

SOLVING THE POISSON-EQUATION IN ONE DIMENSION: TDMA VS LU

FYS3150: COMPUTATIONAL PHYSICS

TRYGVE LEITHE SVALHEIM
SEBASTIAN G. WINTHER-LARSEN
[GITHUB.COM/GREGWINTHER](https://github.com/gregwinther)

ABSTRACT. A one-dimensional version of the Poisson equation with Dirichlet boundary conditions is solved using two different algorithms. The two algorithms employed are the tridiagonal matrix algorithm and the LU decomposition method. We find that a fine-tuned version of the tridiagonal matrix algorithm is 10^5 times faster than the general LU decomposition method and 10 times as precise.

Contents

1. INTRODUCTION

In this study the one-dimensional Poisson dimension will be solved using two different algorithms. The Poisson equation is a partial differential equation of elliptic type with broad application. In physics it is for instance used to describe the potential field of a charge, or mass density distribution. By assuming spherical symmetry, the Poisson equation is reduced to one dimension.

Herein, we devise two algorithms for solving the Poisson equation. The first one is the tridiagonal matrix algorithm, and the second is the LU decomposition method. We see how the tridiagonal matrix algorithm can be tailored to the problem and therefore becomes vastly more efficient than the more general LU method. The algorithms are tested at different precision levels.

Firstly, the physical background for the Poisson equation and how the second derivative can be approximated by a three point formula is presented. Secondly, we provide a thorough elaboration of how the algorithms can be implemented. Thirdly, results of the algorithms performance is presented. Lastly, we provide a discussion of the findings and some concluding remarks.

2. THEORY

2.1. The Poisson Equation. The Poisson equation is a classical equation from electromagnetism. The electrostatic potential Φ is generated by a localized charge distribution $\rho(\mathbf{r})$. In three dimensions the equation reads

$$(1) \quad \nabla^2 \Phi = -4\pi\rho(\mathbf{r})$$

where ∇^2 is the Laplace operator. In three dimensions the Laplace operator can be expressed using spherical coordinates, but in this study we are assuming that Φ and ρ are spherically symmetric, thus reducing the equation to a one-dimensional problem. only dependent on radius r .

$$(2) \quad \nabla^2 = \frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi}{dr} \right)$$

By substituting $\Phi(r) = \phi(r)/r$ the Poisson equation is reduced to

$$(3) \quad \frac{d^2\phi}{dr^2} = -4\pi r\rho(r)$$

and by letting $\phi \rightarrow u$ and $r \rightarrow x$ one is left with the very simple equation

$$(4) \quad -u''(x) = f(x)$$

The inhomogenous term f , or source term, is given by the charge distribution ρ multiplied by r and the constant -4π . In this study, however, the source term will be $f(x) = 100e^{-10x}$ and the results can be compared to the analytical solution $u(x) = 1 - (1 - e^{-10})x - e^{-10x}$.

2.2. Approximation of the Second Derivative. In this study the one-dimensional Poisson equation will be solved with Dirichlet boundary conditions by rewriting it as a set of linear equations. The discretized approximation of u is defined as v_i with grid points $x_i = ih$, step size of $h = \frac{1}{n+1}$, in the interval $x_0 = 0$ to $x_{n+1} = 1$ and with boundary conditions $v_0 = v_{n+1} = 0$. The interior solution $v_i \forall i \in 1, \dots, n$ is to be found. The second order derivative is approximated with the three point formula such that equation ?? becomes

$$(5) \quad -\frac{v_{i+1} - 2v_i + v_{i-1}}{h} = f_i$$

By defining $\mathbf{f} = h^2 f_i$ one can rewrite equation ?? as $-v_{i+1} - 2v_i + v_{i-1} = h^2 f_i$. If we ignore the end points, $i = 0$ and $i = n + 1$, this equation can be represented as a matrix equation.

$$(6) \quad A\mathbf{v} = \mathbf{f}$$

$$(7) \quad \begin{bmatrix} 2 & -1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & -1 & 2 & -1 & \cdots & 0 & 0 & 0 \\ & & \vdots & & \ddots & & \vdots & \\ 0 & 0 & 0 & 0 & \cdots & -1 & 2 & -1 \\ 0 & 0 & 0 & 0 & \cdots & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_{n-1} \\ v_n \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \\ f_{n-1} \\ f_n \end{bmatrix}$$

3. ALGORITHMS

Two main methods are implemented and compared. The first method is gaussian elimination of the tridiagonal matrix A , also known as the *Thomas Algorithm* [?]. This is a simplified form of Gaussian elimination that can be used to solve tridiagonal systems of equations. The method is improved upon in order to take into account the fact that the matrix we are dealing with has the same numbers along the diagonals. The second method is the LU-decomposition method. Code snippets are included in this section, for the full programs, follow the link to github at the first page of this article.

3.1. Tridiagonal Matrix Algorithm. Our tridiagonal system can be represented by

$$(8) \quad a_i x_{i-1} + b_i x_i + c_i x_{i+1} = \mathbf{f}$$

with $a_i = -1, b_i = 2$ and $c_i = -1$, except for $a_1 = 0$ and $c_n = 0$. Row reducing a matrix will reveal how the algorithm functions. Limiting the problem to four dimensions for easier reading and to save the rainforest¹.

$$(9) \quad \left[\begin{array}{cccc|c} b_1 & c_1 & 0 & 0 & f_1 \\ a_2 & b_2 & c_2 & 0 & f_2 \\ 0 & a_3 & b_3 & c_3 & f_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right] \sim \left[\begin{array}{cccc|c} 1 & c_1/b_1 & 0 & 0 & f_1/b_1 \\ 0 & b_2 - \frac{c_1}{b_1}a_2 & c_2 & 0 & f_2 - \frac{f_1}{b_1}a_2 \\ 0 & a_3 & b_3 & c_3 & f_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right]$$

Now let $\beta_1 = b_1$ and $\beta_2 = b_2 - \frac{c_1}{b_1}a_2$. As new elements start to appear in vector \mathbf{f} , right to the vertical bar in the augmented matrix, they are also relabeled to \tilde{f}_i . For example $\tilde{f}_1 = f_1/\beta_1$. One more iteration will reveal the pattern of the algorithm.

¹Writing out a general case will also take up more paper space

(10)

$$\left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & (f_2 - \tilde{f}_2 a_2)/\beta_2 \\ 0 & a_3 & b_3 & c_3 & f_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right] = \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & a_3 & b_3 & c_3 & f_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right]$$

(11)

$$\sim \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & 0 & b_3 - \frac{c_2}{\beta_2} a_3 & c_3 & f_3 - \tilde{f}_2 a_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right] \sim \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & 0 & 1 & c_3/\beta_3 & (f_3 - \tilde{f}_2 a_3)/\beta_3 \\ 0 & 0 & a_4 & b_4 & f_4 \end{array} \right]$$

(12)

$$\sim \dots \sim \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & 0 & 1 & c_3/\beta_3 & \tilde{f}_3 \\ 0 & 0 & 0 & 1 & \tilde{f}_4 \end{array} \right]$$

This pattern can be summarized nicely by the following difference equations. This computation will from now on be referred to as the “forward substitution”.

$$(13) \quad \beta_i = b_i - \frac{a_i c_{i-1}}{\beta_{i-1}}, \quad \tilde{f}_i = f_i - a_i \tilde{f}_{i-1} / \beta_{i-1}, \quad i \in [2, n],$$

In c++, this part of the algorithm is implemented just below. Notice that b_i is overwritten instead of initializing a separate vector for β .

```
// GAUSSIAN ELIMINATION
// Forward substitution
f_tilde[1] = f[1];
for (int i=2; i<n+1; i++){
    b_vec[i] = b_vec[i] - a_vec[i] * c_vec[i-1] / b_vec[i-1];
    f_tilde[i] = f[i] - a_vec[i] * f_tilde[i-1] / b_vec[i-1];
}
```

The Gaussian elimination is not fully complete until a “backward substitution” has been performed as well. This ensures that all elements in the row reduced matrix are pivot elements.

$$(14) \quad \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & 0 & 1 & c_3/\beta_3 & \tilde{f}_3 \\ 0 & 0 & 0 & 1 & \tilde{f}_4 \end{array} \right] \sim \left[\begin{array}{cccc|c} 1 & c_1/\beta_1 & 0 & 0 & \tilde{f}_1 \\ 0 & 1 & c_2/\beta_2 & 0 & \tilde{f}_2 \\ 0 & 0 & 1 & 0 & \tilde{f}_3 - \frac{c_3}{\beta_3} \tilde{f}_4 \\ 0 & 0 & 0 & 1 & \tilde{f}_4 \end{array} \right]$$

The notation is updated again to $v_3 = \tilde{f}_3 - \frac{c_3}{\beta_3} \tilde{f}_4$. By continuing this iterative process all the way to the beginning of the matrix by backward substitution leaves us with a second recurrence relation.

$$(15) \quad v_i = (\tilde{f}_i - c_i v_{i+1}) / \beta_i, \quad i \in \{n-1, n-2, \dots, 1\}$$

The backward substitution is implemented in C++ like so:

```
// Backward substitution:
v[n] = f_tilde[n] / b_vec[n];
for (int i=n-1; i>=1; i--){
    v[i] = (f_tilde[i] - c_vec[i] * v[i+1]) / b_vec[i];
}
```

An ever-important aspect of any algorithm is how efficient the algorithm is. It is therefor worthwhile to consider how many floating point operations the algorithm requires². For the forward substitution there are 6 FLOPS in each iteration and $n - 1$ iterations, resulting in $6(n - 1)$ FLOPS for the entire loop. Additionally, $1(n - 1)$ FLOPS are required for the backward substitution, adding up to

$$(16) \quad N_{tridiag} = 8(n - 1) = \mathcal{O}(8n).$$

3.1.1. Optimizing the algorithm. Our case is special, because the elements in each of the diagonal arrays in the tridiagonal matrix, A , are equal. $a_i = -1, b_i = 2$ and $c_i = -1$, except for $a_1 = 0$ and $c_n = 0$. This is an advantage, because the algorithm can be adjusted, and the number of FLOPS consequently reduced.

The algorithm for the forward substitution, in equation ??, can be simplified to

$$(17) \quad \beta_i = b_i - 1/\beta_{i-1} \quad \tilde{f}_i = f_i - \tilde{f}_{i-1}/\beta_{i-1}$$

and the backward substitution, in equation ??, can be simplified to

$$(18) \quad v_i = (\tilde{f}_i + v_{i+1})/\beta_i$$

This results in some new and sleeker C++ code

```
// GAUSSIAN ELIMINATION
// Forward substitution
f_tilde[1] = f[1];
for (int i=2; i<n+1; i++){
    b_vec[i] -= 1/b_vec[i-1];
    f_tilde[i] = f[i] + f_tilde[i-1]/b_vec[i-1];
}

// Backward substitution:
v[n] = f_tilde[n]/b_vec[n];
for (int i=n-1; i>=1; i--){
    v[i] = (f_tilde[i] + v[i+1])/b_vec[i];
}
```

One can see that this new and improved algorithm only requires $6(n-1)$ FLOPS. In general, a Gaussian row reduction algorithm requires $\frac{2}{3}n^3 + \mathcal{O}(n^2)$ FLOPS [?], a great improvement!

3.2. LU Decomposition. The LU decomposition method is a method where the matrix A in equation ?? ($A\mathbf{v} = \mathbf{f}$) can be rewritten as a product of two matrices L and U , where L is lower triangular and U is upper triangular. This method only works for non-singular, invertible matrices, which the matrix A certainly is.

$$(19) \quad A = LU = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ l_{21} & 1 & 0 & \dots & 0 & 0 \\ l_{31} & l_{32} & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ l_{n-11} & l_{n-12} & l_{n-13} & \dots & 1 & 0 \\ l_{n1} & l_{n2} & l_{n3} & \dots & l_{nn-1} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} & \dots & u_{1n-1} & u_{1n} \\ 0 & u_{22} & u_{23} & \dots & u_{2n-1} & u_{2n} \\ 0 & 0 & u_{33} & \dots & u_{3n-1} & u_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & u_{n-1n-1} & u_{n-1n} \\ 0 & 0 & 0 & \dots & 0 & u_{nn} \end{bmatrix}$$

The equation $A\mathbf{v} = LU\mathbf{v} = \mathbf{f}$ can be solved in two steps

$$(20) \quad L\mathbf{w} = \mathbf{f}, \quad U\mathbf{v} = \mathbf{w}$$

²The number of $+, -, *$ and $/$ operations are counted

Written out, the full set of these linear equations will look like

$$\begin{aligned}
(21) \quad & w_1 = f_1 \\
(22) \quad & l_{21}w_1 + w_2 = f_2 \\
(23) \quad & l_{31}w_1 + l_{32}w_2 + w_3 = f_3 \\
(24) \quad & \vdots \\
(25) \quad & l_{n-11}w_1 + l_{n-12}w_2 + l_{n-13}w_3 + \cdots + w_{n-1} = f_{n-1} \\
(26) \quad & l_{n1}w_1 + l_{n2}w_2 + l_{n3}w_3 + \cdots + l_{nn-1}w_{n-1} + w_n = f_n \\
(27) \quad & u_{11}v_1 + u_{12}v_2 + u_{13}v_3 + \cdots + u_{1n-1}v_{n-1} + u_{1n}v_n = w_1 \\
(28) \quad & u_{22}v_2 + u_{23}v_3 + \cdots + u_{2n-1}v_{n-1} + u_{2n}v_n = w_2 \\
(29) \quad & u_{33}v_3 + \cdots + u_{3n-1}v_{n-1} + u_{3n}v_n = w_3 \\
(30) \quad & \vdots \\
(31) \quad & u_{n-1n-1}v_{n-1} = w_{n-1} \\
(32) \quad & u_{nn}v_n = w_n
\end{aligned}$$

First \mathbf{w} can be determined by iterating forwards through equations ?? to ?. Then \mathbf{v} can be found by iterating backwards from equation ? to equation ?.

In this study the Armadillo library is employed instead of implementing the LU decomposition algorithm from scratch. The LU method is used as a benchmark against which the speed of the tridiagonal method is measured. The number of FLOPS required by the LU is $2n^3/3$, by any algorithm used to implement it[?]. Additionally, n^2 FLOPS is required by the forward and backwards substitution.

$$(33) \quad N_{LU} = \frac{2}{3}n^3 - 2n^2 = \mathcal{O}(\frac{2}{3}n^3)$$

Using LU computation results in a quadratic slowdown compared to the tridiagonal matrix algorithm.

4. RESULTS

The run time for the optimized tridiagonal matrix algorithm and the LU decomposition solver is in table ?. The reason run times are not given for $n > 10000$ for the LU method is because the method uses too much memory for the solver to even function. Regardless, the message of table ? is quite clear, the specialized tridiagonal matrix method is vastly more efficient.

TABLE 1. Elapsed time for increasing n

n	TDMA [s]	LU [s]
10	0.0000035	0.00106083
100	0.0000116	0.0022319
1000	0.000077892	0.0677764
10000	0.000878769	21.9247
100000	0.00757418	n/a
1000000	0.08616075	n/a
10000000	0.76534	n/a

Tables ? and ? show the fit of the tridiagonal matrix solver for $n = 10$ and $n = 100$, respectively. One can easily see that already at $n = 100$ the fit is already very good.

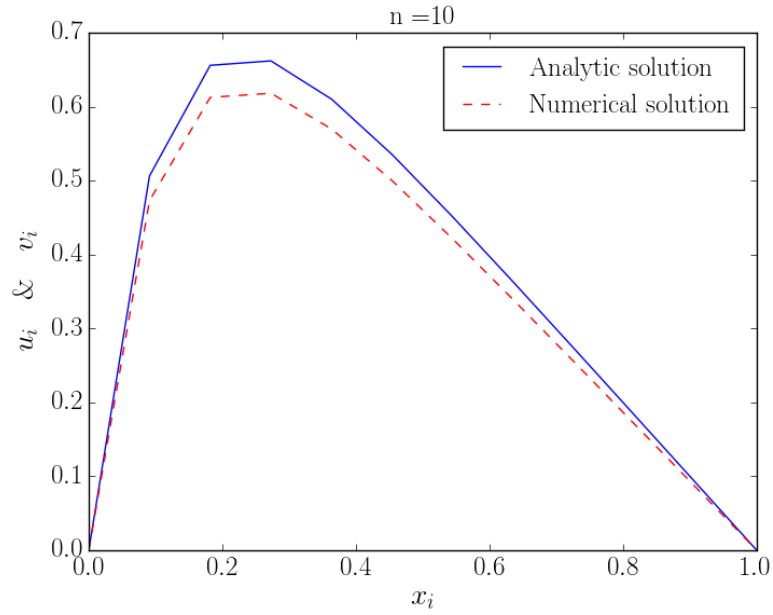


FIGURE 1. Numerical approximation of the tridiagonal matrix algorithm for $n = 10$

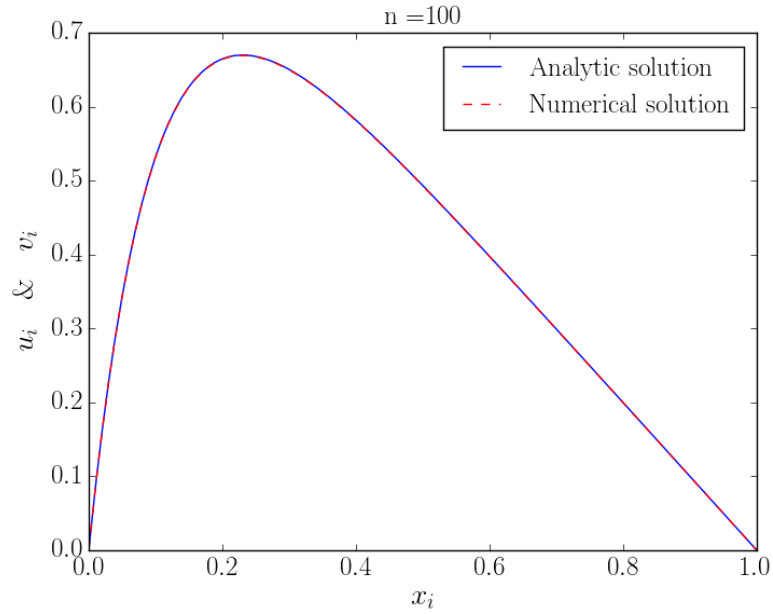


FIGURE 2. Numerical approximation of the tridiagonal matrix algorithm for $n = 100$

Figure ?? shows a plot of the maximum relative error of the tridiagonal solver as a function of step size. One can see that the relative error falls as the step size falls, but only to a certain point. By decreasing the step size further the error increases,

most likely because of decreasing numerical precision because of the computers difficulty in representing very small numbers.

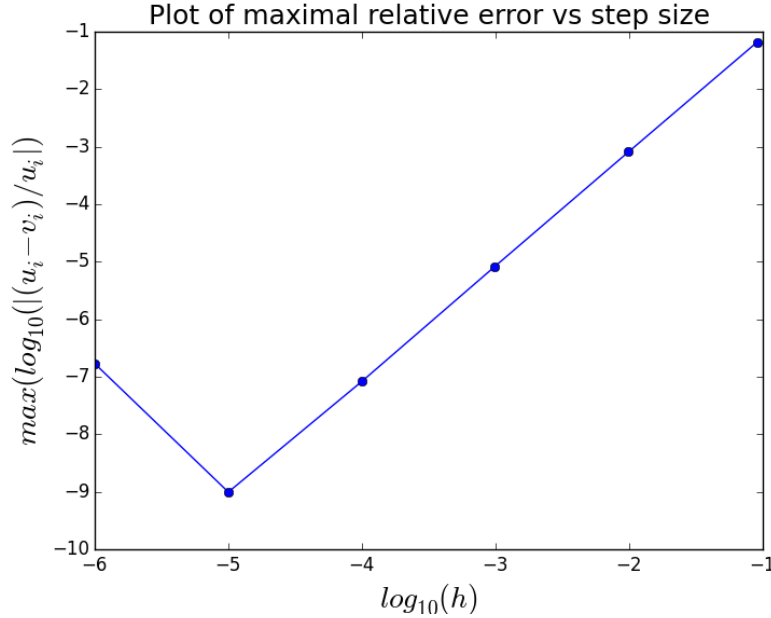


FIGURE 3. Plot of maximum relative error as a function of step size

5. DISCUSSION

When choosing a suitable value for n , one must take into account both precision and speed. Our intuition tells us that a larger n should result in a higher precision for the numerical derivative. We can conclude from figure ?? that this is only true up to a certain point, because of loss in numerical representation accuracy within the computer. An ideal n will be in the magnitude of 10^5 which is equivalent to $h = 10^{-5}$. Increasing this value will not increase run time dratically for the tridiagnoal matrix algorithm. Choosing $n = 10^5$ gives us the ability to calculate with a high certainty of 5 decimals points and with low memory cost.

By looking at the results one is led to believe that there are more reasons to prefer the tridiagonal matrix algorithm for this kind of problem. When using the LU decomposition method the short term memory is quickly filled with nothing but the matrix A . Choosing $n = 10^5$ is equivalent to 8×10^{10} bytes = 80GB of memory! A brand new computer ships with around 8 GB of memory today, which is a tenth of the space required. Obviously, this inhibits our ability to perform the calculation.

6. CONCLUSION

As we have seen, there is significant advantage to fine-tuning and picking the correct algorithm before solving a problem numerically. We have experienced the restrictions in memory when computing, and have been made aware of the problem this poses when doing numerical calculations. The specialised tridiagonal matrix algorithm was able to solve the one-dimensional Poisson equation more than 20000 times faster than the more general LU decomposition method. This can be explained by the fact that the LU decomposition method require $\mathcal{O}(\frac{2}{3}n^3)$ FLOPS

while the specialized tridiagonal matrix algorithm will do with a meager $\mathcal{O}(8n)$ FLOPS.

REFERENCES

- [1] Thomas, L.H. (1949), *Elliptic Problems in Linear Differential Equations over a Network*. Watson Sci. Comput. Lab Report, Columbia University, New York.
- [2] Hjørth-Jensen, M. (2015). *Computational Physics - Lecture Notes 2015*. University of Oslo
- [3] Golub, G.H., van Loan, C.F. (1996). *Matrix Computations* (3rd ed.), Baltimore: John Hopkins.