

Computational Physics III: Report 2
Linear systems solving and diagonalization methods

Due on April 30, 2020

April 26, 2020

Raffaele Ancarola

Contents

Introduction

Solving a system of linear equations

A *linear problem* could be defined as a system of which the describing equations are all linear. Furthermore, a linear system is said to be *determined* if the number of equations N is finite and it corresponds to the number of the unknowns. Such a system defined on a field \mathbb{K} takes the advantage to be written in a matrix form:

$$A \cdot \vec{x} = \vec{b} \quad (1)$$

where A is the describing matrix and $\vec{b} \in \mathbb{K}^N$ the affine component of the system, or the components which are independent with respect to the unknowns contained in \vec{x} . Because the system is *determined*, the condition that A must satisfy is the *invertibility*, then $A \in \mathcal{GL}(N)$ and a solving the system means to find $\vec{x} \in \mathbb{K}^N$ such that eq. (1) is satisfied. There exist various approaches that can reach this attempt, in this report three cases will be analysed: the *Gauss elimination*, the *LU decomposition* and the *diagonalisation*.

Gauss elimination algorithm

The *Gauss elimination* bases to the fact that any square matrix can be decomposed into a finite sequence of elementary operations $\{P_k\}_{1 \leq k \leq M}$, $M \in \mathbb{N}^*$. There are basically three kinds of them:

- Multipling of a row by a scalar factor $\lambda \in \mathbb{K}$
- Switching a row with another
- Adding a row with a multiple of another

The purpose of this method is to reduce the involved matrix A into the identity applying the same operations to the vector \vec{b} , as shown in the equation (2).

$$A = P_1 \cdot \dots \cdot P_M \implies \vec{x} = P_M^{-1} \cdot \dots \cdot P_1^{-1} \cdot \vec{b}, \quad M \in \mathbb{N}^* \quad (2)$$

LU decomposition

The *LU decomposition* is not a direct method which solves a linear system, but it allows to simplify the resolution by decomposing the A matrix into a lower-triangular matrix L and an upper-triangular matrix U . The simplification is due to the major facility to invert the two matrices precedently presented. Once A is decomposed, the process is straight-forward:

$$A \cdot \vec{x} = L \cdot U \cdot \vec{x} = \vec{b} \quad (3)$$

$$L \cdot \vec{y} = \vec{b} \quad (3)$$

$$U \cdot \vec{x} = \vec{y} \quad (4)$$

Both equations (3) and (4) can be solved sequentially using the *Gauss elimination* method.

Diagonalization: introduction

In case A is a symmetric matrix, the spectral theorem \square states that such a matrix is equivalent (definition of equivalence here: \square) to a diagonal matrix D , where the transition matrix P is unitary ($P^{-1} = \bar{P}^T$), then:

$$A = P \cdot D \cdot \bar{P}^T \implies \vec{x} = P \cdot D^{-1} \cdot \bar{P}^T \cdot \vec{b} \quad (5)$$

Generally diagonalization is not used to solve general systems of linear equations, but it's convenient when the problem is related to find the eigen-base related to the eigen-values.

Problem 1

(1) LU decomposition implementation

This algorithm separates the input matrix A into a lower triangular L and an upper triangular U , guaranteeing that $A = L \cdot U$. Nevertheless, not all the invertible square matrices are purely LU decomposable, then it may happen that the output can result ill formed. The code (1) shows at line 23 that a division by the diagonal values is performed, causing eventually a singularity. A possible work-around is to apply the partial pivoting technique in order to swap the problematic lines. In listing (1) is shown a full implementation with partial pivoting.

Listing 1: *LU decomposition* implementation with partial pivoting

```

1  function [L, U, P] = lu_decomposition(A)
2      [Ni, Nj] = size(A);
3      assert(Ni == Nj, "The input must be diagonal");
4
5      N = Nj;
6      assert(N > 0, "The input must non empty");
7
8      L = eye(N); % if zeros doesn't give the same result
9      U = A; % if zeros doesn't give the same result
10     P = eye(N); % identity matrix
11
12     for k=1:(N-1)
13         % pivoting section
14         [Amax,r] = max(abs(U(k:N, k)));
15         r = r + k - 1;
16         % swap rows if it's not the identity swap operation
17         U([k r], :) = U([r k], :);
18         P([k r], :) = P([r k], :);
19         L([k r], 1:k-1) = L([r k], 1:k-1);
20
21         % computing LU
22         for i=(k+1):N
23             L(i,k) = U(i,k) / U(k,k);
24             U(i,:) = U(i,:) - L(i,k) * U(k,:);
25         end
26     end
27 end

```

Partial pivoting

The *LU decomposition* algorithm (presented below in exercise 1.1) can easily run into singularities, especially when A presents zeros as diagonal terms. In order to avoid divergent results, it would better select the rows of which element is not zero in the requested columns and swap them with the current one. More precisely, at the k -th step, select the r -th row such that $A_{rk} = \max_{k \leq i \leq N} |A_{ik}|$, then swap rows at the position k and r . If the pivoting is applied the resulting *LU decomposition* won't be anymore like it was defined in the previous section, but a correction to equation (3) must be applied:

$$P \cdot A = L \cdot U \implies L \cdot \vec{y} = P \cdot \vec{b} \quad (6)$$

where P is the orthogonal matrix that accumulated all row switching applications. The rest of the solving method remains unchanged.

(1.1) Solving a linear system

A linear system can be solved applying the LU decomposition and then a gauss elimination process, as shown in the equations (6) and (4).

For example, the system in equation (7) is determined and can be solved using the `solve.m` script. Additionally the `test_solve.m` script compares with the matlab `x = A \ b` verifying that the solution \vec{x} is given correctly by the `solve.m` script.

$$\begin{cases} 2x_1 + x_2 - x_3 + 5x_4 = 13 \\ x_1 + 2x_2 + 3x_3 - x_4 = 37 \\ x_1 + x_3 + 6x_4 = 30 \\ x_1 + 3x_2 - x_3 + 5x_4 = 19 \end{cases} \implies A = \begin{pmatrix} 2 & 1 & -1 & 5 \\ 1 & 2 & 3 & -1 \\ 1 & 0 & 1 & 6 \\ 1 & 3 & -1 & 5 \end{pmatrix}, \quad \vec{b} = \begin{pmatrix} 13 \\ 37 \\ 30 \\ 19 \end{pmatrix} \implies \vec{x} = A^{-1} \cdot \vec{b} = \begin{pmatrix} 2 \\ 4 \\ 10 \\ 3 \end{pmatrix} \quad (7)$$

(1.3) Decomposition of a matrix

The example taken in equation (8) is a problematic case where a pure LU decomposition doesn't exist. A necessary and sufficient condition to the existence of a pure LU decomposition is that the matrix must be gauss reducible without any row exchange (ref. [?]), that's why if such a decomposition exists, then pivoting matrix P is the identity matrix. So, the form $P \cdot A = L \cdot U$ is obtainable using the pivoting described in the previous section.

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 9 \\ 4 & -3 & 1 \end{pmatrix} \implies L = \begin{pmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0.25 & 0.5 & 1 \end{pmatrix}, \quad U = \begin{pmatrix} 4 & -3 & 1 \\ 0 & 5.5 & 8.5 \\ 0 & 0 & -1.5 \end{pmatrix} \quad P = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad (8)$$

(2) Puzzle board game

The puzzle game problem is described by using the following formalization:

$$b_k = x_k + x_{k-N} + x_{k+N} + x_{k-1} + x_{k+1}, \quad A = T \otimes 1 + 1 \otimes T - 1 \otimes 1 \quad (9)$$

where b_k corresponds to the number of times that the k_{ij} is pressed. Using the `kron` matlab function it's possible to easily construct the matrix A , as shown in the code listing (2).

Listing 2: Construction of the matrix describing the puzzle game problem

```

1 function A = puzzleA()
2     T = diag(ones(N-1,1), -1) + diag(ones(N,1), 0) + diag(ones(N-1,1), 1);
3     I = eye(N);
4     A = kron(T, I) + kron(I, T) - kron(I, I);
5 end

```

The index k describes a remapping of a $N \times N$ grid into a N^2 vector, precisely $K_{ij} = (i-1) \cdot N + j$, $1 \leq i, j \leq N$. Once the A matrix is composed and given the \vec{b} vector, then the solution is simply given by `x = A \ b`. The code is shown in the attached script `puzzle.m`.

(3) Helicopter power formula: dimensional analysis

The helicopter problem is a dimension problem, because knowing that the involved quantities are P , g , L , ρ_h and ρ_a , their relation will only depend on the units of measure expression. Thus, if the second helicopter has $1/3$ of the length with respect to the first one, then, taking the formula in the document [?], its power is given by $P_2 = 3^{-\beta} P_1$.

(3.1) Approaching the problem

The same formula cited above can be expressed in a logarithmic form:

$$\ln(P) = \alpha \cdot \ln(g) + \beta \cdot \ln(L) + \gamma \cdot \ln(\rho_h) + \delta \cdot \ln(\rho_a) \quad (10)$$

Assigning for each quantity its corresponding SI unit of measure [?], or rather, $[P] = \text{kg m}^2/\text{s}^3$, $[g] = \text{m}/\text{s}^2$, $[L] = \text{m}$, $[\rho] = \text{kg}/\text{m}^3$, then the logarithm of m, s and kg can be treated as a vector basis. At this point the equation (10) can be rewritten as:

$$\ln(\text{kg}) \cdot (1 - \gamma - \delta) + \ln(\text{m}) \cdot (2 - \alpha - \beta + 3\gamma + 3\delta) + \ln(\text{s}) \cdot (-3 + 2\alpha) = 0 \quad (11)$$

$$\implies \begin{cases} \gamma + \delta = & 1 \\ \alpha + \beta - 3\gamma - 3\delta = & 2 \\ 2\alpha = & 3 \end{cases} \quad (12)$$

This system is indetermined, thus it cannot be computationally solved using the `solve.m` script, because it's matrix representation is not a square matrix.

(3.2) + (3.3), Adding a constraint

In the case where $\alpha = \gamma$, the equation found in the previous point reduces to a determined system of linear equations, which has a square matrix form A .

$$\begin{cases} \alpha + \delta = 1 \\ 2\alpha - \beta + 3\delta = -2 \\ 2\alpha = 3 \end{cases} \implies A = \begin{pmatrix} 1 & 0 & 1 \\ 2 & -1 & 3 \\ 2 & 0 & 0 \end{pmatrix}, \quad \vec{b} = \begin{pmatrix} 1 \\ -2 \\ 3 \end{pmatrix} \quad (13)$$

Now the system is solveable and the solution is straight forward:

$$\alpha = \gamma = \frac{3}{2}, \quad \beta = \frac{7}{2}, \quad \delta = -\frac{1}{2} \quad (14)$$

So, the output power of the second helicopter $P_2 = 3^{-\frac{7}{2}} \cdot P_1 \approx 0.021P_1$.

The eigenvalue problem and diagonalization

Let \hat{A} be an operator defined over an hilbert space \mathcal{H} . By solving an eigenvalue problem is meant to find all vectors (or functions) $x \in \mathcal{H}$ such that there exists a real (or complex) value λ that satisfies the following condition:

$$\hat{A} \cdot x = \lambda \cdot x, \lambda \in \mathbb{K} \quad (15)$$

In the case of this report, the interest is to computationally solve the eigenvalue problem for finite rank operators, which can be expressed as square matrices. So, let N be rank of a square matrix A and $\vec{v} \in \mathbb{K}^N$, then the equation (15) is equivalent to:

$$A \cdot \vec{v} = \lambda \cdot \vec{v}, \lambda \in \mathbb{K} \quad (16)$$

Power method

The power method bases its functioning on the iterative application of a specific operation T . The principle is that every iteration step tends to minimise of the distance between the old evaluated eigen value λ_{k-1}

and the current λ_k . Given the unitary vector $\vec{v}_k \in \mathbb{K}^N$ with $\|\vec{v}_k\| = 1$ at the iteration step k , the corresponding diagonal value, relative to a square matrix A , is given by the hermitian scalar product (see [?] for the notation):

$$\lambda_k = \langle \vec{v}_k, A\vec{v}_k \rangle, \quad \lambda_k \in \mathbb{K} \quad (17)$$

The operation T mentioned above varies depending on the specific method, which of there are three:

- *Power* method: $T = A$.
- *Inverse power* method: $T = (A - 1 \cdot \tau)^{-1}$, $\tau \in \mathbb{K}$ is a fixed eigenvalue target.
- *Rayleigh quotient* method: $T = (A - 1 \cdot \lambda_{k-1})^{-1}$, $\lambda_{k-1} \in \mathbb{K}$ is the old evaluated eigenvalue, as defined in equation (17).

Notice that for the *Rayleigh quotient* method the application is adapting during each iteration, guaranteeing a faster convergence.

Jacobi method

Problem 2

(1) Power iteration methods implementation

As mentioned in the section (), the all methods involve a matrix application, which means that the temporal complexity is at least N^2 , where N is the size of the involved matrix. Additionally, the number of iterations is not simple to estimate and the worst case is potentially infinite. The same could be claimed for the tolerance that the minimized value $|\lambda_k - \lambda_{k-1}|$ should have. In fact the matlab constant `eps` is often too low as tolerance and for larger values of N the loop doesn't exit. The implementation codes can be found in the scripts `eig_power.m`, `eig_ipower.m` and `eig_rq.m`.

(2) Eigenmodes of a vibrating string

(2.1) Formalization of the problem

The description of a vibrating string follows the wave equation (18). Applying a then separation of variables, the problem reveals to be depend on an additional variable $\lambda \in \mathbb{R}$.

$$\frac{\delta^2 u}{\delta x^2} = \frac{\kappa}{\rho} \frac{\delta^2 u}{\delta t^2} \quad u(0, t) = 0, \quad u(L, t) = 0 \quad u(x, y) = \omega(x) \cdot v(t) \quad \Longrightarrow \quad \begin{cases} \omega''(x) + \lambda \omega(x) = 0 \\ v''(t) + \frac{\lambda \rho}{\kappa} v(t) = 0 \end{cases} \quad (18)$$

The sign of λ will now determine the form of the solution. Analysing first the case of $\lambda = 0$, the solution would be a linear equation $\omega(x) = Cx + B$. Unfortunately, applying the boundary conditions $\omega(0) = B = 0$ and $\omega(L) = CL + B = 0$ the implication is $B = 0$ and $C = 0$, thus $\omega(x) = 0$. The same proof is applicable to the case $\lambda < 0$: the solution would be $\omega(x) = B \exp(\gamma x) + C \exp(-\gamma x)$, $\gamma = \sqrt{-\lambda}$, but $\omega(0) = B + C = 0$ and $\omega(L) = B \exp(\gamma L) + C \exp(-\gamma L) = 0$. These conditions imply that $C = -B$, $2B \sinh(\gamma L) = 0$, which means that $\gamma = 0$ or $B = 0$, in both cases the solution turns to be trivial.

It remains the case where $\lambda > 0$, here the solution takes a sinusoidal form:

$$\begin{cases} \omega(x) = B_x \sin(\gamma x) + C_x \cos(\gamma x), & \gamma = \sqrt{\lambda} \\ v(t) = B_t \sin(2\pi \nu t) + C_t \cos(2\pi \nu t), & \nu = \frac{\sqrt{\lambda \rho}}{2\pi} \end{cases} \quad (19)$$

This solution is not ill formed because $\omega(0) = C_x = 0$, then by the second boundary condition $\omega(L) = B_x \sin(\gamma L) = 0$. Imposing $B_x \neq 0$, the resulting condition is $\sin(\gamma L) = 0 \iff \gamma_n L = \pi n$, $n \in \mathbb{N}^*$. In other words, this means that for each natural number n , there exists a λ_n given by:

$$\lambda_n = \left(\frac{\pi n}{L} \right)^2 \quad (20)$$

Thus, the implicit condition to obtain a non-trivial solution is given by equation (20). Notice, by its expression, that dimensionally λ is $1/\text{m}^2$ in the SI unit system [?].

(2.2) Implementation

Using a finite difference discretization, or rather given a step Δx and setting $x_i = (i - 1) * \Delta x$, the problem is expressed as follow:

$$\frac{-\omega_{i-1} + 2\omega_i - \omega_{i+1}}{\Delta x^2} = \lambda \omega_i \quad \iff \quad A \vec{\omega} = \lambda \vec{\omega}, \quad \vec{\omega} = (\omega_1, \dots, \omega_N) \quad (21)$$

A in equation (21) is the tridiagonal matrix containing -1 in the lower and the upper diagonal and 2 on the diagonal, all the terms divided by the squared step Δx^2 . This matrix is symmetric, then by the spectral theorem [?], it's eigenvalues are real and for each pair of eigenvalues, the associated eigenspaces are orthogonal each other. Furthermore, the matrix is positive definite, thus all the eigenvalues are

strictly positive. Notice that this setup automatically meets the boundary conditions, because the first and the last line of A already discards the boundary terms (setting them implicitly to zero).

(2.3) + (2.4), Find the first four eigenvalues

The first thing to notice is that the matrix A is not degenerate (by its positivity), then each eigenspace's dimension is 1. This clearly means that the first four eigenvalues are different each other. Theoretically, by the expression in equation (20), the eigenvalues for $L = 1$ are $\lambda_1 = \pi^2 \approx 9.8696$, $\lambda_2 = 4\pi^2 \approx 39.4784$, $\lambda_3 = 9\pi^2 \approx 88.8264$ and $\lambda_4 = 16\pi^2 \approx 157.9137$. In order to find the first one, it's enough to use the *inverse power* method setting the target to 0. The other three can be found by using a *range* strategy: let

(3) Jacobi method implementation

(4) Landau levels in a square-lattice model

(a) Grey-scaled image file `stm.png`

(b) Fourier transformed image file `stm.png`

Conclusion

Appendix: matlab codes

Listing 3: Algorithm which solves a system

```

1 function x = solve(A, b)
2     % decompose LU
3     [L, U, P] = lu_decomposition(A);
4
5     y = solve_lower(L, P * b);
6     x = solve_upper(U, y);
7 end
8
9 function y = solve_lower(L, b)
10    [~,N] = size(L);
11    y = b;
12    % start by the top
13    for k=1:N-1 % iteration on columns
14        y(k) = y(k) / L(k,k);
15        for i=k+1:N % iteration on rows
16            y(i) = y(i) - y(k) * L(i,k);
17        end
18    end
19    y(N) = y(N) / L(N,N);
20 end
21
22 function x = solve_upper(U, y)
23    [~,N] = size(U);
24    x = y;
25    % start by the top
26    for k=N:-1:2 % iteration on columns
27        x(k) = x(k) / U(k,k);
28        for i=k-1:-1:1 % iteration on rows
29            x(i) = x(i) - x(k) * U(i,k);
30        end
31    end
32    x(1) = x(1) / U(1,1);
33 end

```

Listing 4: Power iteration method implementation

```

1 function [ vec, val ] = eig_power(inputmatrix)
2     [Ni, Nj] = size(inputmatrix);
3     assert(Ni == Nj || Ni == 0, 'Matrix must be square or non-zero');
4     vec = zeros(Ni, 1);
5     vec(1) = 1;
6     val = ctranspose(vec) * inputmatrix * vec;
7     oldval = 0;
8
9     while abs(oldval - val) > eps
10        vec = inputmatrix * vec;
11        vec = vec / norm(vec);
12        oldval = val;
13        val = ctranspose(vec) * inputmatrix * vec;
14    end

```

```
15 end
```

Listing 5: Inverse power iteration method implementation

```
1 function [ vec, val ] = eig_ipower(inputmatrix, target)
2     [Ni, Nj] = size(inputmatrix);
3     assert(Ni == Nj || Ni == 0, 'Matrix must be square or non-zero')
4     vec = rand(Ni,1);
5     vec = vec / norm(vec);
6     val = ctranspose(vec) * inputmatrix * vec;
7     oldval = target;
8     I_target = eye(Ni) * target;
9
10    while abs(oldval - val) > eps
11        vec = (inputmatrix - I_target) \ vec;
12        vec = vec / norm(vec);
13        oldval = val;
14        val = ctranspose(vec) * inputmatrix * vec;
15    end
16
17    %val = val + target;
18 end
```

Listing 6: Rayleigh quotient iteration method implementation

```
1 function [ vec, val ] = eig_rq(inputmatrix, target)
2     [Ni, Nj] = size(inputmatrix);
3     assert(Ni == Nj || Ni == 0, 'Matrix must be square or non-zero')
4     vec = zeros(Ni, 1);
5     vec(1) = 1;
6
7     val = ctranspose(vec) * inputmatrix * vec;
8     oldval = target;
9     I = eye(Ni);
10
11    while abs(val - oldval) > 1e-10
12        vec = (inputmatrix - I * oldval) \ vec;
13        vec = vec / norm(vec);
14        oldval = val;
15        val = ctranspose(vec) * inputmatrix * vec;
16    end
17    %val = newval - target;
18 end
```

Documentation and sources

[1] https://en.wikipedia.org/wiki/SI_base_unit

[2] <https://math.stackexchange.com/questions/1274373/proof-for-existence-of-lu-decomposition>

[3] https://en.wikipedia.org/wiki/Inner_product_space

[4] https://en.wikipedia.org/wiki/Spectral_theorem