

# Project 1

Antonio Gómez Garrido and David Martínez Hernández

(Dated: September 12, 2022)

*List a link to your github repository here!*

## PROBLEM 1

First, we need to check if the given function satisfies the boundary conditions, which are:  $u(0) = u(1) = 0$

In our case,

$$u(x) = 1 - (1 - e^{-10})x - e^{-10x} \quad (1)$$

So, when computing the boundary conditions, we get:

*For the case of  $u(0)$  :*

$$u(0) = 1 - (1 - e^{-10})0 - e^{-10 \cdot 0}$$

$$u(0) = 1 - e^0$$

$$u(0) = 1 - 1$$

$$u(0) = 0$$

*For the case of  $u(1)$  :*

$$u(1) = 1 - (1 - e^{-10})1 - e^{-10 \cdot 1}$$

$$u(1) = 1 - 1 + e^{-10} - e^{-10}$$

$$u(1) = 0$$

Now that we have verified that the function given in this section satisfies the boundary conditions, we must prove that it also satisfies the Poisson equation.

$$\begin{aligned} -\frac{d^2}{dx^2}(1 - (1 - e^{-10})x - e^{-10x}) &= -\frac{d}{dx}(-1 + e^{-10} + 10e^{-10x}) \\ -\frac{d^2u}{dx^2} &= -(-100e^{-10x}) \\ -\frac{d^2u}{dx^2} &= 100e^{-10x} \end{aligned}$$

## PROBLEM 2

**2.a)**

The values of  $x$  and  $u(x)$  are all stored inside the file 'valores.txt'

**2.b)**

The plot obtain in python from the previously mentioned data is the following:

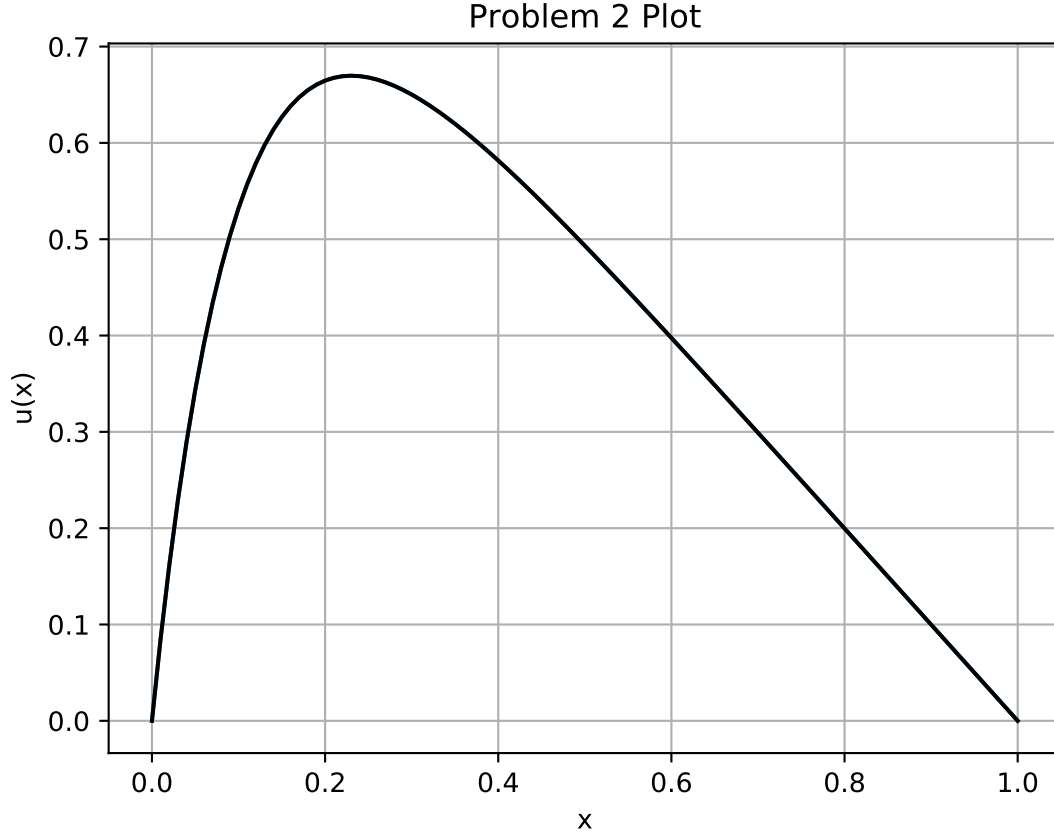


Figure 1.  $u(x)$  against  $x$

### PROBLEM 3

If we Taylor expand our function  $u(x)$  and evaluate it at points  $x_o = x+h$  and  $x_1 = x-h$ , we obtain the expressions:

$$u(x+h) = u(x) + u'(x) \cdot h + \frac{1}{2!}u''(x) \cdot h^2 + \frac{1}{3!}u'''(x) \cdot h^3 + \mathcal{O}(h^4)$$

$$u(x-h) = u(x) - u'(x) \cdot h + \frac{1}{2!}u''(x) \cdot h^2 - \frac{1}{3!}u'''(x) \cdot h^3 + \mathcal{O}(h^4)$$

By summing these two expressions and rearranging the terms, we get the well known formula for the second derivative of a function:

$$\frac{d^2u}{dx^2} = -\frac{u(x+h) - 2u(x) + u(x-h)}{h^2} + \mathcal{O}(h^2) \quad (2)$$

Once we get to this, the next step is to simplify the equation by setting  $h$  to be arbitrarily small and thus disregarding the last term of the sum.

$$\frac{d^2v}{dx^2} = -\frac{v(x+h) - 2v(x) + v(x-h)}{h^2} \quad (3)$$

Here, we have redefined our function to take into account the fact that this is an approximation.

Now we insert this expression in the Poisson formula and discretize the domain of the function. This is, instead of taking a continuous set of values of  $x$ , we choose a finite set of points of this domain and evaluate the function in such points.

Finally, the problem we're trying to solve acquires the form:

$$\frac{v_{i+1} - 2v_i + v_{i-1}}{h^2} = f_i \quad (4)$$

Here, the subindex  $i$  indicates the point of the previously mentioned set of points that we're taking into account.

#### PROBLEM 4

First of all, we need to rearrange the terms in equation 4 to make them of the form:  $-v_{i+1} + 2v_i - v_{i-1} = h^2 f_i$ . Suppose that, in our discretization, we have chosen a set of  $N$  points from the function's domain (without counting the extremes). We, as a result, have  $N$  equations:

$$\begin{array}{lll} -v_0 + 2v_1 - v_2 & \dots & = h^2 f_1 \\ -v_1 + 2v_2 - v_3 & \dots & = h^2 f_2 \\ -v_2 + 2v_3 - v_4 & \dots & = h^2 f_3 \\ \dots & \dots & \dots \\ \dots & -v_{N-1} + 2v_N - v_{N+1} & = h^2 f_N \end{array}$$

But, as both  $v_0$  and  $v_{N+1}$  are already known (they are the boundary conditions in our problem), we can rearrange the terms in the first and last equations in order to have all of the unknowns on one side and all of the given data on the other one. If we now redefine the right sides of the equations (which we will choose to be the one with the problem's data) as follows:  $h^2 f_1 + v_0 \equiv g_1$ ,  $h^2 f_i \equiv g_i \ \forall i = 2, 3, \dots, N-1$ ,  $h^2 f_N + v_{N+1} \equiv g_N$  we get the system of equations:

$$2v_1 - v_2 \quad \dots \quad = h^2 g_1 \quad (5)$$

$$-v_1 + 2v_2 - v_3 \quad \dots \quad = h^2 g_2 \quad (6)$$

$$-v_2 + 2v_3 - v_4 \quad \dots \quad = h^2 g_3 \quad (7)$$

$$\dots \quad \dots \quad \dots \quad (8)$$

$$\dots \quad -v_{N-1} + 2v_N \quad = h^2 g_N \quad (9)$$

Which, if we take into account the rules for the multiplication of matrices, can be written as:

$$A\vec{v} = \vec{g} \quad (10)$$

Where the vector  $\vec{v} = (v_1, v_2, \dots, v_N)$  contains the values of our function  $v(x)$  (which is the unknown of the problem),  $\vec{g} = (g_1, g_2, \dots, g_N)$  contains all of the  $g_i$  defined above and  $A$  represents a tridiagonal matrix with all of its diagonal elements equal to two and all of the elements on its subdiagonal and superdiagonal equal to -1.

#### PROBLEM 5

From the reasoning developed in the previous section, we can easily see that, when setting the equations 5 to 9, two of the elements of  $\vec{v}^*$  ( $v_0^*$  and  $v_{N+1}^*$  to be more precise) are missing. Therefore, the relation between  $m$  and  $n$  is that  $m = n + 2$ . Apart from this, as we've stated previously, the part of the solution that we get from solving equation 10 is the values of the function in every point of its domain except for the extremes.

## PROBLEM 6

### 6.a)

This algorithm is based on the Gaussian elimination method which, in a nutshell, consists on performing a series of row-operations on the matrix of the system ( $A$ ) in order to turn it into a triangular matrix; this is, a matrix for which every element below the diagonal is equal to 0.

Take, for example, the first two rows of our  $A$  matrix (let  $R_1$  be the first, and  $R_2$  the second one). In order to remove the element  $a_{21}$ , which is right below the first element of the diagonal, we must do the operation  $\tilde{R}_2 = R_2 - \frac{a_{21}}{b_1} R_1$ . After that, we do the same for all of the consecutive rows to eliminate every element below the diagonal.

After doing all of these operations, the new elements of the matrix's diagonal and  $\vec{g}$  look like this:

$$\tilde{b}_1 \equiv b_1 \qquad \tilde{g}_1 \equiv g_1 \qquad (11)$$

$$\tilde{b}_i \equiv b_i - \frac{a_i}{\tilde{b}_{i-1}} c_{i-1} \qquad \tilde{g}_i \equiv g_i - \frac{a_i}{\tilde{b}_{i-1}} \tilde{g}_{i-1} \qquad (12)$$

Consequently, the last one of the system's equations acquires the form:

$$\tilde{b}_N \cdot v_N = \tilde{g}_N$$

From which we can obtain the value for  $v_N$ . Once we have this value, we can insert it in the previous equation in order to know  $v_{N-1}$  and so on.

$$v_N = \frac{\tilde{g}_N}{\tilde{b}_N} \qquad (13)$$

$$v_i = \frac{\tilde{g}_i - c_i \cdot v_{i+1}}{\tilde{b}_i} \qquad (14)$$

Once we take all of this into account, the algorithm consists on the following steps:

1. Store the elements of the matrix's subdiagonal, main diagonal and superdiagonal in three different vectors of the type `std::vector<double>`.
2. Create three empty vectors of the same size and type of the matrix's diagonal (one for storing the values of  $\vec{b}$ , another one for the values of  $\vec{g}$ , and one last vector for storing the solution to the problem  $\vec{v}$ )
3. Run a loop for computing the values of  $\tilde{b}_i$  and  $\tilde{g}_i$  given by the expressions in eq: 12 and storing them in the previously mentioned vectors.
4. Run another loop for doing this same operation for the case of the  $v_i$  values which are now given by eq: 14. This can't be done in the same loop of the other two vectors because the values of  $\tilde{b}_i$  and  $\tilde{g}_i$  are needed in order to compute the values of the different  $v_i$ .

In a more technical way, our algorithm would look something like this:

**Algorithm 1** General Algorithm

---

```

std :: vector < double > a = a1, a2, ..., aN - 1; ;
std :: vector < double > b = b1, b2, ..., bN;
std :: vector < double > c = c1, c2, ..., cN - 1
std :: vector < double > g = g1, g2, ..., gN
std :: vector < double > b_tilde(N);
std :: vector < double > g_tilde(N);
std :: vector < double > v(N);
if a.size() != c.size() or b.size() - a.size() != 1 or b.size() != g.size() then
    std :: cout << "the equality is not possible";
    return 0
else
    b_tilde[0] = b[0]; g_tilde[0] = g[0]; v[0] = 0;
    for i=2,3,...,N do
        b_tilde[i] = b[i] -  $\frac{a[i-1]}{b_tilde[i-1]} \cdot c[i-1]$ 
        g_tilde[i] = g[i] -  $\frac{a[i-1]}{b_tilde[i-1]} \cdot g_tilde[i-1]$ 
    v[N] =  $\frac{g_tilde[N]}{b_tilde[N]}$ 
    for i=1,2,...,N-1 do
        int j = N - i
        v[j] = (g_tilde[j] - c[j] * v[j+1]) / (b_tilde[j])

```

---

▷ This is the definition  
 ▷ of the subdiagonal, main diagonal and  
 ▷ superdiagonal vectors respectively  
 ▷ We also define the g vector  
 ▷ Now define empty vectors  
 ▷ where we will store the values of  
 ▷  $\tilde{b}_i, \tilde{g}_i$  and  $v_i$

▷ we run this if to make sure that  
 ▷ the problem's data have been introduced correctly

▷ in our c++ code the indexes will be  
 ▷ different because of the nature of std::vectors

▷ This is the way we've come up with to run the loop "backwards"

**6.b)**

To count the number of FLOPs in this algorithm we need to take a deeper look into equations 11, 12, 13 and 14

Following the procedure used in class, every FLOP will be preceded by  $fl()$  to make it easier to count them. Now, if we look at eq. 12 (we will only be examining the expressions for the values of  $\tilde{b}$  given that they are identical to the ones of  $\tilde{g}$ ):

$$fl \left\{ fl(b_i) - fl \left[ fl(c_{i-1}) \cdot fl \left( \frac{fl(a_{i-1})}{fl(\tilde{b}_{i-1})} \right) \right] \right\} \quad (15)$$

We can count up to seven FLOPs for every one of these operations. If we now take into account the fact that this operation is repeated twice (one time for  $\tilde{b}$  and another one for  $\tilde{g}$ ) and that, on top of that, these operations are repeated  $N - 1$  times in the for loop, we get  $14(N - 1)$  FLOPs. We must now add the 2 FLOPs present in the definitions of  $\tilde{b}_1$  and  $\tilde{g}_1$ .

Now for the case of eq. 14:

$$fl \left\{ \frac{fl[fl(g_i) - fl(fl(c_i) \cdot fl(v_{i+1}))]}{fl(b_i)} \right\} \quad (16)$$

Where we count 7 more FLOPs for each one of the  $N - 1$  iterations in the loop.

Finally, when calculating  $v_N$  we have 3 more FLOPs:

$$fl \left( \frac{fl(\tilde{g}_N)}{fl(\tilde{b}_N)} \right) \quad (17)$$

In conclusion, the total ammount of FLOPs performed in or algorithm is:  $FLOPs = 21(N - 1) + 5$ .

**PROBLEM 7****7.b)**

From the program in section 7.a we have obtained the following solutions:

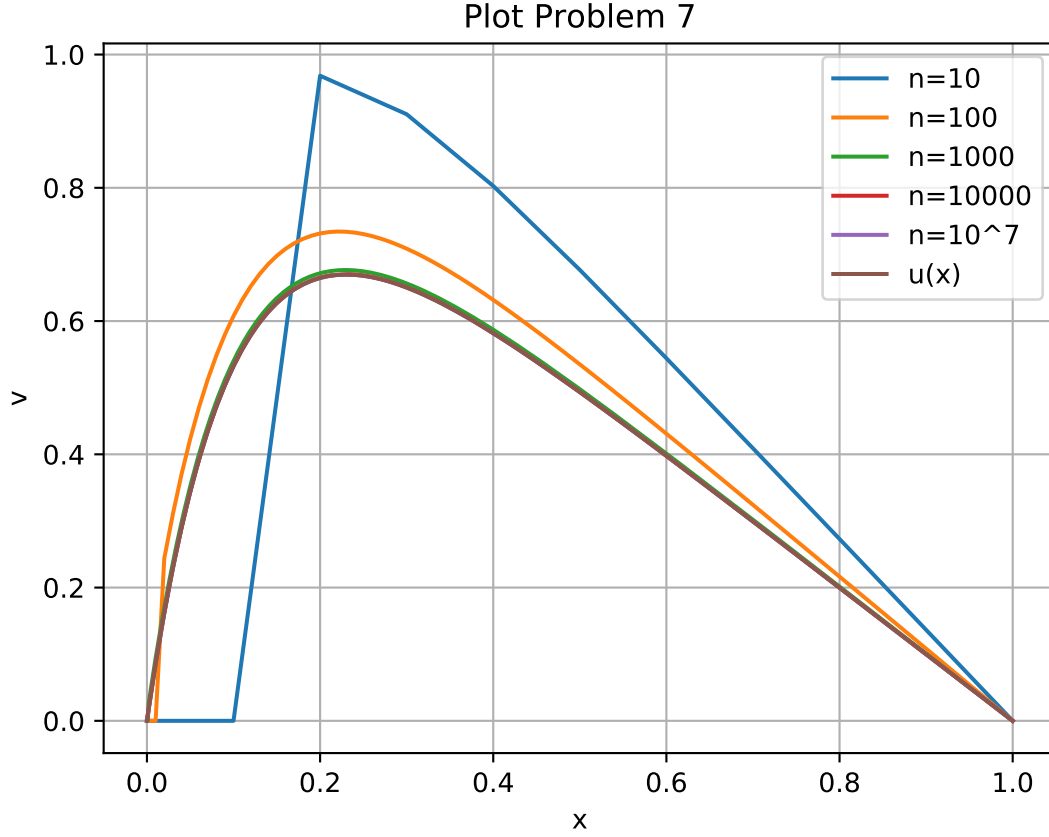


Figure 2. The different colors represent the solutions obtained from different approximations (different values of  $n_{steps}$ )

As we can see, when we increase the number of  $n_{steps}$ , the approximations become closer to the analytic solution.

## PROBLEM 8

## PROBLEM 9

### 9.a)

The theoretical background for this problem is the same as the one for the case discussed in problem 6. This is because the case we're about to treat is just a concrete example of the more general problem discussed in section 6.

Contrary to what was set in such section, now the values for the elements of  $\vec{a}, \vec{b}$  and  $\vec{c}$ , are not only constant, but they are also equal ( $\vec{a} = \vec{c} = (-1, -1, \dots, -1)$ ;  $\vec{b} = (2, 2, \dots, 2)$ )

As a result of this, equations 12 and 14 acquire an easier form, as now the values of  $a_i, b_i$  and  $c_i$  can be replaced by constants. Such form is the following one:

$$\tilde{b}_1 \equiv 2 \qquad \tilde{g}_1 \equiv g_1 \qquad (18)$$

$$\tilde{b}_i \equiv 2 - \frac{-1}{\tilde{b}_{i-1}}(-1) \qquad \tilde{g}_i \equiv g_i - \frac{-1}{\tilde{b}_{i-1}}\tilde{g}_{i-1} \qquad (19)$$

$$\tilde{b}_i \equiv 2 - \frac{1}{\tilde{b}_{i-1}} \qquad \tilde{g}_i \equiv g_i + \frac{1}{\tilde{b}_{i-1}}\tilde{g}_{i-1} \qquad (20)$$

$$v_N = \frac{\tilde{g}_N}{\tilde{b}_N} \quad (21)$$

$$v_i = \frac{\tilde{g}_i - (-1) \cdot v_{i+1}}{\tilde{b}_i} \quad (22)$$

$$v_i = \frac{\tilde{g}_i + v_{i+1}}{\tilde{b}_i} \quad (23)$$

Apart from this, the fact that the three main diagonals of the matrix are constant imply that, for this new algorithm, there is no need to define the three vectors (their values are already reflected in the expressions 20 and 23). Therefore, this new algorithm looks like this:

---

**Algorithm 2** Special Algorithm

---

```

std :: vector<double> g = g1, g2..., gN
std :: vector< double > b_tilde(N);
std :: vector< double > g_tilde(N);
std :: vector< double > v(N);
if a.size() ≠ c.size() or b.size() - a.size() ≠ 1 or b.size() ≠ g.size() then
    std :: cout << "the equality is not possible";
    return 0
else
    b_tilde[0] = 2; g_tilde[0] = g[0]; v[0] = 0;
    for i=2,3,...,N do
        b_tilde[i] = 2 -  $\frac{1}{b_{tilde}[i-1]}$ 
        g_tilde[i] = g[i] +  $\frac{1}{b_{tilde}[i-1]} \cdot g_{tilde}[i-1]$ 
        v[N] =  $\frac{g_{tilde}[N]}{b_{tilde}[N]}$ 
        for i=1,2,...,N-1 do
            int j = N - i
            v[j] =  $\frac{g_{tilde}[j] + v[j+1]}{(b_{tilde}[j])}$ 

```

▷ We replace the a,b and c variables with constants

---

**9.b)**

In order to count the number of FLOPs in this algorithm, we will follow the same method as in section 6.b) but applied to equations 20 and 23.

$$\begin{aligned}
 & fl \left\{ fl(2.0) - fl \left[ \frac{fl(1.0)}{fl(\tilde{b}_{i-1})} \right] \right\} \\
 & fl \left\{ fl(g_i) + fl \left[ fl(\tilde{g}_{i-1}) \cdot fl \left( \frac{fl(1.0)}{fl(\tilde{b}_{i-1})} \right) \right] \right\} \\
 & fl \left\{ \frac{fl[fl(\tilde{g}_i) + fl(v_{i+1})]}{fl(\tilde{b}_i)} \right\}
 \end{aligned}$$

Between the three of these equations there are 17 FLOPs which are repeated  $N - 1$  times inside the for loops. Therefore, the total number of FLOPs in this algorithm is:  $FLOPs = 17(N - 1) + 5$