

# FYS4150 – Project 1

Joachim Falck Brodin, Fredrik Jaibeer Mahal Nordeng  
and Endrias Getachew Asgedom

September 7, 2020

## **Abstract**

Numerical methods for solving the one-dimensional Poisson equation with Dirichlet boundary conditions are investigated. In this report, we show how to discretize a one-dimensional Poisson equation and transform it into a linear tridiagonal set of equations. To solve these equations numerically, three different algorithms are developed (i.e., general tridiagonal Gaussian elimination, special tridiagonal Gaussian elimination, and LU-decomposition) and the results are analysed in terms of their computational cost, memory usage, and numerical error. Comparing our numerical solution with the closed form solution, we observe that the numerical approximation error is proportional to the square of the discretization step size, until we experience the loss of precision, due to round off error, when the discretization step size is less than  $\sim 10^{-5}$ . Moreover, the specialized algorithm outperforms the two other algorithms in terms of computational cost and memory usage.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Theory</b>	<b>3</b>
2.1	Gaussian Elimination for a Tridiagonal Matrix . . . . .	4
2.2	Gaussian Elimination for a Toepliz Tridiagonal Matrix . . . . .	6
2.3	LU-decomposition Method . . . . .	6
<b>3</b>	<b>Method</b>	<b>6</b>
3.1	General Tridiagonal Matrix Algorithm . . . . .	7
3.2	Toepliz Tridiagonal Matrix Algorithm . . . . .	7
3.3	LU-decomposition based Toepliz Tridiagonal Matrix Algorithm .	8
3.4	Error Analysis . . . . .	9
<b>4</b>	<b>Results</b>	<b>9</b>
4.1	Computational Time . . . . .	9
4.2	Algorithm Benchmark and Error Analysis . . . . .	10
<b>5</b>	<b>Discussion</b>	<b>14</b>
<b>6</b>	<b>Conclusion</b>	<b>14</b>
<b>A</b>	<b>Analytical Solution</b>	<b>14</b>

GitHub repository at <https://github.com/endrias34/FYS4150>

# 1 Introduction

Problems in the field of computational sciences is often formulated in terms of differential equations. Unfortunately, most of these differential equations do not have a closed form solution and hence needs to be solved using numerical methods. In this project, we develop algorithms for a numerical method that solves the one-dimensional Poisson equation with Dirichlet boundary conditions. We use central Finite Difference (FD) method to approximate the second-order differentiation and consequently cast the Poisson equation into a tridiagonal matrix equations. This linear system of equations may be solved using a Gaussian elimination method customized for tridiagonal matrices, also known as Thomas algorithm. Two versions of the Thomas algorithm are considered here; the first assumes a general tridiagonal matrix while the second assumes a Toeplitz tridiagonal matrix. As a benchmark, we have also used a more general LU-decomposition method. The accuracy, computational efficiency, and memory usage of the different algorithms are quantitatively analysed.

This report is structured in the following way: First we will briefly present how to discretize the one-dimensional Poisson equation and generate the tridiagonal matrix equations. Second, we will show how to formulate the three different algorithms used to solve the problem. We then analyse the performance of the different algorithms in terms of the number of floating point operations (FLOPs), CPU time, and numerical approximation and precision errors. Finally, we will make a concluding remark.

# 2 Theory

To discretize the one-dimensional Poisson equation and write it in a matrix-vector form, we start with the continuous equation and its corresponding boundary conditions

$$-u''(x) = f(x), \quad x \in (0, 1), \quad u(0) = u(1) = 0, \quad (1)$$

where  $f(x)$  is a known function, called the *source term* and  $u(x)$  is a dimensionless physical quantity defined in a dimension-less domain  $x$ . To solve equation (1) numerically, we first need an approximation of the second derivative of  $u$ . This can be achieved by performing the Taylor expansion of  $u(x)$  around  $x$

$$u(x \pm h) = u(x) \pm hu'(x) + \frac{h^2u''(x)}{2!} \pm \frac{h^3u'''(x)}{3!} + \mathcal{O}(h^4). \quad (2)$$

To find an expression containing only the second derivative, we perform the following operation

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

We can now solve for  $u''$  as

$$u''(x) = \frac{u(x+h) + u(x-h) - 2u(x)}{h^2} + \mathcal{O}(h^2). \quad (4)$$

Notice that the approximation error in equation (4) is  $\mathcal{O}(h^2)$  since we divided (3) by  $h^2$ .

For a discrete domain  $x_i = ih$  with  $x_0 = 0$  and  $x_{n+1} = 1$  as the boundary points, our unknown physical quantity  $u$  is represented as  $v_i$ . Here, the step size is given by  $h = 1/(n+1)$ . The boundary conditions can now be written as  $v_0 = v_{n+1} = 0$ . Finally, the discrete form of equation (1) is thus given by

$$-\frac{v_{i+1} + v_{i-1} - 2v_i}{h^2} = f_i \quad \text{for } i = 1, \dots, n, \quad (5)$$

where  $f_i = f(x_i)$  is the discrete version of  $f(x)$ . To represent equation (5) in a matrix-vector form, we multiply equation (5) with  $h^2$  on both sides and write out the equations for all  $i$  as follows

$$\begin{aligned} 2 * v_1 - 1 * v_2 &= f_1 * h^2 \\ -1 * v_1 + 2 * v_2 - 1 * v_3 &= f_2 * h^2 \\ -1 * v_2 + 2 * v_3 - 1 * v_4 &= f_3 * h^2 \\ &\vdots \\ -1 * v_{n-2} + 2 * v_{n-1} - 1 * v_n &= f_{n-1} * h^2 \\ -1 * v_{n-1} + 2 * v_n &= f_n * h^2. \end{aligned} \quad (6)$$

From equation (6) we can clearly see how equation (5) can be represented in a matrix-vector form as

$$\mathbf{A}\mathbf{v} = \tilde{\mathbf{f}}, \quad (7)$$

where  $\tilde{f}_i = f_i * h^2$ ,

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 & \dots & \dots & 0 \\ -1 & 2 & -1 & 0 & \dots & \dots \\ 0 & -1 & 2 & -1 & 0 & \dots \\ & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & & -1 & 2 & -1 \\ 0 & \dots & & 0 & -1 & 2 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \dots \\ v_{n-1} \\ v_n \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{f}} = \begin{bmatrix} \tilde{f}_1 \\ \tilde{f}_2 \\ \tilde{f}_3 \\ \dots \\ \tilde{f}_{n-1} \\ \tilde{f}_n \end{bmatrix}.$$

Here we can see how each row in the matrix-vector equation corresponds to the sequence of terms as we wrote it out in equation (6). We can also observe that  $\mathbf{A}$  is tridiagonal matrix.

## 2.1 Gaussian Elimination for a Tridiagonal Matrix

Gaussian elimination, or row reduction, is a method for solving a system of linear equations. In the form presented here the procedure is also known as

the *Thomas algorithm*. We rewrite our matrix  $\mathbf{A}$  in terms of one-dimensional vectors  $a, b, c$  and the linear equations can then be presented on the form

$$\begin{bmatrix} b_1 & c_1 & 0 & \dots & \dots & 0 \\ a_2 & b_2 & c_2 & 0 & & \vdots \\ 0 & a_3 & b_3 & c_3 & 0 & \\ \vdots & 0 & a_4 & b_4 & c_4 & 0 \\ & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & a_n & b_n \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \tilde{f}_1 \\ \tilde{f}_2 \\ \tilde{f}_3 \\ \tilde{f}_4 \\ \vdots \\ \tilde{f}_n \end{bmatrix}. \quad (8)$$

Next, we eliminate the  $a_2$  from the second row. Here, the first step is to subtract  $\frac{a_2}{b_1}$  times the first row from the second row and obtain

$$\begin{bmatrix} b_1 & c_1 & 0 & \dots & \dots & 0 \\ 0 & b'_2 & c_2 & 0 & & \vdots \\ 0 & a_3 & b_3 & c_3 & 0 & \\ \vdots & 0 & a_4 & b_4 & c_4 & 0 \\ & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & a_n & b_n \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \tilde{f}_1 \\ \tilde{f}'_2 \\ \tilde{f}_3 \\ \tilde{f}_4 \\ \vdots \\ \tilde{f}_n \end{bmatrix}, \quad (9)$$

where  $b'_2 = b_2 - \frac{a_2}{b_1}c_1$  and  $\tilde{f}'_2 = \tilde{f}_2 - \frac{a_2}{b_1}\tilde{f}_1$ . For the following rows we proceed in the same manner and we can generally express the different rows as

$$b'_i = b_i - \frac{a_i}{b'_{i-1}}c_{i-1}, \quad (10)$$

with  $b'_1 = b_1$ , and

$$\tilde{f}'_i = \tilde{f}_i - \frac{a_i}{b'_{i-1}}\tilde{f}'_{i-1}, \quad (11)$$

where  $\tilde{f}'_1 = \tilde{f}_1$ . Finally, the system of equations become an upper triangular matrix equation of the form

$$\begin{bmatrix} b'_1 & c_1 & 0 & \dots & \dots & 0 \\ 0 & b'_2 & c_2 & 0 & & \vdots \\ 0 & 0 & b'_3 & c_3 & 0 & \\ \vdots & 0 & 0 & b'_4 & c_4 & 0 \\ & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & 0 & b'_n \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \tilde{f}'_1 \\ \tilde{f}'_2 \\ \tilde{f}'_3 \\ \tilde{f}'_4 \\ \vdots \\ \tilde{f}'_n \end{bmatrix}. \quad (12)$$

The process of converting the matrix equation from (8) to (12) is known as *forward substitution*. To find the solution  $v_i$  for the matrix equation (12), we can start from the bottom and solve for  $v_n$ , then continue to solve for  $v_{n-1}$  until we reach  $v_0$ . This way of solving the matrix equation is known as *backward substitution* and it can be generalized to

$$v_i = \begin{cases} \tilde{f}'_n/b'_n, & \text{if } i = n \\ (\tilde{f}'_i - c_i v_{i+1})/b'_i, & \text{otherwise.} \end{cases} \quad (13)$$

## 2.2 Gaussian Elimination for a Toeplitz Tridiagonal Matrix

Matrix  $\mathbf{A}$  in equation (7) has a Toeplitz tridiagonal structure. This means the non-zero diagonal elements have the same value along all the diagonal. In our case, we have  $a_i = c_i = -1$  and  $b_i = 2$ . The fact that the elements of matrix  $\mathbf{A}$  only have two different values greatly simplifies the Gaussian elimination method discussed in the previous section. The forward substitution step in equations (10) and (11) can now be generalized to

$$b'_i = \frac{i+1}{i}, \quad i = 1, \dots, n-1, \quad (14)$$

and

$$\tilde{f}'_i = \tilde{f}_i + \frac{\tilde{f}'_{i-1}}{b'_{i-1}}, \quad i = 1, \dots, n-1. \quad (15)$$

Similarly, the backward substitution step in equation (13) simplifies to

$$v_i = \frac{\tilde{f}'_i + v_{i+1}}{b'_i} \quad i = n-2, n-1, \dots, 0. \quad (16)$$

## 2.3 LU-decomposition Method

The LU-Decomposition method is the most commonly used method to solve a linear system of equation which has a densely populated matrix. The method is a special form of Gaussian elimination method and considers we have a non-singular matrix  $\mathbf{A}$  that can be decomposed into a product of two triangular matrices  $\mathbf{L}$  and  $\mathbf{U}$ , where  $\mathbf{L}$  is a lower-triangular matrix with all diagonal entries equal to 1 and  $\mathbf{U}$  is an upper-triangular matrix. To solve the linear set of equation  $\mathbf{A}\mathbf{v} = \tilde{\mathbf{f}}$  using LU-decomposition, we substitute  $\mathbf{A} = \mathbf{LU}$  such that

$$\mathbf{LUv} = \mathbf{Ly} = \tilde{\mathbf{f}}, \quad (17)$$

where  $\mathbf{y} = \mathbf{Uv}$ . Consequently, the solution  $\mathbf{v}$  can now be obtained by solving the following two triangular linear set of equations

$$\begin{aligned} \mathbf{Ly} &= \tilde{\mathbf{f}} \\ \mathbf{Uv} &= \mathbf{y}. \end{aligned} \quad (18)$$

Solving triangular set of equations is trivial, as the lower triangular equation can be solved using *forward substitution* while the upper triangular equation may be solved by *backward substitution*.

## 3 Method

The three different algorithms required to solve the tridiagonal matrix equation are implemented using C++ and the corresponding source codes can be found in the github repository address linked to this report. In this section, we outline the algorithms structure using pseudo-codes, discuss about the Floating point Operations (FLOPs) counts, and the memory handling of each of the algorithms.

### 3.1 General Tridiagonal Matrix Algorithm

The algorithm we will present here is known as the Thomas algorithm. The main elements of the algorithm are the *forward* and *backward* substitution (see Algorithm (1)). The algorithm takes the vectors  $\mathbf{a} \in \mathbb{R}^{n-1}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{c} \in \mathbb{R}^{n-1}$ , and  $\tilde{\mathbf{f}} \in \mathbb{R}^n$  as input and outputs the solution  $\mathbf{v} \in \mathbb{R}^n$ . Notice that we do not reserve storage for the full  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrix, but rather only store the nonzero tridiagonal elements as vectors. This allows us to utilize this algorithm for solving problems with a very large number of  $n$  without running into a memory allocation problem.

---

#### Algorithm 1 General Tridiagonal Matrix Algorithm

---

**Input:**  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \tilde{\mathbf{f}}$

**Output:**  $\mathbf{v}$

```

1:  $b'_0 = b_0$            // First element along main diagonal
2:  $\tilde{f}'_0 = \tilde{f}_0$        // First element in the RHS
3:
4: // Forward substitution
5: for  $i = 1$  to  $i = n - 1$  do
6:    $b'_i = b_i - a_{i-1} * c_{i-1} / b'_{i-1}$    // Eliminate lower diag
7:    $\tilde{f}'_i = \tilde{f}_i - a_{i-1} * \tilde{f}_{i-1} / b'_{i-1}$  // Change RHS
8: end for
9:
10: // Backward substitution
11:  $v_{n-1} = \tilde{f}'_{n-1} / b'_{n-1}$            // Final element of the solution
12:
13: for  $i = n - 2$  to  $i = 0$  do
14:    $v_i = (\tilde{f}'_i - c_i * v_{i+1}) / b'_i$ 
15: end for
```

---

The number of FLOPs for the Thomas algorithm is composed of  $6(n-1)$  for the forward and  $3(n-1)$  for the backward substitution. Moreover, pre-computing  $a_{i-1}/b'_{i-1}$  (cf. line 6 and 7 in Algorithm 1) in the forward substitution, we can reduce the number of FLOPs by  $(n-1)$ . Therefore, in total the Thomas algorithm requires  $8(n-1)$  FLOPs.

### 3.2 Toeplitz Tridiagonal Matrix Algorithm

The Toeplitz structure of matrix  $\mathbf{A}$  allows us to reduce both the memory demand and the FLOPs of the Thomas algorithm. This specialized form of the Thomas algorithm takes only  $\tilde{\mathbf{f}} \in \mathbb{R}^n$  as input and outputs the solution  $\mathbf{v} \in \mathbb{R}^n$ . Similarly, the FLOPs reduce to  $2(n-1)$  for the forward and another  $2(n-1)$  for the backward substitution. The reduction in the forward substitution is due to the fact that the diagonal elements in line 6 of algorithm 2 can now be fully precomputed since they have analytical expression. Therefore, the total FLOPs for the specialized Thomas algorithm is  $4(n-1)$  which is a reduction by half compared to the general Thomas algorithm.

---

**Algorithm 2** Toeplitz Tridiagonal Matrix Algorithm

---

**Input:**  $\tilde{\mathbf{f}}$ **Output:**  $\mathbf{v}$ 

```
1:  $b'_0 = 2$            // First element along main diagonal
2:  $f'_0 = \tilde{f}_0$        // First element in the RHS
3:
4: // Forward elimination
5: for  $i = 1$  to  $i = n - 1$  do
6:    $b'_i = (i + 2)/(i + 1)$  // Eliminate lower diag
7:    $f'_i = y_i + f'_{i-1}/b'_{i-1}$  // Change RHS
8: end for
9:
10: // Backward elimination
11:  $v_{n-1} = \tilde{f}'_{n-1}/b'_{n-1}$  // Final element in solution
12: for  $i = n - 2$  to  $i = 0$  do
13:    $v_i = (\tilde{f}'_i + v_{i+1})/b'_i$ 
14: end for
```

---

### 3.3 LU-decomposition based Toeplitz Tridiagonal Matrix Algorithm

Here, we used *Armadillo* (a C++ Linear Algebra Library [1]) to perform LU-decomposition and solve the Toeplitz triangular set of linear equations. *Armadillo* has a built in function called *solve* that is a highly optimized function that performs LU-decomposition by default to solve a linear set of equations. However, in our pseudo-code (cf. algorithm 3) we have separately specified the LU-decomposition and solving the two triangular equations for clarity of presentation. Any LU-decomposition based algorithm is composed of three steps; first LU-decomposition of the system matrix, and then triangular forward substitution, and at last triangular backward substitution. The pseudo-code in algorithm 3 takes as input  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\tilde{\mathbf{f}} \in \mathbb{R}^n$ , and outputs the solution  $\mathbf{v} \in \mathbb{R}^n$ . Here, it is pertinent to note that algorithm 3 takes a matrix as input, and hence requires a large amount of memory to store this matrix. Nevertheless, LU-decomposition is a method of choice when solving a linear set of equation with a densely populated matrix. For the case  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the FLOPs for LU-decomposition takes  $O(2/3n^3)$ , and the two substitutions take each  $\sim O(n^2)$  [2]. However, here we are interested in solving a system of equations with a Toeplitz tridiagonal matrix structure. For such a case the LU-decomposition has an analytical expression and takes only  $O(3n)$  FLOPs, the forward substitution takes  $O(2n)$  while the backward substitution takes  $O(3n)$  [3]. In total the number of FLOPs for LU-decomposition based Toeplitz tridiagonal matrix system of equations is  $8n$ .



---

**Algorithm 3** LU-decomposition based Toeplitz Tridiagonal Matrix Algorithm

---

**Input:**  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{f}}$ **Output:**  $\mathbf{v}$ 

```
lu(L, U, A)      // LU-decomposition
y = solve(L, f)   // Forward substitution
v = solve(U, y)   // Backward substitution
```

---

### 3.4 Error Analysis

The error generated by the different algorithms developed in this project are quantified by comparing the results of the algorithm with that of the closed form (analytical) solution using the relative error defined by

$$\epsilon_i = \left| \frac{v_i^{num} - v_i^{ana}}{v_i^{ana}} \right|, \quad i = 1, \dots, n, \quad (19)$$

where  $v_i^{num}$  and  $v_i^{ana}$  are the numerical and analytical solutions at the  $i^{th}$  discretization location, respectively. Moreover,  $|\cdot|$  is the absolute value operator. The analytical solution for the one-dimensional Poisson equation with a known source function is derived in Appendix A.

## 4 Results

### 4.1 Computational Time

To compare the computational efficiency of the three different algorithms discussed in this project we analysed their execution time elapsed in terms of CPU time (cf. Tabel 1). The time is computed by taking an average of 1000 runs of each of the algorithms. Considering the number of FLOPs count for the three different algorithms, one might expect to see the CPU time for LU-decomposition and general Thomas algorithm to be similar while the special Thomas algorithm being twice faster than the two other algorithms. However, what we actually observe from the CPU time is that the special Thomas algorithm is slightly faster than the general Thomas algorithm and the LU-decomposition algorithm is at least an order of magnitude slower than the two Thomas algorithms. Another observation that we can make from the CPU time is that, as the matrix size gets larger the special Thomas algorithm gets more faster than the general Thomas algorithm. Furthermore, the CPU time difference between the two Thomas algorithms and the LU-decomposition reduces as the matrix size gets larger.

One possible reason for explaining the discrepancy between the FLOPs count and the CPU time is that, not all operations consume the same CPU time. For example it is known that division is the most computationally costly operation. Thus, comparing the two Thomas algorithms, it is possible to see that the number of division operations are the same. Therefore, it might be possible that the division operation is the cause for not seeing the special algorithm not being twice as fast as the general algorithm. On the other hand, the very slow LU-decomposition algorithm could be because Armadillo might have treated the matrix  $\mathbf{A}$  as a fully populated matrix.

Matrix size	General [s]	Special [s]	LU-decomp. [s]
$10 \times 10$	$6.970 \cdot 10^{-7}$	$6.910 \cdot 10^{-7}$	$1.139 \cdot 10^{-5}$
$100 \times 100$	$2.114 \cdot 10^{-6}$	$2.047 \cdot 10^{-6}$	$6.882 \cdot 10^{-5}$
$1000 \times 1000$	$1.676 \cdot 10^{-5}$	$1.577 \cdot 10^{-5}$	$2.722 \cdot 10^{-4}$

Table 1: CPU time in seconds for the three different algorithms (i.e., General and special Thomas, and LU-decomposition).

## 4.2 Algorithm Benchmark and Error Analysis

Due to the availability of the analytical solution to the one-dimensional Poisson equation, we have benchmarked our numerical algorithms by computing the relative error. We considered three different cases (i.e.,  $n = 10$ ,  $n = 100$ , and  $n = 1000$ ) and show the results only for the general Thomas algorithm. This is mainly because the two other algorithms provide a similar result as that of the general Thomas algorithm. Figure 1 presents the numerical solution for the three different cases as well as the corresponding analytical solution. Here, we can observe that for the cases  $n = 100$  and  $n = 1000$  the numerical solution is the same as the analytical solution. To quantify the difference between the numerical and analytical solutions, we computed the relative error for the three different cases and shown it in logarithmic scale in Figure 2. Notice that the relative error is the same for all values of  $x$  and that for  $n = 10$ ,  $\epsilon \sim 10^{-1.2}$ , for  $n = 100$ ,  $\epsilon \sim 10^{-3}$ , and for  $n = 1000$ ,  $\epsilon \sim 10^{-5}$ .

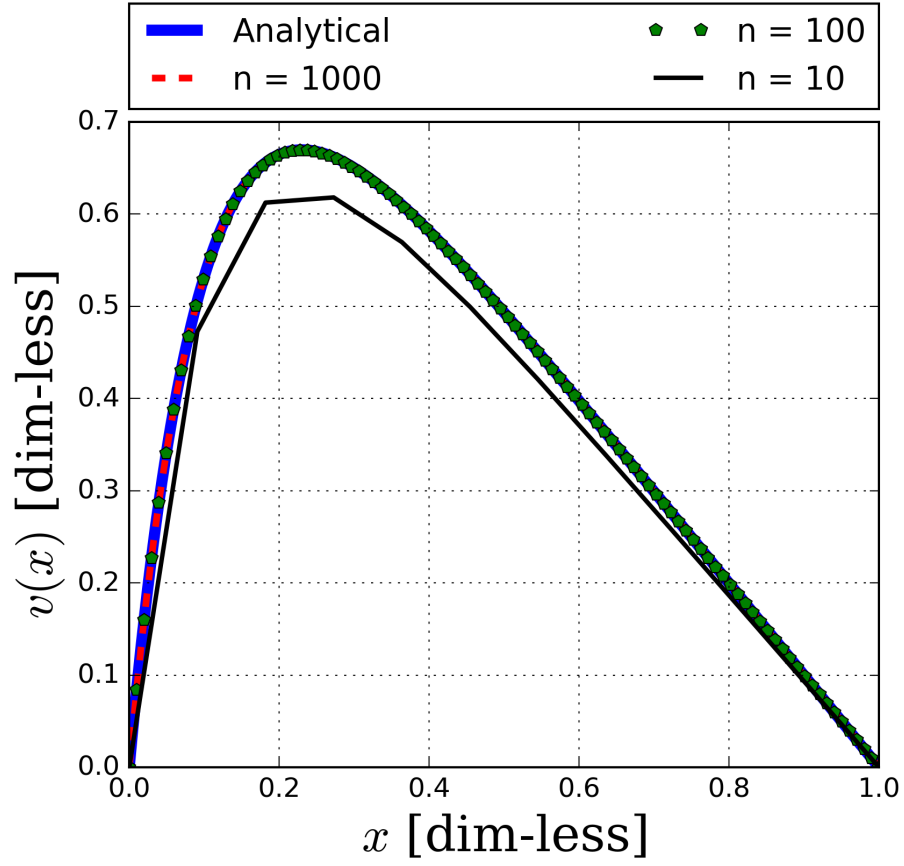


Figure 1: Numerical approximation and analytical solution to the 1D Poisson equation. The numerical solution is shown for different  $n \times n$  matrix sizes. The convergence of the numerical method is apparent as its graph approaches the analytical solution as the number of points  $n$  grows. When  $n \geq 10^2$ , the numerical solution becomes indistinguishable from the analytical solution.

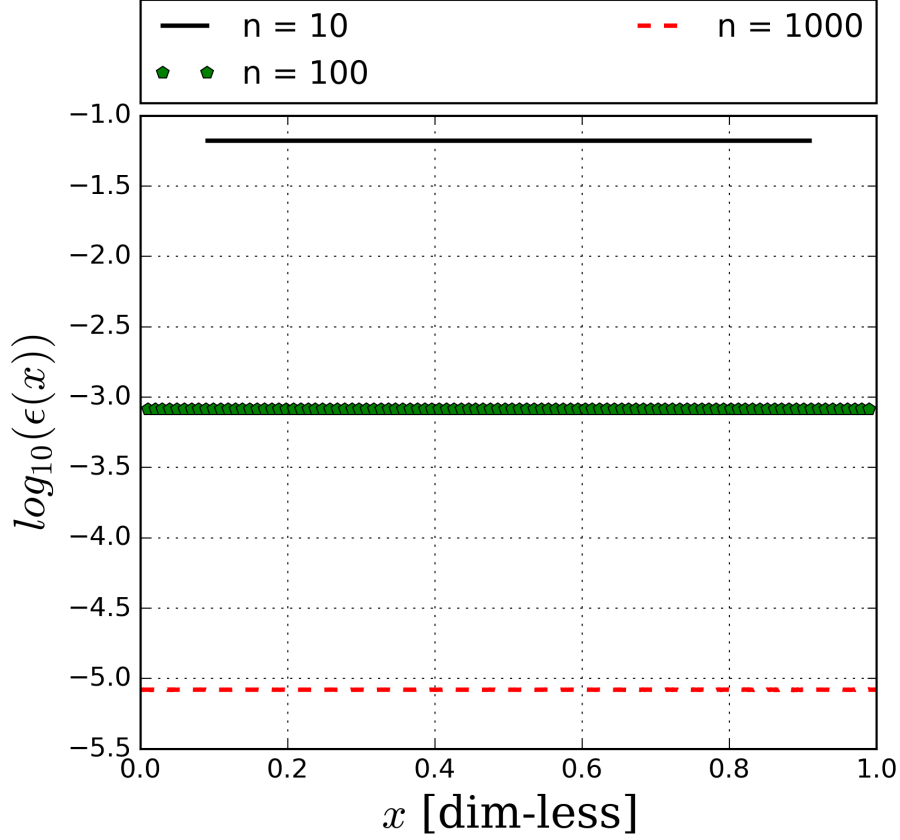


Figure 2: Relative error to the numerical approximations of the 1D Poisson equation. The relative error is shown for different  $n \times n$  matrix sizes.

To analyse the effect of having extremely small discretization step size (or having very large number of grid points  $n$ ) on the precision of the numerical algorithm, we computed the relative error for the general Thomas algorithm for different values of discretization step sizes (cf. Tabel 2 and Figure 3). The logarithmic relative error as a function of the logarithmic discretization step size shows a linear trend and it reduces until it reaches a discretization step size of  $\sim 10^{-5}$  (or  $n = 10^5$ ), where it increases again. This loss of precision is due to round off errors which arises from the fact that computers do not represent numbers with an infinite decimal places and hence need to be rounded off to a number the computer can represent.

Grid points ( $n$ )	Step size ( $\log_{10}(h)$ )	Relative error ( $\log_{10}(\epsilon)$ )
$10^1$	-1.041	-1.179
$10^2$	-2.004	-3.0880
$10^3$	-3.000	-5.0807
$10^4$	-4.000	-7.0791
$10^5$	-5.000	-8.843
$10^6$	-6.000	-6.075
$10^7$	-7.000	-5.525

Table 2: Logarithmic relative error as a function of logarithmic discretization step size (or the number of grid points). Here we have calculated the maximum relative error for each of the discretization step sizes.

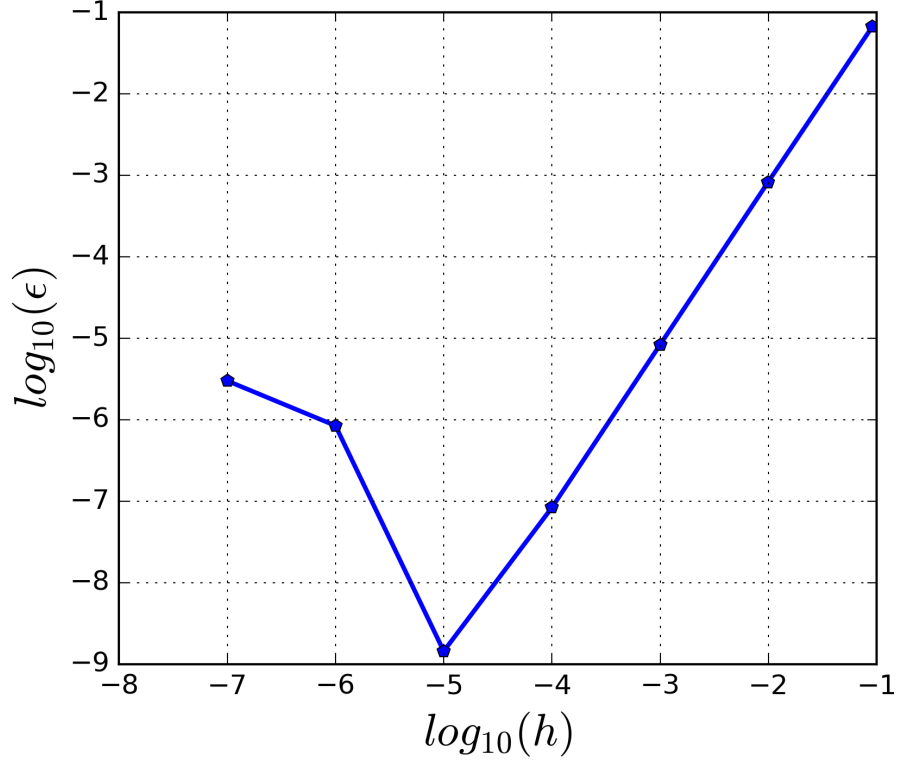


Figure 3: A plot of the logarithmic relative error as a function of the logarithmic discretization step size. Here we can see that the relative error follows our expectation for the truncation error, with a slope value of -2 on a logarithmic scale, up to a point of  $n \sim 10^5$ .

## 5 Discussion

The memory usage of the three different algorithms depends on whether the algorithm uses static or dynamic memory allocation. In addition, it also depends whether the algorithm has to store vectors or matrices. In this project we have allocated all our variables using dynamic memory. It is also important to notice that the LU-decomposition algorithm has to store a matrix while the two Thomas algorithms store only vectors. Consequently, we would run into a memory storage problem when using the LU-decomposition algorithm. For example, consider a problem with the number of grid points  $n = 10^5$ . This implies matrix  $\mathbf{A} \in \mathbb{R}^{10^5 \times 10^5}$ . To store this matrix in memory requires  $10^{10} * 64 / (8 * 10^9) = 80\text{GB}$  for double precision numbers. Unfortunately, this is by far larger than a normal personal computer memory ( $\sim 8\text{GB}$ ).

## 6 Conclusion

Three different numerical algorithms were developed for solving the one dimensional Poisson equation with Dirichlet boundary conditions. The algorithms utilize the fact that the problem can be casted into a matrix equation. The general Thomas algorithm uses the tridiagonal structure of the matrix equation while the special Thomas and LU-decomposition algorithms use the Toeplitz tridiagonal structure for the matrix  $\mathbf{A}$ . This allows the three algorithms to be computationally efficient compared to the general Gaussian elimination algorithm. Moreover, the results of the three algorithms approach the corresponding analytical solution with a small relative error for the number of grid points  $n \geq 100$ . The special Thomas algorithm out performed the two other algorithms by having the fastest computational speed and the smallest memory requirement. The LU-decomposition algorithm has the highest memory requirement and it was even impossible to use the method in a personal computer for  $n > 10^4$ . Finally, we have observed that the numerical precision of the algorithms increases with increasing the  $n$  until it reaches the round off limit  $n \sim 10^5$ .

## References

- [1] C. Sanderson and R. Curtin, *Armadillo: a template-based C++ library for linear algebra.*, Journal of Open Source Software, <http://dx.doi.org/10.21105/joss.00026>, 2016.
- [2] M. Hjorth-Jensen, "Computational physics," University Lecture Notes, 2015.
- [3] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in C*, 2nd ed. Cambridge, USA: Cambridge University Press, 1992.

## A Analytical Solution

The one-dimensional Poisson equation with Dirichlet boundary condition has a closed form solution for a known and twice integrable source function. To

derive this closed form solution, we first rewrite equation 1 for a source function  $f(x) = 100e^{-10x}$  as

$$-u''(x) = 100e^{-10x}, \quad x \in (0, 1), \quad u(0) = u(1) = 0. \quad (20)$$

A closed form solution to equation 20 can now be found by integrating it twice respect to  $x$ . This results

$$u(x) = -e^{-10x} + A + Bx, \quad (21)$$

where  $A$  and  $B$  are integration constants. We now impose the Dirichlet boundary condition to equation 21 and find that  $A = 1$  and  $B = e^{-10} - 1$ . Finally, replacing the values of  $A$  and  $B$  into equation 21 results the analytical solution

$$u(x) = 1 - (1 - e^{-10})x - e^{-10x}. \quad (22)$$