

## Self-study guide

### Week 1

**Keywords:** Introduction, Permutation matrices, Block matrix notation, Gaussian elimination, Back-substitution,  $LU$ -factorisation.

**Homework:** Problems 9, 10, 21 and 26. In addition, solve any additional four problems from 1-27 to gain extra points.

[See outline of Week 1 in Youtube](#)

**Pages:** 5-36.

**Synopsis:** During the first week we prepare for proving the existence of the Cholesky factorisation by discussing permutation matrices,  $LU$ -factorisation, block matrix notation, and recursive definition of matrix algorithms. There is lots of revision material on Gaussian elimination that can be skipped, so do not worry about the large number of pages.

### Week 2

**Keywords:** Cholesky factorisation, fill-in, fill-in reducing permutation, minimum degree ordering.

**Homework:** Problems, 29, 30, 35 and 37. In addition, solve any additional four problems from 28-37 to gain extra points.

[See outline of Week 2 in Youtube](#)

**Pages:** 37-52.

**Synopsis:** The topic of the second week is Cholesky factorisation of sparse matrices. First, we prove existence of the Cholesky factorisation for s.p.d. matrices without taking sparsity into account. Our existence proof uses block matrix notation and induction with respect to dimension of the matrix. Unfortunately, the Cholesky factor of a sparse matrix can be dense. To mitigate this, we discuss methods for predicting location of non-zero entries in the factor without actually computing it. Then we introduce the minimum degree ordering method with the aim of obtaining a sparse factor by permuting the matrix before computing its Cholesky factorisation.

### Week 3

**Keywords:** Numerical stability analysis, Backward error analysis, floating-point representation, floating-point arithmetic model, round-off error,

[See outline of Week 3 in Youtube](#)

**Homework:** Problems, 40, 41, 43 and 47. In addition, solve any additional four problems from 38-47 to gain extra points.

**Pages:** 53-66.

**Synopsis:** A computer can perform billions of arithmetic operations when computing the Cholesky factorisation of a large matrix. When double precision floating-point numbers are used, as often is the case, all of these operations are computed slightly inaccurately. Hence, the computed Cholesky factor is an approximation of the exact factor. During Week 3, we develop tools used to study the accuracy of solutions to linear systems computed using such approximate Cholesky factorisation and back-substitution. We begin by outline, then discuss perturbation theory, derive a model for floating-point arithmetic errors, and develop technical estimates we need later.

### Week 4

**Keywords:** Numerical stability analysis, Backward error analysis, Back-substitution, Cholesky factorisation,  $QR$ -factorisation, Givens Rotation.

[See outline of Week 4 in Youtube](#)

**Homework:** Problems P53, P54, P55, and P56. In addition, solve any additional four problems from 48-57 to gain extra points.

**Pages:** 67-85.

**Synopsis:** During week 4, we give two examples on numerical stability analysis. First, we estimate the error due to solving  $2 \times 2$  - linear system with upper triangular coefficient matrix using the back-substitution method in floating-point representation. Then we study replacing  $A$  in linear system  $A\mathbf{x} = \mathbf{b}$  by  $\hat{L}\hat{L}^T$  where  $\hat{L}$  is the Cholesky factor of  $A$  computed in floating-point representation. In both cases, we formulate a linear system for the floating-point solution and obtain error estimate by perturbation theory. This requires us to bound the relative error due to floating-point arithmetic errors. We also discuss a method for computing numerically stable  $QR$

factorisation.

## Week 5

**Keywords:** Iterative solution method, Fixed-point iteration, Convergence, Conjugate Gradient method, Line search method, Gradient Descend.

**Homework:** Problems P59, P66, P67, and P68. In addition, solve any additional four problems from 58-68 to gain extra points.

[See video related to P66 and P63 in Youtube](#)

**Pages:** 87-102. **Synopsis:** During week 5 we discuss iterative solution

[See outline of Week 5 in Youtube](#)

methods for approximately solving linear system  $A\mathbf{x} = \mathbf{b}$ . Iterative solution method is a process generating a sequence  $\{\mathbf{x}_i\} \subset \mathbb{R}^n$  such that  $\mathbf{x}_i$  converges to the solution  $\mathbf{x}$ . When sufficiently accurate approximation has been obtained, the iteration is stopped. Iterative methods are based on various principles. First, we discuss methods based on fixed point techniques. Then we assume that  $A$  is s.p.d. and show that solving the linear system is equivalent with finding the global minimizer of quadratic functional. We end the week by deriving the Conjugate Gradient method as a line search minimisation iteration applied to this functional.

## Week 6

**Keywords:** Iterative solution method, Conjugate Gradient method, orthogonal projection, error estimate, Krylov subspace.

**Homework:** Problems P70, P72, P74, and P76. In addition, solve any additional four problems from 69-76 to gain extra points.

**Pages:** 103-113.

[See outline of Week 6 in Youtube](#)

**Synopsis:** During week 6 we give an alternative point-of-view to conjugate gradient method. We show that iterates generated by CG are  $A$ -orthogonal projections of the exact solution to certain subspaces of  $\mathbb{R}^n$ . Surprisingly, these projections can be computed without knowledge of the exact solution. We derive an error estimate for CG and discuss how convergence can be improved by using a preconditioner.



# Chapter 1

## Direct solution of sparse linear systems

In this Chapter, we study solution methods for **linear systems**: Find  $\mathbf{x} \in \mathbb{R}^n$  s.t.

$$A\mathbf{x} = \mathbf{b}, \quad (1.1)$$

where  $\mathbf{b} \in \mathbb{R}^n$  and the coefficient matrix  $A \in \mathbb{R}^{n \times n}$  is **large, sparse, symmetric and positive definite (s.p.d.)**. By *sparse matrix*, we mean a matrix with mostly zero entries. If a matrix is not sparse it is called as a *dense matrix*. The definition of positive definiteness will be given at the beginning of Section 1.7.

Large, sparse, s.p.d. coefficient matrices are related, e.g., to solution of partial differential equations (PDEs) using finite element method (FEM) or finite difference method (FDM). For example, application of FDM to two dimensional Laplace operator leads to a coefficient matrix having at most five non-zero entries on every row. If accurate discretisation is required, the dimension of these coefficient matrices can be of the order  $n \approx 10^5 - 10^6$ .

We use the sparse Cholesky factorisation to solve (1.1). In sparse Cholesky factorisation, sparse, s.p.d. matrix  $A \in \mathbb{R}^{n \times n}$  is decomposed as

$$P^T A P = L L^T, \quad (1.2)$$

where  $P \in \mathbb{R}^{n \times n}$  is a *permutation* matrix and  $L \in \mathbb{R}^{n \times n}$  is a lower triangular matrix. As a permutation matrix  $P$  is invertible, and equation (1.1) is equivalent to

$$P^T A P P^{-1} \mathbf{x} = P^T \mathbf{b} \quad \text{and} \quad L L^T P^{-1} \mathbf{x} = P^T \mathbf{b}.$$

Hence, the solution of (1.1) is obtained by solving the auxiliary problems

$$L\mathbf{z} = P^T \mathbf{b}, \quad L^T \mathbf{y} = \mathbf{z}, \quad \text{and setting } \mathbf{x} = P\mathbf{y}.$$

As  $L$  is a lower triangular matrix, the first two equations above are solved using back-substitution.

If  $P = I$  in (1.2), it becomes the Cholesky factorisation of  $A$  that is related to the Gaussian elimination process. Recall that writing the row-operations conducted during the Gaussian elimination process using elimination matrices yields the  $LU$ -factorisation of the coefficient matrix. In  $LU$ -factorisation, matrix  $A$  is written as  $A = LU$  where  $L$  is a lower triangular and  $U$  an upper triangular matrix. The Cholesky factorisation is derived using the same elimination matrices but taking advantage of symmetry and positive definiteness of  $A$ . In sparse Cholesky factorisation, additional permutations are used to obtain a sparse factor  $L$  for a sparse matrix  $A$ .

To convince the reader that sparse matrices appear in practice, we begin this Chapter by application of finite difference method to solution of the Poisson's equation that results in a linear system with a sparse, s.p.d. coefficient matrix. Next, we discuss how sparse matrices are stored in the memory of a computer. Then we prepare to prove existence of the Cholesky factorisation by recalling the Gaussian elimination process and  $LU$ -factorisation. Our existence proof uses block matrix notation that is discussed next. Finally, we show existence of the Cholesky factorisation and introduce minimum degree ordering method for obtaining a sparse factor  $L$  for a sparse matrix  $A$ . We end the section by studying numerical stability or accuracy of solving linear systems using Cholesky factorisation computed using floating-point numbers.

## 1.1 Preliminaries

### 1.1.1 Permutation matrices

See video on permutation matrices in Youtube

In this section, we discuss permutation matrices that encode information on changing the order of rows or the columns of a matrix. Vector  $\mathbf{p} \in \mathbb{R}^n$  is called as a *permutation vector*, if its entries satisfy the conditions:  $p_i \in \{1, \dots, n\}$  and  $p_i \neq p_j$  for all  $i, j \in \{1, \dots, n\}$ ,  $i \neq j$ . This is, a permutation vector is a re-ordering of  $[1 \ \cdots \ n]$ . Matrix  $P \in \mathbb{R}^{n \times n}$  is called as a *permutation matrix*, if

$$P = [\mathbf{e}_{p_1} \ \cdots \ \mathbf{e}_{p_n}] \quad \text{where } \mathbf{p} \in \mathbb{R}^n \text{ is a permutation vector.}$$

As  $P$  has orthonormal columns it is unitary, i.e.,  $P^{-1} = P^T$ .

Let  $P \in \mathbb{R}^{n \times n}$  be a permutation matrix corresponding to permutation vector  $p \in \mathbb{R}^n$  and split  $A, B \in \mathbb{R}^{n \times n}$  into column and row vectors as

$$A = [\mathbf{a}_1 \quad \cdots \quad \mathbf{a}_n] \quad \text{and} \quad B = \begin{bmatrix} \mathbf{b}_1^T \\ \vdots \\ \mathbf{b}_n^T \end{bmatrix}.$$

Recall that  $\mathbf{e}_i^T A$  and  $A \mathbf{e}_i$  are the  $i$ th row and column of a matrix  $A \in \mathbb{R}^{n \times n}$ , respectively. By direct computation

$$AP \mathbf{e}_i = A \mathbf{e}_{p_i} = \mathbf{a}_{p_i} \quad \text{and} \quad \mathbf{e}_i^T P^T B = (P \mathbf{e}_i)^T B = \mathbf{e}_{p_i}^T B = \mathbf{b}_{p_i}^T.$$

Hence, these operations reorder the columns and rows according to permutation vector  $\mathbf{p}$ , this is,

$$AP = [\mathbf{a}_{p_1} \quad \cdots \quad \mathbf{a}_{p_n}] \quad \text{and} \quad P^T B = \begin{bmatrix} \mathbf{b}_{p_1}^T \\ \vdots \\ \mathbf{b}_{p_n}^T \end{bmatrix}.$$

**Example 1.1.** The permutation matrix changing rows 2 and 3 of a  $3 \times 3$ -matrix is related to the permutation vector is  $\mathbf{p} = [1 \quad 3 \quad 2]$  and obtained simply as

$$P = [\mathbf{e}_1 \quad \mathbf{e}_3 \quad \mathbf{e}_2] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

### 1.1.2 Problems

P1. (0.5p) Let

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix}.$$

Find the permutation matrix corresponding to operations

- (a) Swap rows 2 and 3
- (b) Swap column 1 and 4
- (c) Order rows as 3, 2, 1

P2. (0.5p) Prove the claim:

Let  $A \in \mathbb{R}^{n \times n}$  have orthonormal column vectors. Then  $A$  is unitary.

## 1.2 Block matrix notation

*Block matrix notation is extensively used in this lecture note. Hence, this section should be studied with care.*

See [video introduction to block matrices in Youtube](#)

In this section, we introduce block matrix notation which is used to avoid index notation in proofs and derivations. We limit the discussion to  $2 \times 2$  block matrices, which are sufficient for our needs. Block matrices are obtained by splitting entries of a matrix vertically and horizontally into sub-matrices called blocks. In the following, we often divide matrices to  $2 \times 2$  matrix blocks. For example, split  $A \in \mathbb{R}^{n \times k}$  as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \text{where } n = n_1 + n_2, \text{ and } k = p + q.$$

In the above equation, the size of each sub-matrix is written under its symbol.

**Example 1.2.** Consider the block decomposition of  $3 \times 3$  matrix

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

to  $2 \times 2$  block matrix as

$$A = \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix} \quad \text{where } a_{11} = 1, \mathbf{a}_{12} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \mathbf{a}_{21} = \begin{bmatrix} 4 \\ 7 \end{bmatrix}, A_{22} = \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}.$$

This is, we have sliced  $A$  as  $\left[ \begin{array}{c|cc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right]$ .

We proceed to derive  $2 \times 2$  block-matrix-matrix-product formula. Let  $A \in \mathbb{R}^{n \times k}$ ,  $B \in \mathbb{R}^{k \times m}$ , and recall the matrix-matrix product formula

$$AB \in \mathbb{R}^{n \times m} \quad \text{and} \quad (AB)_{ij} = \sum_{l=1}^k a_{il}b_{lj}.$$

Matrices are often written using their column and row vectors as

$$A = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} \quad \text{and} \quad B = [\mathbf{b}_1 \quad \cdots \quad \mathbf{b}_m],$$



where  $\{\mathbf{a}_i\}_{i=1}^n \subset \mathbb{R}^k$  and  $\{\mathbf{b}_i\}_{i=1}^m \subset \mathbb{R}^k$ . Observe, that we use column vectors, hence,  $\mathbf{a}_1^T$  is a row vector. Using row and column vectors, the matrix-matrix product  $AB$  can be written as

$$AB = [\mathbf{A}\mathbf{b}_1 \quad \cdots \quad \mathbf{A}\mathbf{b}_m] = \begin{bmatrix} \mathbf{a}_1^T B \\ \vdots \\ \mathbf{a}_n^T B \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{b}_1 & \cdots & \mathbf{a}_1^T \mathbf{b}_m \\ \vdots & \ddots & \vdots \\ \mathbf{a}_n^T \mathbf{b}_1 & \cdots & \mathbf{a}_n^T \mathbf{b}_m \end{bmatrix}. \quad (1.3)$$

Using the above formula gives a Lemma for computing  $2 \times 2$  block-matrix-matrix-product:

**Lemma 1.1.** Let  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \in \mathbb{R}^{n \times k}$  and  $B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \in \mathbb{R}^{k \times m}$ . Then

$$AB = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}. \quad (1.4)$$

Observe, that the  $2 \times 2$  block-matrix-matrix product  $AB$  is computed similar to the  $2 \times 2$  matrix-matrix product. This holds in general for all block-matrix-matrix-products. The sizes of matrix blocks must match in the sense that all products appearing in (1.4) are well defined. We prove Lemma 1.1 after giving a helper result.

**Lemma 1.2.** Let  $\begin{bmatrix} C & D \end{bmatrix} \in \mathbb{R}^{n \times k}$  and  $\begin{bmatrix} F \\ G \end{bmatrix} \in \mathbb{R}^{k \times m}$  for  $k = p + q$ .

Then

$$\begin{bmatrix} C & D \end{bmatrix} \begin{bmatrix} F \\ G \end{bmatrix} = CF + DG. \quad (1.5)$$

Observe that the sizes of matrix blocks match in the sense that products  $CF$  and  $DG$  are well defined.

*Proof.* Denote the row vectors of  $C, D$  and column vectors of  $F, G$  as

$$C = \begin{bmatrix} \mathbf{c}_1^T \\ \vdots \\ \mathbf{c}_n^T \end{bmatrix}, \quad D = \begin{bmatrix} \mathbf{d}_1^T \\ \vdots \\ \mathbf{d}_n^T \end{bmatrix}, \quad F = [\mathbf{f}_1 \quad \cdots \quad \mathbf{f}_m], \quad \text{and} \quad G = [\mathbf{g}_1 \quad \cdots \quad \mathbf{g}_m].$$

[See video on computing product of  \$2 \times 2\$  matrices in Youtube](#)

[See video on proving the product formula of  \$2 \times 2\$  matrices in Youtube](#)

We proceed to give a formula for computing entries of the product matrix  $\begin{bmatrix} C & D \end{bmatrix} \begin{bmatrix} F \\ G \end{bmatrix} \in \mathbb{R}^{n \times m}$ . The entry  $ij$  of the product matrix is obtained as

$$\mathbf{e}_i^T \begin{bmatrix} C & D \end{bmatrix} \begin{bmatrix} F \\ G \end{bmatrix} \mathbf{e}_j$$

where  $\mathbf{e}_i \in \mathbb{R}^n$  and  $\mathbf{e}_j \in \mathbb{R}^m$  are the  $i$ th and  $j$ th unit vectors. A direct calculation

$$\mathbf{e}_i^T \begin{bmatrix} C & D \end{bmatrix} \begin{bmatrix} F \\ G \end{bmatrix} \mathbf{e}_j = \begin{bmatrix} \mathbf{c}_i^T & \mathbf{d}_i^T \end{bmatrix} \begin{bmatrix} \mathbf{f}_j \\ \mathbf{g}_j \end{bmatrix} = \mathbf{c}_i^T \mathbf{f}_j + \mathbf{d}_i^T \mathbf{g}_j = (CF)_{ij} + (DG)_{ij}$$

gives the formula

$$\begin{bmatrix} C & D \end{bmatrix} \begin{bmatrix} F \\ G \end{bmatrix} = CF + DG. \quad (1.6)$$

□

*Proof of Lemma 1.1.* To prove (1.4) observe that by (1.3)

$$AB = \begin{bmatrix} \begin{bmatrix} A_{11} & A_{12} \end{bmatrix} \begin{bmatrix} B_{11} \\ B_{21} \end{bmatrix} & \begin{bmatrix} A_{11} & A_{12} \end{bmatrix} \begin{bmatrix} B_{12} \\ B_{22} \end{bmatrix} \\ \begin{bmatrix} A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} \\ B_{21} \end{bmatrix} & \begin{bmatrix} A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{12} \\ B_{22} \end{bmatrix} \end{bmatrix}.$$

Application of product formula (1.5) completes the derivation. □

[See video on Example 1.3 in Youtube](#)

**Example 1.3.** *Next, we illustrate how block matrix notation is used in proofs and show that the product of two  $n \times n$  lower triangular matrices is a lower triangular matrix. We formulate an induction proof with respect to the dimension of the lower triangular matrix using suitable  $2 \times 2$  block division.*

**Base step  $n = 1$ :** *Trivially true.*

**Induction assumption:** *Product of two  $k \times k$  lower triangular matrices is lower triangular.*

**Induction step:** Let  $L, T \in \mathbb{R}^{(k+1) \times (k+1)}$  be lower triangular matrices.  
*Split*

$$L = \begin{bmatrix} l_{11} & 0 \\ 1 \times 1 & \\ \mathbf{l}_{21} & L_{22} \\ k \times 1 & k \times k \end{bmatrix} \quad \text{and} \quad T = \begin{bmatrix} t_{11} & 0 \\ 1 \times 1 & \\ \mathbf{t}_{21} & T_{22} \\ k \times 1 & k \times k \end{bmatrix},$$

where  $L_{22}, T_{22}$  lower triangular matrices. Using the  $2 \times 2$  block matrix-matrix product formula gives

$$LT = \begin{bmatrix} l_{11}t_{11} & 0 \\ \mathbf{l}_{21}t_{11} + L_{22}\mathbf{t}_{21} & L_{22}T_{22} \end{bmatrix}.$$

By induction assumption  $L_{22}T_{22}$  is lower triangular matrix, which completes the proof.

### 1.2.1 Problems

P3. (1p) Let

$$A = \begin{bmatrix} A_{11} & 0 \\ n \times n & \\ A_{21} & A_{22} \\ m \times n & m \times m \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} B_{11} & 0 \\ n \times n & \\ B_{21} & B_{22} \\ m \times n & m \times m \end{bmatrix}.$$

- (a) Compute the block-matrix-matrix product  $AB$ .
- (b) Find the inverse matrix of  $A$ . Hint: find  $B_{11}, B_{21}, B_{22}$  such that

$$\begin{bmatrix} A_{11} & 0 \\ n \times n & \\ A_{21} & A_{22} \\ m \times n & m \times m \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ n \times n & \\ B_{21} & B_{22} \\ m \times n & m \times m \end{bmatrix} = \begin{bmatrix} I & 0 \\ n \times n & \\ 0 & I \\ m \times m & \end{bmatrix}.$$

List assumptions (if any) that you have to make on  $A_{11}, A_{21}$ , and  $A_{22}$ .

- (c) Argue that  $\det A = 0$  implies that either  $\det A_{11} = 0$  or  $\det A_{22} = 0$ .

P4. (1p) Let  $E = \begin{bmatrix} 1 & 0 \\ 1 \times 1 & 1 \times n \\ -\mathbf{a}_{21} & I \\ n \times 1 & n \times n \end{bmatrix}.$

- (a) Compute the product  $E \begin{bmatrix} 1 & \mathbf{a}_{12}^T \\ 1 \times n & \\ \mathbf{a}_{21} & A_{22} \\ n \times 1 & n \times n \end{bmatrix}$

- (b) Find the inverse matrix of  $E$  using the formula derived in the previous problem. Check that your inverse is correct by computing the product  $EE^{-1}$ .

P5. (2p)

- (a) Show that

$$\det \begin{bmatrix} I & 0 \\ 0 & A_{22} \\ n \times n & m \times m \end{bmatrix} = \det A_{22}.$$

Hint: recall the Laplace expansion for computing determinants and use induction with respect to parameter  $n$ .

- (b) Modify the proof in (a) to show that

$$\det \begin{bmatrix} I & A_{12} \\ 0 & A_{22} \\ n \times n & n \times m \\ & m \times m \end{bmatrix} = \det A_{22}. \quad (1.7)$$

P6. (0.5p)

- (a) Compute  $\begin{bmatrix} 0 & I \\ I & 0 \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & A_{22} \end{bmatrix} \begin{bmatrix} 0 & I \\ I & 0 \end{bmatrix}$  and  $\begin{bmatrix} 0 & I \\ I & 0 \end{bmatrix} \begin{bmatrix} 0 & I \\ I & 0 \end{bmatrix}$ .
- (b) Use properties of determinant, Problem 3, and (a) to show that  $\det \begin{bmatrix} A_{11} & 0 \\ 0 & I \end{bmatrix} = \det A_{11}$ .

P7. (1p) Consider the block matrix  $A = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \\ n \times n & n \times m \\ & m \times m \end{bmatrix}$ , where  $A_{11}$  and  $A_{22}$  are invertible matrices.

- (a) Compute the product

$$\begin{bmatrix} A_{11} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} I & A_{11}^{-1} A_{12} A_{22}^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & A_{22} \end{bmatrix}.$$

- (b) Use, equation (1.7), Problems 3,4, and decomposition in (a) to show that  $\det A = \det A_{11} \det A_{22}$ .
- (c) Argue by Problem 3 that  $\det A = \det A_{11} \det A_{22}$  even if  $A_{11}$  or  $A_{22}$  are not invertible.

P8. (1p) Let

$$M = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}. \quad (1.8)$$

- (a) Use suitable  $2 \times 2$  block decomposition to compute  $M^2$ .
- (b) Use inverse matrix formula from Problem 3 to compute  $M^{-1}$ .

### 1.2.2 Back-substitution in block matrix notation

This section gives a **recursive definition of the back-substitution** algorithm. Using recursion is necessary to express the algorithm in block matrix notation. This section should be studied with care.

In this section, we use block matrix notation to define the back - substitution algorithm. Our definition is recursive with respect to dimension of the linear system. Using such definition allows simple treatment of matrices with different dimension using the block matrix notation. We use similar techniques to study the  $LU$  and the Cholesky factorisations.

[See video on solution of upper triangular systems in Youtube](#)

Consider the linear system: Find  $\mathbf{x} \in \mathbb{R}^n$  satisfying

$$U\mathbf{x} = \mathbf{b},$$

where the coefficient matrix  $U \in \mathbb{R}^{n \times n}$  is upper triangular and  $\mathbf{b} \in \mathbb{R}^n$ .

**Definition 1.1.** Matrix  $U \in \mathbb{R}^{n \times n}$  is upper triangular, if

$$U_{ij} = 0 \quad \text{for } i > j.$$

This is

$$U = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ & u_{22} & \cdots & u_{2n} \\ & & \ddots & \vdots \\ & & & u_{nn} \end{bmatrix} \quad \text{or} \quad U = \begin{bmatrix} \# & \# & \cdots & \# \\ & \# & \cdots & \# \\ & & \ddots & \vdots \\ & & & \# \end{bmatrix}.$$

Here we use notational convention where the location of non-zero entries in the matrix is indicated by  $\#$  and zero entries are omitted. Such convention

is used when the location of non-zero entries is important but their value is not.

Triangular linear systems are solved using back-substitution algorithm. We use a definition that is recursive with respect to the dimension of the coefficient matrix. The function *triusolve*( $U, \mathbf{b}$ ) returns solution to linear system  $U\mathbf{x} = \mathbf{b}$  for invertible upper triangular matrix  $U \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ .

---

For  $n = 1$ ,  $\text{triusolve}(U, b) = \frac{b}{U}$ .

For  $n > 1$ , we use a recursive definition. First, split the linear system  $U\mathbf{x} = \mathbf{b}$  as

$$\begin{bmatrix} U_{11} & \mathbf{u}_{12} \\ 0 & u_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ b_2 \end{bmatrix} \quad (1.9)$$

where  $x_2, b_2$  ja  $u_{22}$  are scalars,  $U_{11} \in \mathbb{R}^{(n-1) \times (n-1)}$  and  $\mathbf{u}_{12}, \mathbf{x}_1, \mathbf{b}_1 \in \mathbb{R}^{n-1}$ . As  $U$  is invertible,  $u_{22} \neq 0$ ,  $U_{11}$  is invertible<sup>1</sup>, and

$$x_2 = \frac{b_2}{u_{22}}.$$

First equation in (1.9) states  $U_{11}\mathbf{x}_1 = \mathbf{b}_1 - \mathbf{u}_{12}x_2$ . As coefficient matrix  $U_{11} \in \mathbb{R}^{(n-1) \times (n-1)}$  is invertible and upper triangular,  $\mathbf{x}_1$  is obtained recursively as  $\mathbf{x}_1 = \text{triusolve}(U_{11}, \mathbf{b}_1 - \mathbf{u}_{12}x_2)$ . Hence,

$$\text{triusolve}(U, b) = \begin{bmatrix} \mathbf{x}_1 \\ x_2 \end{bmatrix}.$$

See a video on implementing the back substitution algorithm in Youtube

---

An example implementation of the above function is given below.

```
function x = triusolve2(U,b)

n = size(U,2);
x = zeros(n,1);

% Define matrix and vector blocks.
U11 = U(1:(n-1),1:(n-1));
u12 = U(1:(n-1),n);
u22 = U(n,n);

b1 = b(1:(n-1));
b2 = b(n);
```

---

<sup>1</sup>See problem 7 on page 12

```

% solve x2.
x(n) = b2/u22;

if( n > 1 )
% solve x1 using recursive function call.
x(1:(n-1)) = triusolve2(U11,b1-u12*x(n));
end

end

```

Using recursive function calls is not very efficient. A better strategy is to update the vector **b** during the algorithm and use a for-loop to conduct the computation. An example implementation using such *update strategy* is given below.

```

function x = triusolve(U,b)

N = size(U,2);

x = zeros(N,1);

for n=N:-1:1
% Define matrix and vector blocks.

U11 = U(1:(n-1),1:(n-1));
u12 = U(1:(n-1),n);
u22 = U(n,n);

b1 = b(1:(n-1));
b2 = b(n);

% solve x(i).
x(n) = b2/u22;

% update vector b
b(1:(n-1)) = b1 - u12*x(n);
end

```

The above algorithm can be easily modified to solve lower triangular linear systems.

### 1.2.3 Problems

- P9. (2p) Use block matrix notation to give a recursive definition of function *trilsolve*(*L*, **b**) that returns solution of linear system  $L\mathbf{x} = \mathbf{b}$  where *L* is a lower triangular matrix.

P10. (2p)

- (a) Give a recursive implementation of *trilsolve* in Matlab
- (b) Modify recursive implementation in (a) to use the update strategy.

Device a test verifying that both of your implementations are correct.

P11. (1p)

- (a) Compute, how many arithmetic operations are needed to solve a  $N \times N$  - upper triangular system.
- (b) Measure the time required to solve upper triangular linear systems using Matlab backslash, back substitution using recursive implementation, and back substitution using update strategy. Generate random upper triangular matrices with dimension  $N = 10, 50, 100, 200, 300, 400$ , and 500 using commands `rand` and `triu`. For each dimension, compute average solution time for each method from 100 solves. Plot average solution times as a function of  $N$  using a logarithmic scale. Does the result correspond to (a) ?

### 1.3 Finite difference method

*This section gives an example application that leads to linear system with large, sparse and s.p.d coefficient matrix. It is extra material and can be skipped. Or just have a look at the video.*

See [video introduction to finite difference method](#)

Let  $\Omega \subset \mathbb{R}^2$  be a bounded open set with sufficiently regular boundary and recall the definition of the Laplace operator  $\Delta$  in  $\mathbb{R}^2$ ,

$$\Delta := \frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2}.$$

The Poisson's equation in  $\Omega$  is: Find  $u \in C^2(\Omega) \cap C(\overline{\Omega})$  such that

$$\begin{cases} -\Delta u = f & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases} \quad (1.10)$$

where  $f$  is a given function<sup>2</sup>. The Poisson's equation is a simple model problem for other PDEs that appear, e.g., in electrical or mechanical engineering.

---

<sup>2</sup>Here  $C^2(\Omega)$  and  $C(\overline{\Omega})$  are spaces of functions that have two derivatives in open set  $\Omega$  and functions that are continuous in closure of  $\Omega$ , respectively. The differentiability is required for the equation  $-\Delta u = f$  to be well defined, and continuity up to boundary for the boundary condition  $u = 0$  to be meaningful



Several different numerical methods have been developed to find approximate solutions to (1.10). We use the finite difference method, in which one seeks for an approximation to the point-wise values of  $u$ . The first step is to derive the central difference approximation of the Laplace operator.

Let  $h \in \mathbb{R}$ ,  $h > 0$ . The Taylor expansion<sup>3</sup> of the solution  $u$  with respect to the variable  $x_1$  gives

$$\begin{aligned} u(x_1 + h, x_2) &= u(x_1, x_2) + \frac{\partial u}{\partial x_1}(x_1, x_2)h + \frac{1}{2} \frac{\partial^2 u}{\partial x_1^2}(x_1, x_2)h^2 + \frac{1}{6} \frac{\partial^3 u}{\partial x_1^3}(x_1, x_2)h^3 + h.o.t. \\ u(x_1 - h, x_2) &= u(x_1, x_2) - \frac{\partial u}{\partial x_1}(x_1, x_2)h + \frac{1}{2} \frac{\partial^2 u}{\partial x_1^2}(x_1, x_2)h^2 - \frac{1}{6} \frac{\partial^3 u}{\partial x_1^3}(x_1, x_2)h^3 + h.o.t., \end{aligned}$$

where  $h.o.t$  is used to denote higher order terms with respect to  $h$ . Subtracting the two above equations and dividing by  $h^2$  gives

$$\frac{\partial^2 u}{\partial x_1^2}(x_1, x_2) \approx \frac{u(x_1 + h, x_2) - 2u(x_1, x_2) + u(x_1 - h, x_2)}{h^2}. \quad (1.11)$$

Similar computations for the  $x_2$  - component give

$$\frac{\partial^2 u}{\partial x_2^2}(x_1, x_2) \approx \frac{u(x_1, x_2 + h) - 2u(x_1, x_2) + u(x_1, x_2 - h)}{h^2}. \quad (1.12)$$

Combining (1.11) and (1.12) yields the *central difference approximation* of the Laplace operator:

$$(\Delta u)(x_1, x_2) \approx \frac{u(x_1 - h, x_2) + u(x_1 + h, x_2) - 4u(x_1, x_2) + u(x_1, x_2 - h) + u(x_1, x_2 + h)}{h^2}.$$

The accuracy of this approximation depends on  $h$  as well as on the properties of the function  $u$ .

Next, consider the domain  $\Omega = (0, 1)^2$  and a uniform  $N \times N$ -grid composed of points

$$\mathbf{x}_{ij} = \frac{1}{N-1} \begin{bmatrix} i-1 \\ j-1 \end{bmatrix} \quad \text{for } i, j \in \{1, \dots, N\}$$

see Figure 1.1. The distance between grid points is denoted by  $h := \frac{1}{N-1}$  and the value of  $u$  at the grid point  $\mathbf{x}_{ij}$  by  $\mathbf{u}_{ij} := u(\mathbf{x}_{ij})$ .

Observe that the indices of interior grid points  $\mathbf{x}_{ij} \in \Omega$  and boundary grid points  $\mathbf{x}_{ij} \in \partial\Omega$  are

$$I := \{ (i, j) \mid i, j \in \{2, \dots, N-1\} \}$$

---

<sup>3</sup>Observe that the expansion requires additional regularity of  $u$ , i.e  $u \in C^3(\Omega)$ .

and

$$B := \{ (i, j) \mid i, j \in \{1, \dots, N\} \} \setminus I,$$

respectively. At interior grid points, the finite difference approximation states that:

$$\frac{\mathbf{u}_{(i-1)j} + \mathbf{u}_{(i+1)j} + \mathbf{u}_{i(j-1)} + \mathbf{u}_{i(j+1)} - 4\mathbf{u}_{ij}}{h^2} \approx f(\mathbf{x}_{ij}). \quad (1.13)$$

Due to the boundary condition  $u = 0$  on  $\partial\Omega$ ,

$$\mathbf{u}_{ij} = 0 \quad (1.14)$$

at boundary grid points.

In finite difference method, one poses (1.13) as equality and seeks for *approximate point wise values of  $u$  satisfying* the resulting linear system. For notional simplicity, we denote the FD-approximation also by  $\mathbf{u}_{ij}$ . The challenge in solving  $\mathbf{u}_{ij}$  is constructing the coefficient matrix of the linear system (1.13)-(1.14), which requires careful index handling. First, collect the variables  $\mathbf{u}_{ij}$  into the vector  $\mathbf{U} \in \mathbb{R}^{N^2}$  as

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}_{11} \\ \mathbf{u}_{12} \\ \mathbf{u}_{13} \\ \vdots \\ \mathbf{u}_{21} \\ \mathbf{u}_{22} \\ \mathbf{u}_{23} \\ \vdots \end{bmatrix}$$

It is helpful to explicitly define mapping  $\sigma(i, j) = (i - 1)N + j$  that aids in index handling. The value  $\mathbf{u}_{ij}$  resides in the element  $\sigma(i, j)$  of vector  $\mathbf{U}$ . The vector  $\mathbf{U}$  satisfies

$$A\mathbf{U} = \mathbf{b}.$$

The non-zero entries of the coefficient matrix  $A \in \mathbb{R}^{N^2 \times N^2}$  and vector  $\mathbf{b} \in \mathbb{R}^{N^2}$  are:

$$\begin{aligned} a_{\sigma(i,j)\sigma(i-1,j)} &= 1, & a_{\sigma(i,j)\sigma(i+1,j)} &= 1, \\ a_{\sigma(i,j)\sigma(i,j-1)} &= 1, & a_{\sigma(i,j)\sigma(i,j+1)} &= 1, \\ a_{\sigma(i,j)\sigma(i,j)} &= -4, & b_{\sigma(i,j)} &= f(\mathbf{x}_{ij}). \end{aligned}$$

for interior indices  $i, j \in I$  and

$$a_{\sigma(i,j)\sigma(i,j)} = 1, \quad b_{\sigma(i,j)} = 0$$

for boundary indices  $i, j \in B$ . The matrix  $A$  is assembled in the following code.

```

N = 50;
A = sparse( N^2,N^2);
h = 1/(N-1);

ijmap = @(i,j) ( (i-1)*N + j);
active = []; % collect not boundary nodes here.

for i=1:N
    for j=1:N

        x(i,j) = (i-1)/(N-1); y(i,j) = (j-1)/(N-1);

        if( (i > 1) & (i < N) & ( j > 1) & ( j < N))

            % This is the row corresponding to point (i,j)
            I1 = ijmap(i,j);

            active = [active I1];

            A(I1, ijmap(i-1,j)) = -1/h^2;
            A(I1, ijmap(i+1,j)) = -1/h^2;
            A(I1, ijmap(i,j-1)) = -1/h^2;
            A(I1, ijmap(i,j+1)) = -1/h^2;
            A(I1, I1) = 4/h^2;

            b(I1,1) = 1;
        end
    end
end

% system without active rows
A = A(active,active);
b = b(active);

% solve !
u = zeros(N^2,1);
u(active) = A\b;

% visualize u.
U = reshape(u,N,N);
figure;S = surf(x',y',U);
set(S,'facecolor','interp');
```

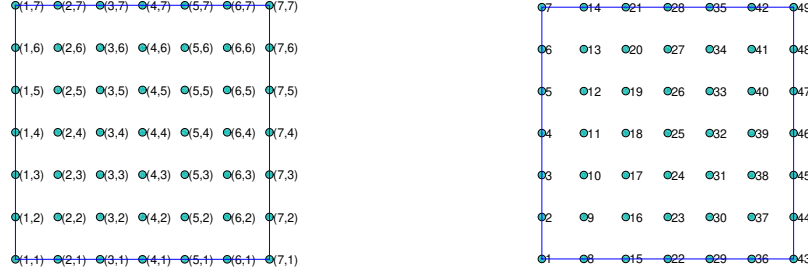


Figure 1.1: Node numbering in  $i, j$  - system vs. node numbering corresponding to vector  $\mathbf{U}$

The rows of  $A$  related to boundary indices are not interesting and they are eliminated. Let  $P \in \mathbb{R}^{N^2 \times N^2}$  be a permutation matrix ordering the rows of  $U$  as

$$P^T \mathbf{U} = \begin{bmatrix} \mathbf{U}_I \\ \mathbf{U}_B \end{bmatrix}$$

where  $\mathbf{U}_I \in \mathbb{R}^{(N-2)^2}$  and  $\mathbf{U}_B \in \mathbb{R}^{4(N-1)}$  are the values of  $\mathbf{u}_{ij}$  related to interior and boundary grid points, respectively. Application of the same splitting to  $A$  and  $\mathbf{b}$  gives

$$P^T A T = \begin{bmatrix} A_{II} & A_{IB} \\ A_{BI} & A_{BB} \end{bmatrix} \quad \text{and} \quad P^T \mathbf{b} = \begin{bmatrix} \mathbf{b}_I \\ \mathbf{b}_B \end{bmatrix}.$$

As  $\mathbf{U}_B = 0$  by (1.14),  $\mathbf{U}_I$  satisfies the system  $A_{II} \mathbf{U}_I = \mathbf{b}_I$  where the matrix  $A_{II}$  depends on the permutation  $P$ . The matrix  $A_{II} \in \mathbb{R}^{N^2 \times N^2}$  is symmetric and has at most five non-zero entries on every column. Its sparsity structure, i.e. location of non-zero entries, generated by the above code is visualized in Figure 1.2 using the Matlab command `spy(A)`. The accuracy of the computed approximate point-wise values depends on  $h$ . If accurate solutions are sought for,  $h$  is small and the number of grid points  $N$  can be large. For example,  $N$  can be of the order  $N = 1000$ , which results to linear system with dimension  $(N - 2)^2 \approx 10^6$ .

### 1.3.1 Problems

P12. (1p)

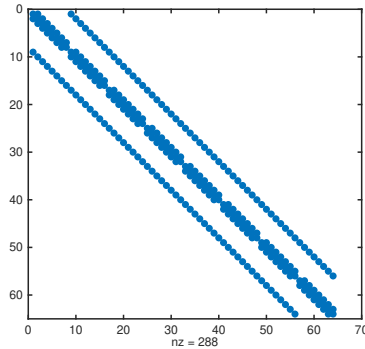


Figure 1.2: Nonzero entries of the matrix  $A_{II}$  related to the linear system given in equation (1.13).

- (a) Derive the finite difference approximation of Laplace operator in 1D.
- (b) Write a Matlab code to solve the 1D Poisson's equation: find  $u(x) \in C^2((0,1)) \cap C([0,1])$  satisfying

$$-u''(x) = 1 \text{ in } (0,1) \quad \text{and} \quad u(0) = u(1) = 0.$$

Plot the solution  $u$ .

P13. (2p) Let  $A \in \mathbb{R}^{2n \times 2n}$ ,  $n > 3$ , satisfy

$$A = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix}. \quad (1.15)$$

- (a) Let  $\mathbf{x}$  satisfy  $A\mathbf{x} = 0$ . Show that  $\mathbf{x}$  also satisfies

$$\begin{bmatrix} x_{i+1} \\ x_{i+2} \end{bmatrix} = C \begin{bmatrix} x_{i-1} \\ x_i \end{bmatrix} \quad \text{for } i \in \{1, \dots, 2n-2\} \quad \text{and} \quad C = \begin{bmatrix} -1 & 2 \\ -2 & 3 \end{bmatrix}$$

- (b) Use the Jordan decomposition of  $C$  to show that

$$\begin{bmatrix} x_{2n-1} \\ x_{2n} \end{bmatrix} = \begin{bmatrix} -2n+1 & 2n \\ -2n & 2n+1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- (c) Show that  $x_2$  and  $x_1$  satisfy

$$\begin{bmatrix} 2 & -1 \\ -2n-1 & 2n+2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0.$$

Use (b) to argue that  $N(A) = \{0\}$  and  $A$  is invertible.

P14. (1p) Consider the matrix  $A$  defined in (1.15).

- (a) Show by direct computation that  $\mathbf{x}^T A \mathbf{x} \geq 0$ , for any  $\mathbf{x} \in \mathbb{R}^n$  i.e.  $A$  is positive semi-definite matrix.
- (b) Argue that any symmetric and positive semi-definite matrix with a trivial null-space is positive definite, i.e. the inequality in part a) is strict. By *trivial null-space* we mean that  $A\mathbf{x} = 0$  if and only if  $\mathbf{x} = 0$ .
- (c) Use (b) and Problem 13 to argue that  $A$  is positive definite.

## 1.4 Compressed column storage format

*This section discusses sparse matrix storage formats used in practical implementation of sparse matrix data types. The aim is to highlight the fact that computational complexity of accessing matrix rows, columns, and elements depends on the chosen storage format. This has to be taken into account when designing high-level matrix algorithms. It also explains why sparse matrix literature gives several alternative ways to compute, e.g., the Cholesky factorisation. This Section is extra material and can be skipped.*

In this section, we discuss how sparse matrices are stored in the memory of a computer. The applied storage format affects the time required to access matrix elements which should be taken into account when designing sparse matrix algorithms.

[See video on CCS storage format in Youtube](#)

A dense matrix is typically stored as a two-dimensional array of numbers, whereas only non-zero entries of a sparse matrix are stored. There are several data structures used for this purpose, the most common ones being compressed row storage (CRS) and compressed column storage (CCS) formats. For example, Matlab uses CCS format to store sparse matrices.

The compressed column storage format uses three arrays:

- **Values:** List of matrix entries ordered column wise.
- **Row indices:** The row index for each of the entries
- **Column pointers:** Index of the first entry of a every column in the values and row index lists.

The CCS format is best illustrated by examples.

**Example 1.4.** *Let*

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}.$$

*In CCS format A is stored as*

$$\begin{aligned} vals &= [a_{11} \ a_{21} \ a_{12} \ a_{22}] \\ row\_ind &= [1 \ 2 \ 1 \ 2] \\ col\_ptr &= [1 \ 3 \ 5] \end{aligned}$$

**Example 1.5.** *Let*

$$B = \begin{bmatrix} -2 & 1 & 0 \\ 1 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix}.$$

In CCS format,  $B$  is stored as

$$\begin{aligned} \text{vals} &= [-2 \ 1 \ 1 \ -2 \ 1 \ 1 \ -2] \\ \text{row\_ind} &= [1 \ 2 \ 1 \ 2 \ 3 \ 2 \ 3] \\ \text{col\_ptr} &= [1 \ 3 \ 6 \ 8] \end{aligned}$$

In the above examples, the column pointer has an extra entry with value  $\text{length}(\text{vals})+1$  that is used to simplify implementation of matrix operations. If the extra entry is used, the column  $i$  is accessed simply as

```
A.col_ptr = [1 3 6 8];
A.rowind = [1 2 1 2 3 2 3];
A.val = [-2 1 1 -2 1 1 -2];

col_i = A.val( A.col_ptr(i):(A.col_ptr(i+1)-1) );
```

The CCS format has constant access time for columns of a matrix. Accessing rows requires looping over the row index array, hence the required time depends linearly on the size of the matrix. Element access is done by first accessing the column and then finding the desired entry. If the row indices are sorted, the desired entry can be sought for using, e.g., bisection search. In this case, the access time for the element  $ij$  has logarithmic dependency on the number of nonzero entries in the column  $j$ .

The access times in Matlab can be studied with the following test code. The resulting times are plotted in Figure 1.3

```
Nlist = floor(linspace(1,1e5,10));
row_timer = []; col_timer = []; ele_timer = [];

for n = Nlist

    e = ones(n,1);
    A = spdiags([e -2*e e], -1:1, n, n);

    I = randi(n,1e3,1); J = randi(n,1e3,1);

    T = tic;
    for j=1:1e3
        x=A(I(j),J(j));
    end
    ele_timer = [ele_timer toc(T)/1e3];

    T = tic;
```



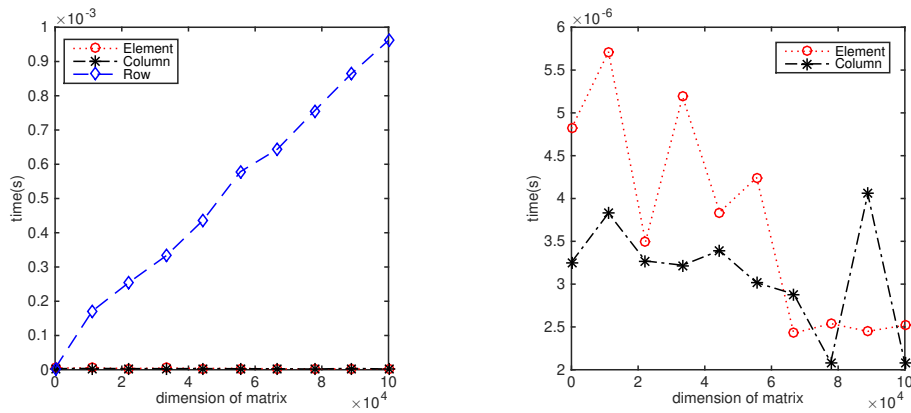


Figure 1.3: Example of access times for elements, rows, and columns of the one dimensional finite difference matrix  $A \in \mathbb{R}^{n \times n}$  in (1.15) as a function of the dimension  $n$ . The test is done in Matlab.

```

for j=1:1e3
    x=A(:,I(j));
end
col_timer = [col_timer toc(T)/1e3];

T = tic;
for j=1:1e3
    x=A(I(j),:);
end
row_timer = [row_timer toc(T)/1e3];

end

figure; plot(Nlist,ele_timer,'ro:',Nlist,col_timer,'k*-.',Nlist,row_timer,'bd--');
legend('Element','Column','Row');
ylabel('time(s)'); xlabel('dimension of matrix');

figure; plot(Nlist,ele_timer,'ro:',Nlist,col_timer,'k*-.');
legend('Element','Column');
ylabel('time(s)'); xlabel('dimension of matrix');
```

### 1.4.1 Additional material

- For more information on sparse matrices in Matlab, see  
John R. Gilbert, Cleve Moler, and Robert Schreiber. Sparse matrices in matlab: Design and implementation. *SIAM Journal on Matrix Analysis and Applications*, 13(1):333–356, 1992

### 1.4.2 Problems

P15. (0.5p) Let

$$A_1 := \begin{bmatrix} 1 & 0 & 2 & 0 \\ 3 & 0 & 4 & 0 \\ 0 & 5 & 0 & 6 \\ 7 & 8 & 9 & 10 \end{bmatrix}. \quad (1.16)$$

and

```
N = 5;
A2 = 2*eye(N) + diag(-ones(N-1,1),1) + diag(-ones(N-1,1),-1)
```

Write  $A_1$  and  $A_2$  using the compressed column storage scheme.

- P16. (1p) Write a Matlab-function `[val,row,col] = mat2ccs(A)` that returns the CCS representation of matrix  $A$ . Test your implementation using matrices  $A_1$  and  $A_2$  defined in Problem 15.
- P17. (1p) Write Matlab functions `coli = ccs_col(val,row,col,i)` and `rowi = ccs_row(val,row,col,i)` that return column and row  $i$  of a matrix represented in CCS format by  $val$ ,  $row$ , and  $col$ -vectors. Repeat the column and row access time test using your own functions.

## 1.5 Gaussian elimination

*This section is a review of the Gaussian elimination process. Read it to refresh your memory, or skip it.*

See [video introduction to Gaussian elimination in Youtube](#) Let  $A \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ , and consider the linear system: Find  $\mathbf{x} \in \mathbb{R}^n$  satisfying

$$A\mathbf{x} = \mathbf{b}. \quad (1.17)$$

Gaussian elimination is an algorithm that transforms (1.17) to the equivalent system: Find  $\mathbf{x} \in \mathbb{R}^n$  satisfying

$$U\mathbf{x} = \tilde{\mathbf{b}}, \quad (1.18)$$

where the coefficient matrix  $U \in \mathbb{R}^{n \times n}$  is upper triangular and  $\tilde{\mathbf{b}} \in \mathbb{R}^n$ . System (1.18) can be easily solved using the back substitution algorithm, see Section 1.2.2.

We proceed by applying the Gaussian elimination to (1.17) in its component form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n &= b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + \dots + a_{3n}x_n &= b_3. \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + a_{nn}x_n &= b_n \end{aligned} \quad (1.19)$$

For simplicity, assume that entry  $a_{11} \neq 0$ . The case  $a_{11} = 0$  is discussed in Section 1.5.1. The variable  $x_1$  is solved from the first equation in (1.19) as

$$x_1 = \frac{b_1}{a_{11}} - \sum_{j=2}^n \frac{a_{1j}}{a_{11}} x_j.$$

Using this expression, we eliminate variable  $x_1$  from equations  $\{2, \dots, n\}$  in (1.19). This yields new linear system for  $\mathbf{x}$ :

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 \\ a_{22}^{(2)}x_2 + a_{23}^{(2)}x_3 + \dots + a_{2n}^{(2)}x_n &= b_2^{(2)} \\ a_{32}^{(2)}x_2 + a_{33}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n &= b_3^{(2)} \\ &\vdots \\ a_{n2}^{(2)}x_2 + a_{n3}^{(2)}x_3 + \dots + a_{nn}^{(2)}x_n &= b_n^{(2)}, \end{aligned} \quad (1.20)$$

with coefficients  $a_{ij}^{(2)}$

$$a_{ij}^{(2)} = a_{ij} - \frac{a_{i1}}{a_{11}} a_{1j} \quad \text{for } i, j \in \{2, \dots, n\}.$$

This is, the transformed system is obtained by multiplying the first equation in (1.19) with  $-a_{i1}a_{11}^{-1}$  and adding it to the equation  $i$  in (1.19). Observe that the resulting equations  $\{2, \dots, n\}$  in (1.20) are independent of  $x_1$ .

The above process is the first step of the Gaussian elimination algorithm. Assuming that  $a_{22}^{(2)} \neq 0$ , the algorithm proceeds by eliminating variable  $x_2$  from the transformed equations  $\{3, \dots, n\}$  in system (1.20). Under assumption  $a_{22}^{(2)} \neq 0$ ,

$$x_2 = \frac{b_{22}^{(2)}}{a_{22}^{(2)}} - \sum_{j=3}^n \frac{a_{2j}^{(2)}}{a_{22}^{(2)}} x_j.$$

Identically, variable  $x_2$  is eliminated from the transformed equations  $\{3, \dots, n\}$  in (1.20). New coefficients are computed as :

$$a_{ij}^{(3)} = a_{ij}^{(2)} - \frac{a_{i2}^{(2)}}{a_{22}^{(2)}} a_{2j} \quad \text{for } i, j \in \{3, \dots, n\}.$$

Assuming  $a_{ii}^{(i)} \neq 0$  for  $i \in \{3, \dots, n\}$ , the above process can be repeated until (1.19) has been transformed to the system

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 \\ a_{22}^{(2)}x_2 + a_{23}^{(2)}x_3 + \dots + a_{2n}^{(2)}x_n &= b_2^{(2)} \\ a_{33}^{(3)}x_3 + \dots + a_{3n}^{(3)}x_n &= b_3^{(3)} \\ &\vdots \\ a_{nn}^{(n)}x_n &= b_n^{(n)} \end{aligned}$$

The matrix elements  $a_{ii}^{(i)}$  for  $i \in \{1, \dots, n\}$  are called pivots. Here and in the following we set  $a_{ij}^{(1)} := a_{ij}$ .

We denote the coefficient matrix of intermediate transformed system on step  $k \in \{1, \dots, n\}$  as  $A^{(k)} \in \mathbb{R}^{n \times n}$ . For  $k = 1$  we define  $A^{(1)} := A$ . The systems  $A^{(2)}\mathbf{x} = \mathbf{b}^{(2)}$  and  $A^{(3)}\mathbf{x} = \mathbf{b}^{(3)}$  are given in (1.19) and (1.20). For  $k \in \{2, \dots, n\}$ , matrix  $A^{(k)}$  has the block structure

$$A^{(k)} = \begin{bmatrix} U^{(k)} & A_{12}^{(k)} \\ 0 & A_{22}^{(k)} \end{bmatrix}.$$

where the matrix  $U^{(k)} \in \mathbb{R}^{(k-1) \times (k-1)}$  is upper triangular.

Example 1.6 demonstrates how the Gaussian elimination algorithm is used in hand calculations.

**Example 1.6.** Consider the linear system

$$\begin{cases} x_1 + x_2 + x_3 &= 0 \\ x_1 + 2x_2 + 4x_3 &= 1 \\ x_1 + 3x_2 + 2x_3 &= 7. \end{cases}$$

In matrix form, the above system is: find  $\mathbf{x} \in \mathbb{R}^3$  satisfying

$$A\mathbf{x} = \mathbf{b}, \quad \text{where} \quad A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 2 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 7 \end{bmatrix}.$$

When running Gaussian elimination algorithm by hand, matrix  $A$  and vector  $\mathbf{b}$  are written in the same table as

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 1 & 2 & 4 & 1 \\ 1 & 3 & 2 & 7 \end{array} \right].$$

The row operations are marked on the left hand side of the table.

$$\begin{array}{l} -Y1 \\ -Y1 \end{array} \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 1 & 2 & 4 & 1 \\ 1 & 3 & 2 & 7 \end{array} \right] \rightarrow \begin{array}{l} \\ -2Y2 \end{array} \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 2 & 1 & 7 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & -5 & 5 \end{array} \right].$$

The resulting linear system is solved using the back-substitution algorithm.

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & -5 & 5 \end{array} \right] \xrightarrow{x_3 = -1} \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & -5 & 5 \end{array} \right] \xrightarrow{x_2 = 4} [1 \mid -3] \rightarrow x_1 = -3.$$

This process yields the solution  $\mathbf{x} = [-3 \ 4 \ -1]^T$ .

### Problems

P18. (0.5p) Solve the linear system

$$\begin{bmatrix} 1 & 0 & 2 & 1 \\ 0 & 1 & 2 & 2 \\ -2 & 1 & 0 & 1 \\ -1 & 0 & -4 & 2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

by hand using Gaussian elimination and back-substitution. Check your solution using Matlab.

- P19. (1p) Let  $A \in \mathbb{R}^{n \times n}$ . Assume, that all pivots during Gaussian elimination are no-zeros. Estimate the total number of arithmetic operations  $\cdot, +, -, /$  in the elimination process of  $A$ .

Use the identity

$$\sum_{x=1}^{n-1} (x + \alpha)^k \leq \int_0^{n-1} (x + \alpha + 1)^k, \quad (1.21)$$

for  $\alpha \in \mathbb{R}$  and  $k \geq 0$  to give a simple upper bound for the number of operations. Identity (1.21) follows from geometric interpretation of the sum, see Figure 1.4.

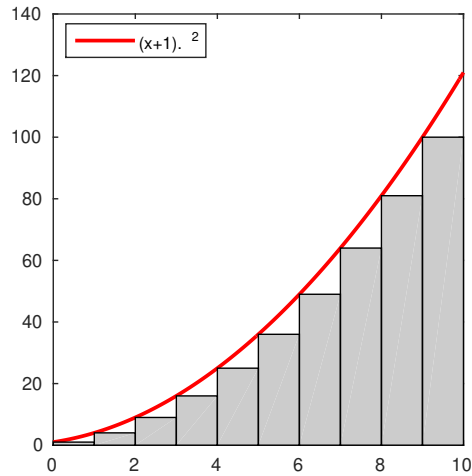


Figure 1.4: Geometry interpretation of estimate (1.21)

### 1.5.1 Pivoting

In this section, we modify Gaussian elimination process to cope with zero pivot elements. If pivot is zero, an additional pivoting step changing the order of equations or unknowns is conducted before the elimination step. Changing the order of rows and/or columns is expressed using permutation matrices.

**Example 1.7.** Consider the linear system

$$\begin{cases} x_1 + x_2 + x_3 &= 0 \\ x_1 + x_2 + 4x_3 &= 3 \\ x_1 + 3x_2 + 2x_3 &= 7. \end{cases}$$

To perform Gaussian elimination by hand, we write the system in a table:

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 1 & 1 & 4 & 3 \\ 1 & 3 & 2 & 7 \end{array} \right]$$

First step of elimination yields:

$$\begin{array}{l} -Y1 \\ -Y1 \end{array} \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 1 & 1 & 4 & 3 \\ 1 & 3 & 2 & 7 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 0 & 0 & 3 & 3 \\ 0 & 2 & 1 & 7 \end{array} \right]$$

Because the pivot  $a_{22}^{(2)} = 0$  we exchange rows two and three. This corresponds to changing the order of equations in the original linear system and does not change the solution. We obtain,

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 0 \\ 0 & 2 & 1 & 7 \\ 0 & 0 & 3 & 3 \end{array} \right]. \quad (1.22)$$

The coefficient matrix has now been transformed to upper triangular one, and  $\mathbf{x}$  is solved using back-substitution.

The permutation vector corresponding to changing rows 2 and 3 is  $\mathbf{p} = [1 \ 3 \ 2]$  and the related permutation matrix

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

*In this example, transformed system (1.22) is obtained by applying Gaussian elimination without pivoting to linear system*

$$P^T A \mathbf{x} = P^T \mathbf{b}.$$

We show in Section 1.6 that changing the order of equations or unknowns during the elimination process does not change the solution of the linear system. Further, identical transformed system is obtained by applying Gaussian elimination without pivoting to the permuted linear system

$$P^T A Q (Q^{-1} \mathbf{x}) = P^T \mathbf{b},$$

where  $P$  and  $Q$  are permutation matrices re-ordering equations and entries of  $\mathbf{x}$ .

When running the Gaussian elimination process by hand, the pivot is chosen so that the resulting computations are as simple as possible. When Gaussian elimination is implemented using a computer, pivoting is applied on every step to improve numerical stability of the algorithm. Numerical stability is discussed later in this course.

Different pivoting strategies on step  $k$  are:

- **Row-pivoting:** Choose entry  $a_{ik}^{(k)}$  for  $i \in \{k, \dots, n\}$  with largest absolute value as pivot
- **Column-pivoting:** Choose entry  $a_{kj}^{(k)}$  for  $j \in \{k, \dots, n\}$  with largest absolute value as pivot
- **Full-pivoting:** Choose entry  $a_{ij}^{(k)}$  for  $i, j \in \{k, \dots, n\}$  with largest absolute value as pivot

### 1.5.2 Problems

P20. (0.5p) Solve the linear system

$$\begin{bmatrix} 1 & 0 & 3 & 4 \\ 2 & 0 & 9 & 9 \\ 0 & 1 & 3 & 2 \\ 0 & 3 & 9 & 8 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

Using Gaussian elimination and back substitution.



### 1.5.3 Elimination matrices and $LU$ -factorisation

In this section, we express row operations conducted during Gaussian elimination process using elimination matrices. This representation allows us to prove equivalence between the original and the transformed linear system. It also yields the  $LU$  factorisation of a matrix  $A$ .

For simplicity, assume that all pivot elements are nonzero. On step  $k$  of the elimination process, row  $k$  is first multiplied with a  $-\frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}$  and then added to row  $i$  for  $i \in \{k+1, \dots, n\}$ . The corresponding linear mapping is

$$f_k(\mathbf{x})_i = \begin{cases} \mathbf{x}_i & i \leq k \\ \mathbf{x}_i - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \mathbf{x}_k & i > k \end{cases}.$$

When pivots  $a_{kk}^{(k)} \neq 0$ , the mapping  $f_k$  is invertible and

$$f_k^{-1}(\mathbf{x})_i = \begin{cases} \mathbf{x}_i & i \leq k \\ \mathbf{x}_j + \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \mathbf{x}_k & i > k \end{cases}.$$

First step of the elimination process can be stated as  $f_1(A\mathbf{x}) = f_1(\mathbf{b})$ . Let  $E_1 \in \mathbb{R}^{n \times n}$  be the matrix representation of the linear mapping  $f_1$ , this is  $f_1(\mathbf{x}) = E_1\mathbf{x}$ . The matrix representation is obtained as  $E_1 = [f_1(\mathbf{e}_1) \ f_1(\mathbf{e}_2) \ \dots \ f_1(\mathbf{e}_n)]$ , where  $\{\mathbf{e}_i\}_{i=1}^n$  are the Cartesian unit vectors. This yields

$$E_1 = \begin{bmatrix} 1 & & & \\ -\frac{a_{21}}{a_{11}} & 1 & & \\ \vdots & & \ddots & \\ -\frac{a_{n1}}{a_{11}} & & & 1 \end{bmatrix}$$

Using the above matrix representation gives the relation

$$A^{(2)} = E_1 A \quad \text{and} \quad \mathbf{b}^{(2)} = E_1 \mathbf{b},$$

where  $A^{(2)}$  is the transformed coefficient matrix obtained from step 1. Transformation of linear system  $A\mathbf{x} = \mathbf{b}$  to upper triangular form corresponds to

$$f(A\mathbf{x}) = f(\mathbf{b}). \quad (1.23)$$

where  $f = f_{n-1} \circ \dots \circ f_1$ . Let  $E_k$  be the matrix representation of the linear mapping  $f_k$ . Then the final transformed system satisfies

$$A^{(n)} = E_{n-1} \dots E_2 E_1 A \quad \text{and} \quad \mathbf{b}^{(n)} = E_{n-1} \dots E_2 E_1 \mathbf{b}. \quad (1.24)$$

As  $A^{(n)}$  is an upper triangular matrix, we denote  $U = A^{(n)}$ . Observe that the structure of elimination matrices changes for every  $k$  making them difficult to write using block matrix notation. This difficulty is addressed in Section 1.6 using recursive definition of the Gaussian elimination process.

Observe, that  $f^{-1} = f_1^{-1} \circ \dots \circ f_{n-1}^{-1}$ . Hence  $f$  has an inverse, and  $f(x) = 0 \Rightarrow \mathbf{x} = 0$ . Thus

$$f(A\mathbf{x} - \mathbf{b}) = 0 \Rightarrow A\mathbf{x} - \mathbf{b} = 0.$$

This is, the solution to transformed linear system produced by Gaussian elimination is also the solution to the original system.

Let  $A \in \mathbb{R}^{n \times n}$  be invertible matrix and assume non-zero pivots. By (1.24) it holds that  $E_{n-1} \dots E_2 E_1 A = U$  where  $U$  is an upper triangular matrix. Inverting the product of elimination matrices yields the  $LU$  factorisation

$$A = LU \quad \text{for} \quad L = E_1^{-1} \dots E_{n-1}^{-1}. \quad (1.25)$$

By Problem 21 on page 34, the matrix  $L$  is lower-triangular. Recall that entries of matrix  $L$  can be obtained directly from the row multipliers used in the elimination process. This fact is tricky to prove using index notation, hence, it is proven in Section 1.6 using block matrix notation.

Linear system

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

is reduced to two sub-problems using  $LU$ -factorisation of  $A = LU$

$$L\mathbf{y} = \mathbf{b} \quad \text{and} \quad U\mathbf{x} = \mathbf{y}.$$

Both sub-problems have triangular coefficient matrices and can be efficiently solved using back-substitution, see Section 1.2.2.

### Problems

- P21. (2p) Show that the inverse of any  $n \times n$  lower triangular matrix is lower triangular. Formulate an induction proof with respect to the dimension  $n$  and use Problem 3 on page 11
- P22. (1p) Let  $A \in \mathbb{R}^{n \times n}$  be invertible matrix. Show that on step  $k \in \{2, \dots, n\}$  of Gaussian elimination there exists a nonzero pivot on column  $k$ . Hint: argue by contradiction and recall the block form of  $A^{(k)}$  and use Problem 7 on page 12.

## 1.6 LU Factorization in block matrix notation

[video on introduction to recursive algorithm for computing the LU decomposition](#)

In this section, we use block matrix notation to define a recursive process that returns the  $LU$  factorisation of a given invertible matrix. Recall that the elimination matrices related to the elimination process all have different structure, and hence, they cannot be easily treated using block matrix notation. This problem is remedied by recursive definition that allows us to formulate the elimination process using only the first elimination matrix. The given process could be easily turned into an existence proof of the  $LU$ -decomposition. It also shows that the matrix  $L$  can be constructed from multipliers related to row operations and there is no need to save or construct elimination matrices  $E_1, \dots, E_{n-1}$  or their inverses during the elimination process. We do not assume non-zero pivots and use row pivoting. In this case, the  $LU$  factorisation of invertible matrix  $A \in \mathbb{R}^{n \times n}$  is

$$P^T A = LU \quad \text{where } P \text{ is a permutation matrix.}$$

Next, we give a recursive definition of  $[P, L, U] = lu(A)$  that returns the  $LU$  factorisation of invertible matrix  $A$ .

[See video on recursive algorithm for computing the LU decomposition](#)

---

For  $n = 1$ ,  $lu(A) = [1, 1, A]$ .

For  $n > 1$ , we use recursive definition. First, we seek the permutation  $P$  such that  $(P^T A)_{11} \neq 0$ . Next, split  $P^T A$  as

$$P^T A = \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix} \quad \text{where } a_{11} \in \mathbb{R}, \mathbf{a}_{12}, \mathbf{a}_{21} \in \mathbb{R}^{(n-1)} \text{ and } A_{22} \in \mathbb{R}^{(n-1) \times (n-1)}.$$

The elimination matrix corresponding to first step of Gauss algorithm is

$$E = \begin{bmatrix} 1 & 0 \\ -\frac{\mathbf{a}_{21}}{a_{11}} & I \end{bmatrix} \quad \text{and} \quad EP^T A = \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ 0 & A_{22} - \frac{\mathbf{a}_{21}\mathbf{a}_{12}^T}{a_{11}} \end{bmatrix}.$$

Let  $[P_2, L_2, U_2] = lu(A_{22} - \frac{\mathbf{a}_{21}\mathbf{a}_{12}^T}{a_{11}})$  so that  $A_{22} - \frac{\mathbf{a}_{21}\mathbf{a}_{12}^T}{a_{11}} = P_2^{-T} L_2 U_2$  and

$$P^T A = \begin{bmatrix} 1 & 0 \\ \frac{\mathbf{a}_{21}}{a_{11}} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & P_2^{-T} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & L_2 \end{bmatrix} \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ 0 & U_2 \end{bmatrix}.$$

By direct computation,

$$\begin{bmatrix} 1 & 0 \\ \frac{\mathbf{a}_{21}}{a_{11}} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & P_2^{-T} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & P_2^{-T} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ P_2^T \frac{\mathbf{a}_{21}}{a_{11}} & I \end{bmatrix}.$$

Thus

$$\begin{bmatrix} 1 & 0 \\ 0 & P_2^T \end{bmatrix} P^T A = \begin{bmatrix} 1 & 0 \\ P_2^T \frac{\mathbf{a}_{21}}{a_{11}} & L_2 \end{bmatrix} \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ 0 & U_2 \end{bmatrix}.$$

And finally

$$lu(A) = \begin{bmatrix} P & \begin{bmatrix} 1 & 0 \\ 0 & P_2 \end{bmatrix} \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ P_2^T \frac{\mathbf{a}_{21}}{a_{11}} & L_2 \end{bmatrix}, \begin{bmatrix} a_{11} & \mathbf{a}_{12}^T \\ 0 & U_2 \end{bmatrix}.$$

We deduce from the above algorithm that the Gaussian elimination with pivoting is Gaussian elimination applied matrix

$$P^T A,$$

where  $P$  collects all row permutations done during the process. Same holds for row- and full-pivoting. The matrix  $L$  is obtained by collecting the multipliers from step  $k$  as

$$L = \begin{bmatrix} 1 & & & & \\ \alpha_{21} & 1 & & & \\ \alpha_{31} & \alpha_{32} & 1 & & \\ \vdots & \vdots & \dots & \ddots & \\ \alpha_{n1} & \alpha_{n2} & \dots & \alpha_{n(n-1)} & 1 \end{bmatrix} \quad \text{where} \quad \alpha_{ij} = \frac{a_{ij}^{(j)}}{a_{jj}^{(j)}}.$$

### Problems

- P23. (0.5p) Write down the elimination matrices used in Example 1.6 and compute the corresponding LU-decomposition
- P24. (0.5p) Write the  $LU$  decomposition corresponding to Example 1.7.
- P25. (2p) Modify the definition of function  $lu$  to use column pivoting instead of row pivoting.
- P26. (2p) Write a recursive implementation of the function  $[P, L, U] = lu(A)$  in Matlab. Device a test verifying that your decomposition is correct.
- P27. (2p) Modify the recursive implementation of function  $lu$  to utilise the update strategy.

## 1.7 Cholesky factorisation

*This section gives existence proof for the Cholesky factorisation, which should be studied with care. The left-looking variant of the Cholesky decomposition is included because it yields a simpler formula for computing entries of the Cholesky factor and can be skipped.*

Symmetric matrix  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$  is also *positive definite* if there exists  $\alpha > 0$  such that

$$\mathbf{x}^T A \mathbf{x} \geq \alpha \|\mathbf{x}\|_2^2 \quad \text{for any } \mathbf{x} \in \mathbb{R}^n. \quad (1.26)$$

In finite dimensional case this is equivalent to

$$\mathbf{x}^T A \mathbf{x} > 0 \quad \text{for any } \mathbf{x} \in \mathbb{R}^n \setminus \{0\}. \quad (1.27)$$

In this section, we prove that every such matrix has a Cholesky decomposition:

**Theorem 1.1.** *Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric and positive definite. Then there exists a lower triangular matrix  $L \in \mathbb{R}^{n \times n}$  such that  $A = LL^T$ .*

The matrix  $L$  is called as the Cholesky factor of  $A$ . We prove Theorem 1.1 using induction with respect to the dimension of the matrix, block matrix notation, and the following technical result:

**Lemma 1.3.** *Let  $F \in \mathbb{R}^{n \times m}$  have a trivial null-space and  $A \in \mathbb{R}^{n \times n}$  be a symmetric and positive definite matrix. Then the  $m \times m$  matrix  $F^T A F$  is positive definite.* [See video proof of this lemma in Youtube](#)

Note that by the rank-nullity Theorem it holds that  $m \leq n$ .

*Proof.* As  $A$  is s.p.d. there exists  $\alpha > 0$  such that

$$\mathbf{x}^T F^T A F \mathbf{x} \geq \alpha \mathbf{x}^T F^T F \mathbf{x} \quad \text{for any } \mathbf{x} \in \mathbb{R}^m. \quad (1.28)$$

As  $F^T F$  is symmetric,  $F^T F = U \Lambda U^T$  where  $U \in \mathbb{R}^{m \times m}$ ,  $U = [\mathbf{u}_1 \ \cdots \ \mathbf{u}_m]$  is unitary and  $\Lambda \in \mathbb{R}^{m \times m}$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$  is a diagonal matrix. Matrix  $F^T F$  has the expansion

$$F^T F = \sum_{i=1}^m \lambda_i \mathbf{u}_i \mathbf{u}_i^T. \quad (1.29)$$

As eigenvectors  $\{\mathbf{u}_i\}$  are orthonormal and  $N(F) = \{0\}$ , it follows that  $\lambda_i = \mathbf{u}_i^T F^T F \mathbf{u}_i = \|F \mathbf{u}_i\|_2^2 > 0$ . Using (1.29) and estimating  $\lambda_i$  from below by  $\lambda_{\min} := \min_{i \in \{1, \dots, m\}} \lambda_i$  gives

$$\mathbf{x}^T F^T F \mathbf{x} = \sum_{i=1}^m \lambda_i (\mathbf{x}^T \mathbf{u}_i)^2 \geq \lambda_{\min} \sum_{i=1}^m (\mathbf{x}^T \mathbf{u}_i)^2 = \lambda_{\min} \mathbf{x}^T \mathbf{x}.$$

Noticing that  $\lambda_{\min} > 0$  and using (1.28) completes the proof.  $\square$

*proof of Theorem 1.1.* The proof proceeds by induction with respect to the dimension  $n$ . **Base step  $n=1$ :**  $A \in \mathbb{R}$ ,  $A > 0$ . Hence,  $L = \sqrt{A}$ .

**Induction Assumption:** The claim holds for  $n = k$

**Induction step:** Let  $A \in \mathbb{R}^{(k+1) \times (k+1)}$  and split

$$A = \begin{bmatrix} a_{11} & \mathbf{a}_{21}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix}$$

where  $a_{11} \in \mathbb{R}$ ,  $\mathbf{a}_{21} \in \mathbb{R}^k$  and  $A_{22} \in \mathbb{R}^{k \times k}$ . Let

$$E = \begin{bmatrix} 1 & 0 \\ -a_{11}^{-1} \mathbf{a}_{21} & I \end{bmatrix}.$$

By direct calculation

$$EAE^T = \begin{bmatrix} a_{11} & 0 \\ 0 & A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T \end{bmatrix}.$$

See video on existence  
proof of Cholesky fac-  
torisation in Youtube

Before applying the induction assumption to the matrix  $A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T$ , we have to show that it is positive definite. Observe that

$$(A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T) = \begin{bmatrix} 0 & I \\ k \times 1 & k \times k \end{bmatrix} EAE^T \begin{bmatrix} 0 & I \\ k \times 1 & k \times k \end{bmatrix}^T \quad (1.30)$$

As both  $\begin{bmatrix} 0 & I \end{bmatrix}^T$  and  $E$  have trivial null-spaces, so does  $F = E^T \begin{bmatrix} 0 & I \end{bmatrix}^T$ . Hence by (1.30) and Lemma 1.3,  $A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T$  is positive definite. Applying the induction assumption gives  $A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T = L_2 L_2^T$ , where  $L_2 \in \mathbb{R}^{k \times k}$  is a lower triangular matrix. Note that  $a_{11} = \mathbf{e}_1^T A \mathbf{e}_1 > 0$ . Hence,

$$EAE^T = \begin{bmatrix} a_{11} & 0 \\ 0 & A_{22} - \mathbf{a}_{21} a_{11}^{-1} \mathbf{a}_{21}^T \end{bmatrix} = \begin{bmatrix} \sqrt{a_{11}} & 0 \\ 0 & L_2 \end{bmatrix} \begin{bmatrix} \sqrt{a_{11}} & 0 \\ 0 & L_2^T \end{bmatrix}.$$

Inverting  $E$  gives

$$L = \begin{bmatrix} 1 & 0 \\ a_{11}^{-1} \mathbf{a}_{21} & I \end{bmatrix} \begin{bmatrix} \sqrt{a_{11}} & 0 \\ 0 & L_2 \end{bmatrix} = \begin{bmatrix} \sqrt{a_{11}} & 0 \\ \frac{\mathbf{a}_{21}}{\sqrt{a_{11}}} & L_2 \end{bmatrix} \quad (1.31)$$

□

The above proof is constructive, this is, it also gives a method for computing  $L$ .

The function  $L = rchol(A)$  returns the Cholesky factorisation of a s.p.d. matrix  $A \in \mathbb{R}^{n \times n}$ .

---

For  $n = 1$ ,  $rchol(A) = \sqrt{A}$ .

For  $n > 1$ , we use recursive definition. Split  $A$  as

$$A = \begin{bmatrix} a_{11} & \mathbf{a}_{21}^T \\ \mathbf{a}_{21} & A_{22} \end{bmatrix} \quad \text{and let } L_2 = rchol(A_{22} - \frac{\mathbf{a}_{21}\mathbf{a}_{21}^T}{a_{11}}).$$

By (1.31) we have

$$rchol(A) = \begin{bmatrix} \sqrt{a_{11}} & 0 \\ \frac{\mathbf{a}_{21}}{\sqrt{a_{11}}} & L_2 \end{bmatrix}$$

---

Similar to functions *triusolve* and *lu*, function *rchol* can be implemented using recursive function calls or using the update strategy. The implementation utilising update strategy is called as the down-looking Cholesky factorisation because the lower right corner is updated on each step of the algorithm.

There exist (at least) two other strategies for computing the Cholesky factorisation. The difference between these variants is the order in which the matrix elements are accessed. One has to choose the best strategy for each sparse matrix storage format and computer architecture. For example, the down-looking variant accesses data column wise and works well with compressed column storage format.

To derive the left-looking Cholesky factorisation, we split

$$L = \begin{bmatrix} \# & & & & \\ \mathbf{l}_i^T & l_{ii} & & & \\ \# & \# & \# & & \\ \mathbf{l}_j^T & l_{ji} & \# & \# & \\ \# & \# & \# & \# & \# \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} \# & & & & \\ \mathbf{a}_i^T & a_{ii} & & & \\ \# & \# & \# & & \\ \mathbf{a}_j^T & a_{ji} & \# & \# & \\ \# & \# & \# & \# & \# \end{bmatrix} \quad \text{sym.}$$

where indices  $i$  and  $j$  refer to rows  $i$  and  $j$  of matrices  $L$  and  $A$ . Computing the matrix product  $LL^T$  gives

$$\begin{bmatrix} \# & & & & \\ \mathbf{l}_i^T & l_{ii} & & & \\ \# & \# & \# & & \\ \mathbf{l}_j^T & l_{ji} & \# & \# & \\ \# & \# & \# & \# & \# \end{bmatrix} \begin{bmatrix} \# & \mathbf{l}_i & \# & \mathbf{l}_j & \# \\ & l_{ii} & \# & l_{ji} & \# \\ & & \# & \# & \# \\ & & & \# & \# \\ & & & & \# \end{bmatrix} = \begin{bmatrix} \# & & & & \\ \# & \mathbf{l}_i^T \mathbf{l}_i + l_{ii}^2 & & & \\ \# & \# & \# & & \\ \# & \mathbf{l}_j^T \mathbf{l}_i + l_{ji} l_{ii} & \# & \# & \\ \# & \# & \# & \# & \# \end{bmatrix} \quad \text{sym.}$$

Using the relation  $A = LL^T$  yields

$$\mathbf{l}_i^T \mathbf{l}_i + l_{ii}^2 = a_{ii} \quad \text{and} \quad \mathbf{l}_j^T \mathbf{l}_i + l_{ji} l_{ii} = a_{ji}. \quad (1.32)$$

Note that the entry  $l_{ii}$  is not uniquely defined by (1.32). The usual choice,  $l_{ii} \in \mathbb{R}, l_{ii} > 0$ , gives the formulas

$$l_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2} \quad \text{and} \quad l_{ji} = \frac{1}{l_{ii}} (a_{ji} - \sum_{k=1}^{i-1} l_{ik} l_{jk}) \quad \text{for } j > i \quad (1.33)$$

We use (1.33) in Section 1.8.2 to the study location of non-zero entries of  $L$ .

### 1.7.1 Additional material

A different inductive existence proof for the Cholesky factorisation is outlined in blog posting [What Is Cholesky Factorisation](#).

A survey on Cholesky factorisation aimed for computer scientist is given in

Nicholas J. Higham. Cholesky factorization. *WIREs Computational Statistics*, 1(2):251–254, 2009

### 1.7.2 Problems

P28. (1p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d.

- Starting from the definition (1.26), show that  $a_{ii} > 0$  and  $A$  is invertible. Hint : Show that system  $Ax = 0$  has only zero solution, i.e.,  $N(A) = \{0\}$ .
- Show that all eigenvalues of  $A$  are positive.
- Assume, that  $A$  also satisfies  $A = F^T F$  for some  $F \in \mathbb{R}^{n \times n}$ . Show that  $F$  is invertible.



P29. (2p)

(a) Using the definition of (1.27), show that if a matrix  $A$  has a decomposition  $A = LL^T$ , it must be positive definite.

(b) Compute by hand the Cholesky decomposition of  $\begin{bmatrix} 1 & 2 & 2 \\ 2 & 8 & 4 \\ 2 & 4 & 15 \end{bmatrix}$ .

(c) Show that the matrix  $\begin{bmatrix} 15 & 2 & 4 \\ 2 & 1 & 2 \\ 4 & 2 & 8 \end{bmatrix}$  is positive definite. Hint: Use Lemma 1.3.

P30. (2p) Write a recursive implementation of the function *rchol*. Demonstrate your implementation with random s.p.d. matrices  $A = F^T F$  where  $F$  is a matrix with random entries. This creates a symmetric positive definite matrix since  $F$  is (almost) always invertible.

P31. (2p) Modify your recursive implementation of *rchol* to use the update strategy.

P32. (1p) Let  $F \in \mathbb{R}^{n \times n}$  and  $A = F^T F$ .

(a) Show that  $\|A\|_2 = \|F\|_2^2$ . Hint: Use the definition of operator norm to obtain the estimates  $\|A\|_2 \leq \|F\|_2^2$  and  $\|F\|_2 \leq \|A\|_2^{1/2}$ .

(b) Validate (a) by numerical examples.

P33. (2p) Let  $A_N \in \mathbb{R}^{N \times N}$  be the 1D-finite difference matrix

$$A_N = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix}.$$

Define matrices  $A_n \in \mathbb{R}^n$  for  $n \in \{N-1, \dots, 1\}$  as follows. Split  $A_n \in \mathbb{R}^{n \times n}$  for  $n \in \{N, \dots, 2\}$  as

$$A_n = \begin{bmatrix} \alpha_n & \mathbf{a}_n^T \\ 1 \times 1 & \hat{A}_n \\ \mathbf{a}_n & \hat{A}_n \\ (n-1) \times 1 & (n-1) \times (n-1) \end{bmatrix}$$

and set  $A_{n-1} = \hat{A}_n - \frac{\mathbf{a}_n \mathbf{a}_n^T}{\alpha_n}$ .

- (a) Compute the block matrix product to verify that  $A_n$  can be factorised as

$$A_n = \begin{bmatrix} \sqrt{\alpha_n} & 0 \\ \frac{\mathbf{a}_n}{\sqrt{\alpha_n}} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & A_{n-1} \end{bmatrix} \begin{bmatrix} \sqrt{\alpha_n} & \frac{\mathbf{a}_n^T}{\sqrt{\alpha_n}} \\ 0 & I \end{bmatrix}$$

- (b) Use induction to show that

$$A_n = \begin{bmatrix} 1 + \frac{1}{(N+1-n)} & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \quad \text{for } n \in \{N, \dots, 1\}$$

- (c) Give a formula for the Cholesky factor of  $A_N$ .

## 1.8 Sparse Cholesky Factorisation

*This section demonstrates that the Cholesky factor of a sparse matrix can be dense and that in some cases sparse factor can be obtained by a suitable symmetric permutation. Core content.*

Let  $L$  be a Cholesky factor of a sparse s.p.d. matrix  $A$ . In this Section, we are study the location of the non-zero entries of  $L$ . Observe that  $L$  is an invertible lower triangular matrix, and thus  $l_{ii} \neq 0$ .  
[See video on sparse Cholesky factorisation in Youtube](#)

Entries  $l_{ij}$  of  $L$  satisfying

$$l_{ij} \neq 0 \quad \text{and} \quad a_{ij} = 0$$

are called *fill-in*. Fill-in increases the amount of memory required to store  $L$  as well as the time required to compute it's entries. To save computational resources, fill-in is reduced by permuting rows and columns of matrix  $A$  before computing it's the Cholesky factorisation. We call the resulting factorisation

$$P^T A P = L L^T$$

where  $P \in \mathbb{R}^{n \times n}$  is a fill-in minimising permutation and  $L \in \mathbb{R}^{n \times n}$  a lower triangular matrix as the *sparse Cholesky factorisation*.

**Example 1.8.** *Consider*

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 10 & 0 & 0 & 0 \\ 1 & 0 & 10 & 0 & 0 \\ 1 & 0 & 0 & 10 & 0 \\ 1 & 0 & 0 & 0 & 10 \end{bmatrix}.$$

*The Cholesky factor of  $A$  is*

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 & 0 \\ 1 & -0.33333 & 2.9814 & 0 & 0 \\ 1 & -0.33333 & -0.37268 & 2.958 & 0 \\ 1 & -0.33333 & -0.37268 & -0.42258 & 2.9277 \end{bmatrix}$$

*Observe, that  $L$  is a full matrix. The fill-in is reduced by permuting the entries of  $A$ . In our example, changing row 1 to row 5 and column 1 to column 5 gives*

$$P^T A P = \begin{bmatrix} 10 & 0 & 0 & 0 & 1 \\ 0 & 10 & 0 & 0 & 1 \\ 0 & 0 & 10 & 0 & 1 \\ 0 & 0 & 0 & 10 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad (1.34)$$

*where  $P$  is the permutation matrix corresponding to permutation vector  $[5 \ 2 \ 3 \ 4 \ 1]$ . The Cholesky factor of  $P^T A P$  is*

$$\tilde{L} = \begin{bmatrix} 3.1623 & 0 & 0 & 0 & 0 \\ 0 & 3.1623 & 0 & 0 & 0 \\ 0 & 0 & 3.1623 & 0 & 0 \\ 0 & 0 & 0 & 3.1623 & 0 \\ 0.31623 & 0.31623 & 0.31623 & 0.31623 & 0.7746 \end{bmatrix}.$$

*The factor  $\tilde{L}$  does not have any fill-in.*

Finding an optimal permutation that minimizes the fill-in is an *NP*-hard problem, hence, heuristics are used instead. In Section 1.8.2, we discuss minimal degree-ordering, which is a method for finding fill-in reducing permutations by utilising an efficient method for determining the location of non-zero entries of  $L$ .

### 1.8.1 Problems

P34. (2p) Let

$$A = \begin{bmatrix} a_{11} & \mathbf{a}_{21}^T \\ \mathbf{a}_{21} & I \end{bmatrix} \quad \text{for } a_{11} \in \mathbb{R}, \mathbf{a}_{21} \in \mathbb{R}^{n-1}.$$

- (a) Show that the matrix  $A$  is positive definite if  $a_{11} > \|\mathbf{a}_{21}\|_2^2$ . Hint: use the definition (1.26) with suitable splitting of  $\mathbf{x}$ .
- (b) Consider the linear system  $A\mathbf{x} = \mathbf{e}_1$ . Decompose  $\mathbf{x} = [x_1 \quad \mathbf{x}_2^T]^T$ , where  $x_1 \in \mathbb{R}$  and  $\mathbf{x}_2 \in \mathbb{R}^{n-1}$ . Show that the solution satisfies

$$(a_{11} - \mathbf{a}_{21}^T \mathbf{a}_{21})x_1 = 1 \quad \text{and} \quad \mathbf{x}_2 = -\mathbf{a}_{21}x_1.$$

### 1.8.2 Non-zero structure of the Cholesky factor

*This section gives tools for computing non-zero entries of  $L$  without knowing their exact values. These tools are then used to construct fill-in reducing permutations. Core content.*

See video on graph associated to matrix in Youtube

The Cholesky factorisation of a sparse matrix is computed in two steps: First, *symbolic factorisation* step constructs a fill-in reducing permutation and finds the location of non-zero entries of the Cholesky factor. The location of nonzero entries is used to set up sparse matrix data structure for storing  $L$ . The entries of the Cholesky factor are then computed in the *numerical factorization* step.

The location of non-zero entries in the Cholesky factor of  $A \in \mathbb{R}^{n \times n}$  is predicted from the undirected graph  $\mathcal{G}(A) = (\mathcal{V}(A), \mathcal{E}(A))$  consisting of a set of vertices  $\mathcal{V}(A) = \{1, \dots, n\}$  and a set of edges

$$\mathcal{E}(A) = \{ (i, j) \mid a_{ij} \neq 0 \quad i, j = 1, \dots, n \quad \text{and} \quad i > j \}.$$

This is, vertices  $i$  and  $j$  of the graph  $\mathcal{G}(A)$  are connected by an edge if the entry  $a_{ij}$  is nonzero.

**Example 1.9.** Let

$$A_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 10 & 0 & 0 & 0 \\ 1 & 0 & 10 & 0 & 0 \\ 1 & 0 & 0 & 10 & 0 \\ 1 & 0 & 0 & 0 & 10 \end{bmatrix} \quad (1.35)$$

and

$$A_2 = \begin{bmatrix} 20 & 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 20 & 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 20 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 20 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 20 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 20 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 20 \end{bmatrix} \quad (1.36)$$

The graphs corresponding to matrices  $A_1$  and  $A_2$  are visualized in Fig. 1.5

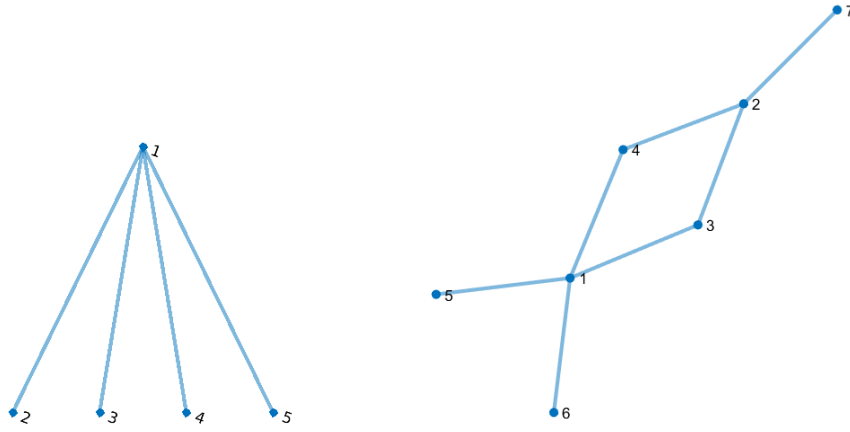


Figure 1.5: Graphs corresponding to the matrices given in (1.35) and (1.36), respectively.

Off-diagonal entries of the Cholesky factor  $L$  are computed using Eq. (1.33) as

$$l_{ij} = \frac{1}{l_{jj}}(a_{ij} - \sum_{k=1}^{j-1} l_{ik}l_{jk}), \quad \text{when } i > j. \quad (1.37)$$

Thus the entry  $l_{ij}$  can be non-zero (Possible numerical cancellations are neglected in the following) if

$$a_{ij} \neq 0 \quad (1.38)$$

or

$$l_{jk} \neq 0 \quad \text{and} \quad l_{ik} \neq 0 \quad \text{for some } k < j. \quad (1.39)$$

Based on equation (1.38), the number of nonzeros in  $L$  will always be greater or equal to the number of nonzeros in  $A$ .

See video on graph notation in Youtube

Before proceeding, we need some notation. We call the ordered set of vertices  $(v_1, v_2, \dots, v_k) \subset \mathcal{V}(A)$  as a *path*, if  $(v_i, v_{i+1}) \in \mathcal{E}(A)$  for  $i \in \{1, \dots, k-1\}$ . Vertex  $x \in \mathcal{V}(A)$  is said to be reachable from vertex  $y \in \mathcal{V}(A)$  via set  $S \subset \mathcal{V}(A)$ , if there exists a path  $(y, v_1, \dots, v_k, x)$  satisfying<sup>4</sup>  $v_i \in S$  for  $i \in \{1, \dots, k\}$ . The reachable set of  $y \in \mathcal{V}(A)$  through  $S \subset \mathcal{E}(A)$  is defined as

$$\text{Reach}(y, S) = \{x \in \mathcal{V}(A) \setminus S \mid x \text{ is reachable from } y \text{ via } S\}. \quad (1.40)$$

Examples of path and reachable set are depicted in Figure 1.6.

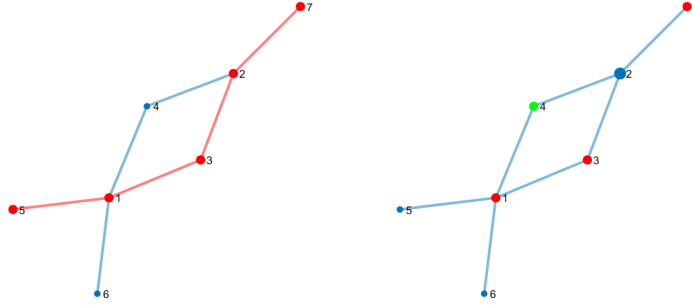


Figure 1.6: Path  $(5, 1, 3, 2, 7)$  is marked in red. Reachable set of vertex 2 via  $S = \{4\}$  is  $\{1, 3, 7\}$ .

See video proof of the following Theorem in Youtube.

The edges of  $\mathcal{G}(L + L^T)$  corresponding to non-zero off-diagonal entries of  $L$  are characterized by the following Theorem.

**Theorem 1.2.** *Let  $A \in \mathbb{R}^{n \times n}$  be a s.p.d. and  $L$  the Cholesky factor of  $A$ . Then*

$$\mathcal{E}(L + L^T) \subset \{ (i, j) \mid i \in \text{Reach}(j, \{1, \dots, j-1\}) \}$$

Recall that diagonal entries of  $L$  are always nonzero as  $L$  is an invertible lower triangular matrix. These entries are not edges of  $\mathcal{G}(L + L^T)$ .

<sup>4</sup>to make the presentation simpler, we abuse notation and use the same notation also for paths  $(y, x)$  and  $(x, v_1, y)$ .

*Proof.* Let  $i > j$  and  $(i, j) \in \mathcal{E}(L + L^T)$ . Then  $l_{ij} \neq 0$ . We proceed by induction with respect to  $j$ .

**Base case:**  $j = 1$  If  $j = 1$ ,  $l_{i1}$  is nonzero iff  $a_{i1} \neq 0$ .

**Induction assumption:** Assume that the claim holds for any  $j < t$ .

**Induction step:** Let  $j = t$ . Then  $l_{ij} \neq 0$  if  $a_{ij} \neq 0$  or there exists index  $k < j$  such that  $l_{ik} \neq 0$  and  $l_{jk} \neq 0$ . By induction assumption, there then exists paths  $(k, v_1, \dots, v_l, i)$  and  $(k, \hat{v}_1, \dots, \hat{v}_m, j)$  satisfying  $v_q < k$  for  $q \in \{1, \dots, l\}$  and  $\hat{v}_{\hat{q}} < k$  for  $\hat{q} \in \{1, \dots, m\}$ . As paths can be "walked" in both directions, there also exists path  $(i, v_l, \dots, v_1, k)$ . Thus,  $(i, v_l, \dots, v_1, k, \hat{v}_1, \dots, \hat{v}_m, j)$  is a path between vertices  $i$  and  $j$  between nodes via vertices with index smaller than  $t$ .  $\square$

A set including edges  $\mathcal{E}(L + L^T)$  is computed by finding the reachable set for ever node of  $\mathcal{V}(A)$ . Such computation can be implemented as a depth-first search (DFS). A naive example implementation is given below.

[See Wikipedia for more information on DFS](#)

```
% Call as : R = my_reach(A, v, S)
%
% A is a matrix, v is the current node, S is a vector of nodes.
%
function [R,visited] = my_reach(A, v, S, R, visited)

    if nargin == 3
        R = [];
        visited(1:size(A,2)) = false;
    end

    visited(v) = true;

    edges = find( abs( A(:,v) ) > 0 );

    if isempty(S)
        R = setdiff(edges,v);
        return;
    end

    for w=edges(:)'
        if not(visited(w))
            if not(ismember(S,w))
                R = [R w];
            else
```

```

[R,visited] = my_reach(A,w,S,R,visited);
    end
end
end
end

```

In the worst case, the cost of computing single reachable set using DFS algorithm is  $O(|N(A)| + |E(A)|)$ . Due to this potentially high cost, more efficient methods have been developed for computing the location of non-zero entries of  $L$ .

**Example 1.10.** Consider the matrix  $A_2$  in (1.36). The off-diagonal non-zero entries of the Cholesky factor are obtained as

See video on this example on Youtube

- off-diagonal non-zeros on column 1 are  $\text{reach}(1, \emptyset) = \{3, 4, 5, 6\}$ .
- off-diagonal non-zeros on column 2 are  $\text{reach}(2, \{1\}) = \{3, 4, 7\}$
- off-diagonal non-zeros on column 3 are  $\text{reach}(3, \{1, 2\}) = \{4, 5, 6, 7\}$
- off-diagonal non-zeros on column 4 are  $\text{reach}(4, \{1, 2, 3\}) = \{5, 6, 7\}$
- off-diagonal non-zeros on column 5 are  $\text{reach}(5, \{1, 2, 3, 4\}) = \{6, 7\}$
- off-diagonal non-zeros on column 6 are  $\text{reach}(6, \{1, 2, 3, 4, 5\}) = \{7\}$

The non-zeros of the computed factor are

$$\begin{bmatrix} \times & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \times & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & 0 & 0 & 0 \\ \times & 0 & \times & \times & \times & 0 & 0 \\ \times & 0 & \times & \times & \times & \times & 0 \\ 0 & \times & \times & \times & \times & \times & \times \end{bmatrix}.$$

### 1.8.3 Problems

P35. (2p) Consider the matrix  $A \in \mathbb{R}^{5 \times 5}$  such that

```

A = zeros(5);
A(1,2) = 1; A(2,3) = 1;
A(2,5) = 1; A(3,4) = 1;
A = 100*eye(5) + A + A';

```



- (a) Draw the graph  $\mathcal{G}(A)$ .
- (b) For each vertex  $i \in \mathcal{V}(A)$  compute the set  $\text{reach}(i, \{1, \dots, i-1\})$ . Use `my_reach.m` to validate your answer.
- (c) Predict the location of non-zero entries in the Cholesky factor of  $A$ .
- (d) Compute the Cholesky factorization of  $A$  in Matlab and validate (c).

P36. (1p) Let s.p.d.  $A \in \mathbb{R}^{n \times n}$  be a banded matrix with bandwidth  $b \in \mathbb{N}$ . This is,

$$a_{ij} = 0 \quad \text{if} \quad i > j + b \quad \text{or} \quad i < j - b.$$

- (a) Let  $n = 10$  and  $b = 2$ . Draw the dependency graph  $\mathcal{G}(A)$ .
- (b) Use  $\mathcal{G}(A)$  to predict the location of nonzero entries of the corresponding Cholesky factor  $L$ .

#### 1.8.4 Minimum degree ordering

*This section outlines how minimum degree ordering is used to construct a fill-in reducing permutation. Core content.*

Minimum degree (MD) ordering is a widely used heuristic for finding a fill-in reducing permutation for the matrix  $A$ . The MD method constructs a permutation vector  $\mathbf{p} \in \mathbb{R}^n$  by choosing entry  $p_i$  from the set of free indices  $\{1, \dots, n\} \setminus \{p_1, \dots, p_{i-1}\}$  so that the number of non-zero entries that appear in the  $i$ th column of  $L$  is minimised. The number of non-zero entries on column  $i$  does not depend on entries  $\{p_{i+1}, \dots, p_n\}$  and can be computed using the `my_reach.m` function. A naive implementation is given below.

[See video on MD on Youtube](#)

```
% Construct a fill-in reducing permutation vector for
% A using minimum degree ordering method. (this is a naive example
% implementation)

function p = my_md(A)
n = size(A,1);
p = 1:n;

for i=1:(n-1)
    i
    % try all remaining entries as entry i
```

```

nnzLi = zeros(1,n);
for j=(i+1):n

    tmp = p; tmp(i) = p(j); tmp(j) = p(i);

    nnzLi(j) = length(unique(my_reach(A(tmp,tmp), i, [1:(i-1)])));
end
% choose permutation minimising nnz in column i.
[~,I] = min(nnzLi((i+1):n));
I = I(1)+i;
pi = p(i); p(i) = p(I(1)); p(I) = pi;
end

```

**Example 1.11.** Consider the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 10 & 0 & 0 \\ 1 & 0 & 10 & 0 \\ 1 & 0 & 0 & 10 \end{bmatrix}.$$

Initially,  $p = [1 \ 2 \ 3 \ 4 \ 5]$ . In the first step of MD-algorithm, we test permutations

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 3 \\ 4 \\ 5 \end{bmatrix}, \begin{bmatrix} 3 \\ 2 \\ 1 \\ 4 \\ 5 \end{bmatrix}, \begin{bmatrix} 4 \\ 2 \\ 3 \\ 1 \\ 5 \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} 5 \\ 2 \\ 3 \\ 4 \\ 1 \end{bmatrix}.$$

The resulting number of non-zeros in  $L(2 : \text{end}, 1)$  is computed using function `my_reach` operator, see Fig. 1.7. In Fig. 1.7 and 1.8, letters  $\{a, b, c, d, e\}$  refer to entries  $\{1, 2, 3, 4, 5\}$  of the original matrix that after permutation have indices larger than 1 and 2, respectively. The alternative choices give 5, 2, 2, 2, 2 - nonzero entries in the first column. According to this,  $p_1 = 2$ . The process is then repeated for  $p_2$  see Fig. 1.8. Different options give 4, 2, 2, 2 - nonzero entries in  $L(3 : \text{end}, 2)$ . Accordingly, we set  $p_2 = 3$ .

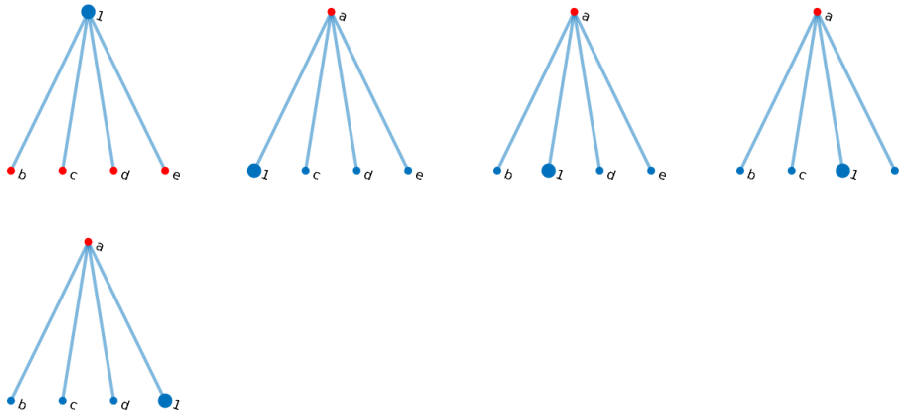


Figure 1.7: The first step of the MD - algorithm

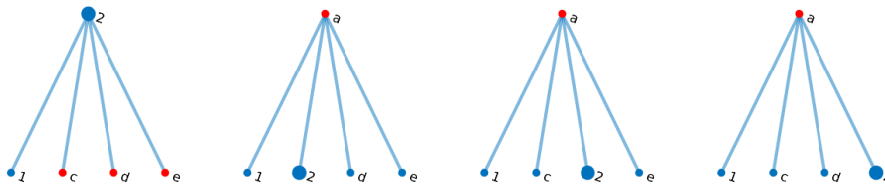


Figure 1.8: The second step of the MD - algorithm

Computing the number of non-zero entries in the column  $i$ , i.e. evaluation of the `my_reach` is the most expensive part of the MD method. This cost is reduced in the approximate minimum degree (AMD) algorithm that approximates the number of non-zeros in the column  $i$ . As AMD is much faster and yields almost as good orderings as MD, latest versions of Matlab only implement it.

### 1.8.5 Problems

P37. (2p)

- (a) Find a fill-in reducing permutation  $P$  for the matrix  $A_2$  in (1.36) using function `my_md`.
- (b) Compute the number of non-zeros in the Cholesky factors of  $A_2$  and  $P^T A_2 P$ .
- (c) Repeat (a) and (b) using permutation generated by Matlab function `amd`.

# Bibliography

- [1] John R. Gilbert, Cleve Moler, and Robert Schreiber. Sparse matrices in matlab: Design and implementation. *SIAM Journal on Matrix Analysis and Applications*, 13(1):333–356, 1992.
- [2] Nicholas J. Higham. Cholesky factorization. *WIREs Computational Statistics*, 1(2):251–254, 2009.

## 1.9 Numerical stability

Most scientific computing is done using double precision floating-point numbers that have a discrete set of possible values  $\mathbb{F} \subset \mathbb{R}$ . The set  $\mathbb{F}$  is not a vector space as it is not closed with respect to addition or multiplication. This is  $x, y \in \mathbb{F}$  does not necessarily imply  $x + y \in \mathbb{F}$ . Thus, the result of arithmetic operations conducted using double precision floating-point representation has to be rounded to the closest element of  $\mathbb{F}$ . In this section, we conduct *numerical stability analysis* and study how the resulting round-off errors affect the accuracy of solving linear systems using the Cholesky factorisation.

[See video introduction to numerical stability in Youtube](#)

In the following, we write

$$fl([expr])$$

when expression  $expr$  is evaluated in floating-point representation. All other expressions are evaluated exactly. If the order of evaluation is important, it will be explicitly stated.

Let  $\widehat{L} \in \mathbb{F}^{n \times n}$  be the (approximate) Cholesky factor of a matrix  $A \in \mathbb{F}^{n \times n}$  computed using floating-point numbers. Due to the inaccurately computed arithmetic operations, there holds that

$$\widehat{L}\widehat{L}^T = A + \delta A, \tag{1.41}$$

where  $\delta A \in \mathbb{R}^{n \times n}$  is a matrix containing round-off errors. Consider the linear system  $A\mathbf{x} = \mathbf{b}$  and let  $L \in \mathbb{R}^{n \times n}$  be the exact Cholesky factor of  $A$ . Recall, that  $\mathbf{x}$  is obtained by solving

$$LL^T \mathbf{x} = \mathbf{b} \quad (1.42)$$

in two steps,  $L\mathbf{y} = \mathbf{b}$ , and  $L^T \mathbf{x} = \mathbf{y}$ . Replacing the exact factor with numerically computed  $\hat{L}$ , as happens in practical computations, we solve

$$\hat{L}\hat{L}^T \hat{\mathbf{x}} = \mathbf{b} \quad (1.43)$$

instead of (1.42). Using the back-substitution method in floating-point representation gives the (approximate) solution  $\tilde{\mathbf{x}} \in \mathbb{F}$  to (1.43). The resulting error satisfies

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| \leq \|\mathbf{x} - \hat{\mathbf{x}}\| + \|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|.$$

Error estimates for factorisation error  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  and back-substitution error  $\|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|$  are given after we have developed sufficient tools. To outline the approach, consider the factorisation error  $\mathbf{x} - \hat{\mathbf{x}}$ . By (1.41),  $\hat{\mathbf{x}}$  satisfies

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

We conduct *backward error analysis* where a bound for the relative error  $\|\mathbf{x} - \hat{\mathbf{x}}\| \|\mathbf{x}\|^{-1}$  is obtained by first estimating the norm of matrix  $\delta A$  and then using perturbation theory of linear systems. Similar approach is used to estimate back-substitution error.

This section is organised as follows. First, we discuss perturbation theory. Then we give a mathematical model for rounding errors related to floating-point arithmetic operations and derive useful technical results. Next, we conduct backward error analysis for back-substitution of  $2 \times 2$ -upper triangular matrices and finally for the Cholesky factorization.

### 1.9.1 Perturbation theory

[See video on operator norms in Youtube](#)

This section is a brief review on perturbation analysis of linear systems. Let  $\|\cdot\|$  be a vector norm and  $\|\cdot\|_{op}$  the induced operator norm

$$\|A\|_{op} := \max_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \neq 0}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \|\mathbf{x}\|=1}} \|A\mathbf{x}\|. \quad (1.44)$$

Following the standard convention in linear algebra, we drop the subscript from (1.44) and denote  $\|\cdot\|_{op} = \|\cdot\|$ . Recall the fundamental properties of operator norms: for any  $A, B \in \mathbb{R}^{n \times n}$  and  $\mathbf{x} \in \mathbb{R}^n$  it holds that

- (i)  $\|A\mathbf{x}\| \leq \|A\|\|\mathbf{x}\|$
- (ii)  $\|AB\| \leq \|A\|\|B\|$  (sub-multiplicativity)
- (iii)  $\|A + B\| \leq \|A\| + \|B\|$  (triangle inequality)
- (iv)  $\|A + B\| \geq \left| \|A\| - \|B\| \right|$  (reverse triangle inequality)

Let  $A \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ , and consider the problem: find  $\mathbf{x} \in \mathbb{R}^n$  such that

$$A\mathbf{x} = \mathbf{b}. \quad (1.45)$$

Assume that the matrix  $A$  is invertible, i.e.,  $A \in \mathbb{R}^{n \times n}$  has an inverse  $A^{-1} \in \mathbb{R}^{n \times n}$ . The perturbed linear equation is

$$(A + \delta A)\hat{\mathbf{x}} = \mathbf{b} + \delta \mathbf{b}, \quad (1.46)$$

See video introduction to perturbation analysis in Youtube.

where perturbations  $\delta A \in \mathbb{R}^{n \times n}$  and  $\delta \mathbf{b} \in \mathbb{R}^n$ . In perturbation analysis, the aim is to relate the error  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  to the size of perturbations  $\|\delta \mathbf{b}\|$  and  $\|\delta A\|$ . The general intuition is that  $\|\delta A\|$  and  $\|\delta \mathbf{b}\|$  are small compared to  $\|A\|$  and  $\|\mathbf{b}\|$ , respectively.

To derive the perturbation estimate, we subtract (1.45) and (1.46), to obtain a linear system that determines the *error*  $\mathbf{e} := \hat{\mathbf{x}} - \mathbf{x}$ ,

$$(A + \delta A)(\hat{\mathbf{x}} - \mathbf{x}) = \delta \mathbf{b} - \delta A\mathbf{x}. \quad (1.47)$$

We begin by stability estimate for this system

**Lemma 1.4.** *Let  $\|\cdot\|$  be a vector norm and use the same notation for the induced operator norm. Let  $A, \delta A \in \mathbb{R}^{n \times n}$  satisfy  $\|A^{-1}\|\|\delta A\| < 1$ . Then the linear system: find  $\mathbf{e} \in \mathbb{R}^n$  satisfying*

$$(A + \delta A)\mathbf{e} = \mathbf{b} \quad (1.48)$$

See video proof of this stability estimate in Youtube

*has a unique solution and*

$$\|\mathbf{e}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|\delta A\|} \|\mathbf{b}\|. \quad (1.49)$$

*Proof.* By rearranging terms and multiplying with  $A^{-1}$  equation (1.48) gives

$$\mathbf{e} = -A^{-1}\delta A\mathbf{e} + A^{-1}\mathbf{b}.$$

Hence  $\|\mathbf{e}\| \leq \|A^{-1}\|\|A\|\|\mathbf{e}\| + \|A^{-1}\|\|\mathbf{b}\|$ . Combining terms related to  $\|\mathbf{e}\|$  and dividing with  $(1 - \|A^{-1}\|\|A\|) > 0$  gives (1.49). Choosing  $\mathbf{b} = 0$  in (1.49) yields  $\mathbf{e} = 0$ . Hence  $N(A + \delta A) = \{0\}$  and (1.48) has a unique solution.  $\square$

The following theorem relates the (relative) error to the relative sizes of the perturbations, i.e.  $\|\delta \mathbf{b}\|/\|\mathbf{b}\|$  and  $\|\delta A\|/\|A\|$ , and the condition number of  $A$ , defined as

$$\kappa(A) := \|A\| \|A^{-1}\|.$$

In other words, the condition number characterizes how much a perturbation affects the error when solving a linear system. Take note that the condition number of a matrix depends on the considered (operator) norm. In Problem 38 it is shown that in the case of the 2-norm the condition number is the ratio of the maximal and minimal singular values of  $A$ .

**Theorem 1.3.** *Suppose the assumptions of Lemma 1.4 are valid. Then it holds that*

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(A)}{1 - \frac{\|\delta A\|}{\|A\|} \kappa(A)} \left( \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta A\|}{\|A\|} \right) \quad (1.50)$$

where  $\kappa(A) = \|A\| \|A^{-1}\|$ .

*Proof.* Application of Lemma 1.4 to (1.47) yields

$$\|\hat{\mathbf{x}} - \mathbf{x}\| \leq \frac{\|A^{-1}\|}{1 - \|\delta A\| \|A^{-1}\|} (\|\delta \mathbf{b}\| + \|\delta A \mathbf{x}\|).$$

Dividing by  $\|\mathbf{x}\|$  and using the estimate  $\|\delta A \mathbf{x}\| \leq \|\delta A\| \|\mathbf{x}\|$  gives

$$\begin{aligned} \frac{\|\hat{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} &\leq \frac{\|A^{-1}\|}{1 - \|\delta A\| \|A^{-1}\|} \left( \frac{\|\delta \mathbf{b}\|}{\|\mathbf{x}\|} + \|\delta A\| \right) \\ &= \frac{\|A\| \|A^{-1}\|}{1 - \frac{\|\delta A\|}{\|A\|} \|A^{-1}\| \|A\|} \left( \frac{\|\delta \mathbf{b}\|}{\|A\| \|\mathbf{x}\|} + \frac{\|\delta A\|}{\|A\|} \right), \end{aligned} \quad (1.51)$$

where the latter step is mere algebraic manipulation. Since  $\|\mathbf{b}\| = \|A \mathbf{x}\| \leq \|A\| \|\mathbf{x}\|$ , we finally obtain

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\|A^{-1}\| \|A\|}{1 - \frac{\|\delta A\|}{\|A\|} \|A^{-1}\| \|A\|} \left( \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta A\|}{\|A\|} \right). \quad (1.52)$$

Substituting the definition of the condition number  $\kappa(A)$  completes the proof.  $\square$

See video proof of this  
perturbation Theorem  
in Youtube



## 1.9.2 Problems

P38. (1p)

- (a) Let  $A \in \mathbb{R}^{n \times n}$  and denote the singular values of  $A$  as  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ . Show that  $\kappa_2(A) = \|A^{-1}\|_2 \|A\|_2 = \frac{\sigma_1}{\sigma_n}$ .
- (b) Let

$$A_1 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad \text{and} \quad A_2 = \begin{bmatrix} \sqrt{2} & 0 \\ \frac{1}{\sqrt{2}} & \sqrt{\frac{3}{2}} \end{bmatrix}.$$

Compute the condition numbers of  $A_1$  and  $A_2$  in 2-norm.

P39. (2p) Let

$$A = \begin{bmatrix} a_{11} & \mathbf{a}_{21}^T \\ \mathbf{a}_{21} & I \end{bmatrix} \quad \text{for} \quad a_{11} \in \mathbb{R}, \mathbf{a}_{21} \in \mathbb{R}^{n-1}.$$

- (a) Let sub-space  $X := \{\mathbf{x} \in \mathbb{R}^{n-1} \mid \mathbf{a}_{21}^T \mathbf{x} = 0\}$  have a basis  $\{\mathbf{v}_1\}_{k=1}^{n-2}$ . Verify that

$$\left\{ \begin{bmatrix} 0 \\ \mathbf{v}_1 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ \mathbf{v}_{n-2} \end{bmatrix} \right\}.$$

are eigenvectors of  $A$  corresponding to eigenvalue 1.

- (b) As  $A$  is symmetric its eigenvectors can be chosen as an orthogonal set. Thus, we choose the two remaining eigenvectors  $\mathbf{u}_1, \mathbf{u}_2$  of  $A$  as  $\mathbf{u}_i = V \mathbf{t}_i$  for

$$V = \begin{bmatrix} 1 & 0 \\ 0 & \frac{\mathbf{a}_{21}}{\|\mathbf{a}_{21}\|} \end{bmatrix} \in \mathbb{R}^{n \times 2}, \quad \mathbf{t}_i \in \mathbb{R}^2, \quad \text{and} \quad i \in \{1, 2\}.$$

Verify that  $\mathbf{u}_i$  satisfies

$$\mathbf{u}_i^T \begin{bmatrix} 0 \\ \mathbf{v}_j \end{bmatrix} = 0 \quad \text{for} \quad i \in \{1, 2\} \quad \text{and} \quad j \in \{1, \dots, n-2\}.$$

- (c) The eigenvalues  $\lambda_1$  and  $\lambda_2$  corresponding to  $\mathbf{u}_1$  and  $\mathbf{u}_2$  can be computed as follows: as  $AV\mathbf{t}_i \in \text{span } V$  there holds that

$$V^T AV \mathbf{t}_i = \lambda_i \mathbf{t}_i \quad \text{for} \quad i \in \{1, 2\}. \quad (1.53)$$

Use matlab to compute  $\lambda_1$  and  $\lambda_2$  from (1.53) when  $n = 10$ ,  $a_{11} = 10$  and  $\mathbf{a}_{21} = [1 \ 1 \ \dots \ 1]^T$ . What is the corresponding condition number?

- P40. (2p) Let  $A \in \mathbb{R}^{n \times n}$  and  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ ,  $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$ . Consider two problems : find  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^n$  such that

$$\begin{aligned}(A + \mathbf{u}\mathbf{v}^T) \hat{\mathbf{x}} &= \mathbf{b} \\ A\mathbf{x} &= \mathbf{b}\end{aligned}$$

- (a) Show that  $\mathbf{w} = \mathbf{x} - \hat{\mathbf{x}}$  satisfies the equation

$$(A + \mathbf{u}\mathbf{v}^T) \mathbf{w} = \mathbf{u}\mathbf{v}^T \mathbf{x}.$$

- (b) Show that  $\mathbf{w} = \alpha A^{-1} \mathbf{u}$  for some  $\alpha \in \mathbb{R}$ .

- (c) Using (a) and (b), show that

$$\alpha = \frac{\mathbf{v}^T A^{-1} \mathbf{b}}{1 + \mathbf{v}^T A^{-1} \mathbf{u}} \quad \text{and} \quad \hat{\mathbf{x}} = A^{-1} \mathbf{b} - \frac{A^{-1} \mathbf{u} \mathbf{v}^T A^{-1} \mathbf{b}}{1 + \mathbf{v}^T A^{-1} \mathbf{u}}.$$

- P41. (2p) Let  $A \in \mathbb{R}^{n \times n}$  be invertible and have a singular value decomposition  $A = U\Sigma V^T$ , where  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$  are the singular values. In addition, let  $\delta A \in \mathbb{R}^{n \times n}$  have the singular value decomposition  $\delta A = -U\delta\Sigma V^T$ , where  $\delta\Sigma = \text{diag}(\delta\sigma_1, \dots, \delta\sigma_n)$ .

- (a) Assume that we are solving  $A\mathbf{x} = \mathbf{b}$  numerically. Show that the error  $\mathbf{e}$  defined in (1.48) satisfies

$$\mathbf{e} = V\hat{\Sigma}U^T \mathbf{b} \quad \text{for } \hat{\Sigma} = \text{diag}(\sigma_1^{-1}(1 - \sigma_1^{-1}\delta\sigma_1)^{-1}, \dots, \sigma_n^{-1}(1 - \sigma_n^{-1}\delta\sigma_n)^{-1}).$$

- (b) For a given matrix  $A$  and assuming that  $|\mathbf{b}| = 1$ , what is the largest value the error  $\mathbf{e}$  can have when measured in 2-norm?

### 1.9.3 Modelling floating point errors

In this section, we discuss floating-point representation and derive a model for arithmetic operations of floating-point numbers. We use the notation

$$(a_n \cdots a_1 a_0 \cdot c_1 c_2 c_3 \cdots)_b$$

for the base- $b$  number

$$(a_n \cdots a_1 a_0 \cdot c_1 c_2 c_3 \cdots)_b = \sum_{k=0}^n a_k b^k + \sum_{k=1}^{\infty} c_k b^{-k}.$$

The symbol  $\cdot$  a called as radix point and it corresponds to the decimal point in the base-10 system. Changing the base of a number is done easily by using the modulo-operation, see the example code below.

[See video on fp. representation in Youtube](#)

```

% p>0 indicates how many numbers there are after the radix point.
% Number is not rounded, but cut-off is used instead.

function str = my_dec2bin(dec,p)

n = max([floor(log2(dec))+1, 0]);
str = '';
for i=1:(n+p)
    i
    remainder = mod(dec,2^(n-i));

    if(remainder == dec)
        bit = '0';
    else
        bit = '1';
    end

    if(i==(n+1))
        str(end+1) = '.';
        str(end+1) = bit;
    else
        str(end+1) = bit;
    end
    dec = remainder;
end

```

Floating-point numbers are based on the (normalised) scientific notation

$$(-1)^S M b^E, \quad (1.54)$$

where  $S \in \{0, 1\}$  is the sign,  $b$  is the base, and  $E$  is the exponent. The term  $M \in [1, b)$  is a base- $b$  number with  $N$ -significant figures called as significand, mantissa, or factor. The term *significant figures* mean the digits that carry information. The rules are (in the following example significant figures are marked in red):

- All nonzero number are significant.
- Leading or trailing zeros are not significant, e.g., 0.00123, and 12300
- All zeros between two nonzero numbers are significant, e.g., 0.120300.

Same rules hold in any base- $b$  system. Observe that  $M \geq 1$ , hence  $M$  has no leading zeros and is at most  $N$  digits long. The exponent has a limited range, but this is not central for our application and it is not discussed in the following. The number zeros has a special representation in floating point system.

Observe that in binary number system, the significand always is of the form

$$(1 \cdot c_1 c_2 c_3 \dots)_2 \quad (1.55)$$

Hence,  $N + 1$  significant figures can be represented by  $N$  bits.

**Example 1.12.** Consider floating point number system with  $b = 10$  and one significant figure for  $M$ . This is,

$$M \in \{1, 2, 3, 4, 5, 6, 7, 8, 9\}.$$

The exponent determines where the decimal point is placed. For example,  $\pi = 3.14 \dots$  is rounded to 3 and represented as  $3 \cdot 10^0$ . Similarly  $0.022 = 2 \cdot 10^{-2}$ .

**Example 1.13.** Consider floating point number system with  $b = 2$  and three significant figures for  $M$ . This is,

$$M \in \{(1.00)_2, (1.01)_2, (1.10)_2, (1.11)_2\}$$

To determine representation of 3.14 in this floating point system, it is first written using binary numbers. There holds that

$$(3.14)_{10} = (11.00100011110101110001 \dots)_2 \approx (11.0)_2$$

As the radix point is shifted one unit to right we obtain  $(3.14)_{10} \approx (1.10)_2 \cdot 2^1$ . Similarly,

$$(0.24)_{10} = (.0011110101110000101000 \dots)_2 \approx (.0100)_2 = (1.00)_2 \cdot 2^{-2}.$$

In scientific presentation (1.54), the resolution between numbers is not constant. Each interval  $[b^E, b^{E+1})$  is divided to sub-intervals according to number of significant figures  $N$  used for the significand  $M$ . In base  $b$ , mantissa has  $(b-1)b^{N-1}$  values, hence on interval  $[b^E, b^{E+1})$  the distance between numbers is

$$\frac{(b^{E+1} - b^E)}{(b-1)b^{N-1}} = ub^E \quad \text{where machine epsilon } u = \frac{1}{b^{N-1}}.$$

See illustration in Figure 1.9. For example, matlab uses double precision floating point numbers where the significand is a binary number with 52 bits, i.e.,  $N = 53$ ,  $b = 2$ , and  $u = 2^{-52}$ .

Roughly speaking, arithmetic operations in floating point number system are conducted by calculating the operation in higher accuracy and then rounding the result to closest floating point number.

See video on Example 1.12 numbers in Youtube

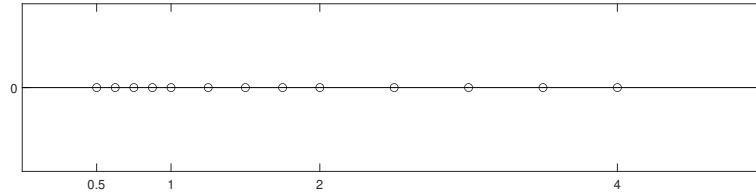


Figure 1.9: Floating point numbers in the system  $b = 2$ ,  $N = 3$  between  $1/2, 1, 2, 4$ . The resolution is different on intervals  $(1/2, 1)$ ,  $(1, 2)$ , and  $(2, 4)$ .

**Example 1.14.** Consider the floating point number system with  $b = 2$  and three significant figures for  $M$ . We compute the sum

$$(1.10)_2 \cdot 2^1 + (1.00)_2 \cdot 2^{-2}.$$

First these numbers are written using the same exponent. Then the mantissa's are added.

$$(1100.00)_2 \cdot 2^{-2} + (1.00)_2 \cdot 2^{-2} = (1101.00)_2 \cdot 2^{-2}$$

We round the mantissa upwards<sup>5</sup> to three significant figures as  $(1101.00)_2 \approx (111)_2$ . Thus, the result is  $(1.11)_2 \cdot 2^1$ .

Let  $a, b$  be two floating-point numbers and  $\odot$  denote some of the operations  $\odot = +, -, *, /$ . The exact value  $a \odot b$  is rounded to the either of the closest two floating-point numbers. For simplicity, assume that  $a \odot b > 0$ . Let  $E$  be such that  $a \odot b \in [b^E, b^{E+1})$ . The two floating point numbers closest to  $a \odot b$  lie within the interval

$$(a \odot b - ub^E, a \odot b + ub^E). \quad (1.56)$$

As exponent  $E$  has somewhat complicated dependency on  $a \odot b$ , we estimate  $b^E \leq a \cdot b$  and use the extended interval

$$(a \odot b - ua \odot b, a \odot b + ua \odot b). \quad (1.57)$$

instead of (1.56). Writing interval (1.57) using an (unknown) parameter  $\delta$ ,  $|\delta| \leq u$  we arrive to our *model for arithmetic operations in floating-point representation*:

$$fl(a \odot b) = (1 + \delta)(a \odot b) \quad \text{where} \quad \odot = +, -, *, /. \quad (1.58)$$

<sup>5</sup>The numbers  $(1100)_2$  and  $(1110)_2$  are equally close to  $(1101)_2$ . Different *tie breaking* rules can be used to choose which one to pick. We round upwards, but one can choose, e.g., to pick the closest even number.

See video on Example 1.14, spacing of fp. numbers, and model for fp. operations in Youtube

This representation is valid independent on the sign of  $a \odot b$ .

Using (1.58) it is straightforward to study rounding errors in evaluation of expressions such as  $fl(x_1y_1 + x_2y_2)$ .<sup>6</sup>

**Example 1.15.** *Consider evaluation of expression  $fl(x_1y_1 + x_2y_2)$ . The model (1.58) gives for the two multiplications*

$$fl(x_1y_1) = (1 + \delta_1)x_1y_1 \quad \text{and} \quad fl(x_2y_2) = (1 + \delta_2)x_2y_2.$$

*and for the summation*

$$fl(x_1y_1 + x_2y_2) = (1 + \delta_1)fl(x_1y_1) + (1 + \delta_2)fl(x_2y_2) = (1 + \delta_3)[(1 + \delta_1)x_1y_1 + (1 + \delta_2)x_2y_2].$$

*Estimate for the width of the interval is obtained as*

$$\begin{aligned} |fl(x_1y_1 + x_2y_2) - (x_1y_1 + x_2y_2)| &\leq |\delta_3 + \delta_1 + \delta_3\delta_1||x_1y_1| + |\delta_3 + \delta_2 + \delta_3\delta_2||x_2y_2| \\ &\leq (2u + u^2)(|x_1y_1| + |x_2y_2|). \end{aligned}$$

#### 1.9.4 Additional material

1. Matlab uses double precision floating-point numbers. A good reference on their working is [Wikipedia](#)
2. Most sources discussing double precision floating-point numbers mention that the *precision* is 53 bits. Precision refers to accuracy of the number system, which is  $b^N/2$ , when proper rounding is used. For more on the topic see [Wikipedia](#)

#### 1.9.5 Problems

P42. (1p) Modify function `my_dec2bin` to change representation of numbers from base-10 system to base- $b$  system for any  $b \in 2, \dots, 10$ .

P43. (2p)

- (a) Write  $x_1 = 345$  and  $x_2 = 1/3$  using floating-point system with  $b = 2$  and  $N = 4$ .
- (b) Compute the absolute error between  $x_i$  and it's floating point representation  $\hat{x}_i$  for  $i \in \{1, 2\}$ .

---

<sup>6</sup>When the order of evaluation for the expression is important, it is explicitly specified, otherwise it is omitted, as is the case with  $x_1y_1 + x_2y_2$ .

- (c) What is the machine epsilon, as defined in this chapter, of this floating-point system?
- (d) Compute the sum  $\hat{x}_1 + \hat{x}_2$ .
- P44. (0.5p) Let  $a \in (1, 3)$ . Find  $a_0 \mathbb{R}$  and smallest  $u$  s.t.  $a = a_0 + \delta$  for some  $|\delta| \leq u$ . Use this expression to determine intervals where the values of the following expressions belong to.
- (a)  $a^2$
- (b)  $a^2 + a$

### 1.9.6 Technical estimates

This section contains technical estimates that are used to simplify expressions related to computations done using the floating-point model (1.58). For example, the term  $(1 + \delta_1)(1 + \delta_2)$ , where  $|\delta_i| \leq u$  for  $i = 1, 2$  satisfies

$$(1 - u)^2 \leq (1 + \delta_1)(1 + \delta_2) \leq (1 + u)^2.$$

However, it is not straightforward to interpret how large the interval  $[(1 - u)^2, (1 + u)^2]$  is. To remedy this, we seek for  $\beta$  satisfying

$$(1 - \beta) \leq (1 - u)^2 \leq (1 + \delta_1)(1 + \delta_2) \leq (1 + u)^2 \leq (1 + \beta),$$

and write  $(1 + \delta_1)(1 + \delta_2) = (1 + \theta)$  for  $|\theta| \leq \beta$ . In the next Lemma, we give a simple expression for  $\theta$ .

[See video proof for Lemma 1.5 in Youtube](#)

**Lemma 1.5.** *Let  $n \in \mathbb{N}, \alpha \in \mathbb{R}, \alpha > 0$ , and  $n\alpha < 1$ . Then there holds that*

$$1 - \frac{n\alpha}{1 - n\alpha} \leq (1 - \alpha)^n \quad \text{and} \quad (1 + \alpha)^n \leq 1 + \frac{n\alpha}{1 - n\alpha}.$$

*Proof.* There holds that

$$(1 + \alpha)^n = 1 + \int_0^\alpha n(1 + t)^{n-1} dt \tag{1.59}$$

Estimating the integral as in Fig. 1.10 gives  $\int_0^\alpha n(1 + \alpha)^{n-1} \leq n\alpha(1 + \alpha)^{n-1}$  so that

$$(1 + \alpha)^n \leq 1 + n\alpha(1 + \alpha)^{n-1}.$$

Rearranging the terms in the equation above leads to

$$(1 + \alpha - n\alpha)(1 + \alpha)^{n-1} \leq 1$$

Which gives

$$(1 + \alpha)^{n-1} \leq \frac{1}{1 - (n-1)\alpha} = 1 + \frac{(n-1)\alpha}{1 - (n-1)\alpha}.$$

Lower follows by observing that the second derivative of the function  $t \mapsto (1+t)^n$  for  $n > 1$  is non-decreasing and using result of P45.  $\square$

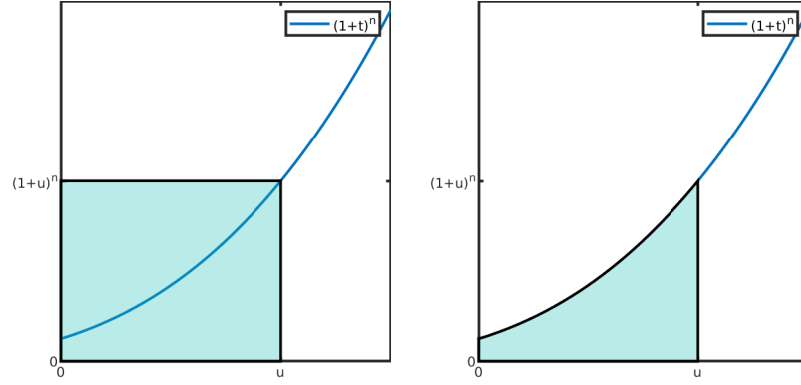


Figure 1.10: The geometric idea in estimating the integral in proof of Lemma 1.5

See [video proof for Lemma 1.6 in Youtube](#) **Lemma 1.6.** Let  $\hat{\delta}_q, \delta_p \in \mathbb{R}$  be such that  $|\delta_p| \leq u$  and  $|\hat{\delta}_q| \leq u$  for all  $p = 1, \dots, n$  and  $q = 1, \dots, \hat{n}$ . Then there holds that

$$\left( \prod_{q=1}^{\hat{n}} (1 + \hat{\delta}_q) \right) \left( \prod_{p=1}^n \frac{1}{1 + \delta_p} \right) = (1 + \theta) \quad \text{where} \quad |\theta| \leq \frac{(n + \hat{n})u}{1 - (n + \hat{n} + 1)u}.$$

*Proof.* Observe, that

$$\frac{1}{1 + \delta_p} = 1 - \frac{\delta_p}{1 + \delta_p} \quad \text{so that} \quad 1 - \frac{u}{1 - u} \leq \frac{1}{1 + \delta_p} \leq 1 + \frac{u}{1 - u}$$

as  $(1 + u) \leq (1 + \frac{u}{1-u})$ , there holds that

$$\left( 1 - \frac{u}{1 - u} \right)^{n + \hat{n}} \leq \left( \prod_{q=1}^{\hat{n}} (1 + \hat{\delta}_q) \right) \left( \prod_{p=1}^n \frac{1}{1 + \delta_p} \right) \leq \left( 1 + \frac{u}{1 - u} \right)^{n + \hat{n}}.$$

Application of Lemma 1.5 completes the proof.  $\square$



**Example 1.16.** Continuing the previous example, Lemma 1.6 gives

$$fl(x_1y_1 + x_2y_2) = (1 + \theta_2^1)x_1y_1 + (1 + \theta_2^2)x_2y_2.$$

For some  $\theta_2^1$  and  $\theta_2^2$  satisfying<sup>7</sup>

$$|\theta_2^i| \leq \frac{2u}{1 - 3u}.$$

The error between exact and floating point number is estimated as

$$|x_1y_1 + x_2y_2 - fl(x_1y_1 + x_2y_2)| \leq \frac{2u}{1 - 3u} (|x_1y_1| + |x_2y_2|).$$

Finally, we arrive to our final technical estimate. This estimate is useful in studying floating point errors related to Cholesky factorisation or back substitution.

**Lemma 1.7.** Let  $b, c \in \mathbb{R}$ ,  $x, y \in \mathbb{R}^n$ . Assume that  $s = \frac{1}{b} (c + \sum_{i=1}^n x_i y_i)$  is evaluated in floating-point arithmetics as

See [outline](#) of [Lemma 1.7](#) in Youtube

```
s = c;
for i=1:n
    s = s + x(i)*y(i)
end
s = s/b
```

Then there holds that

$$b(1 + \theta_{n+1}) fl(s) = c + \sum_{i=1}^n x_i y_i (1 + \theta_i).$$

where  $|\theta_i| \leq \frac{(i+1)u}{1-(i+2)u}$  for  $i = 1, \dots, n+1$ .

*Proof.* Let us first show that computing the sum satisfies

$$fl(\hat{s}) = \left( c \prod_{j=1}^n (1 + \delta_j) + \sum_{i=1}^n x_i y_i (1 + \hat{\delta}_i) \prod_{j=i}^n (1 + \delta_j) \right)$$

where  $|\hat{\delta}_i| \leq u$  and  $|\delta_j| \leq u$ . Denote the partial sums by  $\hat{s}_n$ . This claim is proven using induction with respect to  $n$ .

---

<sup>7</sup>Slightly better estimate can be obtained by application of Lemma 1.5.

**Base case  $n = 1$ :** By (1.58)

$$fl(\hat{s}_1) = \left( c + x_1 y_1 (1 + \hat{\delta}_1) \right) (1 + \delta_1).$$

**Induction assumption:** assume now that the claim holds for  $n = k - 1$ , this is,

$$fl(\hat{s}_{k-1}) = c \prod_{j=1}^{k-1} (1 + \delta_j) + \sum_{i=1}^{k-1} x_i y_i (1 + \hat{\delta}_i) \prod_{j=i}^{k-1} (1 + \delta_j).$$

**Induction step:** Let  $n = k$  and consider computing

$$fl(\hat{s}_k) = \left( s_{k-1} + x_k y_k (1 + \hat{\delta}_k) \right) (1 + \delta_k).$$

Using the induction assumption gives

$$fl(\hat{s}_k) = c \prod_{j=1}^{k-1} (1 + \delta_j) (1 + \delta_k) + \sum_{i=1}^{k-1} x_i y_i (1 + \hat{\delta}_i) \prod_{j=i}^{k-1} (1 + \delta_j) (1 + \delta_k) + x_k y_k (1 + \hat{\delta}_k) (1 + \delta_k).$$

The final division by  $b$  leads to

$$fl(s) = \frac{1}{b} \left( c \prod_{j=1}^{n+1} (1 + \delta_j) + \sum_{i=1}^n x_i y_i (1 + \hat{\delta}_i) \prod_{j=i}^{n+1} (1 + \delta_j) \right)$$

Dividing by  $\prod_{p=1}^j (1 + \delta_p)$  and multiplying with  $b$  gives now

$$b \prod_{p=1}^j (1 + \delta_p)^{-1} fl(s) = c + \sum_{i=1}^n x_i y_i (1 + \hat{\delta}_i) \prod_{j=1}^{i-1} (1 + \delta_j)^{-1}$$

The proof is completed by application of Lemma 1.6. □

### 1.9.7 Problems

P45. (1p) Let  $f : \mathbb{R} \mapsto \mathbb{R}$  be a smooth function such that  $f'$  is increasing function. In addition, let  $x_0, \alpha, \delta_{x_0} \in \mathbb{R}$ ,  $\delta_{x_0} > 0$  be such that

$$f(x_0 + \alpha) \leq f(x_0) + \delta_{x_0}.$$

Show that

$$f(x_0 - \alpha) \geq f(x_0) - \delta_{x_0}.$$

Hint : expand  $f(x_0 + \alpha)$  and  $f(x_0 - \alpha)$  as in (1.59). Argue that

$$\int_{x_0}^{x_0+\alpha} f'(s) ds \geq - \int_{x_0}^{x_0-\alpha} f'(s) ds.$$

P46. (2p) Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Show that

$$fl(\sum x_i y_i) = \sum_{i=1}^n x_i y_i (1 + \hat{\delta}_i) \prod_{k=i}^n (1 + \delta_k).$$

where  $\delta_1 = 0$ ,  $|\delta_k| \leq u$  for  $k = 2, \dots, n$  and  $|\hat{\delta}_i| \leq u$  for  $i=1, \dots, n$ .

P47. (2p) Following the notation and resulting assumptions of problem above;

- (a) Assume that  $\|\mathbf{x}\|_2, \|\mathbf{y}\|_2 = 1$ . Show that  $|fl(\sum x_i y_i) - \sum x_i y_i| \leq \frac{nu}{1-nu}$ . Use the result of problem P46.
- (b) Let  $U, V \in \mathbb{R}^{n \times n}$  be unitary matrices and assume that  $nu < 1$ . Show that  $fl(UV) = UV + E$ , where  $E \in \mathbb{R}^{n \times n}$  is such that  $|E_{ij}| \leq \frac{nu}{1-nu}$ .

### 1.9.8 Backward error analysis of back-substitution method

Let  $U \in \mathbb{F}^{2 \times 2}$ ,  $\mathbf{b} \in \mathbb{F}^2$ , and consider the problem: find  $\mathbf{x} \in \mathbb{R}^2$  satisfying  $U\mathbf{x} = \mathbf{b}$ , i.e.,

$$\begin{bmatrix} u_{11} & u_{12} \\ 0 & u_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}. \quad (1.60)$$

In this section, we conduct backward error analysis for computing an approximate solution  $\hat{\mathbf{x}} \in \mathbb{F}$  to (1.60) using the back-substitution method in floating-point representation. Our aim is to estimate the relative error  $\|\mathbf{x}\|_2^{-1} \|\hat{\mathbf{x}} - \mathbf{x}\|_2$ . As we are not working with any particular problem setting, all estimates are given in the  $\|\cdot\|_2$ -norm ("2-norm", also known as Euclidean norm).

The first step in backward error analysis is to use the model of floating-point arithmetic operations in (1.58) to determine which linear system  $\hat{\mathbf{x}}$  solves exactly. The back-substitution method computes  $\hat{\mathbf{x}}$  as follows:

```

hat_x2 = b2/u22; % solve x2
a = u12*hat_x2; % compute update to the load
hat_b1 = b1 - a; % update the load
hat_x1 = hat_b1/u11; % solve x1

```

[See outline of this section in Youtube](#)

See derivation of equation for  $\hat{\mathbf{x}}$  in Youtube

For clarity, all arithmetic operations in the above script are conducted one-by-one. By model (1.58),

$$\hat{x}_2 = fl\left(\frac{b_2}{u_{22}}\right) = (1 + \delta_1)\frac{b_2}{u_{22}}. \quad (1.61)$$

$$a = fl(u_{12}\hat{x}_2) = (1 + \delta_2)u_{12}\hat{x}_2 \quad (1.62)$$

$$\hat{b}_1 = fl(b_1 - a) = (1 + \delta_3)(b_1 - a) = (1 + \delta_3)(b_1 - (1 + \delta_2)u_{12}\hat{x}_2) \quad (1.63)$$

$$\hat{x}_1 = fl\left(\frac{\hat{b}_1}{u_{11}}\right) = (1 + \delta_4)\frac{\hat{b}_1}{u_{11}} \quad (1.64)$$

$$= \frac{1}{u_{11}}(b_1 - u_{12}\hat{x}_2(1 + \delta_2))(1 + \delta_3)(1 + \delta_4) \quad (1.65)$$

for some  $|\delta_i| \leq u$ ,  $i \in \{1, \dots, 4\}$ . Next, we modify (1.65) and (1.61) to find  $\delta U \in \mathbb{R}^{2 \times 2}$  such that  $\hat{\mathbf{x}}$  satisfies

$$(U + \delta U)\hat{\mathbf{x}} = \mathbf{b}.$$

Multiplying (1.61) by  $(1 + \delta_1)^{-1}u_{22}$  gives

$$\frac{u_{22}}{1 + \delta_1}\hat{x}_2 = b_2. \quad (1.66)$$

Multiplying (1.65) by  $(1 + \delta_3)^{-1}(1 + \delta_4)^{-1}u_{11}$  and rearranging the terms yields

$$\frac{u_{11}}{(1 + \delta_3)(1 + \delta_4)}\hat{x}_1 + u_{12}(1 + \delta_2)\hat{x}_2 = b_1. \quad (1.67)$$

Before proceeding, we simplify expressions (1.66) and (1.67). As  $\frac{1}{1 + \delta_1} = 1 - \frac{\delta_1}{1 + \delta_1}$ ,

$$(1 + \theta_1)u_{22}\hat{x}_2 = b_2 \quad \text{for} \quad |\theta_1| \leq \frac{u}{1 - u}. \quad (1.68)$$

Application of Lemma 1.6 to (1.67) leads to

$$(1 + \theta_2)u_{11}\hat{x}_1 + u_{12}(1 + \delta_2)\hat{x}_2 = b_1 \quad \text{for} \quad |\theta_2| \leq \frac{2u}{1 - 3u}. \quad (1.69)$$

By (1.66)-(1.69),

$$(U + \delta U)\hat{\mathbf{x}} = \mathbf{b} \quad \text{for} \quad \delta U := \begin{bmatrix} \theta_2 u_{11} & \delta_2 u_{12} \\ 0 & \theta_1 u_{22} \end{bmatrix}, \quad (1.70)$$

where

$$|\theta_2| \leq \frac{2u}{1 - 3u}, \quad |\theta_1| \leq \frac{u}{1 - u}, \quad \text{and} \quad |\delta_2| \leq u. \quad (1.71)$$

Thus,  $\hat{\mathbf{x}}$  is the exact solution to the *perturbed linear system* (1.70). Perturbation estimate in Theorem 1.3 gives an upper bound for the relative error,

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{\kappa_2(U)}{1 - \kappa_2(U) \frac{\|\delta U\|_2}{\|U\|_2}} \frac{\|\delta U\|_2}{\|U\|_2}. \quad (1.72)$$

Application of (1.72) requires estimate for the condition number  $\kappa_2(U)$  and the size of the relative perturbation  $\|\delta U\|_2 \|U\|_2^{-1}$ . The condition number depends on the (unknown) matrix  $U$ , hence it is computed numerically when  $U$  is fixed. The size of the relative perturbation is estimated using the following technical result:

[See video on estimating the relative perturbation in Youtube](#)

**Lemma 1.8.** *Let  $A \in \mathbb{R}^{n \times n}$ . Then there holds that*

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2$$

where  $\|A\|_F := \left( \sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2}$  is the Frobenius-norm of  $A$ .

Using the Frobenius-norm allows us to obtain estimates for  $\|U\|_2$  from entry-wise estimates of  $\delta U$ .

*Proof.* Problem P48. □

The Frobenius norm of  $\delta U$  defined in (1.70) satisfies

$$\|\delta U\|_F = \left( |\theta_2|^2 |u_{11}|^2 + |\delta_2|^2 |u_{12}|^2 + |\theta_2|^2 |u_{22}|^2 \right)^{1/2}. \quad (1.73)$$

The coefficients in the RHS of (1.73) are estimated by their maximum as

$$|\theta_1|, |\theta_2|, |\delta_2| \leq \frac{2u}{1 - 3u}.$$

Hence,

$$\|\delta U\|_F \leq \frac{2u}{1 - 3u} (|u_{11}|^2 + |u_{12}|^2 + |u_{22}|^2)^{1/2} = \frac{2u}{1 - 3u} \|U\|_F.$$

Using Lemma 1.8 twice and dividing by  $\|U\|_2$  gives

$$\frac{\|\delta U\|_2}{\|U\|_2} \leq \frac{2u}{1 - 3u} \sqrt{2}. \quad (1.74)$$

Application of Theorem 1.3 gives the relative error estimate

$$\frac{\|\mathbf{x} - fl(\mathbf{x})\|_2}{\|\mathbf{x}\|_2} \leq \frac{\kappa_2(U)}{1 - \frac{2u}{1-3u} \sqrt{2} \kappa_2(U)} \frac{2u}{1 - 3u} \sqrt{2}.$$

Assuming that  $\kappa(U)2u(1-3u)^{-1} \ll 1$  and neglecting the higher-order terms leads to the approximation

$$\frac{\kappa_2(U)}{1 - \frac{2u}{1-3u}\sqrt{2}u\kappa_2(U)} \frac{2u}{1-3u}\sqrt{2} \approx \kappa_2(U)2\sqrt{2}u.$$

The error due to floating-point representation is relative to condition number of matrix  $U$ . This is typical result in numerical stability analysis.

### 1.9.9 problems

P48. (1p) Prove Lemma 1.8

P49. (1p) Let  $A, \delta A \in \mathbb{R}^{n \times n}$  satisfy  $|\delta A_{ij}| \leq \epsilon_{ij}|A_{ij}|$  for  $i, j \in \{1, \dots, n\}$ . In addition, denote  $\epsilon := \max_{i,j \in \{1, \dots, n\}} |\epsilon_{ij}|$ . Show that

- (i)  $\|\delta A\|_F \leq \epsilon \|A\|_F$
- (ii)  $\|\delta A\|_1 \leq \epsilon \|A\|_1$
- (iii)  $\|\delta A\|_\infty \leq \epsilon \|A\|_\infty$ .

Here  $\|\cdot\|_F$  is the Frobenius norm and  $\|\cdot\|_1, \|\cdot\|_\infty$  are the operator norms induced by the vector norms

$$\|\mathbf{x}\|_1 := \sum_{i=1}^n |x_i| \quad \text{and} \quad \|\mathbf{x}\|_\infty := \max_{i \in \{1, \dots, n\}} |\mathbf{x}|_i \quad (1.75)$$

for  $\mathbf{x} \in \mathbb{R}^n$ .

P50. (1p) Let  $A \in \mathbb{R}^{n \times n}$  and denote it's floating-point representation by  $\hat{A} \in \mathbb{R}^{n \times n}$ . The floating point representation satisfies

$$\hat{A}_{ij} = (1 + \delta_{ij})A_{ij}$$

for  $|\delta_{ij}| \leq u$  and  $i, j \in \{1, \dots, n\}$ . Let  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^n$  satisfy

$$A\mathbf{x} = \mathbf{b} \quad \text{and} \quad \hat{A}\hat{\mathbf{x}} = \mathbf{b}$$

for some  $\mathbf{b} \in \mathbb{R}^n$ . Give estimate for the relative error  $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \|\mathbf{x}\|_2^{-1}$ .

### 1.9.10 Numerical stability of Cholesky decomposition

Let  $A \in \mathbb{F}^{n \times n}$  be s.p.d.,  $\mathbf{b} \in \mathbb{F}^n$ , and consider the problem: find  $\mathbf{x} \in \mathbb{R}^n$  satisfying

$$A\mathbf{x} = \mathbf{b}. \quad (1.76)$$

[introduction to numerical stability of Cholesky factorisation](#) [youtube](#) As discussed in Section 1.9, solution  $\mathbf{x}$  can be computed by first calculating the Cholesky factorisation of  $A$  and then applying the back-substitution method twice. The accuracy of a solution computed using this strategy in floating-point representation is studied by separately analysing errors due to Cholesky factorisation and back-substitution method.

Denote the (approximate) Cholesky factor of  $A$  computed in floating-point representation by  $\hat{L} \in \mathbb{F}^{n \times n}$ . In this section, we give backward error analysis for replacing  $A$  in (1.76) by  $\hat{L}\hat{L}^T$ . Let  $\hat{\mathbf{x}} \in \mathbb{R}^n$  satisfy:

$$\hat{L}\hat{L}^T\hat{\mathbf{x}} = \mathbf{b}.$$

Our aim is to bound the relative error  $\|\mathbf{x}\|_2^{-1}\|\mathbf{x} - \hat{\mathbf{x}}\|_2$ . First, we study the entries of  $\delta A = A - \hat{L}\hat{L}^T$ , by estimating the terms

$$\left| a_{ij} - \sum_{k=1}^j \hat{l}_{ik}\hat{l}_{jk} \right|.$$

from above. Recall that there are different strategies for computing the factor  $\hat{L}$ . Here, we consider the left-looking strategy that computes the off-diagonal entries ( $i, j \in \{1, \dots, n\}, i < j$ ) of the factor  $\hat{L} \in \mathbb{F}^{n \times n}$  as

$$\hat{l}_{ij} = fl \left( \frac{1}{\hat{l}_{jj}} \left[ a_{ij} - \sum_{k=1}^{j-1} \hat{l}_{ik}\hat{l}_{jk} \right] \right). \quad (1.77)$$

Diagonal terms are computed in a similar manner and they are not explicitly treated in the following. The floating-point error related to evaluation of the sum in (1.77) is estimated using Lemma 1.7. In the following, denote by  $|\cdot| : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^{n \times n}$  the entry-wise absolute value of a matrix, i.e.,

$$(|B|)_{ij} = |b_{ij}|$$

for  $i, j \in \{1, \dots, n\}$  and  $B \in \mathbb{R}^{n \times n}$ .

**Theorem 1.4.** *Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d. and  $\hat{L} \in \mathbb{F}^{n \times n}$  be the Cholesky factor of  $A$  computed in floating-point representation. In addition, let  $\delta A = A - \hat{L}\hat{L}^T$ . Then for  $i, j \in \{1, \dots, n\}$  there holds that*

[See video on Theorem 1.4 in Youtube](#)

$$|\delta A_{ij}| \leq \gamma_n (|\hat{L}||\hat{L}|^T)_{ij} \quad \text{where} \quad \gamma_n := \frac{(n+1)u}{1 - (n+2)u}.$$

*Proof.* We give the proof only for off-diagonal entries. Observe, that  $\delta A = \delta A^T$ , hence, we assume that  $i < j$ . Proof for diagonal entries follows using similar arguments. Application of Lemma 1.7 to (1.77) gives

$$(1 + \theta_j) \hat{l}_{ij} \hat{l}_{jj} = a_{ij} - \sum_{k=1}^{j-1} \hat{l}_{ik} \hat{l}_{jk} (1 + \theta_k)$$

for  $|\theta_k| \leq \frac{(k+1)u}{1-(k+2)u}$ ,  $k = \{1, \dots, j\}$ . Rearranging the terms gives

$$\sum_{k=1}^j \hat{l}_{ik} \hat{l}_{jk} (1 + \theta_k) = a_{ij},$$

and further

$$\left| a_{ij} - \sum_{k=1}^j \hat{l}_{ik} \hat{l}_{jk} \right| \leq \gamma_k \sum_{k=1}^j |\hat{l}_{ik}| |\hat{l}_{jk}|.$$

where  $\gamma_n$  is the upper bound for the absolute value of coefficients  $|\theta_k|$  for  $k \in \{1, \dots, n\}$ .  $\square$

Identical to Section 1.9.8, an upper-bound for the relative error follows from the perturbation estimate in Theorem 1.3,

$$\|\mathbf{x}\|_2^{-1} \|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq \frac{\kappa_2(A)}{1 - \frac{\|\delta A\|_2}{\|A\|_2} \kappa_2(A)} \frac{\|\delta A\|_2}{\|A\|_2}. \quad (1.78)$$

Application of (1.78) requires estimate for the size of the relative perturbation,  $\|\delta A\|_2 \|A\|_2^{-1}$ . Before proceeding, we need the following technical estimates:

P51. (1p) Let  $B, C \in \mathbb{R}^{n \times n}$ , and  $|B|, |C| \in \mathbb{R}^{n \times n}$ . Show that

$$\|B\|_2 \leq \||B|\|_2 \quad (1.79)$$

$$\|B\|_2 \leq \||C|\|_2 \quad \text{if} \quad |B|_{ij} \leq |C|_{ij} \quad \text{for} \quad i, j \in \{1, \dots, n\} \quad (1.80)$$

$$\|B^T\|_2 \leq \|B\|_2. \quad (1.81)$$

Also, recall Lemma 1.8 and the result proven in problem P32: Let  $B \in \mathbb{R}^{n \times n}$  and  $F \in \mathbb{R}^{n \times n}$  satisfy  $B = FF^T$ . Then

$$\|B\|_2 = \|F\|_2^2.$$



[part 1 of this proof](#) **Lemma 1.9.** *Make the same assumptions and use the same notation as in Theorem 1.4. In addition, assume that  $n\gamma_n < 1$ . Then there holds that*

$$\|\hat{L}\|\hat{L}^T\|_2 \leq \frac{n}{1-n\gamma_n}\|A\|_2 \quad \text{and} \quad \frac{\|\delta A\|_2}{\|A\|_2} \leq \frac{n\gamma_n}{1-n\gamma_n}.$$

*Proof.* We begin by estimating  $\|\hat{L}\|\hat{L}^T\|_2$  from above. Using (1.81) gives

$$\|\hat{L}\|\hat{L}^T\|_2 \leq \|\hat{L}\|_2^2. \quad (1.82)$$

We eliminate the entry-wise absolute value using the norm equivalence given in Lemma 1.8 and the definition of the Frobenius-norm as [See part 2 of this proof in Youtube](#)

$$\|\hat{L}\|_2 \leq \|\hat{L}\|_F = \|\hat{L}\|_F \leq \sqrt{n}\|\hat{L}\|_2. \quad (1.83)$$

As  $A - \delta A = \hat{L}\hat{L}^T$ , problem P32 states that  $\|\hat{L}\|_2^2 = \|A - \delta A\|_2$ . Further, using triangle inequality gives

$$\|\hat{L}\|_2^2 = \|A - \delta A\|_2 \leq \|A\|_2 + \|\delta A\|_2 \quad (1.84)$$

Combining (1.82), (1.83), (1.84) gives the estimate

$$\|\hat{L}\|\hat{L}^T\|_2 \leq n\|A\|_2 + n\|\delta A\|_2 \quad (1.85)$$

By (1.80) and Theorem 1.4,

$$\|\delta A\|_2 \leq \gamma_n \left\| \hat{L}\|\hat{L}^T\|_2 \right\|_2. \quad (1.86)$$

Combining (1.86) and (1.85), rearranging the terms, and dividing by  $(1 - n\gamma_n) > 0$  yields

$$\|\hat{L}\|\hat{L}^T\|_2 \leq \frac{n}{1-n\gamma_n}\|A\|_2$$

Combining the above equation with (1.86) completes the proof.  $\square$

### 1.9.11 Problems

P52. (2p) Let  $U \in \mathbb{F}^{n \times n}$  be an upper triangular matrix and  $\mathbf{b} \in \mathbb{F}^n$ . In floating-point representation, the back-substitution method computes an approximate solution  $\hat{\mathbf{x}} \in \mathbb{F}^n$  to the problem

$$U\mathbf{x} = \mathbf{b}.$$

as

$$\hat{x}_i = fl \left( \frac{1}{u_{ii}} \left( b_i - \sum_{j=i+1}^n u_{ij}\hat{x}_j \right) \right). \quad (1.87)$$

- (a) Let  $\delta U \in \mathbb{R}^{n \times n}$  have entries

$$(\delta U)_{ij} = \theta_{ij} u_{ij},$$

where the scalars  $\theta_{ij}$  have an upper bound

$$|\theta_{ij}| \leq \frac{(n-j+2)u}{1-(n-j+3)u}$$

for  $i, j \in \{1, \dots, n\}$ . Recall that  $u$  is the machine epsilon.

Use Lemma 1.7 to show that  $\hat{\mathbf{x}}$  satisfies  $(U + \delta U)\hat{\mathbf{x}} = \mathbf{b}$ .

- (b) Show that  $\|\delta U\|_2 \leq \frac{(n+1)u}{1-(n+2)u} \|U\|_2$ .  
 (c) Give estimate for the relative error

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2}{\|\mathbf{x}\|_2}.$$

P53. (2p) Consider the matrix  $A$  defined in Problem 34, where

$$\mathbf{a}_{21} = [\sqrt{\pi+1} \quad \sqrt{\pi+2} \quad \dots \quad \sqrt{\pi+(n-1)}]^T$$

and  $a_{11} = \|\mathbf{a}_{21}\|_2^2 + 1$ .

- (a) Compute by hand the value  $\|\mathbf{a}_{21}\|_2$   
 (b) Construct the matrix  $A$  in Matlab using the exact value of  $\|\mathbf{a}_{21}\|$ . Remember to define  $A$  as a sparse matrix.  
 (c) Study the estimate of Theorem 1.4 by computing the Cholesky factor  $\hat{L} \in \mathbb{F}^{n \times n}$  of  $A$  numerically for  $n = 10, 20, 40, 80, 160, 320$ . Plot

$$\max_{ij} \frac{|a_{ij} - (LL^T)_{ij}|}{(|L||L^T|)_{ij}}$$

in logarithmic scale as a function of  $n$ . Compare to  $\gamma_n$ .

P54. (2p) Make same assumptions and use the same notation as in Problem P53

- (a) Let  $n = 2^k, k = 2, \dots, 15$  and consider the linear system  $A\mathbf{x} = \mathbf{e}_1$ . Use formula given in problem P34 to solve  $\mathbf{x}$  using pen and paper.  
 (b) Plot the relative error between the exact solution  $\mathbf{x}$  computed in a) and the approximate solution calculated by Matlab (which in this case uses Cholesky factorization). Remember to define  $A$  again as a sparse matrix. Choose an informative plot. Also plot the condition number of  $A$ .

## 1.10 Stable QR-decomposition

The QR-decomposition of a matrix  $A \in \mathbb{R}^{n \times n}$ ,

$$A = QR \quad \text{where } Q \in \mathbb{R}^{n \times n} \text{ is unitary and } R \in \mathbb{R}^{n \times n} \text{ is upper triangular,}$$

[See introduction to stable QR-factorization in Youtube](#)

is an important ingredient in iterative solution methods, solution of least squares problems, etc.

The simplest way to compute the QR decomposition is by the *Gram–Schmidt orthogonalization process*. However, when the Gram–Schmidt orthogonalization process is implemented in floating-point representation, the computed QR factorisation suffers from *loss of orthogonality*. This is, the computed matrix  $Q$  can be far from an orthogonal matrix.

In this section, we discuss two processes that are used to compute QR factorisation of  $A$  using *unitary elimination matrices*  $\{U_i\}_{i=1}^N \subset \mathbb{R}^{n \times n}$  that transform  $A$  into an upper triangular matrix  $R \in \mathbb{R}^{n \times n}$  as

$$U_N \cdots U_1 A = R. \quad (1.88)$$

As unitary matrices are invertible,

$$A = QR \quad \text{where } Q = U_1 \cdots U_N. \quad (1.89)$$

By direct computation,  $Q^T Q = I$ , hence  $Q$  is a unitary matrix and (1.89) is the QR factorization of  $A$ .

Computing the product of unitary matrices in floating-point representation is very accurate, see Problem 47. Due to this, the loss of orthogonality in  $Q$  is well under control (much better in comparison to the Gram–Schmidt process).

This section is organised as follows. We begin by review of computing QR-factorisation by the Gram–Schmidt process. Then we discuss two processes that generate unitary elimination matrices, Givens rotation and Householder reflection. Both constructions are based on geometric arguments.

### 1.10.1 Gram–Schmidt orthogonalization process

*This section is review material and can be skipped*

Gram–Schmidt orthogonalization process is based on the following Lemma. In what follows, “orthogonality” refers to orthogonality in the sense of the Euclidean inner product.

**Lemma 1.10.** Let  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\} \subset \mathbb{R}^m$ ,  $k < m$ , be a set of orthonormal vectors, i.e.,  $\|\mathbf{q}_i\|_2 = 1$ ,  $i = 1, \dots, k$ , and

$$\mathbf{q}_i^T \mathbf{q}_j = 0 \quad \text{for } i \neq j.$$

In addition, assume that  $\mathbf{a} \in \mathbb{R}^m$  does not belong to  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k)$  and define

$$\mathbf{q}_{k+1} = \frac{\tilde{\mathbf{q}}_{k+1}}{\|\tilde{\mathbf{q}}_{k+1}\|_2}, \quad \text{where } \tilde{\mathbf{q}}_{k+1} = \mathbf{a} - \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i. \quad (1.90)$$

Then  $\{\mathbf{q}_1, \dots, \mathbf{q}_{k+1}\} \subset \mathbb{R}^n$  is a set of orthonormal vectors and

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{k+1}) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a}). \quad (1.91)$$

*Proof.* To begin with note that  $\tilde{\mathbf{q}}_{k+1}$  defined in (1.90) is a nonzero vector because  $\mathbf{a} \notin \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k)$ , i.e.,  $\mathbf{a}$  cannot be given as a linear combination of  $\mathbf{q}_1, \dots, \mathbf{q}_k$ .

To prove the orthonormality of the set  $\{\mathbf{q}_1, \dots, \mathbf{q}_{k+1}\}$  it is enough to prove that  $\mathbf{q}_{k+1}$  is of unit Euclidean length and orthogonal to  $\mathbf{q}_1, \dots, \mathbf{q}_k$ . From the first equation in (1.90), it is obvious that  $\|\mathbf{q}_{k+1}\|_2 = 1$ . Moreover, since  $\mathbf{q}_{k+1}$  and  $\tilde{\mathbf{q}}_{k+1}$  are parallel, it is actually enough to show that  $\tilde{\mathbf{q}}_{k+1}$  is orthogonal to  $\mathbf{q}_1, \dots, \mathbf{q}_k$ : for any  $j = 1, \dots, k$ , we have

$$\tilde{\mathbf{q}}_{k+1}^T \mathbf{q}_j = \left( \mathbf{a} - \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i \right)^T \mathbf{q}_j = \mathbf{a}^T \mathbf{q}_j - \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) (\mathbf{q}_i^T \mathbf{q}_j) = \mathbf{a}^T \mathbf{q}_j - \mathbf{a}^T \mathbf{q}_j = 0$$

due to the orthonormality of  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$ .

Although (1.91) follows straightforwardly from the definition of linear span, let us anyway carefully prove it for the sake of completeness. Assume first that  $\mathbf{x} \in \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{k+1})$ , i.e.,

$$\mathbf{x} = \sum_{i=1}^{k+1} \alpha_i \mathbf{q}_i = \sum_{i=1}^k \alpha_i \mathbf{q}_i + \alpha_{k+1} \mathbf{q}_{k+1}$$

for some  $\alpha \in \mathbb{R}^{k+1}$ . Note that (1.90) can be rewritten in the form

$$\mathbf{q}_{k+1} = \frac{1}{\|\tilde{\mathbf{q}}_{k+1}\|_2} \left( \mathbf{a} - \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i \right).$$

Hence,

$$\mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{q}_i + \frac{\alpha_{k+1}}{\|\tilde{\mathbf{q}}_{k+1}\|_2} \left( \mathbf{a} - \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i \right) = \sum_{i=1}^k \left( \alpha_i - \frac{\alpha_{k+1} \mathbf{a}^T \mathbf{q}_i}{\|\tilde{\mathbf{q}}_{k+1}\|_2} \right) \mathbf{q}_i + \frac{\alpha_{k+1}}{\|\tilde{\mathbf{q}}_{k+1}\|_2} \mathbf{a},$$

which is obviously in  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a})$ , meaning that  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{k+1}) \subset \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a})$ .

On the other hand, if  $\mathbf{x} \in \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a})$ , then for some  $\alpha \in \mathbb{R}^{k+1}$ ,

$$\begin{aligned} \mathbf{x} &= \sum_{i=1}^k \alpha_i \mathbf{q}_i + \alpha_{k+1} \mathbf{a} = \sum_{i=1}^k \alpha_i \mathbf{q}_i + \alpha_{k+1} \left( \|\tilde{\mathbf{q}}_{k+1}\|_2 \mathbf{q}_{k+1} + \sum_{i=1}^k (\mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i \right) \\ &= \sum_{i=1}^k (\alpha_i + \alpha_{k+1} \mathbf{a}^T \mathbf{q}_i) \mathbf{q}_i + \alpha_{k+1} \|\tilde{\mathbf{q}}_{k+1}\|_2 \mathbf{q}_{k+1}, \end{aligned}$$

which clearly belongs to  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{k+1})$ . Hence, also  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a}) \subset \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{k+1})$ , which completes the proof.  $\square$

The intuitive idea of (1.90) is that one first subtracts from  $\mathbf{a}$  its projections onto the one-dimensional subspaces defined by  $\mathbf{q}_1, \dots, \mathbf{q}_n$ , leaving only the component of  $\mathbf{a}$  orthogonal to  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k)$ , and then this component is normalized. In fact, one can write the second equation of (1.90) in the form

$$\tilde{\mathbf{q}}_{k+1} = (I - P_k) \mathbf{a},$$

where  $P_k \in \mathbb{R}^{m \times m}$  is the orthogonal projection matrix onto the subspace  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_k)$ .

Using Lemma 1.10, it is straightforward to compute an orthonormal basis for the subspace

$$R(A) = \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_n) \subset \mathbb{R}^m,$$

assuming the columns  $\mathbf{a}_1, \dots, \mathbf{a}_n$  of the matrix  $A \in \mathbb{R}^{m \times n}$  are linearly independent, i.e., assuming  $N(A) = \{0\}$ . Indeed, such basis  $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$  can be recursively obtained via

$$\mathbf{q}_j = \frac{\tilde{\mathbf{q}}_j}{\|\tilde{\mathbf{q}}_j\|_2}, \quad \text{where} \quad \tilde{\mathbf{q}}_j = \mathbf{a}_j - \sum_{i=1}^{j-1} (\mathbf{a}_j^T \mathbf{q}_i) \mathbf{q}_i, \quad \text{for } j = 1, \dots, n.$$

In other words, one first defines  $\mathbf{q}_1$  by simply normalizing  $\mathbf{a}_1$ , then one computes a unit vector  $\mathbf{q}_2$  that is orthogonal to  $\mathbf{q}_1$  and satisfies  $\text{span}(\mathbf{q}_1, \mathbf{q}_2) =$

$\text{span}(\mathbf{q}_1, \mathbf{a}_2) = \text{span}(\mathbf{a}_1, \mathbf{a}_2)$ , then one continues by computing a unit vector  $\mathbf{q}_2$  that is orthogonal to both  $\mathbf{q}_1$  and  $\mathbf{q}_2$  and satisfies

$$\text{span}(\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3) = \text{span}(\mathbf{q}_1, \mathbf{q}_2, \mathbf{a}_3) = \text{span}(\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3),$$

and so on until  $\mathbf{q}_n$  is computed and it holds that  $\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_n) = \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_n) = R(A)$ .

Take note that one can get the original columns of  $A$  back via

$$\mathbf{a}_j = \|\tilde{\mathbf{q}}_j\|_2 \mathbf{q}_j + \sum_{i=1}^{j-1} (\mathbf{a}_j^T \mathbf{q}_i) \mathbf{q}_i, \quad j = 1, \dots, n, \quad (1.92)$$

which demonstrates that, for any  $j = 1, \dots, n$ , the  $j$ th column  $\mathbf{a}_j$  of  $A$  can be given as a linear combination of  $\mathbf{q}_1, \dots, \mathbf{q}_j$ , i.e., of (only) the first  $j$  orthonormal basis vectors of  $R(A)$  produced by the Gram-Schmidt process. Defining in the standard manner  $Q = [\mathbf{q}_1, \dots, \mathbf{q}_n] \in \mathbb{R}^{m \times n}$  and collecting the coefficients in the linear combinations of (1.92) as columns of an upper triangular matrix  $R \in \mathbb{R}^{n \times n}$ , the equations (1.92) can be written neatly in a matrix form

$$A = QR.$$

To be more precise,  $R$  can be given elementwise as

$$R_{i,j} = \begin{cases} \mathbf{a}_j^T \mathbf{q}_i & \text{if } i < j, \\ \|\tilde{\mathbf{q}}_j\|_2 & \text{if } i = j, \\ 0 & \text{if } i > j. \end{cases}$$

Note also that  $Q^T Q = I \in \mathbb{R}^{n \times n}$  because the columns of  $Q$  are orthonormal. There are two implementations of the Gram-Schmidt procedure. Modified:

```
function [Q,R] = my_gsmith(A)

Q = [];
for i=1:size(A,2)
    q = A(:,i);

    for k=1:size(Q,2)
        R(k,i) = q'*Q(:,k);
        q = q - R(k,i)*Q(:,k);
    end
    R(i,i) = norm(q);
    Q(:,i) = q/R(i,i);
end
```

and the classical:

```
function [Q,R] = my_c_gsmith(A)

Q = [];
for i=1:size(A,2)
    q = A(:,i);

    for k=1:size(Q,2)
        R(k,i) = q'*Q(:,k);
    end

    for k=1:size(Q,2)
        q = q - R(k,i)*Q(:,k);
    end
    R(i,i) = norm(q);
    Q(:,i) = q/R(i,i);
end
```

The two different implementations of the Gram-Schmidt process have very different numerical stability properties. The quality of the factorization is measured by computing error in orthogonality of  $Q$ ,

$$\|I - Q^T Q\|$$

and error in the decomposition,  $\|A - QR\|$ . Orthogonality of  $Q$  is more sensitive to floating point errors than the error in the decomposition. For the modified GS, one can prove numerical stability in both of these measures, where as the classical GS is not numerically stable.

**Example 1.17.** Let  $A \in \mathbb{R}^{4 \times 3}$  be such that

$$A = \begin{bmatrix} \mathbf{1}^T \\ \epsilon I \end{bmatrix},$$

in which  $\mathbf{1} = [1 \ 1 \ \dots 1]^T$  and  $\epsilon > 0$ . Let us measure the orthogonality in the maximum-norm  $\|A\|_{\max} := \max_{ij} |A_{ij}|$ . One obtains,

$$\|I - Q_{MGS}^T Q_{MGS}\|_{\max} = \frac{1}{\sqrt{2}}\epsilon \quad \text{and} \quad \|I - Q_{CGS}^T Q_{CGS}\|_{\max} = \frac{1}{2}.$$

And

$$\|A - Q_{MGS} R_{MGS}\|_{\max} = 0 \quad \text{and} \quad \|A - Q_{CGS} R_{CGS}\|_{\max} = 0$$

Note, that these numbers were computed in floating point arithmetics.

### 1.10.2 Givens rotation

In this section, compute  $QR$ -factorization using Givens rotation matrices. In  $\mathbb{R}^{2 \times 2}$  rotation matrix has the entries

$$\begin{bmatrix} \sin \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix},$$

where  $\theta$  is a given rotation angle. The Givens rotation matrix is constructed from a  $2 \times 2$  rotation matrix that turns a given vector  $\mathbf{a} \in \mathbb{R}^2$  to the direction of  $\mathbf{e}_1$ . Using angles in program code is cumbersome, hence they are avoided in the following. We begin with an example.

**Example 1.18.** Let  $A \in \mathbb{R}^{2 \times 2}$  be such that

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = [\mathbf{a}_1 \quad \mathbf{a}_2]. \quad (1.93)$$

Next, we construct unitary matrix  $U \in \mathbb{R}^{2 \times 2}$  satisfying  $U\mathbf{a}_1 = \alpha\mathbf{e}_1$  for some  $\alpha \in \mathbb{R}$ . As  $U$  is unitary,  $\|U\mathbf{a}_1\|_2 = \|\mathbf{a}_1\|_2$  and  $\alpha = \|\mathbf{a}_1\|$ . Any  $2 \times 2$ -unitary matrix  $U$  satisfies

$$U = \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \end{bmatrix} \quad \text{where} \quad \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij} \quad \text{for} \quad i, j \in \{1, 2\}.$$

Using the condition  $UA = \|\mathbf{a}_1\|\mathbf{e}_1$  gives

$$\mathbf{u}_1^T \mathbf{a}_1 = \|\mathbf{a}_1\|_2, \quad \mathbf{u}_2^T \mathbf{a}_1 = 0, \quad \text{and} \quad \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij} \quad \text{for} \quad i, j \in \{1, 2\}.$$

We choose  $\mathbf{u}_1$  as the unit vector to the direction of  $\mathbf{a}_1$  and  $\mathbf{u}_2$  as a unit vector orthogonal to  $\mathbf{a}_1$ . This is,

$$\mathbf{u}_1 = -\frac{\mathbf{a}_1}{\|\mathbf{a}_1\|_2} \quad \text{and} \quad \mathbf{u}_2 = \frac{1}{\|\mathbf{a}_1\|_2} \begin{bmatrix} -a_{21} \\ a_{11} \end{bmatrix}.$$

Computing the product  $UA$  gives

$$UA = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix} \quad \text{i.e.} \quad A = U^T \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}.$$

which is the  $QR$ -decomposition of  $A$ .

Let of  $A \in \mathbb{R}^{n \times n}$  and  $i, j \in \{1, \dots, n\}$ ,  $i < j$ . We proceed to construct unitary matrix  $G$  such that  $(GA)_{ji} = 0$ . Let  $\hat{\mathbf{a}} = [a_{ii} \quad a_{ji}]^T$  and  $U \in \mathbb{R}^{2 \times 2}$  a unitary matrix satisfying  $U\hat{\mathbf{a}} = \|\hat{\mathbf{a}}\|_2\mathbf{e}_1$ . Suitable matrix  $U$  is constructed

See video on  $2 \times 2$ -  
Givens rotation matrices  
in Youtube

See video on Givens  
rotation matrices in  
Youtube



in Example 1.18.

Consider the linear mapping  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  (see Problem P??) defined as

$$\mathbf{y} = f(\mathbf{x}), \quad \begin{bmatrix} y_i \\ y_j \end{bmatrix} = U \begin{bmatrix} x_i \\ x_j \end{bmatrix} \quad \text{and} \quad y_k = x_k, \text{ for } k \neq i, k \neq j.$$

This is, the matrix  $U$  operates on rows  $i$  and  $j$  of vector  $\mathbf{x}$ , while all other rows are left untouched. In Matlab, the linear mapping  $f$  is evaluated simply as

```
function y = fmap(x,i,j,U)
```

```
y = x; % copy all entries to x
y([i;j]) = U*x([i;j]); % operate to rows i and j by U.
```

The matrix representation of mapping  $f$  is the unitary matrix

$$G := \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & q_{11} & & q_{12} & \\ & & & \ddots & & \\ & & q_{21} & & q_{22} & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix} \quad (1.94)$$

called as the *Givens rotation matrix*. Observe, that  $G$  depends on  $i, j$ , and  $A$ . It is customary to leave these dependencies implicit. The product  $GA$  can be computed as

$$GA = [f(\mathbf{a}_1) \quad \cdots \quad f(\mathbf{a}_n)] \quad \text{where} \quad A = [\mathbf{a}_1 \quad \cdots \quad \mathbf{a}_n].$$

Hence, the multiplication  $GA$  only modifies rows  $i$  and  $j$  of  $A$ . In addition, the entry  $(GA)_{ji}$  satisfies

$$(GA)_{ji} = \mathbf{e}_j^T GA \mathbf{e}_i = \mathbf{e}_j^T f(\mathbf{a}_i).$$

Using the definition of  $f$  gives

$$(GA)_{ji} = [0 \quad 1] U \begin{bmatrix} a_{ii} \\ a_{ji} \end{bmatrix} = 0.$$

To collect,  $G$  is an unitary matrix that only modifies rows  $i, j$  and  $(GA)_{ji} = 0$ . The  $QR$ -factorisation is computed using Givens rotation matrices as

follows. Zeros are introduced one-by-one starting from the first column. When first row has zeros at appropriate locations, the process continues to the second row. The operation  $GA$  will preserve the zeros on the preceding rows. Graphically, the  $QR$  factorisation is computed as follows

$$\begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \\ \times & \times & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \end{bmatrix}$$

$$\begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \mathbf{0} & \times \\ \mathbf{0} & \times & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \mathbf{0} & \times \\ \mathbf{0} & \mathbf{0} & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \mathbf{0} & \times \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

The rows  $i$  and  $j$  are marked with red color. As  $U$  depends on  $i, j, A$ , it is constructed separately in each step in above. The matrix  $G$  is not explicitly computed. The corresponding Matlab code is

```
function [Q,A] = my-givens_qr(A)

Q = eye(max(size(A)));

for i=1:size(A,2)
    for j=(i+1):size(A,1)

        % Construct G
        x = [A(i,i) ; A(j,i)];
        xN = [-x(2) ; x(1)];

        G = [ x'/norm(x) ; xN'/norm(xN) ];

        % Operate with G
        Q([i j],:) = G*Q([i j],:);
        A([i j],:) = G*A([i j],:);

    end
end

Q = Q';
```

P55. (2p) Let  $U \in \mathbb{R}^{2 \times 2}$  be such that  $U^T U = I$ . Consider the mapping  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  such that

$$\mathbf{y} = f(\mathbf{x}), \quad \begin{bmatrix} y_i \\ y_j \end{bmatrix} = U \begin{bmatrix} x_i \\ x_j \end{bmatrix} \quad \text{and} \quad y_k = x_k, \text{ for } k \neq i, k \neq j.$$

This is, the matrix  $U$  operates on rows  $i$  and  $j$  of vector  $\mathbf{x}$ , while all other rows are left untouched.

- (a) Show that  $f$  is a linear mapping
- (b) Show that for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{x})^T f(\mathbf{y}) = \mathbf{x}^T \mathbf{y}. \quad (1.95)$$

- (c) As  $f$  is a linear mapping, there exists  $G \in \mathbb{R}^{n \times n}$  such that  $f(x) = G\mathbf{x}$ . Show that  $G$  is unitary, if  $f$  satisfies Eq. (1.95).

P56. (2p) Let  $A \in \mathbb{R}^{2 \times 2}$  be such that

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}. \quad (1.96)$$

- (a) Construct a unitary matrix  $U \in \mathbb{R}^{2 \times 2}$  such that

$$UA = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}$$

- (b) Construct a unitary matrix  $U_2 \in \mathbb{R}^{3 \times 3}$  s.t.

$$U_2 \begin{bmatrix} 1 & 2 & 5 \\ 6 & 7 & 8 \\ 3 & 4 & 9 \end{bmatrix} = \begin{bmatrix} \times & \times & \times \\ 6 & 7 & 8 \\ 0 & \times & \times \end{bmatrix}.$$

### 1.10.3 Householder reflection

*The QR factorisation using Householder reflections is computed identically to using Givens rotations. The difference lies in the construction of the unitary elimination matrix. Additional material, read it or skip it.*

Fix  $\mathbf{x} \in \mathbb{R}^n$ . Householder reflection is the linear matrix

$$H = I - 2 \frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T \mathbf{u}},$$

where  $\mathbf{u}$  is chosen such that  $H\mathbf{x} = \|\mathbf{x}\|\mathbf{e}_1$ . This transformation is symmetric and unitary for each  $\mathbf{u} \in \mathbb{R}^n$ , so that  $H^T = H$  and  $H^T H = I$ .

The Householder reflection is based on a geometric construction. Let vector  $\mathbf{u} = \|\mathbf{x}\|\mathbf{e}_1 - \mathbf{x}$ . The Householder reflection is a reflection with respect to the hyperplane  $\mathcal{V}$  orthogonal to  $\mathbf{u}$ . Let  $P \in \mathbb{R}^{n \times n}$  be the orthogonal projection to the sub-space  $\text{span}\{\mathbf{u}\}$ , this is

$$P = \frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T \mathbf{u}}.$$

The orthogonal projection  $P$  introduces the splitting

$$\mathbf{x} = (I - P)\mathbf{x} + P\mathbf{x}, \quad (1.97)$$

in which  $(I - P)\mathbf{x} \in \mathcal{V}$  and  $P\mathbf{x} \in \mathcal{V}^\perp$ . Thus, a reflection of  $\mathbf{x}$  with respect to the hyperplane  $\mathcal{V}$  is simply

$$(I - P)\mathbf{x} - P\mathbf{x} = I - 2P. \quad (1.98)$$

Geometrically, it is easy to see that

$$P\mathbf{x} = \frac{\mathbf{u}}{2}.$$

hence, the condition  $Hx = \|\mathbf{x}\|\mathbf{e}_1$  is satisfied by the construction of  $H$ . QR-decomposition is computed using Householder transformation as

```
function [Q,A] = my_house_qr(A)

Q = eye( size(A,1) );

if( size(A,1) > size(A,2) )
    N = size(A,2);
else
    N = min(size(A)-1);
end

for i=1:N
    % Construct H
    x = A(i:end,i);
    u = -x;
    u(1,1) = norm(x)+u(1,1);

    H = eye(length(x)) - 2*u*u'/(u'*u);

    % Operate with H
```

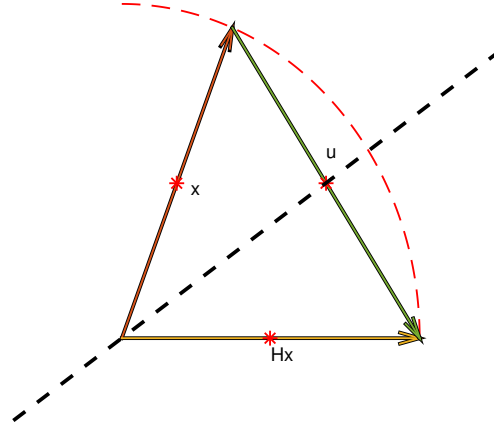


Figure 1.11: Householder transformation in  $2D$ . The transformation is a reflection of given vector  $\mathbf{x}$  with respect to the dotted line to  $H\mathbf{x}$ .

```

A(i:end,:) = H*A(i:end,:);
Q(i:end,:) = H*Q(i:end,:);
end
Q = Q';

```

Graphically, we proceed as follows

$$\begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix} \quad \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{bmatrix}$$

There are two alternative ways to construct the Householder reflection, transform the vector either to the direction of  $\mathbf{e}_1$  or  $-\mathbf{e}_1$ . Reflection with respect to the longest  $\mathbf{u}$  is chosen to avoid division by zero and to guarantee numerical stability. Note that this important feature is not included in the example code given above.

P57. (1p) Consider the matrix

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

- (a) Find a rotation matrix  $Q \in \mathbb{R}^{2 \times 2}$  and a permutation matrix  $P \in \mathbb{R}^{4 \times 4}$  such that  $(UA)_{31} = 0$ , in which

$$U = P^T \begin{bmatrix} Q & 0 \\ 0 & I \end{bmatrix} P.$$

Check that  $U$  is an unitary matrix.

- (b) Find the Householder reflection matrix  $H \in \mathbb{R}^{4 \times 4}$  such that

$$H \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = 2\mathbf{e}_1.$$

Compute  $HA$ .

## Chapter 2

# Iterative solution methods

Let  $A \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ , and consider the linear system

$$A\mathbf{x} = \mathbf{b}. \quad (2.1)$$

Computing the exact solution to (2.1) is often not necessary. This is the case, for example, when the linear system is related to the (approximate) solution of the two-dimensional Poisson's equation  $-\Delta u = f$  using the finite difference method, see Section 1.3. In this case, the entries  $x_i$  are the approximate nodal values of the exact solution  $u$ . Even if (2.1) is solved exactly, the obtained nodal values have *discretization error* due to the finite difference approximation. Hence, it is sufficient to compute an approximate solution  $\hat{\mathbf{x}}$  with error  $\hat{\mathbf{x}} - \mathbf{x}$  of the same order as the discretization error. [See introduction of Chapter 2 \(this is the outline of Week 5\) in Youtube](#)

In this Chapter, we discuss iterative solution methods for finding an approximate solution  $\hat{\mathbf{x}}$  to (2.1). Iterative solution methods are processes that generate a sequence of approximate solutions  $\{\mathbf{x}_i\}$  satisfying  $\mathbf{x}_i \rightarrow \mathbf{x}$ . The process is stopped, when a sufficiently good approximation has been constructed. The quality of approximation  $\mathbf{x}_i$  can be estimated by computing the relative residual  $\|A\mathbf{x} - \mathbf{b}\|_2 \|\mathbf{b}\|_2^{-1}$ .

This Chapter is organised as follows. We begin by discussing iterative methods based on fixed-point techniques. Then we assume that  $A$  is s.p.d. and show that solution of (2.1) is equivalent to solution of a quadratic minimisation problem. We derive *Conjugate Gradient* (CG) method as line search iteration for solving this minimisation problem. Finally, we discuss iterative methods based on subspace projection techniques.

## 2.1 Methods based on fixed-point iteration

Let  $f : \mathbb{R} \mapsto \mathbb{R}$ , and consider the problem: find  $x \in \mathbb{R}$  satisfying

$$x = f(x), \quad (2.2)$$

See video on fixed-point methods in Youtube

We say that equation (2.2) is in *fixed-point form* and call  $x$  as *fixed-point*. Fixed-point iteration generates a sequence  $\{x_i\} \subset \mathbb{R}$  from a given initial guess  $x_0 \in \mathbb{R}$  by

$$x_i = f(x_{i-1}) \quad \text{for } i \in \mathbb{N}.$$

The function  $f : \mathbb{R} \mapsto \mathbb{R}$  satisfying

$$|f(x) - f(y)| \leq L|x - y| \quad \text{for } L < 1 \quad \text{and each } x, y \in \mathbb{R}.$$

is called as *contraction*. If  $f$  is a contraction, the Banach-fixed point theorem guarantees that the sequence  $\{x_i\}$  converges to fixed point  $x$  satisfying  $x = f(x)$ , i.e.,

$$x = \lim_{i \rightarrow \infty} x_i.$$

The convergence rate, i.e., decay rate of  $|x - x_i|$  depends on  $L$ .

Next, we apply fixed-point iteration to solution of (2.1). First, we write the linear system in fixed-point form by decomposing  $A$  as

$$A = B + (A - B), \quad (2.3)$$

where the matrix  $B \in \mathbb{R}^{n \times n}$  is chosen so that the linear system  $B\mathbf{y} = \mathbf{g}$  is computationally inexpensive to solve. In Gauss-Seidel method,  $B$  is the lower triangular part of  $A$ , this is

$$\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}_A = \begin{bmatrix} \times & 0 & 0 \\ \times & \times & 0 \\ \times & \times & \times \end{bmatrix}_B + \begin{bmatrix} 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{bmatrix}_{A-B}.$$

In the Jacobi-iteration  $B = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$ , this is,

$$\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}_A = \begin{bmatrix} \times & 0 & 0 \\ 0 & \times & 0 \\ 0 & 0 & \times \end{bmatrix}_B + \begin{bmatrix} 0 & \times & \times \\ \times & 0 & \times \\ \times & \times & 0 \end{bmatrix}_{A-B}.$$

Here  $\times$  indicates entries of matrix  $A$ . Next, assume that  $B$  is invertible. Using the splitting in (2.3) to (2.1), and inverting  $B$  gives

$$\mathbf{x} = B^{-1}(\mathbf{b} - (A - B)\mathbf{x}). \quad (2.4)$$



The corresponding fixed-point iteration<sup>1</sup>

$$\mathbf{x}_{i+1} = \mathbf{c} + C\mathbf{x}_i, \quad (2.5)$$

where  $\mathbf{c} := B^{-1}\mathbf{b}$  and  $C := -B^{-1}(A - B)$  is called as the *iteration matrix*. Convergence of  $\mathbf{x}_i$  to  $\mathbf{x}$  is studied in the next Lemma.

**Lemma 2.1.** *Let  $\|\cdot\|$  be a vector norm,  $\mathbf{c} \in \mathbb{R}^n$ , and  $C \in \mathbb{R}^{n \times n}$  a matrix satisfying  $N(I - C) = \{0\}$ . In addition, let  $\mathbf{x}_0 \in \mathbb{R}^n$  be a given initial value, and consider the iteration*

$$\mathbf{x}_i = \mathbf{c} + C\mathbf{x}_{i-1} \quad \text{for } i \in \mathbb{N}.$$

Then the error  $\mathbf{e}_i := \mathbf{x} - \mathbf{x}_i$  satisfies

$$\mathbf{e}_i = C^i \mathbf{e}_0 \quad \text{and} \quad \|\mathbf{e}_i\| \leq \|C\|^i \|\mathbf{e}_0\|$$

for any  $i \in \mathbb{N}$ .

The above Lemma states that the fixed-point iteration in (2.5) converges to  $\mathbf{x} \in \mathbb{R}^n$  satisfying  $A\mathbf{x} = \mathbf{b}$  if the iteration matrix satisfies  $C^k \rightarrow 0$ , when  $k \rightarrow \infty$ .

*Proof.* See Problem P59 □

### 2.1.1 Problems

P58. (1p) Consider the vector norm,

$$\|\mathbf{x}\|_\infty := \max_i |x_i|. \quad (2.6)$$

(a) Show that the related matrix norm

$$\|A\|_\infty = \max_{\mathbf{x} \in \mathbb{R}^n} \frac{\|A\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \quad (2.7)$$

can be computed as  $\|A\|_\infty := \max_i \sum_j |a_{ij}|$ .

(b) Show that

$$\|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty \quad \text{and} \quad \|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2.$$

for any  $\mathbf{x} \in \mathbb{R}^n$ .

---

<sup>1</sup>The iteration matrix should  $C$  not be constructed in any practical implementation. Instead  $\mathbf{x}_{i+1}$  is solved from the linear system  $B\mathbf{x}_{i+1} = \mathbf{b} - (A - B)\mathbf{x}_i$ .

P59. (2p) Make same assumptions and use same notation as in Lemma 2.1.

- (a) Show that there exists a fixed point  $\mathbf{x} \in \mathbb{R}^n$  satisfying  $\mathbf{x} = \mathbf{c} + C\mathbf{x}$ .
- (b) Let  $\mathbf{e}_i := \mathbf{x} - \mathbf{x}_i$ . Show that  $\mathbf{e}_i = C^i \mathbf{e}_0$  and further  $\|\mathbf{e}_i\| \leq \|C\|^i \|\mathbf{e}_0\|$ .
- (c) For which  $\alpha \in \mathbb{R}$  does the iteration (2.5) converge, when

$$C = \frac{1}{\alpha} \begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & -1 \\ 0 & -1 & 0 \end{bmatrix} \quad (2.8)$$

The initial guess can be any  $\mathbf{x}_0 \in \mathbb{R}^3$ .

P60. (1p) Let the matrix  $A \in \mathbb{R}^{n \times n}$  be strictly diagonally dominant, this is,

$$\sum_{j=1, j \neq i}^n |a_{ij}| < |a_{ii}| \quad \forall i = 1, \dots, n.$$

In addition, let

$$D = \begin{bmatrix} a_{11} & & & \\ & a_{22} & & \\ & & \ddots & \\ & & & a_{nn} \end{bmatrix} \quad \text{and} \quad N = A - D.$$

Consider solving the linear system  $A\mathbf{x} = \mathbf{b}$  using the fixed point iteration

$$D\mathbf{x}_{i+1} = \mathbf{b} - N\mathbf{x}_i. \quad (2.9)$$

Show that

(a)

$$(D^{-1}N)_{ij} = \begin{cases} a_{ij}a_{ii}^{-1} & i \neq j \\ 0 & i = j \end{cases}$$

(b)

$$\|D^{-1}N\|_{\infty} < 1$$

(c)

$$\|\mathbf{x} - \mathbf{x}_k\|_2 \leq \sqrt{n}\rho^k \|\mathbf{x} - \mathbf{x}_0\|_2 \quad (2.10)$$

where  $\rho = \|D^{-1}N\|_{\infty}$ . Hint: Use Problem P58

P61. (0.5p) Consider the following approximations:

(a)

$$\int_0^1 x^2 dx \approx \frac{1}{N+1} \sum_{i=0}^N \left( \frac{i}{N+1} \right)^2.$$

(b)

$$\frac{1}{1-r} \approx \sum_{i=0}^N r^i, r = \frac{1}{2}.$$

Plot the error in these approximations as a function of  $N$  using commands *plot*, *semilogy* and *loglog*. Which graph is the most informative? What can you say about the relation between error and  $N$  based on the plots?

P62. (1p) Let  $\alpha \in \mathbb{R}$  and

```
N = 10;
A = alpha*eye(N) + diag(-ones(N-1,1),1) + diag(-ones(N-1,1),-1)
b = zeros(N,1); b(end) = 1;
```

(a) For which values of  $\alpha$  is  $A$  diagonally dominant?

(b) Approximately solve  $A\mathbf{x} = \mathbf{b}$  using the Jacobi iteration. Plot the error norm  $\|\mathbf{e}_i\|_2$  as a function of  $i$  for different values of  $\alpha$ . Compare your results to (2.10)

## 2.2 Conjugate gradient method

Conjugate gradient method computes approximate solutions to the linear system

$$A\mathbf{x} = \mathbf{b} \quad \text{for } \mathbf{b} \in \mathbb{R}^n \text{ and s.p.d. } A \in \mathbb{R}^{n \times n}. \quad (2.11)$$

CG is an iterative method that can be understood either as an iteration for finding the minimizer of a quadratic functional related to (2.11) or as a projection method approximately solving (2.11) in a subspace of  $\mathbb{R}^n$ .

In this section, we derive CG as an energy minimisation method. The alternative derivation as an projection method is discussed later. First, we show that solving the linear system  $A\mathbf{x} = \mathbf{b}$  is equivalent to solving a quadratic minimisation problem. Then we discuss line search methods for solving the resulting problem and present the CG method. Throughout this section, we make the following assumptions,

**Assumption 2.1.** Assume that  $A \in \mathbb{R}^{n \times n}$  is s.p.d.,  $\mathbf{b} \in \mathbb{R}^n$ , and define  $J : \mathbb{R}^n \mapsto \mathbb{R}$  as

$$J(\mathbf{u}) := \frac{1}{2} \mathbf{u}^T A \mathbf{u} - \mathbf{b}^T \mathbf{u}.$$

We begin the discussion with the following Lemma.

**Lemma 2.2.** Under Assumptions 2.1, the problems

(P1) Find  $\mathbf{x} \in \mathbb{R}^n$  satisfying  $A\mathbf{x} = \mathbf{b}$ .

and

(P2) Find  $\mathbf{y} \in \mathbb{R}^n$  satisfying  $J(\mathbf{y}) < J(\mathbf{y} + \mathbf{v})$  for any  $\mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{v} \neq 0$ .

are equivalent.

See video proof of  
Lemma 2.3 in Youtube

As problems (P1) and (P2) are equivalent, we can find the global minimizer of  $J$  instead of solving the linear system. Observe that the following proof does not rely on any previous results characterising minimum of  $J$ .

*Part 1 of the proof, solution to (P1) is solution to (P2).* Let  $\mathbf{x} \in \mathbb{R}^n$  satisfy  $A\mathbf{x} = \mathbf{b}$ . Our aim is to prove that  $J(\mathbf{x}) < J(\mathbf{x} + \mathbf{v})$  for all  $\mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{v} \neq 0$ . Expanding  $J(\mathbf{x} + \mathbf{v})$  gives

$$J(\mathbf{x} + \mathbf{v}) = J(\mathbf{x}) + \frac{1}{2} \mathbf{v}^T A \mathbf{v} + \mathbf{x}^T A \mathbf{v} - \mathbf{b}^T \mathbf{v}. \quad (2.12)$$

The last two terms are written as  $\mathbf{x}^T A \mathbf{v} - \mathbf{b}^T \mathbf{v} = (A\mathbf{x} - \mathbf{b})^T \mathbf{v}$ . By assumption,  $A\mathbf{x} - \mathbf{b} = 0$ . Hence, (2.12) becomes

$$J(\mathbf{x} + \mathbf{v}) = J(\mathbf{x}) + \frac{1}{2} \mathbf{v}^T A \mathbf{v}. \quad (2.13)$$

The proof is completed by observing that  $A$  is s.p.d., this is,  $\mathbf{v}^T A \mathbf{v} > 0$  for  $\mathbf{v} \neq 0$ .

□

*Part 2 of proof, solution to (P2) is solution to (P1).* Let  $\mathbf{y} \in \mathbb{R}^n$  satisfy  $J(\mathbf{y}) < J(\mathbf{y} + \mathbf{v})$  for all  $\mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{v} \neq 0$ . Our aim is to show that  $A\mathbf{y} = \mathbf{b}$ . We argue by contradiction, and assume that  $A\mathbf{y} \neq \mathbf{b}$ . Expanding  $J(\mathbf{y} + \mathbf{v})$  gives

$$J(\mathbf{y} + \mathbf{v}) = J(\mathbf{y}) + \frac{1}{2} \mathbf{v}^T A \mathbf{v} + \mathbf{y}^T A \mathbf{v} - \mathbf{b}^T \mathbf{v}.$$

As  $\mathbf{y}$  is the global minimizer of  $J$ ,

$$\frac{1}{2} \mathbf{v}^T A \mathbf{v} + \mathbf{y}^T A \mathbf{v} - \mathbf{b}^T \mathbf{v} > 0 \quad (2.14)$$

for all  $\mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{v} \neq 0$ . We choose  $\mathbf{v} = -t(\mathbf{A}\mathbf{y} - \mathbf{b})$  and define  $p : \mathbb{R} \rightarrow \mathbb{R}$  as

$$p(t) := \frac{1}{2} \mathbf{r}^T \mathbf{A} \mathbf{r} t^2 - \|\mathbf{r}\|_2^2 t \quad \text{where} \quad \mathbf{r} = \mathbf{A}\mathbf{y} - \mathbf{b}.$$

By (2.14) and assumption  $\mathbf{r} \neq 0$ ,  $p(t) > 0$  for every  $t \neq 0$ . As  $A$  is s.p.d. and  $\mathbf{r} \neq 0$ ,  $p(t)$  is an upwards opening parabola with minimum at point  $t_{min} \in \mathbb{R}$  satisfying  $p'(t_{min}) = 0$ . Direct calculation gives

$$p(t_{min}) = -\frac{1}{2} \frac{\|\mathbf{r}\|_2^4}{\mathbf{r}^T \mathbf{A} \mathbf{r}} < 0 \quad \text{if} \quad \|\mathbf{r}\|_2 \neq 0,$$

which is a contradiction with (2.14). Hence,  $\mathbf{r} = 0$  □

Let  $\hat{\mathbf{x}} \in \mathbb{R}^n$  be an approximate solution to (P1). When working with s.p.d. matrices, the error  $\mathbf{x} - \hat{\mathbf{x}}$  is measured in the  $A$ -weighted norm,  $\|\cdot\|_A : \mathbb{R}^n \mapsto \mathbb{R}$  induced by the  $A$ -weighted inner product  $\langle \cdot, \cdot \rangle_A : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ ,

$$\langle \mathbf{v}, \mathbf{w} \rangle_A := \mathbf{v}^T \mathbf{A} \mathbf{w} \quad \text{as} \quad \|\mathbf{v}\|_A = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle_A}.$$

By direct computation, the solution  $\mathbf{x}$  to (P1) satisfies

$$J(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} = -\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} = -\frac{1}{2} \|\mathbf{x}\|_A^2.$$

By (P2),  $J(\mathbf{v}) \geq -\frac{1}{2} \|\mathbf{x}\|_A^2$  for any  $\mathbf{v} \in \mathbb{R}^n$ . The following Lemma relates error in  $\|\mathbf{x} - \hat{\mathbf{x}}\|_A$  to the difference of  $J(\mathbf{x})$  and  $J(\hat{\mathbf{x}})$ .

**Lemma 2.3.** *Under assumptions 2.1*

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_A^2 = 2(J(\hat{\mathbf{x}}) - J(\mathbf{x})).$$

for any  $\hat{\mathbf{x}} \in \mathbb{R}^n$ .

*Proof.* See Problem P64. □

### 2.2.1 Problems

P63. (2p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d.,  $\hat{\mathbf{x}}, \mathbf{b} \in \mathbb{R}^n$ , and  $\mathbf{x} \in \mathbb{R}^n$  satisfy  $A\mathbf{x} = \mathbf{b}$ . Show that

- (a)  $\|A\mathbf{x}\|_2 \leq \|A\|_2^{1/2} \|\mathbf{x}\|_A$
- (b)  $\|\mathbf{x} - \hat{\mathbf{x}}\|_A^2 \leq \|\mathbf{x} - \hat{\mathbf{x}}\|_A \|L^{-1}(A\hat{\mathbf{x}} - \mathbf{b})\|_2$  where  $L \in \mathbb{R}^{n \times n}$  is the Cholesky factor of  $A$ .

(c)

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_A}{\|\mathbf{x}\|_A} \leq \sqrt{\kappa(A)} \frac{\|A\hat{\mathbf{x}} - \mathbf{b}\|_2}{\|\mathbf{b}\|_2}.$$

This result relates the *relative residual* and error in the  $A$  norm. The relative residual can be computed easily and is used to determine accuracy of solution in iterative methods. Hint: Have a look at Problem 65.

P64. (1p) Prove Lemma 2.3.

### 2.2.2 Line search method

In this section, we briefly discuss line search methods and their application to minimisation of  $J$ . Let  $f : \mathbb{R}^n \mapsto \mathbb{R}$  be a given function. A line search method generates a sequence  $\{\mathbf{x}_i\} \subset \mathbb{R}^n$  from a given initial guess  $\mathbf{x}_0 \in \mathbb{R}^n$  and the set of search directions  $\{\mathbf{p}_i\} \subset \mathbb{R}^n$  as follows:

See video on line search methods in Youtube

**Definition 2.1** (line search method). *Let function  $f : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$ , and  $\{\mathbf{p}_i\} \subset \mathbb{R}^n$ . A line search method generates terms of the sequence  $\{\mathbf{x}_i\} \subset \mathbb{R}^n$  in two steps:*

(i) Find  $\alpha_i \in \mathbb{R}$  satisfying<sup>2</sup>

$$f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) < f(\mathbf{x}_i + (\alpha_i + t) \mathbf{p}_i) \quad \text{for all } t \in \mathbb{R}, t \neq 0. \quad (2.15)$$

(ii) Set  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{p}_i$ .

When the function  $f$  behaves sufficiently nicely, the sequence  $\{\mathbf{x}_i\}$  converges to the global minimum  $\mathbf{x}$  of  $f$ , i.e.,  $\mathbf{x} = \lim_{i \rightarrow \infty} \mathbf{x}_i$ . This is the case if  $f$  is the energy functional  $J$ .

In general,  $\alpha_i$  in (2.15) is computed approximately, e.g., using the bisection search. If the line search method is used to find the global minimum of the energy functional  $J$ ,  $\alpha_i$  is computed exactly. Define the function  $p : \mathbb{R} \mapsto \mathbb{R}$  as  $p(t) := J(\mathbf{x}_i + t \mathbf{p}_i)$ , i.e.

$$p(t) = \frac{1}{2} \mathbf{p}_i^T A \mathbf{p}_i t^2 + (\mathbf{A} \mathbf{x}_i - \mathbf{b})^T \mathbf{p}_i t + \frac{1}{2} \mathbf{x}_i^T A \mathbf{x}_i.$$

Observe that  $p(t)$  is second order polynomial with respect to  $t$ . As  $A$  is s.p.d., the coefficient  $\mathbf{p}_i^T A \mathbf{p}_i > 0$  and  $p(t)$  is an upwards opening parabola. Hence, the parameter  $\alpha_i$  is obtained by solving  $p'(\alpha_i) = 0$  as

$$\alpha_i = \frac{\mathbf{p}_i^T \mathbf{r}_i}{\mathbf{p}_i^T A \mathbf{p}_i}, \quad (2.16)$$

---

<sup>2</sup>We assume that  $f$  and  $\mathbf{p}_i$  are chosen such that  $\alpha_i$  satisfying (2.15) exists.

where  $\mathbf{r}_i = \mathbf{b} - A\mathbf{x}_i$  is the residual on step  $i$ .

### 2.2.3 Gradient descend

Using the line search method requires a process generating the search directions  $\mathbf{p}_i$ . In this section, we assume that function  $f$  is differentiable and discuss *gradient descend* that is a line search method using search directions  $\mathbf{p}_i = -(\nabla f)(\mathbf{x}_i)$ . This is,  $\mathbf{p}_i$  points to the direction of greatest decay of function  $f$  at point  $\mathbf{x}_i$ .

The gradient vector  $-\nabla f$  can be computed directly as

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}. \quad (2.17)$$

However, often the easiest way to compute derivatives of vector or matrix valued functions is to use the definition of directional derivative. Recall, that gradient describes the change in the value of function  $f$  for an infinitesimal change in its argument. By definition

$$(\nabla f)(\mathbf{x})^T \mathbf{v} = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} \quad (2.18)$$

Next, we apply gradient descend to find the minimum of  $J$ . Computing the gradient using (2.17) or (2.18) (the latter one is much easier !) gives

$$(\nabla J)(\mathbf{y}) = A\mathbf{y} - \mathbf{b}. \quad (2.19)$$

Hence, the gradient of  $J$  at  $\mathbf{x}_i$  is equivalent to the negative residual  $-\mathbf{r}_i = A\mathbf{x}_i - \mathbf{b}$ . If gradient descend is applied to minimisation of  $J$ ,  $\mathbf{p}_i = -(\nabla J)(\mathbf{x}_i) = \mathbf{r}_i$  and we obtain the iteration

**Definition 2.2** (Gradient descend for minimisation of  $J$ ). *Make Assumptions 2.1. In addition, let  $\mathbf{x}_0 \in \mathbb{R}^n$ . Gradient descend method applied to minimisation of  $J$  generates the terms of the sequence  $\{\mathbf{x}_i\} \subset \mathbb{R}^n$  for  $i \in \mathbb{N}$  in three steps:*

[See video on gradient descend in Youtube](#)

$$\mathbf{p}_i = \mathbf{b} - A\mathbf{x}_i \quad (2.20)$$

$$\alpha_i = \frac{\mathbf{p}_i^T \mathbf{r}_i}{\mathbf{p}_i^T A \mathbf{p}_i} \quad \text{for } \mathbf{r}_i := \mathbf{b} - A\mathbf{x}_i \quad (2.21)$$

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{p}_i \quad (2.22)$$

Note that the search direction  $\mathbf{p}_i$  and the residual  $\mathbf{r}_i$  are the same due to our choice of  $\mathbf{p}_i$  in gradient descend. Iteration given in Definition 2.2 satisfies the relations

$$\mathbf{r}_{i+1} = \mathbf{r}_i - \alpha_i A \mathbf{p}_i \quad \text{and} \quad \mathbf{r}_{i+1} = \mathbf{r}_0 - \sum_{j=0}^i \alpha_j A \mathbf{p}_j, \quad (2.23)$$

Where the latter equation is obtained by repeating the first one.

**Example 2.1.** *Let*

$$A = \begin{bmatrix} 5 & -3 \\ -3 & 5 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}.$$

*Next, we apply gradient descend method given in Definition 2.2 to solve the problem  $A\mathbf{x} = \mathbf{b}$  by minimising the related energy  $J(\mathbf{y})$ . The exact solution  $\mathbf{x} = [1 \ 1]^T$ . The iterates generated by gradient descend starting from the initial guess  $\mathbf{x}_0 = [-1 \ -4]^T$  are depicted in Figure 2.1 along with level sets of  $J(\mathbf{y})$ .*

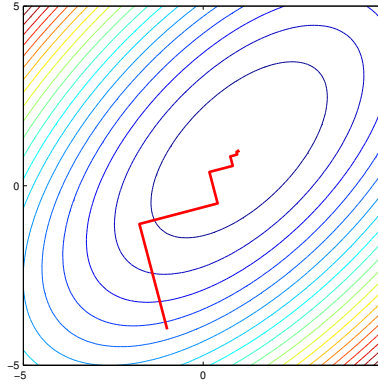


Figure 2.1: Level sets of  $J(\mathbf{y})$  and the ten first iterates of the gradient descend method for solving the problem in Example 2.1.

## 2.2.4 Problems

P65. (1p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d. and  $L$  the Cholesky factor of  $A$ . Prove the following identities

(a)  $\|A^{-1}\|_2 = \|L^{-1}\|_2^2$  Hint: Proceed identical to Problem P32



- (b)  $\|A\mathbf{x}\|_2 \geq \|A^{-1}\|_2^{-1/2} \|\mathbf{x}\|_A$  for any  $\mathbf{x} \in \mathbb{R}^n$ . Hint: start from  $\|\mathbf{x}\|_A$ .
- (c)  $\|A\mathbf{x}\|_A \leq \|A\|_2 \|\mathbf{x}\|_A$  for any  $\mathbf{x} \in \mathbb{R}^n$ .

P66. (2p) Consider the minimisation of  $J$  using the gradient descend method given in Definition 2.2.

- (a) Define the error  $\mathbf{e}_k := \mathbf{x} - \mathbf{x}_k$ . Show that

$$\mathbf{e}_{i+1}^T A \mathbf{p}_i = 0, \quad (2.24)$$

and further

$$\|\mathbf{e}_{i+1}\|_A^2 = \|\mathbf{e}_i\|_A^2 - \alpha_i^2 \mathbf{p}_i^T A \mathbf{p}_i.$$

- (b) Recall that  $\mathbf{r}_i = A\mathbf{e}_i$ ,  $\mathbf{p}_i = A\mathbf{e}_i$ , and show that

$$\alpha_i^2 \mathbf{p}_i^T A \mathbf{p}_i \geq \frac{1}{(\kappa(A))^2} \|\mathbf{e}_i\|_A^2$$

Hint: you need the results given in Problem 65. Start by expanding  $\alpha_i$  and remember the definition of  $\kappa(A)$ .

- (c) Show that  $\|\mathbf{e}_i\|_A^2 \leq (1 - \kappa(A)^{-2})^i \|\mathbf{e}_0\|_A^2$  for  $i \in \mathbb{N}$ .

P67. (2p) Let  $\lambda_1$  and  $\lambda_2$  be given. Generate random s.p.d. matrices in  $\mathbb{R}^{2 \times 2}$  with eigenvalues  $\lambda_1$  and  $\lambda_2$  using the following script.

```
L1 = 1; L2 = 5;
B = rand(2);
[Q,R] = qr(B);
A = Q' * [L1 0; 0 L2] * Q;
```

- (a) Draw a contour plot of the level sets of the energy functional  $J(\mathbf{y})$  for different matrices  $A$  having condition numbers 1.1, 2, and 10. How does the condition number change the function  $J$ ?
- (b) Validate the estimate given in Problem P66 using numerical examples. Use the gradient descend to find minimum of  $J(\mathbf{y})$ . Plot the error in  $A$ -weighted norm using semilogarithmic scale. For reference, draw the predicted rate  $(1 - \kappa^{-2}(A))^i$ . Try several different matrices  $A$  and vectors  $\mathbf{b}$ .

### 2.2.5 $A$ -orthogonal search directions

We call the set of search directions  $\{\mathbf{p}_i\}$  as  $A$ -orthogonal, if

$$\langle \mathbf{p}_i, \mathbf{p}_j \rangle_A = 0, \quad \text{for } i \neq j.$$

In this section, we prove the following Lemma.

**Lemma 2.4.** *Make Assumptions 2.1 and consider the line search method in Definition 2.1 for  $f = J$ . Assume that the search directions  $\{\mathbf{p}_i\} \subset \mathbb{R}^n$  are  $A$ -orthogonal, this is,  $\mathbf{p}_i^T A \mathbf{p}_j = 0$  for  $i \neq j$ . Then*

$$(i) \quad J(\mathbf{x}_i) < J(\mathbf{x}_i + \mathbf{v})$$

and

$$(ii) \quad \|\mathbf{x} - \mathbf{x}_i\|_A < \|\mathbf{x} - \mathbf{x}_i + \mathbf{v}\|_A$$

for any  $\mathbf{v} \in \text{span}\{\mathbf{p}_k\}_{k=0}^{i-1}$ ,  $\mathbf{v} \neq 0$ . Here  $\mathbf{x} \in \mathbb{R}^n$  satisfies  $A\mathbf{x} = \mathbf{b}$ .

Lemma 2.4 is important as it states that the approximate solution  $\mathbf{x}_i$  to (2.11) computed using line search method is *best* in the  $A$ -norm from a the set  $\mathbf{x}_0 + \text{span}\{\mathbf{p}_k\}_{k=0}^{i-1}$ . This result is later used to derive error estimate for solution produced by CG. Before proving it we illustrate the result by an example and give a technical result.

**Example 2.2.** *Consider minimizing the functional  $J$  in the space  $\mathbf{x} + \text{span}\{\mathbf{p}_0, \mathbf{p}_1\}$ , i.e., computing the minimizer to  $F : \mathbb{R}^2 \mapsto \mathbb{R}$ ,*

$$F(s, t) := J(\mathbf{x}_0 + s\mathbf{p}_0 + t\mathbf{p}_1). \quad (2.25)$$

The function  $F$  is quadratic, so it's minimum is attained at point  $(\beta_1, \beta_2) \in \mathbb{R}^2$  satisfying

$$\begin{cases} \frac{\partial F}{\partial s}(\beta_1, \beta_2) = 0 \\ \frac{\partial F}{\partial t}(\beta_1, \beta_2) = 0 \end{cases} \quad (2.26)$$

As  $F$  is quadratic (2.26) is equivalent to

$$\begin{bmatrix} \mathbf{p}_0^T A \mathbf{p}_0 & \mathbf{p}_0^T A \mathbf{p}_1 \\ \mathbf{p}_1^T A \mathbf{p}_0 & \mathbf{p}_1^T A \mathbf{p}_1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \mathbf{p}_0^T \mathbf{b} - \mathbf{p}_0^T A \mathbf{x}_0 \\ \mathbf{p}_1^T \mathbf{b} - \mathbf{p}_1^T A \mathbf{x}_0 \end{bmatrix} \quad (2.27)$$

If search directions  $\mathbf{p}_0$  and  $\mathbf{p}_1$  are  $A$ -orthogonal,  $\mathbf{p}_0^T A \mathbf{p}_1 = 0$  and the linear system in (2.27) reduces to a diagonal one:

$$\begin{bmatrix} \mathbf{p}_0^T A \mathbf{p}_0 & 0 \\ 0 & \mathbf{p}_1^T A \mathbf{p}_1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \mathbf{p}_0^T \mathbf{b} - \mathbf{p}_0^T A \mathbf{x}_0 \\ \mathbf{p}_1^T \mathbf{b} - \mathbf{p}_1^T A \mathbf{x}_0 \end{bmatrix} \quad (2.28)$$

Thus the minimizer is

$$\mathbf{x}_0 + \frac{\mathbf{p}_0^T \mathbf{r}_0}{\mathbf{p}_0^T A \mathbf{p}_0} \mathbf{p}_0 + \frac{\mathbf{p}_1^T \mathbf{r}_0}{\mathbf{p}_1^T A \mathbf{p}_1} \mathbf{p}_1 \quad (2.29)$$

Next, we verify that the line search method using  $A$ -orthogonal search directions  $\mathbf{p}_0$  and  $\mathbf{p}_1$  produces an identical solution. By Definition 2.1,  $\mathbf{r}_1 = \mathbf{r}_0 + \alpha_0 A \mathbf{p}_0$ . Using  $A$ -orthogonality gives

$$\mathbf{p}_1^T \mathbf{r}_1 = \mathbf{p}_1^T \mathbf{r}_0.$$

Hence, (2.29) is

$$\mathbf{x}_0 + \frac{\mathbf{p}_0^T \mathbf{r}_0}{\mathbf{p}_0^T A \mathbf{p}_0} \mathbf{p}_0 + \frac{\mathbf{p}_1^T \mathbf{r}_1}{\mathbf{p}_1^T A \mathbf{p}_1} \mathbf{p}_1$$

which is  $\mathbf{x}_2$  computed by the line search method.

We proceed by giving an orthogonality result.

**Lemma 2.5.** *Make Assumptions 2.1 and consider the line search method in Definition 2.1 for  $f = J$ . In addition, assume that the search directions are  $A$ -orthogonal, this is,  $\mathbf{p}_i^T A \mathbf{p}_j = 0$  for  $i \neq j$ . Then*

[See proof of Lemma 2.5 in Youtube](#)

$$\mathbf{p}_i^T \mathbf{r}_j = 0 \quad \text{for any } i < j.$$

*Proof.* It follows from Definition 2.1 that

$$\mathbf{x}_i = \mathbf{x}_0 + \sum_{k=0}^{i-1} \alpha_k \mathbf{p}_k.$$

,

$$\mathbf{r}_j = \mathbf{r}_i - \sum_{k=i}^{j-1} \alpha_k A \mathbf{p}_k.$$

Taking inner product of the above equation with  $\mathbf{p}_i$  and using  $A$ -orthogonality gives

$$\mathbf{p}_i^T \mathbf{r}_j = \mathbf{p}_i^T \mathbf{r}_i - \alpha_i \mathbf{p}_i^T A \mathbf{p}_i.$$

Using the relation  $\alpha_i = \frac{\mathbf{p}_i^T \mathbf{r}_i}{\mathbf{p}_i^T A \mathbf{p}_i}$  completes the proof.  $\square$

*Proof of Lemma 2.4.* By Lemma 2.3,  $\|\mathbf{x} - \mathbf{x}_i\|_A^2 = 2(J(\mathbf{x}_i) - J(\mathbf{x}))$ . Hence, the statements (i) and (ii) are equivalent. Next, we give a proof for the statement (i). Denote  $\mathbf{e}_i := \mathbf{x} - \mathbf{x}_i$ . Then

[See proof of Lemma 2.4 in Youtube](#)

$$\|\mathbf{e}_i + \mathbf{v}\|_A^2 = \|\mathbf{e}_i\|_A^2 + 2\langle \mathbf{e}_i, \mathbf{v} \rangle_A + \|\mathbf{v}\|_A^2. \quad (2.30)$$

Observe that  $A\mathbf{e}_i = \mathbf{r}_i$ . By assumption,  $\mathbf{v} = \sum_{k=0}^{i-1} \beta_k \mathbf{p}_k$  for some  $\beta_k$ . Hence,

$$\langle \mathbf{e}_i, \mathbf{v} \rangle_A = \sum_{k=0}^{i-1} \beta_k \mathbf{r}_i^T \mathbf{p}_k.$$

By Lemma 2.6,  $\mathbf{r}_i^T \mathbf{p}_k = 0$  for  $k \in \{0, \dots, i-1\}$ . Thus, (2.30) becomes

$$\|\mathbf{e}_i + \mathbf{v}\|_A^2 = \|\mathbf{e}_i\|_A^2 + \|\mathbf{v}\|_A^2. \quad (2.31)$$

Observing that  $\|\mathbf{v}\|_A^2 > 0$  for  $\mathbf{v} \neq 0$  completes the proof.  $\square$

### 2.2.6 Conjugate gradient method

See introduction on CG  
in Youtube

In this section, we define the Conjugate Gradient method by modifying the gradient descend method given in Definition 2.2. Namely, CG uses the  $A$ -orthogonal search directions  $\mathbf{p}_i$  that are computed from  $\mathbf{r}_i$  using the Gram-Schmidt process as

$$\mathbf{p}_i = \begin{cases} \mathbf{r}_i & i = 0 \\ \mathbf{r}_i - \sum_{k=0}^{i-1} \frac{\mathbf{r}_i^T A \mathbf{p}_k}{\mathbf{p}_k^T A \mathbf{p}_k} \mathbf{p}_k & \text{otherwise} \end{cases} \quad (2.32)$$

We show that almost every term in the above sum has value zero, and the search direction  $\mathbf{p}_i$  can be computed from  $\mathbf{p}_{i-1}$ . Thus, the CG iteration only stores vectors  $\mathbf{p}_{i-1}$ ,  $\mathbf{p}_i$ , and  $\mathbf{x}_i$ . As most terms in the sum are zero, finding a search direction  $\mathbf{p}_i$  is computationally inexpensive. Memory use, small computational cost, and  $A$ -optimality property stated in Lemma 2.4 make CG a popular method for solving linear systems with s.p.d. coefficient matrices.

Let  $\beta_{i,k} := \frac{\mathbf{r}_i^T A \mathbf{p}_k}{\mathbf{p}_k^T A \mathbf{p}_k}$  for  $k \in \{0, \dots, i-1\}$ . Equation (2.32) becomes

$$\mathbf{p}_i = \mathbf{r}_i - \sum_{k=0}^{i-1} \beta_{i,k} \mathbf{p}_k. \quad (2.33)$$

We proceed to show that  $\beta_{i,k} = 0$  for  $k \in \{0, \dots, i-2\}$ , this is,  $\mathbf{r}_i^T A \mathbf{p}_k = 0$  for  $k \in \{0, \dots, i-2\}$ .

**Lemma 2.6.** *Make Assumptions 2.1. Consider the line search method give in 2.1. Let  $f = J$  and define the search directions  $\{\mathbf{p}_i\}$  as in (2.32). Then there holds that*

$$(i) \quad \mathbf{r}_i^T \mathbf{r}_j = 0 \quad \text{for } i, j \in \{0, \dots, n\}, i \neq j$$

See video proof of  
Lemma 2.6 in Youtube

and

$$(ii) \mathbf{r}_i^T A \mathbf{p}_k = 0, \quad \text{for } i \in \mathbb{N}, k = \{0, \dots, i-2\}.$$

*Proof.* We begin by proving identity (i). Without loss of generality, we can assume that  $j > i$ . Taking inner product of (2.33) with  $\mathbf{r}_j$  gives

$$\mathbf{r}_j^T \mathbf{p}_i = \mathbf{r}_j^T \mathbf{r}_i - \sum_k^{i-1} \beta_{i,k} \mathbf{r}_j^T \mathbf{p}_k$$

Applying Lemma 2.5 proves (i). We  $f = J$ , using Defintion 2.1 gives

$$\mathbf{r}_{k+1} - \mathbf{r}_k = -\alpha_k A \mathbf{p}_k.$$

so that

$$\mathbf{r}_{i+1}^T A \mathbf{p}_k = -\alpha_k^{-1} (\mathbf{r}_{i+1}^T \mathbf{r}_{k+1} - \mathbf{r}_{i+1}^T \mathbf{r}_k)$$

By the orthogonality of residual vectors in (i),

$$\mathbf{r}_{i+1}^T A \mathbf{p}_k = \begin{cases} 0 & k < i \\ -\alpha_i^{-1} \mathbf{r}_{i+1}^T \mathbf{r}_{i+1} & k = i \end{cases}$$

Which completes the proof.  $\square$

The orthogonality of the residuals leads to the elimination of the dependence on  $\mathbf{p}_i$  on all but the previous step. Hence, the equation (2.32) simplifies to

$$\mathbf{p}_i = \mathbf{r}_i + \alpha_{i-1}^{-1} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_{i-1}^T A \mathbf{p}_{i-1}} \mathbf{p}_{i-1}.$$

This expression is further simplified to reduce the computational cost of the algorithm. By the definition of  $\alpha_{i-1}$ ,

$$\alpha_{i-1}^{-1} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_{i-1}^T A \mathbf{p}_{i-1}} \mathbf{p}_{i-1} = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_{i-1}^T \mathbf{p}_{i-1}} \mathbf{p}_{i-1}.$$

By (2.32) and proof of Lemma 2.5,

$$\mathbf{r}_{i-1}^T \mathbf{p}_{i-1} = \mathbf{r}_{i-1}^T \mathbf{r}_{i-1},$$

which gives the CG method

$$\alpha_i = \frac{\mathbf{p}_i^T \mathbf{r}_i}{\mathbf{p}_i^T A \mathbf{p}_i} \quad (2.34)$$

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{p}_i \quad (2.35)$$

$$\mathbf{r}_{i+1} = \mathbf{r}_i - \alpha_i A \mathbf{p}_i \quad (2.36)$$

$$\mathbf{p}_{i+1} = \mathbf{r}_{i+1} + \frac{\mathbf{r}_{i+1}^T \mathbf{r}_{i+1}}{\mathbf{r}_i^T \mathbf{r}_i} \mathbf{p}_i \quad (2.37)$$

In Matlab, a simple implementation of the above CG method looks like

```
% a simple implementation of cg.
%
function x = my_cg(A,b,x0,N,tol)

x = x0;
r = b - A*x0;
p = r;

for i=1:N

    alpha = (p'*r)/(p'*A*p);

    x = x + alpha*p;

    rold = r;

    r = r - alpha*A*p;
    p = r + r'*r/(rold'*rold)*p;

    if(norm(r) < tol)
        break;
    end

end
```

### 2.2.7 Problems

- P68. (2p) Implement the CG as line search method given in Definition 2.1 with  $f = J$ . On each step generate the search directions using (2.33). Store coefficients  $\beta_{i,k}$  and validate numerically the orthogonality result given in Lemma 2.6.

## 2.3 Orthogonal Projection Matrices

In this section, we discuss orthogonal projection matrices that are later used to interpret conjugate gradient iteration as a subspace method. Before proceeding, recall the one-to-one correspondence between positive definite matrices and inner products.

[See video on orthogonal projection matrices in Youtube](#)

**Lemma 2.7.** *For any inner product  $\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , there exists a s.p.d. matrix  $A \in \mathbb{R}^{n \times n}$  such that*

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^T A \mathbf{x}. \quad (2.38)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . On the other hand, the formula (2.38) defines an inner product for any symmetric positive definite  $A \in \mathbb{R}^{n \times n}$ .

Let  $\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  be an inner product and  $A \in \mathbb{R}^{n \times n}$  the s.p.d. matrix corresponding to  $\langle \cdot, \cdot \rangle$ , i.e.,  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^T A \mathbf{x}$ . In addition, let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^n$ . The orthogonal complement of a subspace  $\mathcal{V}$  in  $\langle \cdot, \cdot \rangle$  inner product is defined as follows:

**Definition 2.3.** *Let  $\mathcal{V}$  be subspace of  $\mathbb{R}^n$  and  $\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  some innerproduct. The orthogonal complement of  $\mathcal{V}$  is the subspace*

$$\mathcal{V}^\perp := \{\mathbf{w} \in \mathbb{R}^n \mid \langle \mathbf{w}, \mathbf{v} \rangle = 0 \text{ for all } \mathbf{v} \in \mathcal{V}\}. \quad (2.39)$$

Orthogonal projections are related to the splitting of a given vector  $\mathbf{x} \in \mathbb{R}^n$  as

$$\mathbf{x} = \mathbf{v} + \mathbf{v}^\perp, \quad (2.40)$$

where  $\mathbf{v} \in \mathcal{V}$  and  $\mathbf{v}^\perp \in \mathcal{V}^\perp$ . As the spaces  $\mathcal{V}$  and  $\mathcal{V}^\perp$  are orthogonal, the splitting given in (2.40) is unique and well defined. The projection matrix  $P_{\mathcal{V}} \in \mathbb{R}^{n \times n}$  is defined via

$$P_{\mathcal{V}} \mathbf{x} = \mathbf{v}$$

where  $\mathbf{v} \in \mathcal{V}$  is the first component in the RHS of (2.40). An immediate consequence of this definition is that  $P_{\mathcal{V}}^2 = P_{\mathcal{V}}$ . In other words, the projection matrix  $P_{\mathcal{V}}$  is an identity map in  $R(P_{\mathcal{V}})$ . This is a fundamental property that can be taken as the definition of a projection (matrix).

**Definition 2.4.** *A matrix  $P \in \mathbb{R}^{n \times n}$  is a projection if it satisfies*

$$P^2 = P. \quad (2.41)$$

Next, we find explicit representation for  $P_{\mathcal{V}}$ . Let  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  be a basis for  $\mathcal{V}$ . That is, the set of vectors is linearly independent and

$$\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}.$$

We stack the basis as columns of the matrix  $V \in \mathbb{R}^{n \times k}$  as:

$$V = [\mathbf{v}_1, \dots, \mathbf{v}_k].$$

Note once again that  $\mathcal{V} = R(V) = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ . By definition (2.39), the columns of  $V$  are orthogonal to  $\mathcal{V}^\perp$  in the inner product corresponding to  $A \in \mathbb{R}^{n \times n}$ , meaning that

$$V^T A \mathbf{v}^\perp = 0 \quad \text{for any } \mathbf{v}^\perp \in \mathcal{V}^\perp. \quad (2.42)$$

On the other hand, the splitting (2.40) is

$$\mathbf{x} = V \mathbf{z}_V + \mathbf{v}^\perp \quad (2.43)$$

for some yet-to-be-defined  $\mathbf{z}_V \in \mathbb{R}^k$ . Multiplying from the left by  $V^T A$  and using the orthogonality property (2.42), we get

$$V^T A \mathbf{x} = V^T A V \mathbf{z}_V.$$

Because the columns of  $V \in \mathbb{R}^{n \times k}$  are linearly independent and  $A$  is positive definite,  $N(V^T A V) = N(V)$  is trivial,  $V^T A V \in \mathbb{R}^{k \times k}$  is invertible, and  $\mathbf{z}_V$  can be solved as

$$\mathbf{z}_V = (V^T A V)^{-1} V^T A \mathbf{x}.$$

Thus,

$$P_{\mathcal{V}} = V(V^T A V)^{-1} V^T A \in \mathbb{R}^{n \times n}. \quad (2.44)$$

Observe that computing the orthogonal projection does not require a basis for  $\mathcal{V}^\perp$ . Orthogonal projections satisfy some useful identities. First of all, since inversion and transposition commute, we have

$$P_{\mathcal{V}}^T A = A P_{\mathcal{V}},$$

and hence also,

$$P_{\mathcal{V}}^T A (I - P_{\mathcal{V}}) = A P_{\mathcal{V}} (I - P_{\mathcal{V}}) = A (P_{\mathcal{V}} - P_{\mathcal{V}}^2) = A (P_{\mathcal{V}} - P_{\mathcal{V}}) = 0. \quad (2.45)$$

After transposition, this yields

$$(I - P_{\mathcal{V}})^T A P_{\mathcal{V}} = 0. \quad (2.46)$$

See video on best approximation property of orthogonal projection matrices in Youtube

We are interest in projection matrices because of the following *best approximation property*.



**Lemma 2.8.** Let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^n$ ,  $\langle \cdot, \cdot \rangle$  an inner product,  $\|\cdot\| := \langle \cdot, \cdot \rangle^{1/2}$ , and  $P_{\mathcal{V}} \in \mathbb{R}^{n \times n}$  the  $\langle \cdot, \cdot \rangle$ -orthogonal projection to  $\mathcal{V}$ . Then

$$\|\mathbf{x} - P_{\mathcal{V}}\mathbf{x}\| < \|\mathbf{x} - P_{\mathcal{V}}\mathbf{x} + \mathbf{v}\| \quad \text{for any } \mathbf{v} \in \mathcal{V}, \mathbf{v} \neq 0.$$

*Proof.* By direct computation,

$$\|\mathbf{x} - P_{\mathcal{V}}\mathbf{x} + \mathbf{v}\|^2 = \|(I - P_{\mathcal{V}})\mathbf{x} + \mathbf{v}\|^2 = \|(I - P_{\mathcal{V}})\mathbf{x}\|^2 + 2\langle (I - P_{\mathcal{V}})\mathbf{x}, \mathbf{v} \rangle + \|\mathbf{v}\|^2.$$

Let the s.p.d. matrix  $A \in \mathbb{R}^{n \times n}$  correspond to the inner product  $\langle \cdot, \cdot \rangle$ . As  $\mathbf{v} \in \mathcal{V}$ ,  $\mathbf{v} = P_{\mathcal{V}}\mathbf{v}$ , and

$$\langle (I - P_{\mathcal{V}})\mathbf{x}, \mathbf{v} \rangle = \mathbf{x}^T (I - P_{\mathcal{V}})^T A \mathbf{v} = \mathbf{x}^T (I - P_{\mathcal{V}})^T A P_{\mathcal{V}} \mathbf{v} = 0$$

by (2.46). Hence,  $\|\mathbf{x} - P_{\mathcal{V}}\mathbf{x} + \mathbf{v}\|^2 = \|(I - P_{\mathcal{V}})\mathbf{x}\|^2 + \|\mathbf{v}\|^2$ . □

### 2.3.1 Problems

P69. (0.5p) Show that  $\mathcal{V}^{\perp}$  in Definition 2.42 is a subspace.

## 2.4 CG as a subspace method

Next, we interpret the conjugate gradient iteration for solution of  $A\mathbf{x} = \mathbf{b}$  as a subspace method. We show that the iterate  $\mathbf{x}_i$  is the  $A$ -orthogonal projection of the exact solution  $\mathbf{x}$  to a *Krylov subspace*. Further, the iterate  $\mathbf{x}_i$  can be computed without knowledge of the exact solution  $\mathbf{x}$ .

Recall that CG iteration is a line search method for minimisation of the quadratic functional  $J(\mathbf{u}) = \frac{1}{2}\mathbf{u}^T A \mathbf{u} - \mathbf{u}^T \mathbf{b}$ , see Definition 2.1. It generates a sequence of solutions  $\{\mathbf{x}_i\}$  from initial guess  $\mathbf{x}_0$ . CG uses  $A$ -orthogonal search directions  $\{\mathbf{p}_i\}$  that are constructed from the residuals  $\mathbf{r}_i = \mathbf{b} - A\mathbf{x}_i$  using the (slightly modified) Gram-Schmidt process. By (2.32) and Lemma 2.6, the search directions generated by CG satisfy

$$\mathbf{p}_i = \mathbf{r}_i - \beta_i \mathbf{p}_{i-1} \quad \text{for some } \beta_i \in \mathbb{R}.$$

Make a standing assumption that the initial guess  $\mathbf{x}_0 = 0$ . (For treatment of nonzero initial guess, see Problem P71). Under this assumption search directions and iterates generated by CG are elements of the Krylov subspaces.

P70. (2p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d. and  $\mathbf{b} \in \mathbb{R}^n$ . In addition let  $\{\mathbf{x}_i\}$  and  $\{\mathbf{p}_i\}$  be the iterates and search directions of conjugate gradient method with

See video on interpreting CG as a subspace method in Youtube

initial guess  $\mathbf{x}_0 = 0$ . Define the family of Krylov subspaces  $\{\mathcal{K}_i(A, \mathbf{b})\}$  associated to  $A$  and  $\mathbf{b}$  as

$$\mathcal{K}_i(A, \mathbf{b}) = \text{span}\{\mathbf{b}, A\mathbf{b}, A^2\mathbf{b}, \dots, A^{i-1}\mathbf{b}\} \quad \text{for } i \in \{1, \dots, n\}. \quad (2.47)$$

Show that  $\mathbf{p}_{i-1}$  and  $\mathbf{x}_i \in \mathcal{K}_i(A, \mathbf{b})$  for  $i \in \{1, \dots, n\}$ . Hint: formulate an induction proof.

Lemma 2.4 states that iterates  $\{\mathbf{x}_i\}$  generated by any line search method using  $A$ -orthogonal search directions to minimise  $J$  satisfy

$$\|\mathbf{x} - \mathbf{x}_i\|_A < \|\mathbf{x} - \mathbf{x}_i + \mathbf{v}\|_A$$

for all  $\mathbf{v} \in \text{span}\{\mathbf{p}_0, \dots, \mathbf{p}_{i-1}\}, \mathbf{v} \neq 0$ , and  $i \in \mathbb{N}$ . This is, each  $\mathbf{x}_i$  is the best approximation of  $\mathbf{x}$  from subspace spanned by search directions  $\{\mathbf{p}_0, \dots, \mathbf{p}_{i-1}\}$ . Further, by Problem 70, iterate  $\mathbf{x}_i$  computed using conjugate gradient method is the best approximation of  $\mathbf{x}$  from the Krylov subspace  $\mathcal{K}_i(A, b)$ . The best approximation of  $\mathbf{x}$  from  $\mathcal{K}_i(A, b)$  is characterised in Lemma 2.8 as

$$\mathbf{x}_i = P_{\mathcal{K}_i(A, b)} \mathbf{x},$$

where  $P_{\mathcal{K}_i(A, b)} \in \mathbb{R}^{n \times n}$  is the  $A$ -orthogonal projection to  $\mathcal{K}_i(A, b)$ . Let the columns of matrix  $V_i \in \mathbb{R}^{n \times i}$  be a basis of the space  $\mathcal{K}_i(A, b)$ . By (2.44),

$$P_{\mathcal{K}_i(A, b)} := V_i (V_i^T A V_i)^{-1} V_i^T A.$$

Using  $A\mathbf{x} = \mathbf{b}$  yields

$$\mathbf{x}_i = V_i (V_i^T A V_i)^{-1} V_i^T \mathbf{b}. \quad (2.48)$$

Observe that (2.48) can be solved *without knowledge on  $\mathbf{x}$*  in two steps:

$$\begin{aligned} \text{Find } \mathbf{q}_i \in \mathbb{R}^n \text{ satisfying } & V_i^T A V_i \mathbf{q}_i = V_i^T \mathbf{b} \\ \text{Set } \mathbf{x}_i & = V_i \mathbf{q}_i. \end{aligned} \quad (2.49)$$

Rest of this section is organised as follows. We begin by giving a numerically stable method for computing a basis of  $\mathcal{K}_i(A, b)$ . Next, we use the best approximation property to give error estimate for  $\|\mathbf{x} - \mathbf{x}_i\|_A$ . We end this section by discussing how convergence of CG is improved using preconditioning.

### 2.4.1 Problems

P71. (1p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d. and  $\mathbf{x}_0, \mathbf{b} \in \mathbb{R}^n$ . Show that the following processes yield the same solutions:

- (1) Solve  $A\mathbf{c} = \mathbf{b} - A\mathbf{x}_0$  using CG, starting from initial guess  $\mathbf{c}_0 = 0$ , and read solution as  $\mathbf{x}_i = \mathbf{c}_i + \mathbf{x}_0$ .
- (2) Solve  $A\mathbf{x} = \mathbf{b}$  using CG starting from initial guess  $\mathbf{x}_0$ .

P72. (2p) Let  $S$  be a subspace of  $\mathbb{R}^n$ ,  $\{\mathbf{q}_1, \dots, \mathbf{q}_k\} \subset \mathbb{R}^n$  be a basis of  $S$ , and  $Q = [\mathbf{q}_1 \ \dots \ \mathbf{q}_k]$ . In addition, let  $A \in \mathbb{R}^{n \times n}$  be s.p.d.,  $\mathbf{b} \in \mathbb{R}^n$ , and  $\mathbf{x} \in \mathbb{R}^n$  satisfy  $A\mathbf{x} = \mathbf{b}$ . Consider the problem: find  $\hat{\mathbf{x}} \in S$  that minimises

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_A. \quad (2.50)$$

- (a) A standard least squares problem in 2-norm is of the form  $\|G\mathbf{a} - \mathbf{b}\|_2$  with  $G \in \mathbb{R}^{n \times k}$ ,  $\mathbf{a} \in \mathbb{R}^k$  and  $\mathbf{b} \in \mathbb{R}^n$ . Use  $Q$  and Cholesky decomposition to reduce (2.50) to a standard least squares problem in 2-norm.
- (b) The standard least squares problem in 2-norm can be minimized by solving the *normal equations*  $G^T G\mathbf{a} = G^T \mathbf{b}$ . Solve the minimization problem (2.50).
- (c) Show that the solution  $\hat{\mathbf{x}}$  is equivalent to the  $A$ -orthogonal projection of  $\mathbf{x}$  to the subspace  $S$ .

### 2.4.2 Krylov Subspace

The family of Krylov subspaces  $\{\mathcal{K}_i(A, b)\}$  associated to  $A \in \mathbb{R}^{n \times n}$  and  $\mathbf{b}$  is defined as

$$\mathcal{K}_i(A, b) = \text{span}\{\mathbf{b}, A\mathbf{b}, \dots, A^{i-1}\mathbf{b}\}.$$

Next, we give a method for constructing a basis for the space  $\mathcal{K}_i(A, b)$ . The basis is used, e.g., to compute the best approximation of  $\mathbf{x}$  from  $\mathcal{K}_i(A, b)$ .

Finding a basis for  $\mathcal{K}_i(A, b)$  requires care. The trivial basis

$$\{\mathbf{b}, A\mathbf{b}, A^2\mathbf{b}, \dots, A^{i-1}\mathbf{b}\}$$

is not numerically stable and cannot be used. This is demonstrated in the following example.

[See video on computing a basis for the Krylov subspace in Youtube](#)

**Example 2.3.** *Let*

$$A = Q \begin{bmatrix} 10 & & \\ & 1 & \\ & & 10^{-1} \end{bmatrix} Q^T.$$

*The vector  $A^{10}\mathbf{b}$  is then*

$$A^{10}\mathbf{b} = Q \begin{bmatrix} 10^{10} & & \\ & 1 & \\ & & 10^{-10} \end{bmatrix} Q^T \mathbf{b}.$$

*When  $A^i\mathbf{b}$  is computed by repeatedly multiplying  $\mathbf{b}$  by  $A$ , information on the lowest eigenvalue is lost due to rounding-off errors and the resulting vector points to the direction of the eigenvector corresponding to largest eigenvalue.*

Due to stability issues, basis for  $\mathcal{K}_i(A, b)$  is computed using the *Arnoldi iteration* that is based on the Gram-Schmidt process.

**Definition 2.5** (Arnoldi iteration). *Let  $A \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ . Arnoldi iteration computes basis  $\{\mathbf{q}_1, \dots, \mathbf{q}_i\}$  for  $\mathcal{K}(A, b)$  in three steps*

$$\begin{aligned} \text{Step 1} \quad \hat{\mathbf{q}}_{i+1} &= \begin{cases} \mathbf{b} & \text{if } i = 0 \\ A\mathbf{q}_i & \text{otherwise} \end{cases} \\ \text{Step 2} \quad \tilde{\mathbf{q}}_{i+1} &= \hat{\mathbf{q}}_{i+1} - \sum_{k=1}^i \hat{\mathbf{q}}_{i+1}^T \mathbf{q}_k \mathbf{q}_k \\ \text{Step 3} \quad \mathbf{q}_{i+1} &= \|\tilde{\mathbf{q}}_{i+1}\|_2^{-1} \tilde{\mathbf{q}}_{i+1}. \end{aligned}$$

Intuitively speaking, the numerical stability of Arnoldi iteration is due to orthogonalisation step that eliminates directions  $\mathbf{q}_1, \dots, \mathbf{q}_{i-1}$  from  $\mathbf{q}_i$ . This prevents  $A\mathbf{q}_i$  from turning to the direction of the largest eigenvector of  $A$  as happens in Example 2.3. An example implementation in Matlab is given below.

```
function [Q,R] = my_arnoldi(A,b,N)

Q = [];
q = b;
for i=1:N

    for k=1:size(Q,2)
        R(k,i) = q'*Q(:,k);
        q = q - R(k,i)*Q(:,k);
    end
end
```

```

end
R(i,i) = norm(q);
Q(:,i) = q/R(i,i);

q = A*Q(:,i);

end

```

Basis computed by Arnoldi iteration has the following properties.

P73. (1p) Let  $A \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ , and generate  $\{\mathbf{q}_i\} \subset \mathbb{R}^n$  using the Arnoldi Iteration in Definition 2.5. Show that

- (i)  $\mathcal{K}_i(A, \mathbf{b}) = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_i\}$ . Hint: use induction.
- (ii)  $\mathbf{q}_i^T A \mathbf{y} = 0 \quad \forall \mathbf{y} \in \mathcal{K}_{i-2}(A, \mathbf{b})$ .

### 2.4.3 Problems

P74. (2p) Generate a random symmetric, positive definite  $3 \times 3$  matrix using the snippet:

```
[Q,R] = qr(rand(3)); A = Q*diag([10 0.5 0.1])*Q'.
```

- (a) Without using Matlab, what are the eigenvalues and vectors of  $A$ ?
- (b) Compute  $A^i \mathbf{b}$  as  $\mathbf{z}_0 = \mathbf{b}$ ,  $\mathbf{z}_i = A \mathbf{z}_{i-1}$ , with  $\mathbf{b}$  being a random vector in  $\mathbb{R}^3$ . On each step, find  $\alpha_{i1}, \alpha_{i2}, \alpha_{i3}$  s.t.

$$\mathbf{z}_i = \alpha_{i1} \mathbf{q}_1 + \alpha_{i2} \mathbf{q}_2 + \alpha_{i3} \mathbf{q}_3.$$

where  $\mathbf{q}_i$  are the eigenvectors of  $A$ . Explain mathematically how  $\alpha_{ij}$  should behave when  $i$  grows.

- (c) Plot  $|\alpha_{ij}|$  using a semilogarithmic plot. When does the method fail? Give a hypothesis for the reason behind this failure.

### 2.4.4 Error estimate

Next, we estimate the error  $\|\mathbf{x} - \mathbf{x}_i\|_A$ . Lemma 2.4 states that the iterate  $\mathbf{x}_i$  computed by CG is the best possible approximation of  $\mathbf{x}$  from the subspace  $\mathcal{K}_i(A, \mathbf{b})$ . We begin by formulating this result in a different form. First, we recall some properties of symmetric matrices.

Any symmetric real matrix  $A \in \mathbb{R}^{n \times n}$  is unitary diagonalisable, this is,  $A = Q^T \Lambda Q$ , where  $Q \in \mathbb{R}^{n \times n}$  is a unitary matrix (i.e  $Q^T Q = I$ ) and  $\Lambda \in \mathbb{R}^{n \times n}$  is a diagonal matrix. [See video on deriving error estimate for CG in Youtube](#)

$\mathbb{R}^{n \times n}$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  is a diagonal matrix, and  $\{\lambda_k\} \subset \mathbb{R}$  eigenvalues of matrix  $A$  on the diagonal. Let  $p : \mathbb{R} \mapsto \mathbb{R}$  be a degree  $n$ -polynomial, i.e.,

$$p(t) = \sum_{k=0}^n \alpha_k t^k$$

where the coefficients  $\alpha_k \in \mathbb{R}$  for  $k \in \{0, \dots, n\}$ . Recall that  $p(A) \in \mathbb{R}^{n \times n}$  is defined as

$$p(A) = \sum_{k=0}^n \alpha_k A^k.$$

We make use the following result.

P75. (2p) Let  $A \in \mathbb{R}^{n \times n}$  be s.p.d.,  $q(A)$  be any polynomial of  $A$ , and  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  the orthonormal eigenbasis of  $A$ , i.e.

$$A\mathbf{v}_i = \lambda_i \mathbf{v}_i \quad \text{and} \quad \mathbf{v}_i^T \mathbf{v}_j = \delta_{ij}.$$

Show that

$$(a) \quad \|\mathbf{b}\|_A^2 = \sum_{i=1}^n \lambda_i (\mathbf{b}^T \mathbf{v}_i)^2 \quad \text{for any } \mathbf{b} \in \mathbb{R}^n.$$

$$(b) \quad \|q(A)\mathbf{b}\|_A^2 = \sum_{i=1}^n \lambda_i q(\lambda_i)^2 (\mathbf{b}^T \mathbf{v}_i)^2 \quad \text{for any } \mathbf{b} \in \mathbb{R}^n.$$

**Lemma 2.9.** Let  $A \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ , and  $\mathbf{x} \in \mathbb{R}^n$  satisfy  $A\mathbf{x} = \mathbf{b}$ . In addition, let  $\{\mathbf{x}_i\}$  be the iterates computed by CG and  $\mathbf{e}_i = \mathbf{x} - \mathbf{x}_i$  for  $i \in \{0, 1, \dots\}$ . Then there holds that

$$\|\mathbf{e}_i\|_A \leq \min_{\substack{q \in \mathcal{P}_i \\ q(0)=1}} \max_k |q(\lambda_k)| \|\mathbf{e}_0\|_A.$$

where  $\mathcal{P}_i$  is the space of degree  $i$  polynomials.

This estimate is important as it relates errors  $\mathbf{x} - \mathbf{x}_i$  to approximation properties of polynomials in the maximum norm.

*Proof.* As  $\mathbf{x}_i \in \mathcal{K}_i(A, \mathbf{b})$  it can be written as  $p(A)\mathbf{b}$  where  $p \in \mathcal{P}^{i-1}$ . Hence,

$$\|\mathbf{x} - \mathbf{x}_i\|_A = \|\mathbf{e}_0 - p(A)\mathbf{b}\|_A.$$

Using relation  $\mathbf{b} = A\mathbf{x} = A\mathbf{e}_0$  gives

$$\|\mathbf{x} - \mathbf{x}_i\|_A = \|\mathbf{e}_0 - p(A)A\mathbf{e}_0\|_A = \|(I - p(A)A)\mathbf{e}_0\|_A.$$

Application of Lemma 2.4 gives

$$\min_{\substack{q \in \mathcal{P}_i \\ q(0)=1}} \|q(A)\mathbf{e}_0\|_A = \|(I - p(A)A)\mathbf{e}_0\|_A$$

Hence, we obtain

$$\|\mathbf{e}_i\|_A = \min_{\substack{q \in \mathcal{P}_i \\ q(0)=1}} \|q(A)\mathbf{e}_0\|_A.$$

Using identities given in Problem P75 yields

$$\|q(A)\mathbf{e}_0\|_A^2 = \sum_{i=1}^n \lambda_i q(\lambda_i)^2 (\mathbf{e}_0^T \mathbf{v}_i)^2,$$

Taking the maximum of  $q(\lambda_i)^2$  as common factor and using Problem P75 gives

$$\|q(A)\mathbf{e}_0\|_A^2 \leq \max_i q(\lambda_i)^2 \sum_{i=1}^n \lambda_i (\mathbf{e}_0^T \mathbf{v}_i)^2 = \max_i q(\lambda_i)^2 \|\mathbf{e}_0\|_A^2.$$

□

The minimisation problem

$$\min_{\substack{q \in \mathcal{P}_i \\ q(0)=1}} \max_{t \in [\lambda_{min}, \lambda_{max}]} |q(t)|$$

can be solved analytically using Chebyshev polynomials. This leads to the convergence estimate

$$\min_{\substack{q \in \mathcal{P}_i \\ q(0)=1}} \max_k |q(\lambda_k)| \leq \left( \frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right)^i,$$

where  $\kappa_2(A)$  is the condition number of  $A$  in the 2-norm, i.e.,

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\lambda_{max}}{\lambda_{min}}$$

### 2.4.5 Problems

P76. (2p) Generate a random symmetric, positive definite  $n \times n$  matrix using the snippet:

```
n=20; [Q,R] = qr(rand(n)); A = Q*diag(linspace(1,L,n))*Q' .
```

- (a) Let  $\mathbf{b} \in \mathbb{R}^{20}$  be a random vector. Compute a basis for the Krylov subspace  $K_i(A, \mathbf{b})$  using function `my_arnoldi.m`. Compute the best solution  $\mathbf{x}_i$  in the  $A$ -norm to

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

from  $K_i(A, \mathbf{b})$

- (b) Plot the error  $\|\mathbf{x} - \mathbf{x}_i\|_A$  for  $L = 1, 10, 100, 1000$  and  $i = 5, \dots, 15$ . How does the error depend on  $L$ ? How about  $i$ ?

### 2.4.6 Preconditioning

If the condition number  $\kappa_2(A)$  is large, computing a sufficiently accurate approximation  $\mathbf{x}_i$  using CG can be too time consuming. For example, the condition number of the finite difference matrix behaves as  $h^{-2}$ .

The convergence of iterative solution methods is improved by using a preconditioner  $B \in \mathbb{R}^{n \times n}$ . The system  $A\mathbf{x} = \mathbf{b}$  is transformed as

$$\text{Right preconditioner } AB\mathbf{y} = \mathbf{b} \quad \mathbf{x} = B\mathbf{y}$$

$$\text{Left preconditioner } BA\mathbf{x} = B\mathbf{b}$$

$$\text{Split preconditioner, } B \text{ is s.p.d. } B^{1/2}AB^{1/2}\mathbf{y} = B^{1/2}\mathbf{b} \quad \mathbf{x} = B^{1/2}\mathbf{y}.$$

If  $N(B) = \{0\}$ , all three transformed systems in above are equivalent to  $A\mathbf{x} = \mathbf{b}$ . However, only the split preconditioned system is symmetric and can be used with CG. When the split preconditioner is used together with CG, operations only with  $B$  are required and it is not necessary to explicitly construct  $B^{1/2}$ . This important feature of CG.

The convergence properties of the CG iteration applied to the split-preconditioned system are related to the eigenvalues of the matrix

$$B^{1/2}AB^{1/2}.$$

Due to the symmetry, the matrix  $B^{1/2}AB^{1/2}$  is unitary diagonalizable. The eigenvalues satisfy

See video on preconditioning in Youtube



$$B^{1/2}AB^{1/2}\mathbf{v} = \lambda\mathbf{v}$$

now, make a change of variables to  $\mathbf{q} = B^{1/2}\mathbf{v}$  and multiply with  $B^{1/2}$ . This gives

$$BA\mathbf{q} = \lambda\mathbf{q}$$

In addition, the system is typically multiplied with  $A$ , so that one obtains

$$ABA\mathbf{q} = \lambda A\mathbf{q}.$$

This is a generalized eigenvalue problem with same eigenvalues as the system  $B^{1/2}AB^{1/2}$ . This form is especially suitable for convergence analysis of PCG. The maximal and minimal eigenvalues can be computed via Rayleigh-Ritz quotients as

$$\lambda_{min} = \min_{\mathbf{x} \in \mathbb{R}^N} \frac{\mathbf{x}^T ABA\mathbf{x}}{\mathbf{x}^T A\mathbf{x}} \quad \text{and} \quad \lambda_{max} = \max_{\mathbf{x} \in \mathbb{R}^N} \frac{\mathbf{x}^T ABA\mathbf{x}}{\mathbf{x}^T A\mathbf{x}}$$

A good preconditioned  $B$  approximates  $A^{-1}$  in the sense that  $\lambda_{min}$  and  $\lambda_{max}$  are close to one.

# THE END !