

Chapter 2

Direct linear solvers

2.1 Introduction

A linear system is a collection of m linear equations involving n unknowns, with $n, m \in \mathbb{N}$, taking the form

$$\sum_{k=1}^n a_{j,k} x_k = b_j \quad \forall j = 1 \dots m. \quad (2.1)$$

In these equations $\mathbf{x} = (x_k)_{k=1}^n$ is the unknown vector, the $a_{j,k}$ are the entries (or coefficients) of the matrix of the linear system, and $\mathbf{b} = (b_j)_{j=1}^m$ is the right hand side. In what follows the matrix of the linear system will be denoted $\mathbf{A} = (a_{j,k})$, so that this linear system rewrites as

$$\mathbf{A} \cdot \mathbf{x} = \mathbf{b}. \quad (2.2)$$

In this course we shall only consider invertible linear systems, which corresponds to the situations where the matrix \mathbf{A} admits an inverse denoted \mathbf{A}^{-1} . According to the rank-nullity theorem, a consequence of this assumption is that linear systems are square $m = n$ so that the matrix \mathbf{A} shall admit as many rows as columns.

Effective solution methods for problems of the form (2.1) are called linear solvers. There exists essentially two families of linear solvers

- **Direct solvers** where the solution \mathbf{x} is computed in a finite (but potentially large) number of operations.
- **Iterative solvers** where the linear system is reformulated as a fixed point problem $\mathbf{x} = \Phi(\mathbf{x})$. The solution \mathbf{x} is then approximated by means of an algorithm of the form $\mathbf{x}^{(p+1)} = \Phi(\mathbf{x}^{(p)})$.

In this chapter we will focus on direct solvers. Iterative solvers are the subject of subsequent chapters of this course.

2.2 Triangular systems

Firstly let us consider the case of a lower triangular linear system, that is the case where $a_{j,k} = 0$ for $j < k$. In this case the matrix A admits the form

$$A = \begin{bmatrix} a_{1,1} & 0 & \cdots & 0 \\ a_{2,1} & a_{2,2} & \ddots & \vdots \\ \vdots & * & \ddots & 0 \\ a_{n,1} & a_{n,2} & \cdots & a_{n,n} \end{bmatrix}$$

The algorithm to solve the problem $A\mathbf{x} = \mathbf{b}$, forward substitution method, is quite natural. One first observes that the j -th row (j -th equation) can be re-written

$$x_j = \frac{1}{a_{j,j}} \left(b_j - \sum_{k=1}^{j-1} a_{j,k} x_k \right) \quad (2.3)$$

From this identity, the value of x_j is deduced from the values of $x_{j-1}, x_{j-2}, \dots, x_1$. One thus starts with the first row that writes $x_1 = b_1/a_{1,1}$, then x_2, x_3, \dots, x_n are successively obtained by means of the recurrence relations (2.3).

Case of an upper triangular system In the case of an upper triangular system, that is when $a_{j,k} = 0$ for $j > k$, the solution principle is the same. This time the method is called backward substitution method. The solution procedure starts from the last row, then the following recursion formula is applied

$$x_j = \frac{1}{a_{j,j}} \left(b_j - \sum_{k=j+1}^n a_{j,k} x_k \right) \quad (2.4)$$

Algorithmic complexity Let us examine the number of operations required for the forward substitution method with respect to n , when solving a lower triangular system. If we do not make any particular assumption on the matrix A , the solution procedure that we have just discussed requires

- n divisions,
- $n - 1$ subtractions,
- $\sum_{i=2}^n \sum_{j=1}^{i-1} 1 = \sum_{i=2}^n (i-1) = n(n-1)/2$ multiplications,
- $\sum_{i=3}^n \sum_{j=1}^{i-2} 1 = \sum_{i=2}^n (i-2) = (n-1)(n-2)/2$ additions.

This leads to a total cost of $n + (n-1) + n(n-1)/2 + (n-1)(n-2)/2 = n^2$ elementary operations ("floating point operation" = flop). The algorithm that we have just described is thus said to admit an algorithmic complexity of $\mathcal{O}(n^2)$.

We assume that the matrix A admits only $\mathcal{O}(1)$ nonzero term in each row (for example it admits band structure with a band of fixed width), then the solution to this triangular system only costs $\mathcal{O}(n)$, which is the cost of a matrix-vector product.

2.3 Gaussian elimination

We come back to the case of an arbitrary linear system. The gaussian elimination method consists in reducing the initial linear system to an equivalent upper triangular system (obviously, with right hand side modified accordingly). In the sequel $\mathbf{e}_j, j = 1 \dots n$ shall refer to the canonical basis of $\mathbb{C}^{n \times n}$.

We know that $\mathbf{A}\mathbf{x} = \mathbf{b} \iff \mathcal{L}^{(1)}\mathbf{A}\mathbf{x} = \mathcal{L}^{(1)}\mathbf{b}$ if the matrix $\mathcal{L}^{(1)} \in \mathbb{C}^{n \times n}$ is invertible. The matrix of the transformed linear system takes the form $\mathbf{A}^{(1)} := \mathcal{L}^{(1)}\mathbf{A} = (a_{j,k}^{(1)})_{j,k=1 \dots n} \in \mathbb{C}^{n \times n}$, and $\mathbf{b}^{(1)} := \mathcal{L}^{(1)}\mathbf{b}$. Let us choose $\mathcal{L}^{(1)}$ so as to make sure that $a_{j,1}^{(1)} = 0$ for $j = 2, \dots, n$ i.e. so that $\mathbf{A}^{(1)}$ takes the (block lower triangular) form

$$\mathbf{A}^{(1)} = \begin{bmatrix} a_{1,1}^{(1)} & * & \cdots & * \\ 0 & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ 0 & * & \cdots & * \end{bmatrix}$$

It suffices to choose

$$\mathcal{L}^{(1)} := \text{Id} - \boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top$$

where $\boldsymbol{\ell}_1^\top := (0, a_{2,1}/a_{1,1}, \dots, a_{n,1}/a_{1,1})$.

Note that $\boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top$ is indeed a rank 1 matrix of size $n \times n$ not to be confused with $\mathbf{e}_1^\top \cdot \boldsymbol{\ell}_1$ (which is a simple scalar value). Besides we have $\mathbf{e}_1^\top \cdot \boldsymbol{\ell}_1 = 0$. Let us verify that $\mathcal{L}^{(1)}$ is invertible. We have $(\text{Id} - \boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top) \cdot (\text{Id} + \boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top) = \text{Id} - (\boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top)^2 = \text{Id} - \boldsymbol{\ell}_1 \cdot (\mathbf{e}_1^\top \cdot \boldsymbol{\ell}_1) \cdot \mathbf{e}_1^\top = \text{Id}$. We have just proved that

$$(\mathcal{L}^{(1)})^{-1} = \text{Id} + \boldsymbol{\ell}_1 \cdot \mathbf{e}_1^\top.$$

Note that $\mathbf{e}_1^\top \cdot \mathbf{A}^{(1)} = \mathbf{e}_1^\top \cdot \mathbf{A}$. To conclude, note the particular role played by the coefficient $a_{1,1}$ in this procedure. This coefficient is called first *pivot*. For the matrix $\mathcal{L}^{(1)}$ to be properly defined, this pivot needs to be non-zero.

Second iteration We can reiterate the procedure that we have just described. Indeed we have $\mathbf{A}^{(1)}\mathbf{x} = \mathbf{b}^{(1)} \iff \mathcal{L}^{(2)}\mathbf{A}^{(1)}\mathbf{x} = \mathcal{L}^{(2)}\mathbf{b}^{(1)}$ provided that $\mathcal{L}^{(2)}$ is invertible. Set $\mathbf{A}^{(2)} = \mathcal{L}^{(2)}\mathbf{A}^{(1)} = (\mathcal{L}^{(2)} \cdot \mathcal{L}^{(1)}) \cdot \mathbf{A} = (a_{j,k}^{(2)})_{j,k=1 \dots n}$ and $\mathbf{b}^{(2)} = \mathcal{L}^{(2)}\mathbf{b}^{(1)}$. We will choose $\mathcal{L}^{(2)}$ so as to make sure that $\mathbf{A}^{(2)}$ takes the form

$$\mathbf{A}^{(2)} = \begin{bmatrix} \mathbf{U}_2 & * \\ 0 & * \end{bmatrix} \quad \text{où} \quad \mathbf{U}_2 = \begin{bmatrix} a_{1,1}^{(2)} & * \\ 0 & a_{2,2}^{(2)} \end{bmatrix}$$

Here the matrix \mathbf{U}_2 is of size 2×2 . To obtain the form above, it suffices then to define $\mathcal{L}^{(2)}$ as a lower triangular matrix given by

$$\mathcal{L}^{(2)} := \text{Id} - \boldsymbol{\ell}_2 \cdot \mathbf{e}_2^\top$$

avec $\boldsymbol{\ell}_2^\top = (0, 0, a_{3,2}^{(1)}/a_{2,2}^{(1)}, \dots, a_{n,2}^{(1)}/a_{2,2}^{(1)})$.

Again the matrix $\mathcal{L}^{(2)}$ is properly defined provided that the second pivot verifies $a_{2,2}^{(1)} \neq 0$ and, in this case, $\mathcal{L}^{(2)}$ is invertible with inverse given by $(\mathcal{L}^{(2)})^{-1} = \text{Id} + \boldsymbol{\ell}_2 \cdot \mathbf{e}_2^\top$. Then we have $\mathbf{e}_1^\top \mathbf{A}^{(2)} = \mathbf{e}_1^\top \mathbf{A}^{(1)} = \mathbf{e}_1^\top \mathbf{A}$ (in particular $a_{1,1}^{(2)} = a_{1,1}^{(1)} = a_{1,1}$), and $\mathbf{e}_2^\top \mathbf{A}^{(2)} = \mathbf{e}_2^\top \mathbf{A}^{(1)}$.

General iteration We can generalize the procedure above to up any order. Assume that $A^{(p-1)}$ has been defined for $p = 2 \dots n$, and that it is upper triangular on its $p-1$ first columns i.e. it admits the following form

$$A^{(p-1)} = \begin{bmatrix} U_{p-1} & * \\ 0 & * \end{bmatrix}$$

where U_{p-1} is an upper triangular matrix of size $(p-1) \times (p-1)$. We then define $A^{(p)} = \mathcal{L}^{(p)} \cdot A^{(p-1)} = (a_{j,k}^{(p)})_{j,k=1 \dots n}$ and $\mathbf{b}^{(p)} = \mathcal{L}^{(p)} \mathbf{b}^{(p-1)}$, where $\mathcal{L}^{(p)} \in \mathbb{C}^{n \times n}$ is a lower triangular matrix given by the formula

$$\begin{aligned} \mathcal{L}^{(p)} &:= \text{Id} - \boldsymbol{\ell}_p \cdot \mathbf{e}_p^\top \\ \text{avec } \boldsymbol{\ell}_p^\top &= (0, \dots, 0, a_{p+1,p}^{(p-1)}/a_{p,p}^{(p-1)}, \dots, a_{n,p}^{(p-1)}/a_{p,p}^{(p-1)}). \end{aligned} \quad (2.5)$$

A matrix of the form (2.5) is called a Gauss transformation. In the definition of $\boldsymbol{\ell}_p$, the first p coefficients are zero, which writes $\mathbf{e}_j^\top \cdot \boldsymbol{\ell}_p = 0 \forall j = 1 \dots p$. This implies $\mathbf{e}_j^\top \mathcal{L}^{(p)} = \mathbf{e}_j^\top, \forall j = 1 \dots p$. This construction of $A^{(p)}$ hence guarantees that $\mathbf{e}_j^\top A^{(p)} = \mathbf{e}_j^\top \mathcal{L}^{(p)} A^{(p-1)} = \mathbf{e}_j^\top A^{(p-1)}$ for all $j = 1 \dots p$ and, recursively, we deduce that

$$\mathbf{e}_j^\top A^{(p)} = \mathbf{e}_j^\top A^{(j-1)} \quad \forall j = 1 \dots p$$

setting $A^{(0)} := A$. Here again, the matrix $\mathcal{L}^{(p)}$ is properly defined only if the p -th pivot satisfies $a_{p,p}^{(p-1)} \neq 0$, and its inverse is given by $(\mathcal{L}^{(p)})^{-1} = \text{Id} + \boldsymbol{\ell}_p \cdot \mathbf{e}_p^\top$. To conclude, the construction above guarantees that $A^{(p)}$ admits the form

$$A^{(p)} = \begin{bmatrix} U_p & * \\ 0 & * \end{bmatrix}$$

where U_p is an upper triangular matrix of size $p \times p$. From what precedes, we see that the first p rows and columns remain unchanged when transforming the system from $A^{(p-1)}$ to $A^{(p)}$. We obtain the following formulas expressing $a_{j,k}^{(p)}$ with respect to $a_{j,k}^{(p-1)}$,

$$\begin{cases} a_{j,k}^{(p)} = a_{j,k}^{(p-1)} & \text{for } 1 \leq j \leq p \text{ or } 1 \leq k < p \\ a_{j,k}^{(p)} = 0 & \text{if } p < j \leq n \text{ and } k = p \\ a_{j,k}^{(p)} = a_{j,k}^{(p-1)} - a_{j,p}^{(p-1)} a_{p,k}^{(p-1)} / a_{p,p}^{(p-1)} & \text{if } p < j \leq n \text{ and } p < k \leq n \end{cases} \quad (2.6)$$

Conclusion of the algorithm The gaussian elimination algorithm terminates whenever $p = n - 1$, since the matrix $A^{(n-1)}$ is then completely upper triangular. We have obtained that $A\mathbf{x} = \mathbf{b} \iff A^{(n-1)}\mathbf{x} = \mathbf{b}^{(n-1)}$. The later system being upper triangular (by construction...), it can be solved by a backward substitution algorithm as explained in the previous section 2.2.

Algorithmic complexity One can show that the gaussian elimination algorithm requires $2n^3/3 + n^2/3 - n$ operations (this is left as an exercise).

Remarks on the pivots The gaussian elimination algorithm is well defined provided that $a_{p,p}^{(p-1)} \neq 0$ for all p i.e. no pivot equals zero. There are classes of matrices for which this condition is systematically fulfilled. Here are three examples:

- symmetric positive definite matrices,
- row-wise diagonally dominant matrices: $|a_{j,j}| > \sum_{k \neq j} |a_{j,k}| \forall j = 1 \dots n$,
- column-wise diagonally dominant matrices: $|a_{j,j}| > \sum_{k \neq j} |a_{k,j}| \forall j = 1 \dots n$.

2.4 LU decomposition

The gaussian elimination method leads to a factorization of the matrix under the form $A = L \cdot U$ where L is lower triangular with $L_{j,j} = 1, j = 1 \dots n$ and U is upper triangular. Indeed, let us start by setting $U = A^{(n-1)}$ which is the matrix obtained after the $n - 1$ -th step of the Gauss method. We thus have $\mathcal{L}^{(n-1)} \mathcal{L}^{(n-2)} \dots \mathcal{L}^{(1)} A = U \iff A = L \cdot U$ with

$$L := (\mathcal{L}^{(1)})^{-1} (\mathcal{L}^{(2)})^{-1} \dots (\mathcal{L}^{(n-1)})^{-1}$$

Let us study in more details the matrix L . It appears in factorized form. We will develop this product, taking account of the elementary property $\mathbf{e}_j^\top \boldsymbol{\ell}_p = 0$ for $j = 1 \dots p$. As a particular case of this property, we see that $\mathbf{e}_1^\top \boldsymbol{\ell}_2 = 0$, which leads to the conclusion that

$$\begin{aligned} (\mathcal{L}^{(1)})^{-1} (\mathcal{L}^{(2)})^{-1} &= (\text{Id} + \boldsymbol{\ell}_1 \mathbf{e}_1^\top) (\text{Id} + \boldsymbol{\ell}_2 \mathbf{e}_2^\top) \\ &= \text{Id} + \boldsymbol{\ell}_1 \mathbf{e}_1^\top + \boldsymbol{\ell}_2 \mathbf{e}_2^\top + \boldsymbol{\ell}_1 (\mathbf{e}_1^\top \boldsymbol{\ell}_2) \mathbf{e}_2^\top \\ &= \text{Id} + \boldsymbol{\ell}_1 \mathbf{e}_1^\top + \boldsymbol{\ell}_2 \mathbf{e}_2^\top \end{aligned}$$

We then proceed by recurrence on p to show that $(\mathcal{L}^{(1)})^{-1} (\mathcal{L}^{(2)})^{-1} \dots (\mathcal{L}^{(p)})^{-1} = \text{Id} + \sum_{j=1}^p \boldsymbol{\ell}_j \mathbf{e}_j^\top$. We have just proved that this property holds for $p = 1$. Assume that it holds for p , and let us prove that it holds for $p + 1$ as well. We have

$$\begin{aligned} (\mathcal{L}^{(1)})^{-1} (\mathcal{L}^{(2)})^{-1} \dots (\mathcal{L}^{(p)})^{-1} (\mathcal{L}^{(p+1)})^{-1} &= (\text{Id} + \sum_{j=1}^p \boldsymbol{\ell}_j \mathbf{e}_j^\top) (\text{Id} + \boldsymbol{\ell}_{p+1} \mathbf{e}_{p+1}^\top) \\ &= \text{Id} + \left(\sum_{j=1}^p \boldsymbol{\ell}_j \mathbf{e}_j^\top \right) + \boldsymbol{\ell}_{p+1} \mathbf{e}_{p+1}^\top + \left(\sum_{j=1}^p \boldsymbol{\ell}_j \mathbf{e}_j^\top \right) \boldsymbol{\ell}_{p+1} \mathbf{e}_{p+1}^\top \\ &= \text{Id} + \sum_{j=1}^{p+1} \boldsymbol{\ell}_j \mathbf{e}_j^\top + \left(\sum_{j=1}^p \boldsymbol{\ell}_j \mathbf{e}_j^\top \right) \boldsymbol{\ell}_{p+1} \mathbf{e}_{p+1}^\top \\ &= \text{Id} + \sum_{j=1}^{p+1} \boldsymbol{\ell}_j \mathbf{e}_j^\top + \sum_{j=1}^p \boldsymbol{\ell}_j (\mathbf{e}_j^\top \boldsymbol{\ell}_{p+1}) \mathbf{e}_{p+1}^\top \end{aligned}$$

Since $\mathbf{e}_j^\top \boldsymbol{\ell}_{p+1} = 0$ for $j = 1 \dots p$, the last term in the right hand side above is zero, so that the property that we want to prove holds for $p + 1$. We finally obtain, for $p = n - 1$,

$$L = \text{Id} + \sum_{j=1}^{n-1} \boldsymbol{\ell}_j \mathbf{e}_j^\top$$

Let us examine the j -th column of this matrix. We have $L \cdot \mathbf{e}_j = \mathbf{e}_j + \boldsymbol{\ell}_j$, with $(\mathbf{e}_j + \boldsymbol{\ell}_j)^\top = (0, \dots, 0, 1, a_{j+1,j}^{(j-1)}/a_{j,j}^{(j-1)}, \dots, a_{n,j}^{(j-1)}/a_{j,j}^{(j-1)})$. We see that L is indeed upper triangular, with 1's on the diagonal, and the expression that we have just obtained for its columns suggests that this matrix should be assembled on-the-fly during the gaussian elimination algorithm.

2.5 Factorization algorithm

Let us now examine the algorithmic details of an LU decomposition. Before going into the description of the algorithm itself, we introduce a few simple notations. If J and K are two subsets of $\llbracket 1, n \rrbracket := \{1, \dots, n\}$, we will denote $A_{J,K} \in \mathbb{C}^{|J| \times |K|}$ the submatrix obtained out of A by extracting the columns $k \in K$ and the rows $j \in J$ (here we denote $|J|$ the cardinal of J , and accordingly for $|K|$). For a $j \in \llbracket 1, n \rrbracket$ and a subset $K \subset \llbracket 1, n \rrbracket$, we denote $A_{j,K}$ instead of $A_{\{j\},K}$ and, similarly, for $k \in \llbracket 1, n \rrbracket$ and $J \subset \llbracket 1, n \rrbracket$ we will denote $A_{J,k}$ instead of $A_{J,\{k\}}$.

As suggested by the algorithm of the gaussian elimination described above, an LU decomposition method involves $n - 1$ for a square matrix $A \in \mathbb{C}^{n \times n}$. Iteration p consists in an update of the matrix $A^{(p-1)}$ to obtain the matrix $A^{(p)}$. On the theoretical side, this update takes the form of a left-multiplication by the matrix $\mathcal{L}^{(p)} := \text{Id} - \ell_p \cdot \mathbf{e}_p^\top$, but this is not the actual way this update takes place in practice: this would be unnecessarily costly. Here are a few simple observations concerning this update:

- it does not modify the rows $1, \dots, p$
- it does not modify the columns $1, \dots, p - 1$
- in column p , it cancels out the elements $a_{j,p}^{(p-1)}$ for $j = p + 1, \dots, n$

As a consequence, to change $A^{(p-1)}$ into $A^{(p)}$, it suffices to modify the block $A_{J,J}^{(p-1)}$ for $J = \llbracket p + 1, n \rrbracket$. Besides, the matrix $\mathcal{L}^{(p)}$ is entirely determined by ℓ_p . These elementary remarks together with §2.3, lead to Algorithm 1 that follows

Algorithm 1

```

function LU_NAIF(A)
  L = Id
  A(0) = A
  for  $p = 1 \dots n - 1$  do
    J =  $\llbracket p + 1, n \rrbracket$ 
    LJ,p = AJ,p(p-1) / Ap,p(p-1)
    AJ,J(p) = AJ,J(p-1) - LJ,p Ap,J(p-1)
  end for
  U = A(n-1)
  return (L, U)
end function

```

Algorithm 2

```

function LU_DECOMPOSE(A)
  T = A
  for  $p = 1 \dots n - 1$  do
    J =  $\llbracket p + 1, n \rrbracket$ 
    TJ,p = TJ,p / Tp,p
    TJ,J = TJ,J - TJ,p Tp,J
  end for
  return (T)
end function

```

In fact Algorithm 1 makes use of many unnecessary intermediate variables. One may store L (resp. U) in the lower (resp. upper) triangular part of a single matrix $T \in \mathbb{C}^{n \times n}$. During this construction, it is also possible to store the matrices $A^{(p)}$ inside T . This leads to Algorithm 2 that only involves a single additional matrix.

Once the LU decomposition of the matrix A is available et stored in the matrix T , we can solve the linear system $Ax = b$ by means of the successive application of a backward and a forward substitution method based on the upper and lower triangular parts of T (the diagonal coefficients of the lower triangular system must equal 1 hence do not need to be stored). The effective solution of a linear system making use of the a priori knowledge of the LU decomposition then takes the form of Algorithm 3 below.

Algorithm 3

```

function LU_SOLVE(A,b)
  T = LU_DECOMPOSE(A)
  //descente
  for  $j = 1 \dots n$  do
     $v_j = b_j$ 
    for  $k = 1 \dots j - 1$  do
       $v_j = v_j - T_{j,k} v_k$ 
    end for
  end for
  //remontee
  for  $p = 1 \dots n$  do
     $j = n - p + 1$ 
     $u_j = v_j$ 
    for  $k = j + 1 \dots n$  do
       $u_j = u_j - T_{j,k} u_k$ 
    end for
     $u_j = u_j / T_{j,j}$ 
  end for
  return ( $u$ )
end function

```

Algorithm 4

```

function LU_PIVOT(A)
  T = A, P = Id
  for  $q = 1 \dots n - 1$  do
    //pivoting
     $r = \operatorname{argmax}_{j=q, \dots, n} |T_{j,q}|$ 
     $P = \tau(q, r) \cdot P$ ,  $T = \tau(q, r) \cdot T$ 
    //gaussian elimination
     $J = \llbracket p + 1, n \rrbracket$ 
     $T_{J,p} = T_{J,p} / T_{p,p}$ 
     $T_{J,J} = T_{J,J} - T_{J,p} T_{p,J}$ 
  end for
  return (T, P)
end function

```

2.6 Partial pivoting

As we saw, the gaussian elimination algorithm stops if one of the pivots is zero i.e. $a_{q,q}^{(q-1)} = 0$. One way to circumvent this issue relies on a so-called partial pivoting strategy. It consists in a preliminary step taking place at each iteration q , where the rows q and r are swapped with r chosen so that

$$|a_{r,q}^{(q-1)}| = \max_{j=q \dots n} |a_{j,q}^{(q-1)}| \quad (2.7)$$

This is equivalent to a left multiplication of the matrix A^{q-1} by a permutation matrix P_q defined by: $P_q \cdot e_q = e_r$, $P_q \cdot e_r = e_q$ and $P_q \cdot e_j = e_j$ if $j \neq q, r$. It is important to note

that $(P_q)^2 = \text{Id}$ and that $P_q^\top = P_q$. In addition observe that the coefficient $a_{r,q}^{(q-1)}$ in (2.7) necessarily satisfies $a_{r,q}^{(q-1)} \neq 0$ otherwise the matrix A would not be invertible.

With the pivoting strategy described above, the result of $n-1$ iterations of the gaussian elimination algorithm writes $\mathcal{L}^{(n-1)}P_{n-1} \cdots P_2 \mathcal{L}^{(1)}P_1 \cdot A = U$. This can be rewritten equivalently as follows:

$$\begin{aligned} A &= P_1(\mathcal{L}^{(1)})^{-1}P_2 \cdots P_{n-1}(\mathcal{L}^{(n-1)})^{-1} \cdot U \\ \iff PA &= P \cdot P_1(\mathcal{L}^{(1)})^{-1}P_2 \cdots P_{n-1}(\mathcal{L}^{(n-1)})^{-1} \cdot U \\ &\text{with } P = P_{n-1}P_{n-2} \cdots P_1 \\ \iff PA &= \mathcal{T}_1 \cdot \mathcal{T}_2 \cdots \mathcal{T}_{n-1} \cdot U \\ &\text{where } \mathcal{T}_q = P_{n-1} \cdots P_{q+1} \cdot (\mathcal{L}^{(q)})^{-1} \cdot P_{q+1} \cdots P_{n-1} \end{aligned}$$

In the calculation above, we have considered that $\mathcal{T}_{n-1} = (\mathcal{L}^{(n-1)})^{-1}$. Let us check that the matrices \mathcal{T}_q are lower triangular. Recall that $(\mathcal{L}^{(q)})^{-1} = \text{Id} + \ell_q \cdot e_q^\top$, where $\ell_q \in \mathbb{C}^{n \times n}$, verifies $e_j^\top \cdot \ell_q = 0$ for $j = 1 \dots q$. As a consequence, we have

$$\begin{aligned} \mathcal{T}_q &= P_{n-1} \cdots P_{q+1} \cdot (\text{Id} + \ell_q \cdot e_q^\top) \cdot P_{q+1} \cdots P_{n-1} \\ &= \text{Id} + (P_{n-1} \cdots P_{q+1} \cdot \ell_q) \cdot (e_q^\top \cdot P_{q+1} \cdots P_{n-1}) \\ &= \text{Id} + (P_{n-1} \cdots P_{q+1} \cdot \ell_q) \cdot (P_{n-1} \cdots P_{q+1} \cdot e_q)^\top \\ &= \text{Id} + \ell'_q \cdot e_q^\top \quad \text{with } \ell'_q := P_{n-1} \cdots P_{q+1} \cdot \ell_q \end{aligned}$$

In the last step of the calculation above, we have used the fact that $P_q e_j = e_j$ whenever $q > j$. Using this same property, we also see that, for all $j = 1 \dots q$, we have $e_j^\top \cdot \ell'_q = e_j^\top \cdot P_{n-1} \cdots P_{q+1} \cdot \ell_q = (P_{q+1} \cdots P_{n-1} \cdot e_j)^\top \cdot \ell_q = e_j^\top \cdot \ell_q = 0$. This proves that \mathcal{T}_q is lower triangular. Setting this time $L = \mathcal{T}_1 \cdot \mathcal{T}_2 \cdots \mathcal{T}_{n-1}$, we have obtained

$$P \cdot A = L \cdot U$$

where L is lower triangular, and U is upper triangular. Such a decomposition holds as soon as A is invertible.

A modified version of Algorithm 2 taking account of the row-wise partial pivoting is given in Algorithm 4 above. In this algorithm $\tau(q, r) \in \mathbb{C}^{n \times n}$ is the transposition matrix such that $\tau(q, r)e_k = e_k$ if $k \neq q, r$, $\tau(q, r)e_q = e_r$ and $\tau(q, r)e_r = e_q$. Obviously, to compute $\tau(q, r) \cdot P$ and $\tau(q, r) \cdot T$, no need to perform a full matrix-matrix product (a costly operation with $\mathcal{O}(n^3)$ algorithmic complexity a priori); it suffices to swap rows q and r of matrices P and T (a fast operation with $\mathcal{O}(n)$ algorithmic complexity). At the end of the day, the solution algorithm proceeds as in Algorithm 3, applying the permutation P on the right hand side as a preliminary step.

2.7 Theoretical results regarding LU

Apart from the construction and practical details of the LU decomposition that we have presented, theoretical questions naturally arise concerning this factorization.

Théorème 2.1.

Given an invertible matrix $A \in \mathbb{C}^{n \times n}$, we have $\det(A_{J,J}) \neq 0$ for $J = \llbracket 1, j \rrbracket$ and for all $j = 1 \dots n-1$ if and only if there exists a unique pair $L, U \in \mathbb{C}^{n \times n}$ satisfying $A = LU$ with U upper triangular and L lower triangular with $L_{j,j} = 1 \forall j = 1 \dots n$.

Proof:

We proceed by recurrence on the dimension n . The result for $n = 1$ is obvious. Assume that the result holds for $A \in \mathbb{C}^{m \times m}$ invertible with $m = 1 \dots n$, and let us prove that this result still holds $A \in \mathbb{C}^{(n+1) \times (n+1)}$ invertible.

Now assume for a moment the existence/uniqueness of the decomposition $A = LU$. Then we necessarily have $A_{J,J} = L_{J,J}U_{J,J}$ for all $J = \llbracket 1, k \rrbracket$ and all $k = 1 \dots n$. On the other hand, $0 \neq \det(A) = \det(U) = \prod_{j=1}^{n+1} U_{j,j} \Rightarrow 0 \neq \prod_{j=1}^k U_{j,j} = \det(A_{J,J})$. Now let us set $B := A_{J,J}$ for $J = \llbracket 1, n \rrbracket$, so that

$$A = \begin{bmatrix} B & \mathbf{a} \\ \mathbf{b}^* & \alpha \end{bmatrix} \quad (2.8)$$

where $\mathbf{a}, \mathbf{b} \in \mathbb{C}^n$ and $\alpha \in \mathbb{C}$. Now we assume that $\det(A_{J,J}) \neq 0$ for $J = \llbracket 1, j \rrbracket$ and for all $j = 1, \dots, n$, as well as $\det(A) \neq 0$. According to the recurrence hypothesis, there is existence/uniqueness of the factorization $B = \tilde{L}\tilde{U}$. As a consequence we have the decomposition $A = LU$ if and only if there exist $\mathbf{u}, \mathbf{v} \in \mathbb{C}^n$ and $\beta \in \mathbb{C} \setminus \{0\}$ such that

$$L = \begin{bmatrix} \tilde{L} & 0 \\ \mathbf{v}^* & 1 \end{bmatrix} \quad U = \begin{bmatrix} \tilde{U} & \mathbf{u} \\ 0 & \beta \end{bmatrix} \quad \text{et} \quad \begin{array}{l} \tilde{L}\mathbf{u} = \mathbf{a} \\ \tilde{U}^*\mathbf{v} = \mathbf{b} \\ \beta + \mathbf{v}^*\mathbf{u} = \alpha \end{array} \quad (2.9)$$

Besides $\det(\tilde{L}) = 1$ and $\det(\tilde{U}) = \det(\tilde{B}) \neq 0$. Hence this leads to existence and uniqueness of $\mathbf{u}, \mathbf{v}, \beta$ solution to (2.9) and, from the existence/uniqueness of \tilde{L}, \tilde{U} we deduce existence/uniqueness of L, U . \square

The criterium on the $\det(A_{J,J})$'s mentionned above is not very easy to verify in practice. However there exist certain classes of matrices for which this criterium is systematically satisfied. This is the case of diagonally dominant matrices. Recall that a matrix $A = (a_{j,k}) \in \mathbb{C}^{n \times n}$ is said row-wise diagonally dominant if $|a_{j,j}| > \sum_{k \neq j} |a_{j,k}| \forall j = 1 \dots n$. Similarly it is said column-wise diagonally dominant when $|a_{j,j}| > \sum_{k \neq j} |a_{k,j}| \forall j = 1 \dots n$.

Proposition 2.2.

If the matrix $A \in \mathbb{C}^{n \times n}$ is either row-wise or column-wise diagonally dominant, then it admits a unique LU-decomposition.

Démo:

Let us assume row-wise diagonal dominance, since the case of a column-wise diagonal dominance proceed very similarly. According to Theorem 2.1, it suffices to show that $\det(A_{J,J}) \neq 0$ for $J = \llbracket 1, k \rrbracket$ for all $k = 1 \dots n$. It appears clear that if A is row-wise diagonally dominant, this is also the case for $A_{J,J}$. Hence there remains to verify that a row-wise diagonally dominant matrix is necessarily invertible.

Let us proceed by contradiction and consider a row-wise diagonally dominant matrix $B = (b_{j,k}) \in \mathbb{C}^{n \times n}$ such that $B\mathbf{u} = 0$ for a certain $\mathbf{u} = (u_j) \in \mathbb{C}^n \setminus \{0\}$. Let $p \in \{1 \dots n\}$ such that $|u_p| = \max_{j=1 \dots n} |u_j| = |\mathbf{u}|_\infty$. Considering \mathbf{u}/u_p instead of \mathbf{u} if necessary, we may assume

that $u_p = 1$ and thus $|u_j| \leq 1$ for all $j = 1 \dots n$. Then we have $0 = b_{p,p} + \sum_{j \neq p} b_{p,j} u_j$ and thus $|b_{p,p}| = |\sum_{j \neq p} b_{p,j} u_j| \leq \sum_{j \neq p} |b_{p,j}|$ which contradicts the diagonal dominance. \square

The next result shows that if the matrix A has a band structure with a certain width, then each of the factors of the LU decomposition admits the same band structure.

Proposition 2.3.

Let $A \in \mathbb{C}^{n \times n}$ a matrix admitting a unique LU factorization. Let us assume in addition that there exists $p \geq 0$ such that $A_{j,k} = 0$ if $|j - k| > p$. Then we also have $L_{j,k} = U_{j,k} = 0$ for $|j - k| > p$.

Démo:

Again we proceed by recurrence on the band width, and assume that the result holds for all matrices uniquely LU-factorizable with a band width of m with $1 \leq m \leq n$. We pick a matrix $A \in \mathbb{C}^{(n+1) \times (n+1)}$ satisfying the assumptions of the proposition we are seeking to prove. Coming back to the notations of the proof of Theorem 2.1, we must in particular have (2.9) with $\tilde{L}_{j,k} = \tilde{U}_{j,k} = 0$ whenever $|j - k| > p$ according to the recurrence hypothesis.

Let us also take the notation $\mathbf{u} = (u_j), \mathbf{v} = (v_j), \mathbf{a} = (a_j), \mathbf{b} = (b_j)$ for the vectors coming into play in (2.9) and (2.8). To prove that L and U admit the same band structure, that is $L_{j,k} = U_{j,k} = 0$ for $|j - k| > p$, there remains to verify that $u_j = v_j = 0$ for $|j - n| > p \iff j < n - p$. We already know that $a_j = b_j = 0$ for $j < n - p$ according to the band structure satisfied by A . On the other hand \tilde{L} and \tilde{U}^* are lower triangular, hence $(\tilde{L})^{-1}$ and $(\tilde{U}^*)^{-1}$ also, and we have $\mathbf{u} = (\tilde{L})^{-1} \mathbf{a}$ and $\mathbf{v} = (\tilde{U}^*)^{-1} \mathbf{b}$, which indeed implies $u_j = v_j = 0$ for $j < n - p$. \square

2.8 Cholesky factorization

Recall that a matrix $A \in \mathbb{C}^{n \times n}$ is hermitian whenever $A = A^*$, and it is called positive when $\mathbf{x}^* A \mathbf{x} \in (0, +\infty)$ for all $\mathbf{x} \in \mathbb{C}^n$. Finally, when it is positive, the matrix A is called definite if $\mathbf{x}^* A \mathbf{x} = 0 \Rightarrow \mathbf{x} = 0$. In the case where the matrix $A \in \mathbb{C}^{n \times n}$ is hermitian positive definite (HPD), its LU-factorization simplifies, and we can obtain more explicit formulas.

Théorème 2.4.

Let $A = (a_{j,k})_{j,k=1 \dots n} \in \mathbb{C}^{n \times n}$ be hermitian positive definite. Then there exists a lower triangular matrix $H \in \mathbb{C}^{n \times n}$ whose diagonal terms are strictly real positive and such that the following so-called Cholesky factorization holds

$$A = HH^*.$$

Démo:

Again we proceed by recurrence on the size n . In the case $n = 1$, we have $A = (a_{1,1})$. Since A is hermitian positive definite, we have $a_{1,1} > 0$. We can thus define $H = (h_{1,1})$ where $h_{1,1} = \sqrt{a_{1,1}}$. Let us assume that the Cholesky factorization uniquely exists for any $A \in \mathbb{C}^{p \times p}$ HPD, for $1 \leq p \leq n$, and let us prove that this property still holds $A \in \mathbb{C}^{(n+1) \times (n+1)}$ HPD. Pick an arbitrary HPD matrix $A \in \mathbb{C}^{(n+1) \times (n+1)}$. We can write it in the following form

$$A = \begin{bmatrix} B & \mathbf{a} \\ \mathbf{a}^* & \alpha \end{bmatrix} \quad (2.10)$$

where $B \in \mathbb{C}^{n \times n}$ is hermitian, $\mathbf{a} \in \mathbb{C}^n$ and $\alpha \in \mathbb{R}_+^*$. Let us prove that B is positive definite, that is $\mathbf{y}^* B \mathbf{y} > 0$, for all $\mathbf{y} \in \mathbb{C}^n \setminus \{0\}$. Pick an arbitrary $\mathbf{y} \in \mathbb{C}^n \setminus \{0\}$, and define $\mathbf{x} \in \mathbb{C}^{n+1}$ by $\mathbf{x}^\top = (\mathbf{y}^\top, 0)$. Since A is HPD, we have $0 < \mathbf{x}^* A \mathbf{x} = \mathbf{y}^* B \mathbf{y}$ hence B is HPD. According to the recurrence hypothesis, there exists a matrix $M \in \mathbb{C}^{n \times n}$, $M = (m_{i,j})_{i,j=1}^n$, such that $m_{i,j} = 0$ if $j > i$, $m_{i,i} > 0$ and $B = MM^*$. We look for H under the form:

$$H = \begin{bmatrix} M & 0 \\ \mathbf{b}^* & \beta \end{bmatrix} \quad (2.11)$$

with $\mathbf{b} \in \mathbb{C}^n$, $\beta \in \mathbb{R}_+^*$ and such that $HH^* = A$. To determine \mathbf{b} and β , let us compute HH^* from (2.11) and let us identify with A :

$$HH^* = \begin{bmatrix} MM^* & M\mathbf{b} \\ (M\mathbf{b})^* & |\mathbf{b}|^2 + \beta^2 \end{bmatrix} \quad (2.12)$$

Comparing the above identity with (2.10), we deduce that the following equations must be satisfied

$$M\mathbf{b} = \mathbf{a} \quad \text{et} \quad |\mathbf{b}|^2 + \beta^2 = \alpha.$$

Let us note that M is invertible since $\det(M) = m_{1,1} \cdot m_{2,2} \cdots m_{n,n} > 0$. We can thus set $\mathbf{b} := M^{-1}\mathbf{a}$ as a definition. In practice we can determine \mathbf{b} by a simple forward substitution method since M is lower triangular. Plugging $B = MM^*$ into the second equality yields $\mathbf{a}^* B^{-1} \mathbf{a} + \beta^2 = \alpha$. It thus suffices to set $\beta := (\alpha - \mathbf{a}^* B^{-1} \mathbf{a})^{1/2}$. However such a definition is valid only if

$$\alpha - \mathbf{a}^* B^{-1} \mathbf{a} > 0 \quad (2.13)$$

Let us prove that this condition is indeed satisfied by (2.10). Let us consider the vector $\mathbf{z} \in \mathbb{C}^{n+1} \setminus \{0\}$ defined by $\mathbf{z}^\top = ((B^{-1}\mathbf{a})^\top, -1)$. Then we have $0 < \mathbf{z}^* A \mathbf{z} = \alpha - \mathbf{a}^* B^{-1} \mathbf{a}$. This leads to the conclusion. \square

A very explicit algorithm can be proposed for Cholesky factorization. By directly expressing the matrix product $A = HH^*$, coefficient by coefficient, we obtain $a_{j,j} = \sum_{k=1}^j |h_{j,k}|^2$ and $a_{j,k} = \sum_{p=1}^k h_{j,p} \bar{h}_{k,p}$ for $k < j$. From there we deduce the following formulas

$$\begin{aligned} h_{j,j} &= (a_{j,j} - \sum_{k=1}^{j-1} |h_{j,k}|^2)^{1/2}, \\ h_{j,k} &= (a_{j,k} - \sum_{p=1}^{k-1} h_{j,p} \bar{h}_{k,p}) / \bar{h}_{k,k} \quad \text{pour } k < j. \end{aligned}$$

The construction of the matrix H then proceeds for j growing from 1 to n and for k growing from 1 to j . These formulas yield Algorithm 5 below.

Algorithm 5

```

function CHOLESKY(A)
  H = 0
  for  $j = 1 \dots n$  do
    for  $k = 1 \dots j$  do
       $H_{j,k} = A_{j,k}$ 
      for  $p = 1 \dots k - 1$  do
         $H_{j,k} = H_{j,k} - H_{j,p} \bar{H}_{k,p}$ 
      end for
      if  $k < j$  then
         $H_{j,k} = H_{j,k} / H_{k,k}$ 
      else
         $H_{j,j} = \sqrt{H_{j,j}}$ 
      end if
    end for
  end for
  return (H)
end function

```
