrices are orthogonal matrices and real hermitian matrices are symmetric matrices.

2.2 Polar Decomposition

Polar decomposition is a matrix decomposition of the form

$$\mathbf{A} = \mathbf{U}\mathbf{P}, \quad (9)$$

where \mathbf{U} is an unitary matrix and \mathbf{P} is a positive-semidefinite Hermitian matrix. The decomposition always exists, and if \mathbf{A} is invertible then the decomposition is unique and \mathbf{P} positive definite. The matrix \mathbf{P} is given by

$$\mathbf{P} = \sqrt{\mathbf{A}^* \mathbf{A}}, \quad (10)$$

where \mathbf{A}^* is a conjugate transpose of \mathbf{A} . A positive definite Hermitian matrix has a unique positive square root. The matrix \mathbf{U} is then given as

$$\mathbf{U} = \mathbf{A}\mathbf{P}^{-1}. \quad (11)$$

In terms of the singular value decomposition, we have

$$\begin{aligned} \mathbf{P} &= \mathbf{V}\mathbf{\Sigma}\mathbf{V}^* \\ \mathbf{U} &= \mathbf{W}\mathbf{V}^* \end{aligned}, \quad (12)$$

¹ Isomorphism is a bijective map f such that both f and f^{-1} are *homomorphisms* (i.e. structure preserving maps - $\phi(u \bullet v) = \phi(u) \circ \phi(v)$)

which confirms that \mathbf{P} is positive-semidefinite and \mathbf{U} unitary.

Determinant of \mathbf{A} can be expressed, according to (9), as

$$\det \mathbf{A} = \det \mathbf{P} \det \mathbf{U} = r e^{i\theta}, \quad (13)$$

which gives the corresponding polar decomposition of the determinant of \mathbf{A} , since $\det \mathbf{P} = r = |\det \mathbf{A}|$ and $\det \mathbf{U} = e^{i\theta}$.

A matrix can also be decomposed as

$$\mathbf{A} = \mathbf{P}' \mathbf{U}, \quad (14)$$

where \mathbf{U} is the same as before and

$$\mathbf{P}' = \mathbf{U} \mathbf{P} \mathbf{U}^{-1} = \sqrt{\mathbf{A} \mathbf{A}^*} = \mathbf{W} \mathbf{\Sigma} \mathbf{W}^* \quad (15)$$

The matrix \mathbf{A} is normal if and only if $\mathbf{P}' = \mathbf{P}$. Then $\mathbf{U} \mathbf{\Sigma} = \mathbf{\Sigma} \mathbf{U}$ and it is possible to diagonalize \mathbf{U} with a unitary similarity matrix \mathbf{S} that commutes with $\mathbf{\Sigma}$, giving $\mathbf{S} \mathbf{U} \mathbf{S}^* = \mathbf{\Phi}^{-1}$, where $\mathbf{\Phi}$ is a diagonal unitary matrix of phase $e^{i\phi}$. Putting, $\mathbf{Q} = \mathbf{V} \mathbf{S}^*$, the polar decomposition can be re-written as

$$\mathbf{A} = (\mathbf{Q} \mathbf{\Phi} \mathbf{Q}^*) (\mathbf{Q} \mathbf{\Sigma} \mathbf{Q}^*), \quad (16)$$

so \mathbf{A} also has a spectral decomposition

$$\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^* \quad (17)$$

with complex eigenvalues so that $\mathbf{\Lambda} \mathbf{\Lambda}^* = \mathbf{\Sigma}^2$ and a unitary matrix of complex eigenvectors \mathbf{Q} .

2.3 Eigenvalues and eigenvectors

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$. $\lambda \in \mathbb{C}$ is called an eigenvalue of \mathbf{A} if there exists a non-zero (non-null) vector $\mathbf{x} \in \mathbb{C}^n$ such that $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$. \mathbf{x} is an *eigenvector* associated with λ . The set of eigenvalues of \mathbf{A} is called the *spectrum* of \mathbf{A} denoted by $\sigma(\mathbf{A})$. \mathbf{x} and \mathbf{y} are the *left and right eigenvectors* of \mathbf{A} associated with λ , respectively, if

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}, \quad \mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H. \quad (18)$$

The eigenvalue corresponding to the eigenvector \mathbf{x} can be determined by computing the *Raileigh quotient* $\lambda = \mathbf{x}^H \mathbf{A} \mathbf{x} / (\mathbf{x}^H \mathbf{x})$. Eigenvalue λ is solution of a *characteristic equation*

$$p_A(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I}) = 0, \quad (19)$$

where $p_A(\lambda)$ is the *characteristic polynomial*. Because this is a polynomial of degree n , there exist n eigenvalues of \mathbf{A} that are not necessarily distinct. The following is true:

$$\begin{aligned} \det(\mathbf{A}) &= \prod_{i=1}^n \lambda_i, \quad \text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i \\ \sigma(\mathbf{A}) &= \sigma(\mathbf{A}^T), \quad \sigma(\mathbf{A}^H) = \sigma(\bar{\mathbf{A}}) \end{aligned} \quad (20)$$

From the first relation we see that a matrix is *singular* if it has at least one zero eigenvalue, since $p_A(0) = \det \mathbf{A} = \prod_{i=1}^n \lambda_i$.

If \mathbf{A} has *real entries*, then coefficients of $p_A(\lambda)$ are real and therefore *complex eigenvalues occur in complex conjugate pairs*.

The maximum module (absolute value) of the eigenvalues of \mathbf{A} is called the *spectral radius* of \mathbf{A} and denoted by

$$\rho(\mathbf{A}) = \max_{\lambda \in \sigma(\mathbf{A})} |\lambda| \quad (21)$$

λ is an eigenvalue of \mathbf{A} iff $\bar{\lambda}$ is an eigenvalue of \mathbf{A}^H . Therefore $\rho(\mathbf{A}) = \rho(\mathbf{A}^H)$. Also $\rho(\alpha \mathbf{A}) = |\alpha| \rho(\mathbf{A})$ and $\rho(\mathbf{A}^k) = (\rho(\mathbf{A}))^k \quad \forall k \in \mathbb{N}$

Let \mathbf{A} be a block triangular matrix

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1k} \\ \mathbf{0} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2k} \\ \vdots & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{A}_{kk} \end{bmatrix}. \quad (22)$$

Because $p_A(\lambda) = p_{A_{11}}(\lambda) + p_{A_{22}}(\lambda) + \dots + p_{A_{kk}}(\lambda)$, the spectrum of \mathbf{A} is union of the spectra of each diagonal block. As a consequence, if \mathbf{A} is triangular then its eigenvalues are its diagonal elements.

For each eigenvalue of matrix \mathbf{A} the set of eigenvectors associated with λ , together with the null vector, identifies a subspace in \mathbb{R}^n which is called the *eigenspace* associated with λ and corresponds by definition to $\ker(\mathbf{A} - \lambda\mathbf{I})$. The dimension of the eigenspace is

$$\dim(\ker(\mathbf{A} - \lambda\mathbf{I})) = n - \text{rank}(\mathbf{A} - \lambda\mathbf{I}), \quad (23)$$

and is called *geometric multiplicity* of the eigenvalue λ . It can not be greater than the *algebraic multiplicity* of λ , which is the multiplicity of λ as a root of the characteristic polynomial. Eigenvalues that have geometric multiplicity strictly less than algebraic multiplicity are called *defective*.

The eigenspace associated with an eigenvalue of the matrix \mathbf{A} is invariant with respect to \mathbf{A} in the sense of the following definition:

A subspace $S \subseteq \mathbb{R}^n$ is called invariant with respect to a square matrix \mathbf{A} if $\mathbf{A}S \subseteq S$, where $\mathbf{A}S$ is S transformed through \mathbf{A} .

2.4 Spectral theorem

2.5 Similarity transformations, Schur Decomposition & Spectral Decomposition

Similarity transformation is transformation of the form

$$\mathbf{A} \rightarrow \mathbf{C}^{-1} \mathbf{A} \mathbf{C} \quad (24)$$

where \mathbf{C} is a square nonsingular matrix having the same order as \mathbf{A} . We say that matrixes \mathbf{A} and $\mathbf{C}^{-1} \mathbf{A} \mathbf{C}$ are similar matrixes. If \mathbf{C} is unitary then matrixes are unitary similar. Two similar matrixes have the same spectrum and the same characteristic polynomial.

We can easily check that if (λ, \mathbf{x}) are eigenvalue-eigenvector pair for \mathbf{A} then $(\lambda, \mathbf{C}^{-1} \mathbf{x})$ are eigenvalue-eigenvector pair for $\mathbf{C}^{-1} \mathbf{A} \mathbf{C}$ since $(\mathbf{C}^{-1} \mathbf{A} \mathbf{C}) \mathbf{C}^{-1} \mathbf{x} = \mathbf{C}^{-1} \mathbf{A} \mathbf{x} = \lambda \mathbf{C}^{-1} \mathbf{x}$

Matrices \mathbf{AB} and \mathbf{BA} , $\mathbf{A} \in \mathbb{C}^{n \times m}$ and $\mathbf{B} \in \mathbb{C}^{m \times n}$ are not similar, but satisfy the following property:

$$\sigma(\mathbf{AB}) \setminus \{0\} = \sigma(\mathbf{BA}) \setminus \{0\}, \quad (25)$$

i.e. \mathbf{AB} and \mathbf{BA} share the same spectrum apart from null eigenvalues, so that $\rho(\mathbf{AB}) = \rho(\mathbf{BA})$.

2.5.1 Schur Decomposition theorem

Given $\mathbf{A} \in \mathbb{C}^{n \times n}$, there exists an *unitary* matrix \mathbf{U} such that

$$\mathbf{U}^{-1} \mathbf{A} \mathbf{U} = \mathbf{U}^H \mathbf{A} \mathbf{U} = \begin{bmatrix} \lambda_1 & b_{12} & \cdots & b_{1n} \\ 0 & \lambda_2 & & b_{2n} \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & \lambda_n \end{bmatrix} = \mathbf{T} \quad (26)$$

where λ_i are eigenvalues of \mathbf{A} .

It follows that every square matrix \mathbf{A} is unitary similar to an upper triangular matrix. Matrices \mathbf{T} and \mathbf{U} are not necessarily unique.

Among the others, the Schur decomposition theorem gives rise to the following results:

Every Hermitian matrix is unitary similar to a diagonal real matrix. Every Schur decomposition of a Hermitian matrix is diagonal:

$$\mathbf{A}^H = \mathbf{A} \Rightarrow \mathbf{U}^{-1} \mathbf{A} \mathbf{U} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \quad (27)$$

It turns that $\mathbf{A} \mathbf{U} = \mathbf{U} \Lambda$, i.e. $\mathbf{A} \mathbf{u}_i = \lambda_i \mathbf{u}_i$ for $i=1, \dots, n$, i.e. column vectors of \mathbf{U} are eigenvectors of \mathbf{A} . Since eigenvectors are mutually orthogonal, a hermitean matrix has a set of orthogonal vectors that generate the whole space \mathbb{C}^n . Furthermore, it can be shown that a matrix \mathbf{A} of order n is similar to a diagonal matrix \mathbf{D} iff the eigenvectors of \mathbf{A} form a basis in \mathbb{C}^n .

Spectral decomposition (or eigendecomposition) of a normal matrix:

A matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is normal iff it is unitary similar to a diagonal matrix. As a consequence, a normal matrix admits the following *spectral decomposition*:

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^H, \quad (28)$$

with \mathbf{U} being unitary and $\mathbf{\Lambda}$ diagonal.

Let \mathbf{A} and \mathbf{B} be two normal and commutative matrices. Then the generic eigenvalue μ_i of $\mathbf{A} + \mathbf{B}$ is given by the sum $\mu_i = \lambda_i + \xi_i$, where λ_i and μ_i are eigenvalues of \mathbf{A} and \mathbf{B} associated with the same eigenvector.

There are nonsymmetric matrices that are unitary similar to diagonal matrices, but they are not unitary similar.

2.5.1.1 Canonical Jordan form

Schur decomposition can be improved as follows.

Let \mathbf{A} be any square matrix. Then, there exists a nonsingular matrix \mathbf{X} which transforms \mathbf{A} into a block diagonal matrix \mathbf{J} such that

$$\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = \mathbf{J} = \text{diag}(\mathbf{J}_{k_1}(\lambda_1), \mathbf{J}_{k_2}(\lambda_2), \dots, \mathbf{J}_{k_n}(\lambda_n)), \quad (29)$$

which is called a *Jordan canonical form*. λ_j are eigenvalues of \mathbf{A} and $\mathbf{J}_k(\lambda) \in \mathbb{C}^{k \times k}$ is a Jordan block of the form $\mathbf{J}_1(\lambda) = \lambda$ for $k=1$ and

$$\mathbf{J}_k(\lambda) = \begin{bmatrix} \lambda & 1 & 0 & \dots & 0 \\ 0 & \lambda & 1 & \vdots & \vdots \\ \vdots & \ddots & \ddots & 1 & 0 \\ \vdots & & \ddots & \lambda & 1 \\ 0 & \dots & \dots & 0 & \lambda \end{bmatrix}. \quad (30)$$

If an eigenvalue is defective, the size of the corresponding Jordan block is greater than one. Therefore, a matrix can be diagonalized by a similarity transform iff it is nondefective. Nondefective matrices are therefore called *diagonalizable*. Normal matrices are diagonalizable.

2.6 Singular Value Decomposition

The Singular value decomposition (SVD) can be seen as generalization of the spectral theorem to arbitrary matrices (not necessarily square). For a m -by- n matrix over the field K of real or complex numbers \mathbf{A} there exists a factorization of the form

$$\mathbf{A}_{(m \times n)} = \mathbf{U}_{(m \times m)} \mathbf{\Sigma}_{(m \times n)} \mathbf{V}_{(n \times n)}^T, \quad (31)$$

$$\mathbf{U}^T \mathbf{U} = \mathbf{I}; \mathbf{V}^T \mathbf{V} = \mathbf{I}; \Sigma_{ii} = \sigma_i; \Sigma_{ij} = 0; i \neq j$$

where \mathbf{U} is a m -by- m unitary matrix over K , $\mathbf{\Sigma}$ is m -by- n trapezoid matrix with nonnegative numbers on the diagonal and zeros off the diagonal and \mathbf{V} is a n -by- n unitary matrix over K . Such decomposition is called a *singular-value decomposition* of \mathbf{M} . A common convention is to order the values $\sigma_{ii} = \Sigma_{ii}$ in non-increasing fashion. In this case, the diagonal matrix is uniquely determined by \mathbf{M} (but the matrices \mathbf{U} and \mathbf{V} are not).

The matrix \mathbf{V} contains a set of orthonormal “input” vectors for “analysis” basis vector directions for \mathbf{A} , the matrix \mathbf{U} contains a set of orthonormal “output” basis vector directions for \mathbf{A} , the matrix $\mathbf{\Sigma}$ contains the singular values, which can be understood as scalar “gain controls” by which each corresponding input is multiplied to give a corresponding output.

A non-negative number σ is a *singular value* of \mathbf{A} only if there exists unit-length vectors \mathbf{u} in K^m and \mathbf{v} in K^n such that

$$\mathbf{A} \mathbf{v} = \sigma \mathbf{u} \text{ and } \mathbf{A}^* \mathbf{u} = \sigma \mathbf{v} . \quad (32)$$

The vectors \mathbf{u} and \mathbf{v} are called **left-singular** and **right-singular vectors** for σ , respectively. In any singular value decomposition, the diagonal entries of $\mathbf{\Sigma}$ are equal to the singular values of \mathbf{A} . The columns of \mathbf{U} and \mathbf{V} are left and right singular vectors for the corresponding singular values. The theorem states that

- An $m \times n$ matrix \mathbf{A} has at most $p = \min(m, n)$ distinct singular values.
- It is always possible to find a unitary basis for K^m consisting of left-singular vectors of \mathbf{A} .

- It is always possible to find a unitary basis for k^n consisting of right-singular vectors for \mathbf{A} .

A singular value for which we can find two left (or right) singular vectors that are not linearly dependent is called *degenerate*. Non-degenerate singular values always have unique left and right singular vectors, up to multiplication by a unit phase factor $e^{i\phi}$ (for the real case, up to sign). Consequently, if all singular values of M are non-degenerate and non-zero, then its singular value decomposition is unique, up to multiplication of a column of \mathbf{U} by a unit phase factor and simultaneous multiplication of the corresponding column of \mathbf{V} by the same unit phase factor.

Degenerate singular values have non-unique singular vectors (by definition). If \mathbf{u}_1 and \mathbf{u}_2 are two left-singular vectors which both correspond to the singular value σ , then any normalized linear combination of the two vectors is also a left singular vector corresponding to σ (and similar for the right singular vectors). Singular decomposition is then not unique.

When \mathbf{A} is a Hermitian matrix which is positive semi-definite (all eigenvalues are real and non-negative), then the singular values and singular vectors of \mathbf{A} coincide with the eigenvalues and eigenvectors of \mathbf{M} ,

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^* . \quad (33)$$

More generally, given a SVD of \mathbf{M} , the following two relations hold:

$$\begin{aligned} \mathbf{A}^* \mathbf{A} &= \mathbf{V} \mathbf{\Sigma}^* \mathbf{U}^* \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* = \mathbf{V} (\mathbf{\Sigma}^* \mathbf{\Sigma}) \mathbf{V}^* \\ \mathbf{A} \mathbf{A}^* &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* \mathbf{V} \mathbf{\Sigma}^* \mathbf{U}^* = \mathbf{U} (\mathbf{\Sigma} \mathbf{\Sigma}^*) \mathbf{U}^* \end{aligned} \quad (34)$$

The right hand sides of these relations describe the eigenvalue decompositions of the left hand sides. Consequently, the squares of the non-zero singular values of \mathbf{A} are equal to the non-zero eigenvalues of $\mathbf{A}^* \mathbf{A}$. The columns of \mathbf{U} (left singular vectors) are eigenvectors for $\mathbf{A} \mathbf{A}^*$ and the columns of \mathbf{V} (right singular vectors) are eigenvectors of $\mathbf{A}^* \mathbf{A}$.

2.7 Examples in physics

2.7.1 Finite deformation tensor

Finite deformation tensors are used when the deformation of a body is sufficiently large that the assumptions in small strain theory are not valid.

We denote \mathbf{X} the position vector of a particle in the initial (undeformed) state of a body relative to some coordinate basis. The position of the particle in the deformed state is denoted \mathbf{x} . If $d\mathbf{X}$ is a line segment that joins two nearby particles in the undeformed state and $d\mathbf{x}$ is the line segment joining the same two particles in the deformed state, then the linear transformation between the two segments is given by

$$d\mathbf{x} = \mathbf{F} d\mathbf{X} . \quad (35)$$

The quantity \mathbf{F} is called the *deformation gradient*, and is given by:

$$\mathbf{F} = \nabla_{\mathbf{X}} \mathbf{x} = \frac{\partial \mathbf{x}}{\partial \mathbf{X}} , \quad (36)$$

or:

$$F_{ij} = \frac{\partial x_i}{\partial X_j} . \quad (37)$$

\mathbf{F} is a second order tensor and contains information about the stretch and rotation of the body.

Polar decomposition:

$$\mathbf{F} = \mathbf{R}\mathbf{U} = \mathbf{V}\mathbf{R} , \quad (38)$$

where \mathbf{R} is a proper orthogonal tensor representing rotation, and \mathbf{U} and \mathbf{V} are positive definite symmetric tensors that represent stretches. \mathbf{U} is called the right stretch and \mathbf{V} is called the left stretch tensor.

Spectral decompositions of \mathbf{U} and \mathbf{V} are

$$\mathbf{U} = \sum_{i=1}^3 \lambda_i \mathbf{N}_i \otimes \mathbf{N}_i \quad (39)$$

and

$$\mathbf{V} = \sum_{i=1}^3 \lambda_i \mathbf{n}_i \otimes \mathbf{n}_i , \quad (40)$$

where λ_i are the principal stretches, and \mathbf{N}_i and \mathbf{n}_i are the directions of the principal stretches (principal directions). The principal directions are related by

$$\mathbf{n}_i = \mathbf{R} \mathbf{N}_i . \quad (41)$$

Rotation independent deformation measures are introduced because rotation should not induce any stress in a deformable body. Rotation is excluded by multiplying \mathbf{R} by its transpose. In this way we obtain the **Right Cauchy-Green Tensor**

$$\begin{aligned} \mathbf{C} &= \mathbf{F}^T \mathbf{F} = \mathbf{U}^T \mathbf{U} \\ C_{ij} &= \sum_{k=1}^3 \frac{\partial x_k}{\partial X_i} \frac{\partial x_k}{\partial X_j} \end{aligned} \quad (42)$$

and the left Cauchy-Green tensor

$$\begin{aligned} \mathbf{B} &= \mathbf{F} \mathbf{F}^T = \mathbf{V} \mathbf{V}^T = \mathbf{V}^2 \\ B_{ij} &= \sum_{k=1}^3 \frac{\partial x_i}{\partial X_k} \frac{\partial x_j}{\partial X_k} . \end{aligned} \quad (43)$$

The spectral decompositions are

$$\begin{aligned} \mathbf{C} &= \sum_{i=1}^3 \lambda_i^2 \mathbf{N}_i \otimes \mathbf{N}_i \\ \mathbf{B} &= \sum_{i=1}^3 \lambda_i^2 \mathbf{n}_i \otimes \mathbf{n}_i \end{aligned} \quad (44)$$

3 LU DECOMPOSITION

\mathbf{L} – lower triangular

\mathbf{U} – upper triangular with 1 on diagonal

$$L_{ij} = \begin{cases} L_{ij}; & i > j \\ 1; & i = j \\ 0; & i < j \end{cases} \quad (3.0)$$

$$U_{ij} = \begin{cases} U_{ij}; & i \leq j \\ 0; & i > j \end{cases} \quad (3.0)$$

$$[\mathbf{LU}]_{ij} = \sum_{k=1}^n L_{ik} U_{kj}$$

$$\left\{ \begin{array}{l} i < j: \quad \underbrace{\sum_{k=1}^{i-1} L_{ik} U_{kj}}_{k < i < j} + \underbrace{L_{ii}}_{=1} U_{ij} + \underbrace{\sum_{k=i+1}^{j-1} L_{ik} U_{kj}}_{i < k < j} + \underbrace{L_{ij}}_{=0} U_{jj} + \underbrace{\sum_{k=j+1}^n L_{ik} U_{kj}}_{i < j < k} \\ i = j: \quad \underbrace{\sum_{k=1}^{i-1} L_{ik} U_{kj}}_{k < i = j} + \underbrace{L_{ii}}_{=0} U_{ii} + \underbrace{\sum_{k=i+1}^n L_{ik} U_{kj}}_{i = j < k} \\ i > j: \quad \underbrace{\sum_{k=1}^{j-1} L_{ik} U_{kj}}_{k < j < i} + \underbrace{L_{ij}}_{=0} U_{jj} + \underbrace{\sum_{k=j+1}^{i-1} L_{ik} U_{kj}}_{j < k < i} + \underbrace{L_{ii}}_{=0} U_{ij} + \underbrace{\sum_{k=i+1}^n L_{ik} U_{kj}}_{j < i < k} \end{array} \right.$$

$$[\mathbf{LU}]_{ij} = \begin{cases} i < j: & \sum_{k=1}^i L_{ik} U_{kj} \\ i = j: & \sum_{k=1}^{i-1} L_{ik} U_{kj} \\ i > j: & \sum_{k=1}^{j-1} L_{ik} U_{kj} \end{cases} \quad (3.45)$$

4 LDL^T DECOMPOSITION

4.1 Introduction: General on Matrix Products

$$[\mathbf{LU}]_{ij} = \sum_{k=1}^n L_{ik} U_{kj}$$

$$\left\{ \begin{array}{l} i < j: \underbrace{\sum_{k=1}^{i-1} L_{ik} U_{kj}}_{k < i < j} + L_{ii} U_{ij} + \underbrace{\sum_{k=i+1}^{j-1} \underbrace{L_{ik}}_{=0} U_{kj}}_{i < k < j} + L_{ij} U_{jj} + \underbrace{\sum_{k=j+1}^n \underbrace{L_{ik}}_{=0} \underbrace{U_{kj}}_{=0}}_{i < j < k} \\ i = j: \underbrace{\sum_{k=1}^{i-1} L_{ik} U_{kj}}_{k < i = j} + L_{ii} U_{ii} + \underbrace{\sum_{k=i+1}^n \underbrace{L_{ik}}_{=0} \underbrace{U_{kj}}_{=0}}_{i = j < k} \\ i > j: \underbrace{\sum_{k=1}^{j-1} L_{ik} U_{kj}}_{k < j < i} + L_{ij} U_{jj} + \underbrace{\sum_{k=j+1}^{i-1} L_{ik} U_{kj}}_{j < k < i} + \underbrace{L_{ii} U_{ij}}_{=0} + \underbrace{\sum_{k=i+1}^n \underbrace{L_{ik}}_{=0} \underbrace{U_{kj}}_{=0}}_{j < i < k} \end{array} \right.$$

$$[\mathbf{LU}]_{ij} = \sum_{k=1}^{\min(i,j)} L_{ik} U_{kj} \quad (4.0)$$

$$\begin{aligned}
a_{ij} &= [\mathbf{LDL}^T]_{ij} = \sum_{k=1}^n l_{ik} [\mathbf{DL}^T]_{kj} = \sum_{k=1}^n l_{ik} [\mathbf{L}^T]_{kj} d_{kk} = \sum_{k=1}^n l_{ik} l_{jk} d_{kk} = \\
&= \sum_{k=1}^{\min(i,j)} l_{ik} l_{jk} d_{kk} = \begin{cases} i < j: & \sum_{k=1}^i l_{ik} l_{jk} d_{kk} \\ i = j: & \sum_{k=1}^{i-1} l_{ik} l_{ik} d_{kk} + \underbrace{l_{ii}^2}_{=1} d_{ii} \\ i > j: & \sum_{k=1}^j l_{ik} l_{jk} d_{kk} \end{cases} = \begin{cases} i < j: & \sum_{k=1}^{i-1} l_{ik} l_{jk} d_{kk} + \underbrace{l_{ii} l_{ji}}_{=1} d_{ii} \\ i = j: & \sum_{k=1}^{i-1} l_{ik}^2 d_{kk} + \underbrace{l_{ii}^2}_{=1} d_{ii} \\ i > j: & \sum_{k=1}^{j-1} l_{ik} l_{jk} d_{kk} + \underbrace{l_{ij} l_{jj}}_{=1} d_{jj} \end{cases} \\
i < j: & a_{ij} = \sum_{k=1}^{i-1} l_{ik} l_{jk} d_{kk} + \underbrace{l_{ii} l_{ji}}_{=1} d_{ii} \quad a_{ji} = \sum_{k=1}^{i-1} l_{jk} l_{ik} d_{kk} + \underbrace{l_{ji} l_{ii}}_{=1} d_{ii} \\
\Rightarrow & a_{ij} = a_{ji} \Rightarrow \mathbf{A} = \mathbf{A}^T
\end{aligned}$$

- LDL^T decomposition is possible **only for symmetric matrices**.

Let \mathbf{A} be any square matrix and \mathbf{D} any diagonal matrix. Then

$$[\mathbf{AD}]_{ij} = A_{ij} D_{jj}, \quad (4.0.)$$

$$[\mathbf{DA}]_{ij} = A_{ij} D_{ii} \quad (4.0.)$$

and

$$[\mathbf{D}^2]_{ij} = [\mathbf{DD}^T]_{ij} = [\mathbf{D}^T \mathbf{D}]_{ij} = d_{ij}^2$$

$$[\mathbf{AA}^T]_{ij} = \sum_{k=1}^n a_{ik} [\mathbf{A}^T]_{kj} = \sum_{k=1}^n a_{ik} a_{jk}$$

$$[\mathbf{A}^T \mathbf{A}]_{ij} = \sum_{k=1}^n [\mathbf{A}^T]_{jk} a_{ki} = \sum_{k=1}^n a_{ki} a_{kj}$$

$$(\mathbf{AA}^T)^T = (\mathbf{A}^T)^T \mathbf{A}^T = \mathbf{AA}^T$$

from which follows that \mathbf{AA}^T is symmetric.

4.1.1 Eigenvalues and Eigenvectors

Let's have a square matrix \mathbf{A} and let's there exist some non-zero vector \mathbf{v} so that

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}, \quad (4.0)$$

where λ is a scalar. Then \mathbf{v} is an eigenvector and λ an eigenvalue of matrix \mathbf{A} . We also see that if \mathbf{v} is an eigenvector of \mathbf{A} , then also $k\mathbf{v}$ is eigenvector of \mathbf{A} , where k is any scalar different than zero. An invertible square matrix has n eigenvectors where n is matrix dimension (i.e. number of rows or columns):

$$\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i, \quad i = 1, \dots, n. \quad (4.0)$$

Positive definiteness:

Matrix \mathbf{A} is positive-definite, if

$$\mathbf{u}^T \mathbf{A} \mathbf{u} > 0 \quad \forall \mathbf{u} \neq 0. \quad (4.0)$$

Matrix is positive-semidefinite, if

$$\mathbf{u}^T \mathbf{A} \mathbf{u} \geq 0 \quad \forall \mathbf{u} \neq 0, \quad (4.0)$$

negative-definite, if

$$\mathbf{u}^T \mathbf{A} \mathbf{u} < 0 \quad \forall \mathbf{u} \neq 0, \quad (4.0)$$

negative-semidefinite if the relationship includes allows the equality sign, and indefinite otherwise.

Matrix is positive-definite, if and only if all its eigenvalues are greater than zero (which can be seen if the matrix components are written in its eigensystem, i.e. a coordinate system the basis of which form matrix eigenvectors).

If a square matrix \mathbf{A} is invertible (i.e. has a full rank), then matrices $\mathbf{A}\mathbf{A}^T$, $\mathbf{A}^T\mathbf{A}$ and \mathbf{A}^2 are all positive definite.

Let's denote such matrix \mathbf{B} . The relation (4.0) can be verified if we write vector \mathbf{u} as a linear combination of eigenvectors of \mathbf{A} .

4.1.2 General (asymmetric) matrices

4.1.3 Symmetric Real Matrices

For symmetric matrices,

$$\mathbf{A} = \mathbf{A}^T \quad (4.0.)$$

We can perform the LDLT decomposition of such a matrix to a product of a lower triangular matrix \mathbf{L} , diagonal matrix \mathbf{D} and transpose of \mathbf{L} :

$$\mathbf{A} = \mathbf{LDL}^T \quad (4.0.)$$

where \mathbf{L} is a general lower triangular matrix with unit diagonal elements,

$$\mathbf{L}_{ij} = \begin{cases} l_{ij}; & i \geq j \\ 0; & i < j \end{cases}, \quad (4.0.)$$

and \mathbf{D} is diagonal:

$$\mathbf{D}_{ij} = \begin{cases} d_{ij}; & i = j \\ 0; & i \neq j \end{cases} \quad (4.0.)$$

$$\begin{aligned}
a_{ij} &= [\mathbf{LDL}^T]_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} [\mathbf{DL}^T]_{kj} = \sum_{k=1}^{\min(i,j)} l_{ik} [\mathbf{L}^T]_{kj} d_{kk} = \sum_{k=1}^{\min(i,j)} l_{ik} l_{jk} d_{kk} = \\
&= \left\{ \begin{array}{l} i < j: \sum_{k=1}^{i-1} l_{ik} l_{jk} d_{kk} + l_{ii} l_{ji} d_{ii} \\ i = j: \sum_{k=1}^{i-1} l_{ik}^2 d_{kk} + l_{ii}^2 d_{ii} \\ i > j: \sum_{k=1}^{j-1} l_{ik} l_{jk} d_{kk} + l_{ij} l_{jj} d_{jj} \end{array} \right\} \\
i < j: a_{ij} &= \sum_{k=1}^i l_{ik} l_{jk} d_{kk} \quad a_{ji} = \sum_{k=1}^i l_{jk} l_{ik} d_{kk} \\
\Rightarrow a_{ij} &= a_{ji} \Rightarrow \mathbf{A} = \mathbf{A}^T
\end{aligned}$$

- LDL^T decomposition is possible **only for symmetric matrices**.

$$a_{11} = d_{11}$$

$$a_{12} = l_{21} d_{11} + l_{21} d_{11}$$

$$a_{21} = l_{21} d_{11} + l_{21} d_{11}$$

4.2 LDLT Decomposition

$$\mathbf{A} = \mathbf{A}^T \quad (4.0)$$

We can decompose symmetric matrices in the form

$$\mathbf{A} = \mathbf{LDL}^T \quad (4.0)$$

where \mathbf{L} is a lower triangular matrix with unit diagonal elements,

$$\mathbf{L} = \begin{cases} L_{ij}; & i > j \\ 1; & i = j \\ 0; & i < j \end{cases}, \quad (4.0)$$

and \mathbf{D} is diagonal:

$$\mathbf{D}_{ij} = \begin{cases} D_{ij}; & i = j \\ 0; & i \neq j \end{cases} \quad (4.0)$$

$$\begin{aligned} a_{ij} &= [\mathbf{LDL}^T]_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} [\mathbf{DL}^T]_{kj} = \sum_{k=1}^{\min(i,j)} l_{ik} [\mathbf{L}^T]_{kj} d_{kk} = \sum_{k=1}^{\min(i,j)} l_{ik} l_{jk} d_{kk} = \\ &= \begin{cases} i \geq j: \sum_{k=1}^j l_{ik} l_{jk} d_{kk} \\ i < j: a_{ji} \end{cases} \end{aligned}$$

$$l_{ii} = 1 \quad \forall i$$

$$i < j: a_{ij} = \sum_{k=1}^i l_{ik} l_{jk} d_{kk} \quad a_{ji} = \sum_{k=1}^i l_{jk} l_{ik} d_{kk}$$

$$\Rightarrow a_{ij} = a_{ji} \Rightarrow \mathbf{A} = \mathbf{A}^T$$

$$a_{11} = \underbrace{l_{11} l_{11}}_{=1} d_{11}$$

$$a_{21} = \underbrace{l_{21} l_{11}}_{=1} d_{11}$$

$$a_{22} = \underbrace{l_{21} l_{21}}_{=1} d_{11} + \underbrace{l_{22} l_{22}}_{=1} d_{22}$$

$$a_{31} = \underbrace{l_{31} l_{11}}_{=1} d_{11}$$

$$a_{32} = \underbrace{l_{31} l_{21}}_{=1} d_{11} + \underbrace{l_{32} l_{22}}_{=1} d_{22}$$

$$a_{33} = \underbrace{l_{31} l_{31}}_{=1} d_{11} + \underbrace{l_{32} l_{32}}_{=1} d_{22} + \underbrace{l_{33} l_{33}}_{=1} d_{33}$$

$$a_{41} = \underbrace{l_{41} l_{11}}_{=1} d_{11}$$

$$a_{42} = \underbrace{l_{41} l_{21}}_{=1} d_{11} + \underbrace{l_{42} l_{22}}_{=1} d_{22}$$

$$a_{43} = \underbrace{l_{41} l_{31}}_{=1} d_{11} + \underbrace{l_{42} l_{32}}_{=1} d_{22} + \underbrace{l_{43} l_{33}}_{=1} d_{33}$$

$$a_{44} = \underbrace{l_{41} l_{41}}_{=1} d_{11} + \underbrace{l_{42} l_{42}}_{=1} d_{22} + \underbrace{l_{43} l_{43}}_{=1} d_{33} + \underbrace{l_{44} l_{44}}_{=1} d_{44}$$

4.2.1 LDLT Decomposition Using Lower Triangle

A general formula is then

$$\begin{aligned}
 & i = 1, \dots, n : \\
 & l_{ij} = \frac{1}{d_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} d_{kk} \right), \quad j = 1, \dots, i-1 \\
 & d_{ii} = a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 d_{kk}
 \end{aligned} \tag{4.0}$$

Instead by lines we can also evaluate factors by columns:

$$\begin{aligned}
 & i = 1, \dots, n : \\
 & d_{ii} = a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 d_{kk} \\
 & l_{ji} = \frac{1}{d_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} l_{ik} l_{jk} d_{kk} \right), \quad j = i+1, i+2, \dots, n
 \end{aligned} \tag{4.0}$$

(note that $a_{ij} = a_{ji}$ since \mathbf{A} is symmetric).

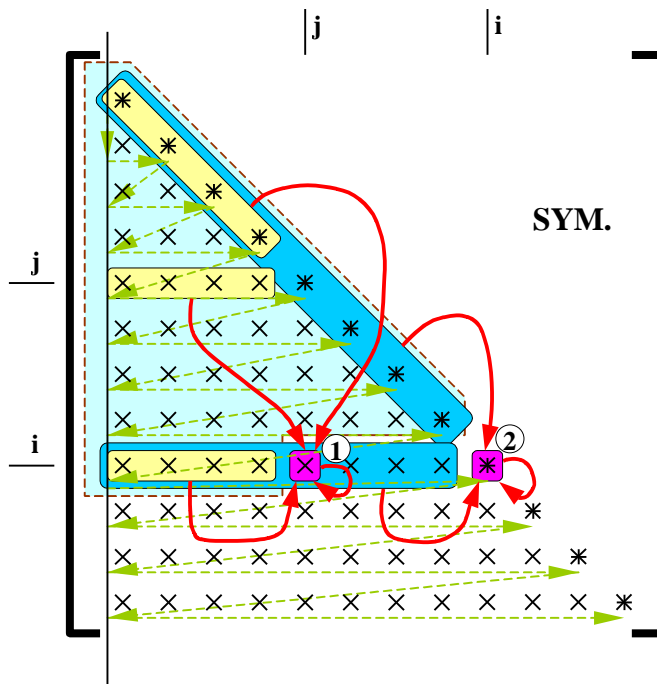


Figure 4.1: Evaluation of elements of \mathbf{D} and \mathbf{L} (i.e. the lower triangle) of the LDLT decomposition from lower triangle of the original matrix according to equation (4.0).

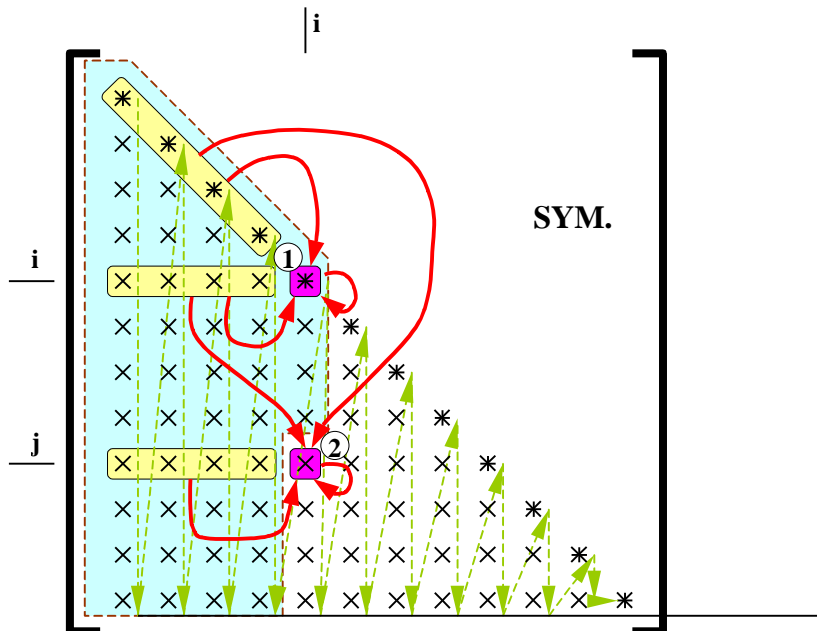


Figure 4.2: Evaluation of elements of \mathbf{D} and \mathbf{L} (i.e. the lower triangle) of the LDLT decomposition from lower triangle of the original matrix according to equation (4.0).

4.2.2 LDLT Decomposition Using Only Upper Triangle:

Equations can be obtained from (4.0 and (4.0 taking into account relations $a_{ij} = a_{ji}$ and $u_{ij} = l_{ji}$, where $\mathbf{U} = \mathbf{L}^T$ is transpose of the lower triangular factor (i.e the upper triangular factor) of the decomposition. Formulas are then (derived from (4.0):

$$\begin{aligned}
 & i = 1, \dots, n: \\
 & u_{ji} = \frac{1}{d_{jj}} \left(a_{ji} - \sum_{k=1}^{j-1} u_{ki} u_{kj} d_{kk} \right), \quad j = 1, \dots, i-1 \\
 & d_{ii} = a_{ii} - \sum_{k=1}^{i-1} u_{ki}^2 d_{kk}
 \end{aligned} \tag{4.0}$$

or (derived from (4.0):

$$\begin{aligned}
 & i = 1, \dots, n: \\
 & d_{ii} = a_{ii} - \sum_{k=1}^{i-1} u_{ki}^2 d_{kk} \\
 & u_{ij} = \frac{1}{d_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} u_{ki} u_{kj} d_{kk} \right), \quad j = i+1, i+2, \dots, n
 \end{aligned} \tag{4.0}$$

Evaluation order for terms in both equations is shown graphically in Figure 4.4 and Figure 6.1.

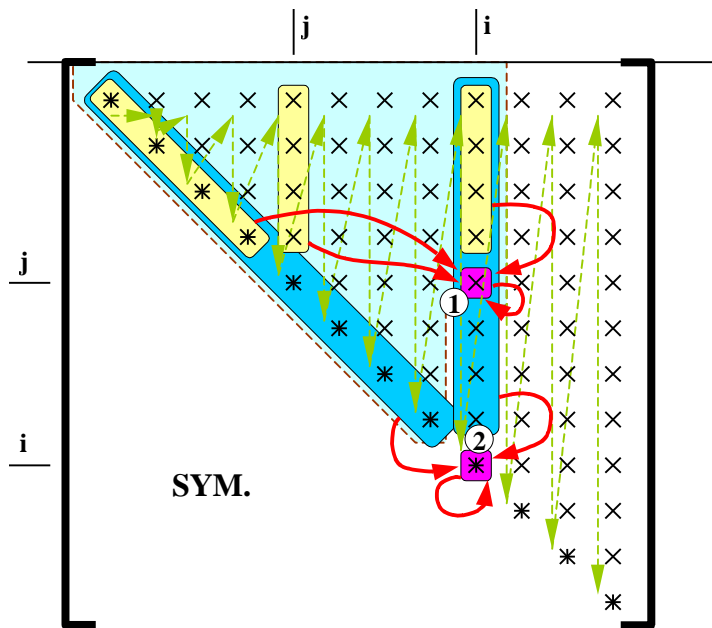


Figure 4.3: Evaluation of elements of \mathbf{D} and \mathbf{L}^T (i.e. the upper triangle) of the LDLT decomposition from upper triangle of the original matrix according to equation (4.0).

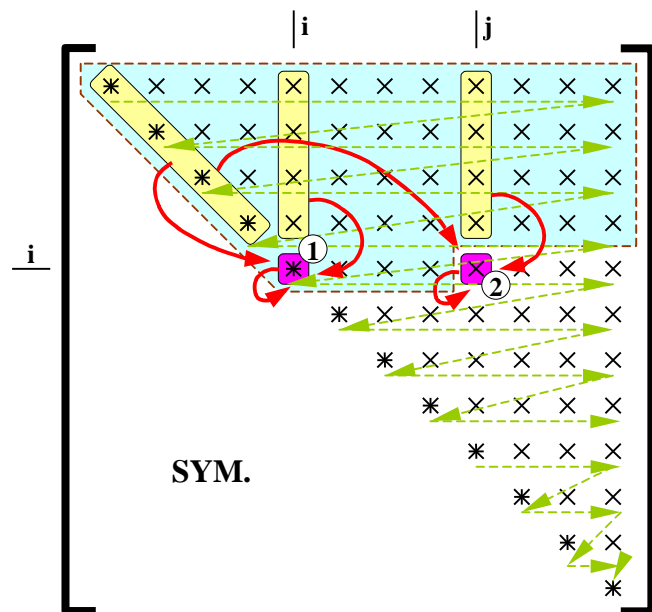


Figure 4.4: Evaluation of elements of \mathbf{D} and \mathbf{L}^T (i.e. the upper triangle) of the LDLT decomposition from upper triangle of the original matrix according to equation (4.0).

5 CHOLESKY DECOMPOSITION

This decomposition is of form

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T, \quad (5.0)$$

where \mathbf{L} is a lower triangular matrix

$$\mathbf{L} = \begin{cases} l_{ij}; & i \geq j \\ 0; & i < j \end{cases}. \quad (5.0)$$

Note that diagonal elements of \mathbf{L} are in general different than 1, unlike in the LDL^T decomposition. \mathbf{A} must be a symmetric positive definite matrix.

$$\begin{aligned} a_{ij} &= [\mathbf{L}\mathbf{L}^T]_{ij} = \sum_{k=1}^n l_{ik} [\mathbf{L}^T]_{kj} = \sum_{k=1}^i l_{ik} [\mathbf{L}^T]_{kj} = \sum_{k=1}^i l_{ik} l_{jk} = \sum_{k=1}^{\min(i,j)} l_{ik} l_{jk} = \\ &= \sum_{k=1}^{\min(i,j)} l_{ik} l_{jk} = \\ &= \begin{cases} i \geq j: \sum_{k=1}^j l_{ik} l_{jk} \\ i < j: a_{ji} \end{cases} \end{aligned}$$

$$a_{11} = \underline{\underline{l_{11}l_{11}}}$$

$$a_{21} = \underline{\underline{l_{21}l_{11}}}$$

$$a_{22} = l_{21}l_{21} + \underline{\underline{l_{22}l_{22}}}$$

$$a_{31} = \underline{\underline{l_{31}l_{11}}}$$

$$a_{32} = l_{31}l_{21} + \underline{\underline{l_{32}l_{22}}}$$

$$a_{33} = l_{31}l_{31} + l_{32}l_{32} + \underline{\underline{l_{33}l_{33}}}$$

$$a_{41} = \underline{\underline{l_{41}l_{11}}}$$

$$a_{42} = l_{41}l_{21} + \underline{\underline{l_{42}l_{22}}}$$

$$a_{43} = l_{41}l_{31} + l_{42}l_{32} + \underline{\underline{l_{43}l_{33}}}$$

$$a_{44} = l_{41}l_{41} + l_{42}l_{42} + l_{43}l_{43} + \underline{\underline{l_{44}l_{44}}}$$

Computing the LL^T factors:

$$l_{11} = \sqrt{a_{11}}$$

$$l_{21} = \frac{1}{l_{11}}a_{21}$$

$$l_{22} = \sqrt{a_{22} - l_{21}l_{21}}$$

$$l_{31} = \frac{1}{l_{11}}a_{31}$$

$$l_{32} = \frac{1}{l_{22}}(a_{32} - l_{31}l_{21})$$

$$l_{33} = \sqrt{(a_{33} - l_{31}l_{31} - l_{32}l_{32})}$$

$$l_{41} = \frac{1}{l_{11}}a_{41}$$

$$l_{42} = \frac{1}{l_{22}}(a_{42} - l_{41}l_{21})$$

$$l_{43} = \frac{1}{l_{33}}(a_{43} - l_{41}l_{31} - l_{42}l_{32})$$

$$l_{44} = \sqrt{a_{44} - l_{41}l_{41} - l_{42}l_{42} - l_{43}l_{43}}$$

General formula:

$$\begin{aligned}
 & i = 1, \dots, n : \\
 & l_{ij} = \frac{1}{l_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} \right), \quad j = 1, \dots, i-1 \\
 & l_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2}
 \end{aligned} \tag{5.0}$$

In this case we evaluate elements of \mathbf{L} by rows, i.e. in the following order (for a 5×5 matrix):

$$\begin{bmatrix}
 1:l_{11} & 0 & 0 & 0 & 0 \\
 2:l_{21} & 3:l_{22} & 0 & 0 & 0 \\
 4:l_{31} & 5:l_{32} & 6:l_{33} & 0 & 0 \\
 7:l_{41} & 8:l_{42} & 9:l_{43} & 10:l_{44} & 0 \\
 11:l_{51} & 12:l_{52} & 13:l_{53} & 14:l_{54} & 15:l_{55}
 \end{bmatrix}$$

Instead of by rows, we can also evaluate factors by columns:

$$\begin{aligned}
 & i = 1, \dots, n : \\
 & l_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2} \\
 & l_{ji} = \frac{1}{l_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} l_{ik} l_{jk} \right), \quad j = i+1, i+2, \dots, n
 \end{aligned} \tag{5.46}$$

Evaluation order is then the following for a 5×5 matrix):

$$\begin{bmatrix}
 1:l_{11} & 0 & 0 & 0 & 0 \\
 2:l_{21} & 6:l_{22} & 0 & 0 & 0 \\
 3:l_{31} & 7:l_{32} & 10:l_{33} & 0 & 0 \\
 4:l_{41} & 8:l_{42} & 11:l_{43} & 13:l_{44} & 0 \\
 5:l_{51} & 9:l_{52} & 12:l_{53} & 14:l_{54} & 15:l_{55}
 \end{bmatrix}$$

Outline of an algorithm for computing the LL^T factors:

A nice property of the algorithm is that we don't need to store the original matrix throughout the algorithm. When a specific element of \mathbf{L} is computed, it can immediately replace the corresponding element of \mathbf{A} since this is not needed any more. The algorithm assumes that diagonal

and below diagonal elements of \mathbf{A} are stored in the input matrix, while elements beyond diagonal do not need to be stored because \mathbf{A} is symmetric.

6 GRAM-SCHMIDT (GS) ORTHOGONALIZATION

In this Section, a procedure of obtaining a set of mutually orthogonal (or orthonormal) vectors from an arbitrary set of linearly independent vectors by the process called *Gram-Schmidt orthogonalization* is described. This process can be used for orthogonal decomposition of a matrix, which includes the known QR factorization described in Section 6.2, where a matrix is factorized as a product of an orthogonal and upper triangular matrix. This can be used for calculating eigenvalues and eigenvectors of the matrix (by application of an iterative procedure known as the QR iteration) or for solution of systems of equations, which is a method of choice when the system matrix is ill conditioned.

After the QR factorization, a non-standard form of factorization (67) is described in Section 6.3, where one obtains an orthogonal matrix (whose columns are orthogonalized columns of the original matrix) expressed as a product of the original and upper triangular matrix. The factorization will not be used in practice for solution of equations (although this is possible, as described in Section 6.3.2), but is described for instructive purposes. The reader can skip this section without harm.

We have vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$, where $m \leq n$ (n is the dimension of the vector space). We want to construct m mutually orthogonal vectors \mathbf{q}_i which are obtained from the original set of vectors, i.e. are linear combinations of vectors \mathbf{v}_i .

6.1.1 Basic Idea

Let us have vectors \mathbf{q}_1 and \mathbf{v}_2 . We want to construct a vector \mathbf{q}_2 which will be a linear combination of these two vectors and will be orthogonal to the vector \mathbf{q}_1 . The most obvious way is to subtract from \mathbf{v}_2 its orthogonal projection on \mathbf{q}_1 (see Figure 6.1). This projection has direction of \mathbf{q}_1 and size

$$\|\mathbf{v}_{2\mathbf{q}_1}\|_2 = \|\mathbf{v}_2\|_2 \cos \phi = \frac{\langle \mathbf{q}_1, \mathbf{v}_2 \rangle}{\|\mathbf{q}_1\|_2}, \quad (6.0)$$

since

$$\langle \mathbf{q}_1, \mathbf{v}_2 \rangle = \|\mathbf{q}_1\|_2 \|\mathbf{v}_2\|_2 \cos \phi.$$

The projection is then

$$\mathbf{v}_{2\mathbf{q}_1} = \|\mathbf{v}_{2\mathbf{q}_1}\|_2 \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|_2} = \frac{\langle \mathbf{q}_1, \mathbf{v}_2 \rangle}{\|\mathbf{q}_1\|_2^2} \mathbf{q}_1 = \frac{\langle \mathbf{q}_1, \mathbf{v}_2 \rangle}{\langle \mathbf{q}_1, \mathbf{q}_1 \rangle} \mathbf{q}_1 \quad (6.0.)$$

and the constructed orthogonal vector \mathbf{q}_2 is

$$\mathbf{q}_2 = \mathbf{v}_2 - \mathbf{v}_{2\mathbf{q}_1} = \mathbf{v}_2 - \frac{\langle \mathbf{v}_2, \mathbf{q}_1 \rangle}{\langle \mathbf{q}_1, \mathbf{q}_1 \rangle} \mathbf{q}_1 \quad (6.0.)$$

It can be easily shown that this vector is orthogonal to \mathbf{q}_1 : dot product of the above equation with \mathbf{q}_1 yields

$$\langle \mathbf{q}_2, \mathbf{q}_1 \rangle = \left\langle \mathbf{v}_2 - \frac{\langle \mathbf{v}_2, \mathbf{q}_1 \rangle}{\langle \mathbf{q}_1, \mathbf{q}_1 \rangle} \mathbf{q}_1, \mathbf{q}_1 \right\rangle = \langle \mathbf{v}_2, \mathbf{q}_1 \rangle - \frac{\langle \mathbf{v}_2, \mathbf{q}_1 \rangle}{\langle \mathbf{q}_1, \mathbf{q}_1 \rangle} \langle \mathbf{q}_1, \mathbf{q}_1 \rangle = 0.$$

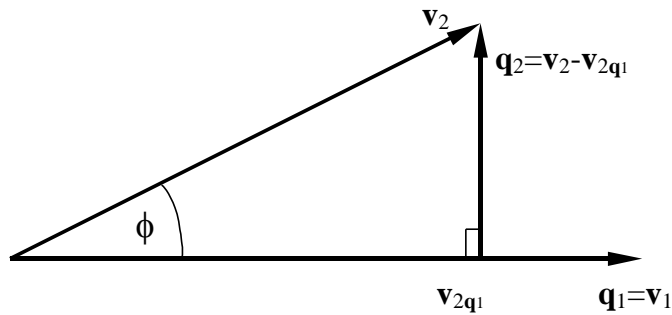


Figure 6.1: Construction of the next orthogonal vector.

The construction described above can be easily extended for a set of n vectors. We set

$$\underline{\underline{\mathbf{q}_1 = \mathbf{v}_1}} \quad (6.0.)$$

and the subsequent orthogonal vectors are formed as

$$\underline{\underline{\mathbf{q}_j = \mathbf{v}_j - \sum_{l=1}^{j-1} \frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \mathbf{q}_l}}, \quad j = 2, 3, \dots, n. \quad (6.0.)$$

We can see by total induction that each \mathbf{q}_j generated in this way is a linear combination of vectors $\mathbf{v}_k, k \leq j$ and that it is orthogonal to all $\mathbf{q}_k, k < j$. At the end of the procedure we therefore obtain a set of m mutually orthogonal vectors \mathbf{q}_i , which are all linear combinations of original vectors and are non-zero if the original vectors are linearly independent.

6.2 Modified Gram-Schmidt Orthogonalization and QR Factorization

The Gram-Schmidt orthogonalization procedure described by (6.0) is numerically less stable because the generated vectors tend to lose their linear independency because of the numerical errors. A modified procedure is therefore used in practice, which is quite similar to the original one, but the order is a bit different.

We will use a bit different notation and will denote the original vectors by $\mathbf{a}_i, i = 1, \dots, n$. We want to obtain an orthonormal basis of vectors \mathbf{q}_i from \mathbf{a}_i such that

$$\begin{aligned} \langle \mathbf{q}_i, \mathbf{q}_i \rangle &= 1 \quad \forall i = 1, \dots, n \\ \langle \mathbf{q}_i, \mathbf{q}_j \rangle &= 0 \quad \forall i \neq j \end{aligned} \quad (6.0)$$

The original procedure is performed by the following steps:

$$\begin{aligned} \mathbf{q}_1 &= \frac{\mathbf{a}_1}{\|\mathbf{a}_1\|} \\ k &= 2, \dots, n : \\ \tilde{\mathbf{q}}_k &= \mathbf{a}_k - \sum_{i=1}^{k-1} \langle \mathbf{q}_i, \mathbf{a}_k \rangle \mathbf{q}_i, \\ \mathbf{q}_k &= \frac{\tilde{\mathbf{q}}_k}{\|\tilde{\mathbf{q}}_k\|_2} \end{aligned} \quad (6.0)$$

In the first part of each step, a new vector is produced from subtracting from the original vector all orthogonal projections of this vector to already computed orthogonal and normal vectors. In this way a vector is produced that is normal to all previously calculated vector, and in the second step this vector is normalized with respect to the Euclidean norm. If vectors were not normalized, then we should perform division by $\langle \mathbf{q}_i, \mathbf{q}_i \rangle$ in order to obtain the correct projection.

The **modified Gram-Schmidt procedure** is performed in such a way that projections are gradually subtracted from the original vector and such reduced vectors are projected on already obtained orthonormal vectors:

$$\begin{aligned}
 \mathbf{q}_1 &= \frac{\mathbf{a}_1}{\|\mathbf{a}_1\|}; \\
 j &= 2, \dots, n: \\
 \mathbf{q}_j^{(0)} &= \mathbf{a}_j \\
 k &= 1, \dots, j-1 \\
 \mathbf{q}_j^{(k)} &= \mathbf{q}_j^{(k-1)} - \langle \mathbf{q}_j^{(k-1)}, \mathbf{q}_k \rangle \mathbf{q}_k \\
 \mathbf{q}_j &= \frac{\mathbf{q}_j^{(j-1)}}{\|\mathbf{q}_j^{(j-1)}\|_2}
 \end{aligned} \tag{6.47}$$

Vectors obtained by this process are the same since, taking into account orthonormality of vectors \mathbf{q}_i , we have

$$\begin{aligned}
 \mathbf{q}_k^{k-1} &= \mathbf{a}_k - \langle \mathbf{q}_1, \mathbf{a}_k \rangle \mathbf{q}_1 - \langle \mathbf{q}_2, \langle \mathbf{q}_1, \mathbf{a}_k \rangle \mathbf{q}_1 \rangle \mathbf{q}_2 - \dots = \\
 \mathbf{a}_k &- \sum_{j=1}^{k-1} \langle \mathbf{q}_j, \mathbf{a}_k \rangle \mathbf{q}_j
 \end{aligned}$$

We can use the Gram-Schmidt procedure for computing the QR factorization of a matrix. Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \geq n$) admits a QR decomposition if there exist an orthogonal matrix (columns orthogonal and normalized) $\mathbf{Q} \in \mathbb{R}^{m \times m}$ and an upper trapezoidal matrix $\mathbf{R} \in \mathbb{R}^{m \times n}$, such that $\mathbf{A} = \mathbf{QR}$.

Reduced QR factorization:

If $\mathbf{A}_{m \times n}$ is of rank n (i.e. has full rank) for which a QR factorization is known, then there exists a unique factorization of \mathbf{A} of the form

$$\mathbf{A}_{m \times n} = \tilde{\mathbf{Q}}_{m \times n} \tilde{\mathbf{R}}_{n \times n} \tag{48}$$

where $\tilde{\mathbf{Q}}$ is orthogonal and $\tilde{\mathbf{R}}$ upper triangular, and these are submatrices of \mathbf{Q} and \mathbf{R} (left upper corner). $\tilde{\mathbf{Q}}$ has orthonormal vectors of columns and coincides with the *Cholesky factor* \mathbf{H} of the symmetric positive definite matrix $\mathbf{A}^T \mathbf{A}$, i.e.

$$\mathbf{A}^T \mathbf{A} = \tilde{\mathbf{R}}^T \tilde{\mathbf{R}} \tag{49}$$

This matrix is positive definite because every matrix of the form $\mathbf{A}^T \mathbf{A}$ is positive semidefinite, and if \mathbf{A} has full rank then it is positive definite.

If \mathbf{A} has full rank n then the column vectors of \mathbf{A} form an orthonormal basis for the vector space $\text{range}(\mathbf{A})$. The QR factorization is therefore a procedure of generating an orthonormal basis for a given set of vectors. If \mathbf{A} has rank $r < n$, then the QR factorization does not necessarily yield an orthonormal basis for $\text{range}(\mathbf{A})$. However, it is possible to obtain a factorization of the form

$$\mathbf{Q}^T \mathbf{A} \mathbf{P} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},$$

where \mathbf{Q} is orthogonal, \mathbf{P} is a permutation matrix and \mathbf{R}_{11} is a nonsingular upper triangular matrix of order r .

$$\mathbf{A}_{(m \times n)} = \mathbf{Q}_{(m \times n)} \mathbf{R}_{(n \times n)}^{\text{upper}}$$

Figure 6: The reduced QR factorization.

If the QR decomposition is performed by the modified Gram-Schmidt procedure, then the columns of \mathbf{A} are vectors \mathbf{a}_i and columns of \mathbf{Q} are vectors \mathbf{q}_i of the equation (47). Matrix \mathbf{R} is obtained by left multiplication of \mathbf{A} by \mathbf{Q}^T (from equation (48), since \mathbf{Q} is orthogonal):

$$\underline{\underline{\mathbf{R} = \mathbf{Q}^T \mathbf{A}}}. \quad (50)$$

6.2.1.1 Case $m=n$

$$\mathbf{A}_{(n \times n)} = \mathbf{Q}_{(n \times n)} \mathbf{R}_{(n \times n)}^{\text{upper}} \quad (51)$$

If \mathbf{A} is non-singular and \mathbf{R} has positive diagonal elements (by agreement), then the factorization (51) is unique.

6.2.2 Remark: Extended QR decomposition

In fact, some QR algorithms produce the expanded trapezoid matrix $\mathbf{R}_{(m \times n)}$. In this way we have

$$\mathbf{A}_{m \times n} = \mathbf{Q}_{(m \times m)} \mathbf{R}_{(m \times n)}, \quad (52)$$

where \mathbf{R} is upper trapezoidal,

$$r_{ij} = \begin{cases} r_{ij} & ; i \leq n \\ 0 & ; i > n \end{cases}, \quad (53)$$

and \mathbf{A} is upper triangular, i.e. $a_{ij} = 0; i > j$. The orthogonal matrix \mathbf{Q} is orthogonal in the whole space $\mathbb{R}^{n \times n}$, therefore its orthonormal basis could not obtain only from columns of \mathbf{A} , but was expanded. We will therefore call the procedure **expanded QR decomposition with Gram-Schmidt orthogonalization**.

While \mathbf{R} is simply supplemented by zeros below the m -th row, there are different possibilities of how to expand \mathbf{Q} . The most obvious possibility is to successively take trial vectors of some basis of \mathbb{R}^n , orthogonalize them with respect to already calculated columns of \mathbf{Q} , and if linear dependence on already obtained orthonormal vectors occurs then reject the current trial vector and skip to the next one (since at least $m-k$ vectors of the basis will be linearly independent on all k vectors that are already there). Orthonormal basis vectors $\mathbf{e}_i = [0, \dots, 1, 0, \dots]$ can be most conveniently taken¹.

For a measure of linear dependency, we can simply take the norm of what remains after subtraction of projections on already calculated \mathbf{q}_k .² Some lower bound must be set on this norm, e.g. 0.1. The norm must be such that there for sure exists such basis vector that after subtraction of all projections on any set of lower dimensional basis vectors, the norm of the remaining vector is greater than this lower bound. Besides, the norm should not be too small because this would increase numerical errors (because we would allow acceptance of vectors that are almost linearly dependent on the currently available basis vectors).

If we would also need to extend \mathbf{A} with additional columns such that it would have a full rank in \mathbb{R}^n , then we can proceed as follows: We first extend \mathbf{Q} as described above. Then we extend \mathbf{R} to the dimension $m \times m$ by diagonal elements taking the value 1 (or maybe an average absolute value of already calculated diagonal elements obtained by the reduced QR factorization, in order to improve scaling) and out of diagonal elements taking the value 0. Then we calculate the extended \mathbf{A} by simply calculating the product \mathbf{QR} by extended matrices. The extended \mathbf{A} would in this way have the full rank m if the original \mathbf{A} had full rank n (because both \mathbf{Q} and \mathbf{R} would have a full rank). Besides, the first n columns of \mathbf{A} would be the same as with the original because of the zeros in \mathbf{R} below the upper-left $n \times n$ block.

6.2.3 Solution of Systems of Equations with QR decomposition

We are solving the system

¹ In special cases when there is a reason to suspect that columns of \mathbf{A} are close to the first m of these vectors, we can reverse the order in which these vectors are taken, in order to minimize the possibility of rejecting vectors.

² In fact, the ratio between this norm and norm of \mathbf{e}_i should be taken, but norm of \mathbf{e}_i is 1.

$$\mathbf{A}_{(m \times n)} \mathbf{x}_{(n \times 1)} = \mathbf{b}_{(m \times 1)} \quad (54)$$

We assume we have factored matrix \mathbf{A} as a product of an orthogonal matrix \mathbf{Q} and an upper trapezoid matrix \mathbf{R} :

$$\mathbf{A}_{(m \times n)} = \mathbf{Q}_{(m \times m)} \mathbf{R}_{(m \times n)} \quad (55)$$

where \mathbf{Q} has a full rank m , and

$$i > j \Rightarrow r_{ij} = 0 . \quad (56)$$

If the $m > n$ then the system

$$\mathbf{A}_{(m \times n)} \mathbf{x}_{(n \times 1)} = \mathbf{b}_{(m \times 1)} \quad (57)$$

is **overdetermined** (i.e. there are more equations) and can therefore be solved in the least square sense.

6.2.3.1 Case $m=n$

We first limit ourselves on the case $m=n$, i.e. the number of equations is the same as the number of variables. We solve the system by first setting (since $\mathbf{Q}^{-1} = \mathbf{Q}^T$)

$$\mathbf{y} = \mathbf{Q}^T \mathbf{b} . \quad (58)$$

and then solving the system

$$\mathbf{R} \mathbf{x} = \mathbf{y} \quad (59)$$

Equation (59) is a system with an upper triangular matrix whose solution is described in Section 8.4.1.

6.2.3.2 Case $m > n$

In this case we have an overdetermined system with more equations than unknowns. Solution of such systems is described in Section 7.

6.3 Non-standard factorization by using the GS

In this chapter a non-standard form of a factorization obtained derived from the Gram-Schmidt orthogonalization is used. A standard form of factorization known as QR decomposition is described in Section 6.2. The Section 6.1.1 explains the idea of Gram-Schmidt orthogonalization that is used for derivation of both forms, while the procedure described in this Section is specific for non-standard form.

In standard form, we perform the QR factorization in such a way that the original matrix is a product of an orthogonal and upper triangular matrix produced by the factorization. The orthogonal matrix is calculated first and the upper triangular matrix is obtained by left multiplying the original matrix by its transpose (since the inverse of the orthogonal matrix equals its transpose).

In the non-standard factorization described here, we express the orthogonal matrix that is obtained from the original one, as a product of the original and upper triangular matrix. Again, columns of the orthogonal matrix are orthogonalized vectors obtained from columns of the original matrix. The orthogonal matrix is obtained by a similar Gram-Schmidt process as in QR decomposition, while the upper triangular matrix is obtained simultaneously with the orthogonal matrix by performing equivalent operations as on the orthogonal matrix, on the matrix that is first set to identity matrix (and these operations transform it to the upper triangular form such that its left product with the original matrix equals the orthogonal matrix). The method will not be often used in practice, but it is described here because it is instructive.

After the Gram-Schmidt orthogonalization, each orthogonal vectors \mathbf{q}_j is a linear combination of the original vectors. We can therefore write

$$\underline{\underline{\mathbf{Q}_{(m \times n)}}} = \mathbf{V}_{(m \times n)} \mathbf{A}_{(n \times n)}, \quad (6.60)$$

where \mathbf{A} is a matrix of coefficients of these linear combinations¹, \mathbf{V} is the matrix whose columns are the original vectors \mathbf{v}_i and \mathbf{Q} is the matrix whose columns are orthogonal vectors \mathbf{q}_j . The above equation can be written by columns, which gives orthogonal vectors expressed as linear combination of the original ones:

Remark: the following derivation should be checked. The final procedure is correct, however, because the functions described in Section 6.3.3 were all verified.

$$\mathbf{q}_j = \sum_{k=1}^n a_{kj} \mathbf{v}_k, \quad j = 1, 2, \dots, m. \quad (6.0)$$

¹ Beware of different notation with respect to the usual notation. Here \mathbf{V} denotes the original matrix and \mathbf{A} denotes an upper triangular matrix, while in most commonly, \mathbf{A} is used for the original matrix and \mathbf{R} is used for the upper triangular matrix.

This comes from

$$(\mathbf{q}_j)_i = q_{ij} = \sum_{k=1}^n v_{ik} a_{kj} = \sum_{k=1}^n a_{kj} v_{ik}, \quad \mathbf{q}_j = \sum a_{kj} (\mathbf{v}^T)_k.$$

It is seen from (6.0 and (6.0 that matrix \mathbf{A} is upper triangular with unit diagonal elements:

$$\mathbf{A} = \begin{bmatrix} 1 & * & \dots & * \\ 0 & 1 & \dots & * \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (6.0.$$

Evaluation of orthogonal vectors \mathbf{q}_j can be accompanied by simultaneous evaluation of coefficients of \mathbf{A} . The complete procedure is then the following:

$j = 1$:

set

$$\mathbf{q}_1 = \mathbf{v}_1 \quad (6.0.$$

and

$$a_{11} = 1, \quad a_{k1} = 0, \quad k = 2, 3, \dots, m \quad (6.0.$$

for $j = 2, 3, \dots, m$:

set

$$\mathbf{q}_j = \mathbf{v}_j + \sum_{l=1}^{j-1} -\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \mathbf{q}_l, \quad (6.0.$$

and

$$a_{jj} = 1, \quad a_{kj} = \sum_{l=k}^{j-1} -\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} a_{kl}, \quad k = 1, 2, \dots, j-1. \quad (6.0.$$

Equation (6.0 follows from (6.0 and (6.0, which give

$$\mathbf{q}_j = \mathbf{v}_j + \sum_{l=1}^{j-1} -\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \mathbf{q}_l = \mathbf{v}_j + \sum_{l=1}^{j-1} -\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \sum_{k=1}^l a_{kl} \mathbf{v}_k = \sum_{k=1}^j a_{kj} \mathbf{v}_k, \quad (6.0)$$

where (6.0) was taken into account twice. Equating terms at \mathbf{v}_k gives (6.0), where it is taken into account that $a_{kl} = 0$ for $l < k$.

In the algorithm the coefficient $-\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \mathbf{q}_l$ should be evaluated only once for each (j, l) , therefore individual contributions to all a_{kj} in (6.0) are evaluated when specific coefficient is available. The algorithm step ((6.0), (6.0)) is therefore actually the following:

$j = 1$:

set

$$\underline{\mathbf{q}_1 = \mathbf{v}_1} \quad (6.0)$$

and

$$\underline{a_{11} = 1, a_{k1} = 0, k = 2, 3, \dots, m} \quad (6.0)$$

for $j = 2, 3, \dots, m$:

set

$$\underline{a_{jj} = 1, a_{kj} = 0, k = 1, 2, \dots, j-1} \quad (6.0)$$

and

$$\underline{\mathbf{q}_j = \mathbf{v}_j}, \quad (6.0)$$

for $l = 2, 3, \dots, j-1$

$$\underline{\mathbf{q}_j = \mathbf{q}_j + -\frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} \mathbf{q}_l}, \quad (6.0)$$

and

$$\underline{a_{kj} = a_{kj} - \frac{\langle \mathbf{v}_j, \mathbf{q}_l \rangle}{\langle \mathbf{q}_l, \mathbf{q}_l \rangle} a_{kl}, \quad k = 1, 2, \dots, l.} \quad (6.0.)$$

6.3.1 Orthogonalization with normalization

This operation generates vectors which are mutually orthogonal and whose Euclidian norm equals 1. The procedure is practically equivalent except for the normalization factor that is applied at the end of calculation of each column of \mathbf{Q} . The procedure follows those described in (6.0 to (6.0, except that orthogonal vectors are normalised immediately after they are constructed and because of that some of the coefficient in the above equations become 1. Orthogonalisation with normalisation is approximately as fast as without normalisation, but is more convenient for solution of systems of equations. Matrix \mathbf{A} does not have units in diagonal any more:

$$\mathbf{A} = \begin{bmatrix} * & * & \dots & * \\ 0 & * & \dots & * \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & * \end{bmatrix} \quad (6.0.)$$

The procedure is the following:

$j = 1$:

set

$$\underline{\mathbf{q}_1 = \frac{\mathbf{v}_1}{\langle \mathbf{v}_1, \mathbf{v}_1 \rangle}} \quad (6.0.)$$

and

$$\underline{a_{11} = 1, \quad a_{k1} = 0, \quad k = 2, 3, \dots, m} \quad (6.0.)$$

for $j = 2, 3, \dots, m$:

set

$$\underline{a_{jj} = 1, \quad a_{kj} = 0, \quad k = 1, 2, \dots, j-1} \quad (6.0.)$$

and

$$\underline{\mathbf{q}_j = \mathbf{v}_j} . \quad (6.0.)$$

for $\underline{l = 2, 3, \dots, j-1}$

set

$$\underline{\mathbf{q}_j = \mathbf{q}_j - \langle \mathbf{v}_j, \mathbf{q}_l \rangle \mathbf{q}_l} , \quad (6.0.)$$

and

$$\underline{a_{kj} = a_{kj} - \langle \mathbf{v}_j, \mathbf{q}_l \rangle a_{kl}, \quad k = 1, 2, \dots, l} . \quad (6.0.)$$

for $k = 1, 2, \dots, j$

set

$$a_{kj} = \frac{a_{kj}}{\sqrt{\langle \mathbf{q}_j, \mathbf{q}_j \rangle}} . \quad (6.0.)$$

set

$$\mathbf{q}_j = \frac{\mathbf{q}_j}{\sqrt{\langle \mathbf{q}_j, \mathbf{q}_j \rangle}} \quad (6.0.)$$

6.3.2 Solution of Systems of Equations with Gram-Schmidt Orthogonalisation (non-standard form)

In Gram-Schmidt Orthogonalisation we construct from matrix of original vectors (columns) \mathbf{V} an orthogonal matrix \mathbf{Q} and upper trapezoid matrix of coefficients \mathbf{A} so that

$$\mathbf{Q}_{(n \times n)} = \mathbf{V}_{(n \times m)} \mathbf{A}_{(m \times m)} ; m \leq n . \quad (6.0.)$$

For solution of system of equations only square matrices will be considered here, i.e. $m=n$.

We can solve the system

$$\underline{\mathbf{V}\mathbf{x} = \mathbf{b}} \quad (6.0)$$

by introducing a new variable $\mathbf{y} = \mathbf{A}^{-1}\mathbf{x}$. The above system then falls to one orthogonal system and matrix multiplication:

$$\underline{\mathbf{Q}\mathbf{y} = \mathbf{b}} \quad (6.0)$$

and

$$\underline{\mathbf{x} = \mathbf{A}\mathbf{y}} \quad (6.0)$$

Solution of (6.0 is described in Section 8.2 while (6.0 is just a simple matrix multiplication.

6.3.3 Implementation of GS (non-standard form)

In IOptLib the following functions are used:

```
void GSortplain(matrix V, matrix Q, matrix A)
```

- From original matrix $\mathbf{V}_{(m \times n)}$ ($m \geq n$), a matrix with mutually orthogonal (but not normal) columns $\mathbf{Q}_{m \times n}$ and upper triangular $\mathbf{A}_{(n \times n)}$ are calculated such that $\mathbf{Q} = \mathbf{V}\mathbf{A}$. All matrices must be allocated and of appropriate dimensions.

```
void GSortnormplain(matrix v, matrix q, matrix a)
```

- Similar as GSortplain, only that \mathbf{Q} is orthogonal (i.e. columns are not only mutually orthogonal but also normalized).

```
void GSort0(matrix V, matrix *Q, matrix *A)
```

and

```
int GSortnorm0(matrix v, matrix *q, matrix *a)
```

are comfortable forms of GSortplain and GSortnormplain. Matrices of results are allocated or reallocated with the appropriate dimensions if necessary.

```
vector solvGS0(matrix M, vector b, vector *x)
```

Solves a system of equations $\mathbf{M}\mathbf{x} = \mathbf{b}$ in such a way that the Gram-Schmidt orthogonalization is first performed by GSortnormplain and then the system is solved according to (0 and (0. The intermediate storage necessary for operations is allocated and then released within the function. Because of this, function is not efficient and should not be used for important computations.

The systems can instead be solved by three calls – first to orthogonalization (`GSortnormplain0(M, &Q, &A)`) and then by `solvortnorm0(Q, b, &x)` to solve the orthogonal system (0 followed by `matprodvec0(A, x, &x)` to perform matrix multiplication (0. Solution of the orthogonal system can also be performed by `mattranspprodvec(Q, b, &x)`).

```
matrix solvmatGS0(matrix M, matrix B, matrix *X)
```

- This function is similar to `solvGS0`, but it solves several systems of equations at once, where **B** contains right-hand sides as its columns and solutions are stored in columns of **X**.

7 OVERDETERMINED SYSTEMS

In this Section we describe solution of **over determined systems** of equations where we have more equations than unknowns:

$$\mathbf{A}_{(m \times n)} \mathbf{x}_{(n \times 1)} = \mathbf{b}_{(m \times 1)}, \quad m > n. \quad (61)$$

We assume that rank of **A** is n . In general this means that we can not simultaneously satisfy all the equations. We therefore search for the solution in a **least squares sense**, i.e. we are searching for such \mathbf{x}_0 for which the sum of squares of differences between left-hand and right-hand sides is minimal:

$$\mathbf{x}_0 = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2. \quad (62)$$

The function to be minimized is expressed component wise as

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 = \sum_{i=1}^m \left(\sum_{j=1}^n (a_{ij}x_j) - b_i \right)^2.$$

A unique solution \mathbf{x}_0 exists (under assumption $\text{rank } \mathbf{A} = n$). This is exactly the solution of equation

$$\underline{\underline{\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}}}, \quad (63)$$

which is obtained by left multiplying the equation by \mathbf{A}^T . This system is called the *normal system* of equations.

We can verify that the solution of (63) is also the solution of (62). We write

$$\mathbf{B} = \mathbf{A}^T \mathbf{A}, \quad \mathbf{c} = \mathbf{A}^T \mathbf{b}. \quad (64)$$

Then we have

$$\| \mathbf{A} \mathbf{x} - \mathbf{b} \|_2^2 = \langle \mathbf{A} \mathbf{x} - \mathbf{b}, \mathbf{A} \mathbf{x} - \mathbf{b} \rangle = (\mathbf{B} \mathbf{x} - \mathbf{c})^T \mathbf{B}^{-1} (\mathbf{B} \mathbf{x} - \mathbf{c}) - \mathbf{c}^T \mathbf{B}^{-1} \mathbf{c} + \mathbf{b}^T \mathbf{b} .$$

Matrix \mathbf{B} is positive definite: $\forall \mathbf{x}, \mathbf{x}^T \mathbf{B} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = \| \mathbf{A} \mathbf{x} \|_2^2 \geq 0$. This is 0 only when $\mathbf{A} \mathbf{x} = \mathbf{0}$, which is possible only when $\mathbf{x} = \mathbf{0}$ because matrix \mathbf{A} has a full rank n .

If \mathbf{B} is positive definite then also \mathbf{B}^{-1} is positive definite. The smallest value of the quadratic form therefore equals 0 and this is when equation (63) holds, while other terms are independent of \mathbf{x} and represent the sum of squares of differences in the solution. Because $\mathbf{A}^T \mathbf{A} = \mathbf{B}$ is not singular, the solution is uniquely defined.

System (63) can be calculated by the Cholesky decomposition of \mathbf{B} ,

$$\underline{\underline{\mathbf{A}^T \mathbf{A} = \mathbf{B} = \mathbf{V}^T \mathbf{V}}} \quad (65)$$

where \mathbf{V} is upper triangular, which is followed by solution the lower triangular system

$$\underline{\underline{\mathbf{V}^T \mathbf{y} = \mathbf{c} \quad (= \mathbf{A}^T \mathbf{b})}} \quad (66)$$

and the upper triangular system

$$\underline{\underline{\mathbf{V} \mathbf{x}_0 = \mathbf{y}}} . \quad (67)$$

The described classical way is not recommendable because sensitivity of the system can be strongly increased at explicit calculation of \mathbf{B} (i.e. the matrix can become very ill conditioned). Let us suppose that eigenvalues of \mathbf{B} arranged in decreasing order are σ_i^2 . Spectral sensitivity of \mathbf{A} is then σ_1 / σ_n , while spectral sensitivity of \mathbf{B} is $\left(\sigma_1 / \sigma_n \right)^2$. More recommendable is solution through orthogonal transformations.

7.1 Orthogonal methods

The basis of the orthogonal methods is the QR factorization of the matrix \mathbf{A} , to a product of an **orthogonal matrix \mathbf{Q}** and **upper trapezoid matrix \mathbf{U}** :

$$\underline{\underline{\mathbf{A}_{(m \times n)} = \mathbf{Q}_{(m \times m)} \mathbf{U}_{(m \times n)}}} , \quad (68)$$

where

$$\begin{aligned} \mathbf{Q}^T \mathbf{Q} &= \mathbf{I}_m \\ i > j &\Rightarrow u_{ij} = 0 \end{aligned} \quad (69)$$

We calculate the solution of the orthogonal system

$$\underline{\underline{\mathbf{z} = \mathbf{Q}^T \mathbf{b}}}, \quad (70)$$

and we write \mathbf{U} and \mathbf{z} in block form,

$$\underline{\underline{\mathbf{U} = \begin{bmatrix} \mathbf{V}_{(n \times n)} \\ \mathbf{0}_{((m-n) \times n)} \end{bmatrix}}}, \quad \underline{\underline{\mathbf{z} = \begin{bmatrix} \mathbf{y}_n \\ \mathbf{w}_{m-n} \end{bmatrix}}}. \quad (71)$$

Taking into account the decomposition (68) and block form (71), we have

$$\mathbf{B} = \mathbf{A}^T \mathbf{A} = (\mathbf{Q}\mathbf{U})^T (\mathbf{Q}\mathbf{U}) = \mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V}. \quad (72)$$

This means that upper triangular matrix \mathbf{V} is the Cholesky factor of the normal matrix $\mathbf{A}^T \mathbf{A}$. This means that with the orthogonal methods, calculation of the normal matrix $\mathbf{A}^T \mathbf{A}$ and its factorization is avoided. Matrices \mathbf{B} in (71) and (65) are the same if we take care in the orthogonal decomposition that all diagonal elements of \mathbf{V} are positive.

We can further calculate

$$\mathbf{c} = \mathbf{A}^T \mathbf{b} = (\mathbf{Q}\mathbf{U})^T \mathbf{b} = \mathbf{U}^T \mathbf{Q}^T \mathbf{b} = \mathbf{U}^T \mathbf{z} = \mathbf{V}^T \mathbf{y}. \quad (73)$$

This means that \mathbf{y} , which is the upper part of the transformed vector \mathbf{z} , solution of the lower triangular system (66).

The least squares solution is obtained by solution of the upper triangular system

$$\underline{\underline{\mathbf{V} \mathbf{x}_0 = \mathbf{y}}}. \quad (74)$$

Vector \mathbf{w} , which is the lower part of \mathbf{z} , also has its meaning – square of its Euclidean norm is the sum of squares of differences:

$$\|\mathbf{A}x - \mathbf{b}\|_2^2 = \|\mathbf{w}\|_2^2. \quad (75)$$

This is derived as follows:

$$\| \mathbf{A}x - \mathbf{b} \|_2^2 = \| \mathbf{Q} \mathbf{U} \mathbf{x}_0 - \mathbf{Q} \mathbf{z} \|_2^2 = \| \mathbf{U} \mathbf{x}_0 - \mathbf{z} \|_2^2 = \| \mathbf{V} \mathbf{x}_0 - \mathbf{y} \|_2^2 + \| \mathbf{w} \|_2^2 = \| \mathbf{w} \|_2^2 .$$

Remark on QR decomposition:

The QR decomposition is often calculated in the **reduced form**, such that

$$\mathbf{A}_{(m \times n)} = \mathbf{Q}_{(m \times n)} \mathbf{U}_{(n \times n)} , \quad (76)$$

In this case \mathbf{Q} must be supplemented by orthogonal columns up to full rank m .

8 SPECIAL SYSTEMS OF EQUATIONS

8.1 Orthogonal systems

By the term orthogonal system, we will refer to a systems with *orthogonal matrix*, i.e. a matrix with **orthonormal columns** (and consequently rows). A more general class of systems with matrices whose columns are mutually orthogonal but not normalized is described in Section 8.2.

For orthogonal real matrices¹ the following equation is valid:

$$\mathbf{Q} \mathbf{Q}^T = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}; \quad \mathbf{Q}^{-1} = \mathbf{Q}^T . \quad (77)$$

Therefore, the system of equations with an orthogonal matrix

$$\mathbf{Q} \mathbf{x} = \mathbf{b} \quad (78)$$

is solved simply by multiplying the vector of right hand sides by the transpose of the system matrix:

$$\mathbf{x} = \mathbf{Q}^T \mathbf{b} \quad (79)$$

¹ Equivalent for matrices defined over complex field is unitary.

8.2 Systems with Matrices that have Orthogonal Columns

Let's have a system

$$\underline{\mathbf{Q}\mathbf{x} = \mathbf{b}}, \quad (8.0)$$

where \mathbf{Q} is orthogonal, i.e. such that its columns are mutually orthogonal (dot product zero). The above system can be written in the form

$$\mathbf{b} = x_1\mathbf{Q}_1 + x_2\mathbf{Q}_2 + \dots + x_n\mathbf{Q}_n, \quad (8.0)$$

where \mathbf{Q}_i denotes the i -th column of \mathbf{Q} and x_i denotes the i -th component of vector \mathbf{x} . If the equation is dot-multiplied by \mathbf{Q}_i , we obtain

$$\langle \mathbf{b}, \mathbf{Q}_i \rangle = x_i \langle \mathbf{Q}_i, \mathbf{Q}_i \rangle, \quad (8.0)$$

which follows from orthogonality of \mathbf{Q}_j and \mathbf{Q}_i for each $j \neq i$. It follows that

$$x_i = \frac{\langle \mathbf{b}, \mathbf{Q}_i \rangle}{\langle \mathbf{Q}_i, \mathbf{Q}_i \rangle}. \quad (8.0)$$

If columns of \mathbf{Q} are also normed in the Euclidian norm, then we simply have

$$x_i = \langle \mathbf{b}, \mathbf{Q}_i \rangle, \quad (8.0)$$

which simplifies to

$$\underline{\mathbf{x} = \mathbf{Q}^T \mathbf{b}}. \quad (8.0)$$

If columns of \mathbf{Q} are orthonormal then \mathbf{Q} is an orthogonal matrix for which $\mathbf{Q}^{-1} = \mathbf{Q}^T$.

8.2.1 Inverse of a Matrix with Orthogonal Columns

We search for an unknown matrix \mathbf{X} such that

$$\mathbf{Q}\mathbf{X} = \mathbf{I}. \quad (8.0)$$

For the i -th column of \mathbf{X} we have

$$\mathbf{Q}\mathbf{X}_i = \mathbf{e}_i, \quad (8.0)$$

where \mathbf{e}_i is the i -th unit vector (i.e. component i is 1 and all others are 0). It follows from (8.0 that

$$x_{ji} = [\mathbf{X}_i]_j = \frac{\langle \mathbf{e}_i, \mathbf{Q}_j \rangle}{\langle \mathbf{Q}_j, \mathbf{Q}_j \rangle} = \frac{\mathbf{q}_{ji}}{\langle \mathbf{Q}_j, \mathbf{Q}_j \rangle} \quad (8.0)$$

or finally

$$\underline{\underline{[\mathbf{Q}^{-1}]_{ij}}} = \frac{\mathbf{q}_{ji}}{\langle \mathbf{Q}_i, \mathbf{Q}_i \rangle} \quad (8.0)$$

If columns of \mathbf{Q} are also normalised with respect to the Euclidian norm, then we have simply

$$\underline{\underline{[\mathbf{Q}^{-1}]_{ij}}} = \mathbf{q}_{ji}. \quad (8.81)$$

8.2.2 Component-wise verification

We verify some of the equations from Section 8.2 in component-wise notation. Let's have a system

$$\mathbf{A}\mathbf{x} = \mathbf{b}. \quad (8.82)$$

$$\mathbf{A}_{(m \times n)} \mathbf{x}_{(n \times 1)} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \dots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix} = x_1 \begin{bmatrix} a_{11} \\ a_{21} \\ \dots \\ a_{m1} \end{bmatrix} + x_2 \begin{bmatrix} a_{12} \\ a_{22} \\ \dots \\ a_{m2} \end{bmatrix} + \dots + x_n \begin{bmatrix} a_{1n} \\ a_{2n} \\ \dots \\ a_{mn} \end{bmatrix} \quad (8.83)$$

$$= x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_n \mathbf{a}_n$$

With \mathbf{a}_i we denote the i -th column of \mathbf{A} .

For product of two matrices it holds

$$\left[\mathbf{A}_{(m \times n)} \mathbf{X}_{(n \times p)} \right]_j = \begin{bmatrix} a_{1k} x_{kj} \\ a_{2k} x_{kj} \\ \dots \\ a_{mk} x_{kj} \end{bmatrix} = \mathbf{A} \mathbf{x}_j, \quad (8.84)$$

i.e. the j -th column of the product equals product of the left matrix with the j -th column of the right matrix. In the above equation, the Einstein summation rule was applied, i.e. we have summation over indices that are doubled.

Let's have a matrix \mathbf{Q} whose columns are mutually orthogonal:

$$\mathbf{q}_i \mathbf{q}_j = \sum_{k=1}^n q_{ki} q_{kj} = \delta_{ij} \|\mathbf{q}_i\|^2, \quad (8.85)$$

where $\delta_{ij} = \begin{cases} 1; i = j \\ 0; i \neq j \end{cases}$.

Product $\mathbf{Q}^T \mathbf{Q} = \mathbf{D}$ is a diagonal matrix because

$$[\mathbf{Q}^T \mathbf{Q}]_{ij} = \langle \mathbf{q}_i, \mathbf{q}_j \rangle = \delta_{ij} \|\mathbf{q}_i\|^2, \quad (8.86)$$

where \mathbf{q}_i is the i -th column vector of matrix \mathbf{Q} .

A special case are **orthogonal matrices** which in addition to (8.0) have normed columns, such that

$$\langle \mathbf{q}_i, \mathbf{q}_j \rangle = \delta_{ij} \quad (8.87)$$

For such a matrix the following is true:

$$\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}. \quad (8.88)$$

8.3 Systems of Equations with Matrix with Orthogonal Columns

Let's have a system

$$\underline{\mathbf{Q}\mathbf{x} = \mathbf{b}}, \quad (8.89)$$

where \mathbf{Q} is orthogonal matrix. It follows that

$$\mathbf{b} = x_1\mathbf{q}_1 + x_2\mathbf{q}_2 + \dots + x_n\mathbf{q}_n = \sum_{k=1}^n x_k\mathbf{q}_k, \quad (8.90)$$

where \mathbf{q}_i is the i -th column of \mathbf{Q} . After multiplication of the above equation with \mathbf{q}_i^T from left we obtain

$$\langle \mathbf{b}, \mathbf{q}_i \rangle = x_i \langle \mathbf{q}_i, \mathbf{q}_i \rangle \quad (8.91)$$

and so

$$x_i = \frac{\langle \mathbf{b}, \mathbf{q}_i \rangle}{\|\mathbf{q}_i\|^2}. \quad (8.92)$$

Solution of the system can be obtained very simply. This can be used for formulae for inverse of an orthogonal matrix:

$$\mathbf{Q}\mathbf{X} = \mathbf{I}, \quad [\mathbf{Q}\mathbf{x}]_i = \mathbf{Q}\mathbf{x}_i = \mathbf{e}_i, \quad (8.93)$$

using (8.0 we have

$$[\mathbf{x}_i]_j = x_{ji} = \frac{\langle \mathbf{e}_i, \mathbf{q}_j \rangle}{\|\mathbf{q}_j\|^2} = \frac{q_{ij}}{\|\mathbf{q}_j\|^2} \quad (8.94)$$

or

$$[\mathbf{Q}^{-1}]_{ij} = \frac{q_{ji}}{\|\mathbf{q}_j\|^2} = \frac{q_{ji}}{\sum_{k=1}^n q_{ki}q_{ki}}. \quad (8.95)$$

If columns of \mathbf{Q} are also normed with respect to the Euclidian norm, then the denominator of the above equation equals 1.

8.4 Triangular Systems

8.4.1 Upper Triangular Systems

We solve the system

$$\mathbf{Ax} = \mathbf{b}, \quad (8.96)$$

where \mathbf{A} is upper triangular, i.e. elements below diagonal are zero, see (3.0). The system looks like this:

$$\begin{array}{c} \mathbf{A} \\ \left[\begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ 0 & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & a_{nn} \end{array} \right] \end{array} \begin{array}{c} \mathbf{x} \\ \left[\begin{array}{c} x_1 \\ x_2 \\ \dots \\ x_n \end{array} \right] \end{array} = \begin{array}{c} \mathbf{b} \\ \left[\begin{array}{c} b_1 \\ b_2 \\ \dots \\ b_n \end{array} \right] \end{array} \quad (8.97)$$

The above equations written in reversed order are:

$$\begin{array}{rcccccl} & & & & + a_{nn}x_n & = b_n \\ & & & & + a_{n-1n}x_n & = b_{n-1} \\ & & & + a_{n-1n-1}x_{n-1} & + a_{n-1n}x_n & = b_{n-1} \\ + a_{n-2n-2}x_{n-2} & + a_{n-2n-1}x_{n-1} & + a_{n-2n}x_n & = b_{n-2} & & (8.98) \\ \dots & \dots & \dots & \dots & \dots & \\ a_{11}x_1 & + \dots & + a_{1n-1}x_{n-1} & + a_{1n}x_n & = b_1 \end{array}$$

The solution is evaluated backwards:

$$x_n = \frac{b_n}{a_{nn}}, \quad (8.99)$$

$$x_i = \left(b_i - \sum_{k=i+1}^n a_{ik} x_k \right) / a_{ii}, \quad i = n-1, n-2, \dots, 1 \quad (8.100)$$

8.4.2 Lower Triangular Systems

We solve the system

$$\mathbf{Ax} = \mathbf{b}, \quad (8.101)$$

where \mathbf{A} is lower triangular, i.e. elements below diagonal are zero, see (3.0). The system looks like this:

$$\begin{matrix} & \mathbf{A} & & \mathbf{x} & & \mathbf{b} \\ \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ a_{21} & a_{22} & \dots & \dots \\ \dots & \dots & \dots & 0 \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} & \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} & = & \begin{bmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{bmatrix} \end{matrix} \quad (8.102)$$

The solution is evaluated forwards:

$$x_1 = \frac{b_1}{a_{11}}, \quad (8.103)$$

$$x_i = \left(b_i - \sum_{k=1}^{i-1} a_{ik} x_k \right) / a_{ii}, \quad i = 2, 3, \dots, n \quad (8.104)$$

9 NUMBER OF OPERATIONS FOR STANDARD MATRIX OPERATIONS

9.1 Basic Operations

9.1.1 Matrix and vector multiplications:

Matrix multiplications:

$$\mathbf{C}_{m \times n} = \mathbf{A}_{m \times p} \mathbf{B}_{p \times n} : N = 2mn p \text{ (vsak element približno } p \text{ množenj in } p \text{ seštevanj)}$$

$$\underline{\underline{\mathbf{C}_{n \times n} = \mathbf{A}_{n \times n} \mathbf{B}_{n \times n} : N = 2n^3}}$$

Matrix times vector:

$$\mathbf{c}_{m \times 1} = \mathbf{A}_{m \times n} \mathbf{b}_{n \times 1} : N = 2mn \text{ (vsak element približno } n \text{ množenj in } n \text{ seštevanj)}$$

$$\underline{\underline{\mathbf{c}_{n \times 1} = \mathbf{A}_{n \times n} \mathbf{b}_{n \times 1} : N = 2n^2 \text{ (vsak element približno } n \text{ množenj in } n \text{ seštevanj)}}}$$

Scalar product:

$$c = \mathbf{a}_{n \times 1}^T \mathbf{b}_n : N = 2n \text{ (} n \text{ množenj in } n \text{ seštevanj)}$$

9.2 Special systems of equations

9.2.1.1 Lower triangular system

$$\underline{\underline{\mathbf{L}_{n \times n} \mathbf{y}_{n \times 1} = \mathbf{b}_{n \times 1}, i < j \Rightarrow l_{ij} = 0, l_{ii} = 1; N = \frac{1}{2}n^2 - \frac{1}{2}n}}$$

this follows from $N = \sum_{r=1}^{n-1} (n-r) = \frac{1}{2}n^2 - \frac{1}{2}n$

9.2.1.2 Upper triangular system

$$\underline{\underline{\mathbf{U}_{n \times n} \mathbf{y}_{n \times 1} = \mathbf{b}_{n \times 1}, i > j \Rightarrow u_{ij} = 0; N = \frac{1}{2}n^2 + \frac{1}{2}n}}$$

this follows from $N = \sum_{i=1}^n (n-i+1) = \frac{1}{2}n^2 + \frac{1}{2}n$; The difference with respect to lower triangular system is that diagonal elements are different than 1, which adds n divisions.

The total number of operations requires $\frac{1}{2}n^2 + \frac{1}{2}n$ multiplications or divisions and $\frac{1}{2}n^2 - \frac{1}{2}n$ summations or subtractions. This is four times as much multiplications as Gaussian elimination and twice as much summations, beside the $n(n-1)/2$ square roots.

9.2.1.3 Orthogonal system

$$\mathbf{Q}_{n \times n} \mathbf{y}_{n \times 1} = \mathbf{b}_{n \times 1}, \mathbf{Q}^T \mathbf{Q} = \mathbf{I} : N = 3n^2 - 3n,$$

i.e. $2n^2 - 2n$ multiplications and $n^2 - n$ summations

9.3 Factorizations

9.3.1.1 LU factorization

$$\underline{\underline{\mathbf{A}_{n \times n} = \mathbf{L}_{n \times n} \mathbf{U}_{n \times n}}}; \underline{\underline{N = \frac{1}{3}n^3 - \frac{1}{3}n}}$$

Solution of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$: $N = \frac{1}{3}n^3 + n^2 - \frac{1}{3}n$

(factorization + lower triangular + upper triangular system)

Solution of m systems of equations: $N = \frac{1}{3}n^3 + mn^2 - \frac{1}{3}n$

Determinant of a matrix $\det \mathbf{A}_{n \times n}$: $\frac{1}{3}n^3 + \frac{2}{3}n - 1$

(factorization + n multiplications of diagonal terms of \mathbf{U})

Inverse matrix: n^3 operations

(by solution of $\mathbf{A}\mathbf{X} = \mathbf{I}$, solution with the s -th unit vector on the right is the s -th column of

\mathbf{X}) Decomposition: $N = \frac{1}{3}n^3 - \frac{1}{3}n$ operations, lower triangular with \mathbf{e}_i as right hand sides (less operations): $\frac{1}{6}n^3 - \frac{1}{2}n^2 + \frac{1}{3}n$ upper triangular: $\frac{n^2}{2} + \frac{n^2}{2}$, toal: n^3

9.3.1.2 QR Factorization

9.3.1.2.1 Reduced factorization with the modified Gram-Schmidt method:

$$\mathbf{A}_{m \times n} = \tilde{\mathbf{Q}}_{m \times n} \tilde{\mathbf{R}}_{n \times n}; \tilde{\mathbf{Q}}^T \mathbf{Q} = \mathbf{I}, i > j \Rightarrow r_{ij} = 0: N = 2mn^2 .$$

Solution of a system of equations:

$N = 2mn^2$ for factorization, $3m^2 - 3m$ for solution of the orthogonal system, and n^2 operations for upper triangular system, together.

9.3.1.2.2 QR factorization with Givens method

$$\mathbf{A}_{n \times n} = \mathbf{Q}_{n \times n} \mathbf{R}_{n \times n}; \mathbf{Q}^T \mathbf{Q} = \mathbf{I}, i > j \Rightarrow r_{ij} = 0 : N = \left(\frac{1}{2}n^2 - \frac{1}{2}n\right) + \left(\frac{4}{3}n^3 - \frac{4}{3}n\right) + \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n$$

$$N = 2n^3 - 2n$$

$\frac{1}{2}n^2 - \frac{1}{2}n$ square roots, $\frac{4}{3}n^3 - \frac{4}{3}n$ divisions or multiplications and $\frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n$ summations

Solution of a system of equations:

$N = 2n^3 - 2n$ for factorization, $3n^2 - 3n$ for solution of the orthogonal system, and n^2 operations for upper triangular system, together $2n^2 + n^2 - 5n$ operations.

9.3.1.2.3 QR factorization with the Housholder method

$\mathbf{A}_{n \times n} = \mathbf{Q}_{n \times n} \mathbf{R}_{n \times n}$; $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}, i > j \Rightarrow r_{ij} = 0$: $N = (n-1) + \left(\frac{2}{3}n^3 + O(n^2)\right) + \left(\frac{2}{3}n^3 + O(n^2)\right)$ - twice less multiplications as the Givens method and twice as much as Gauss elimination. Total number of operations is for 1/3 less than at Givens method and twice greater than with Gauss elimination.

$n-1$ square roots, $\frac{2}{3}n^3 + O(n^2)$ divisions or multiplications and $\frac{2}{3}n^3 + O(n^2)$ summations.

Solution of a system of equations:

$\frac{4}{3}n^3 + O(n^2)$ for factorization, $3n^2 - 3n$ for solution of the orthogonal system, and n^2 operations for upper triangular system.

9.3.1.3 LDMT factorization

$\mathbf{A}_{n \times n} = \mathbf{L}_{n \times n} \mathbf{D}_{n \times n} \mathbf{M}_{n \times n}^T$; $i < j \Rightarrow l_{ij} = 0 \wedge m_{ij} \neq 0, i \neq j \Rightarrow d_{ij} = 0$: $N = \frac{1}{3}n^3$

9.3.1.4 LDLT factorization

$\mathbf{A}_{n \times n} = \mathbf{L}_{n \times n} \mathbf{D}_{n \times n} \mathbf{L}_{n \times n}$; $i < j \Rightarrow l_{ij} = 0, i \neq j \Rightarrow d_{ij} = 0$: $N = \frac{1}{3}n^3$

9.3.1.5 Cholesky factorization

$$N = \frac{1}{6}n^3 + O(n^2)$$

References:

- [1] Alfio Quarteroni, Riccardo Sacco, Fausto Saleri, *Numerical Mathematics*. Texts in Applied Mathematics 37, Springer-Verlag, New York, 2000.
- [2] Richard L. Burden, J. Douglas Faires, *Numerical Analysis*, 6th Edition, Brooks/Cole Publishing Company, 1997.
- [3] James W. Demmel, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- [4] Lloyd N. Trefethen, David Bau, *Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- [5] I. N. Bronstein, K. A. Smendljajew, G. Musiol, H. Mühlig, *Taschenbuch des Mathematik (second edition - in German)*, Verlag Harri Deutsch, Frankfurt am Main, 1995.
- [6] I. Kuščer, A. Kodre, H. Neunzert, *Mathematik in Physik und Technik (in German)*, Springer - Verlag, Heidelberg, 1993.
- [7] I. Kuščer, A. Kodre, *Matematika v fiziki in tehniki (in Slovene)*, Društvo matematikov, fizikov in astronomov Slovenije, Ljubljana, 1994.
- [8] W.H. Press, S.S. Teukolsky, V.T. Vetterling, B.P. Flannery, *Numerical Recipes in C – the Art of Scientific Computing*, Cambridge University Press, Cambridge, 1992.
- [9] K.J. Bathe, *Finite Element Procedures*, p.p. 697-745, Prentice Hall, New Jersey, 1996.
- [10] Z. Bohte, *Numerične metode (in Slovene)*, Društvo matematikov, fizikov in astronomov SRS, Ljubljana, 1987.
- [11] Zvonimir Bohte, *Numerično reševanje sistemov linearnih enačb (in Slovene)*, Društvo fizikov, matematikov in astronomov Slovenije, Ljubljana, 1994.
- [12] Zvonimir Bohte, *Numerično reševanje nelinearnih enačb (in Slovene)*, i Društvo fizikov, matematikov in astronomov Slovenije, Ljubljana, 1993.
- [13] Egon Zakrajšek, *Matematično modeliranje (in Slovene)*, DMFA – založništvo, 2004.
- [14] Jože Petrišič, *Reševanje enačb (in Slovene)*, Univerza v Ljubljani – Fakulteta za strojništvo, 2006.
- [15] Jože Petrišič, *Interpolacija in osnove računalniške grafike (in Slovene)*, Univerza v Ljubljani – Fakulteta za strojništvo, 1999.
- [16] Juan Restrepo: Numerical Linear Algebra, electronic document at <http://www.physics.arizona.edu/~restrepo/475A/Notes/sourcea/node60.html>
- [17] Igor Grešovnik: Optimization shell Inverse, <http://www2.arnes.si/~ljc3m2/inverse/>
- [18] Igor Grešovnik: Investigative Optimization Library, <http://www2.arnes.si/~ljc3m2/igor/ioplib/>
- [19] Igor Grešovnik: Investigative Generic Library, <http://www2.arnes.si/~ljc3m2/igor/iglib/>

References

- [20] Igor Grešovnik: The Use of Moving Least Squares for a Smooth Approximation of Sampled Data, Journal of Mechanical Engineering 53 (2007) 9, 582-598 (UDK 517.518.8:519.243, 2007),
http://www2.arnes.si/~ljc3m2/igor/doc/papers/07_09_Gresovnik_Mech_Eng_MLS_Sampled_Data_Optimization.pdf