

# Two methods of solving linear second-order differential equations with Dirichlet boundary conditions

Anna Lina P. Sjur and Jan-Adrian H. Kallmyr

September 9, 2018

## Abstract

We compare the methods of Gaussian elimination (GE) and LU-decomposition (LU) for solving linear second order differential equations on the form  $\frac{d^2u}{dx^2} = f(x)$  where  $u(0) = u(1) = 0$ . Comparing CPU-times we find that a highly specialised GE algorithm is most efficient with 100% – 40% faster CPU-times for  $10 - 10^7$  integration steps ( $n$ ) over a more general GE algorithm, as well as  $10^6$  times faster than the LU algorithm for  $n = 10^4$ . The relative error is in the same order of magnitude for all algorithms at around  $10^{-1} - 10^{-7}$  for  $n \in [10, 10^4]$ , and for our specialised GE algorithm, we see that the error starts to increase again at around  $n = 10^6$ . We therefore find that the specialised GE algorithm is the most reliable for this problem, however, a good compromise between generality and efficiency is found in the general GE algorithm.

## 1 Introduction

Physics is a field concerned with the behaviour of nature, and nature is everchanging. It is therefore no surprise that differential equations appear everywhere in physics. From global climate dynamics to statistical mechanics, what we find is that differential equations, often many and coupled, are required to explain or model the phenomena. For such large models, efficiency is important, as we would, for example, like to have timely weather forecasts. One way to make a model more efficient is by using an efficient algorithm for solving differential equations.

In this article, we compare two different numerical methods of solving linear second-order differential equations with the Dirichlet boundary conditions. To do this, we will solve the one-dimensional Poisson's equation:

$$\frac{d^2\phi}{dr^2} = -4\pi r\rho(r). \quad (1)$$

In the methods section, we develop an approximation for the 2nd derivative to the 2nd order. We will then solve eq. 1 numerically, using gaussian elimination and lower-upper decomposition. As for the former, we will further specialise it to solve eq. 1 more efficiently. Next, in the results section we present the comparison between our numerical solutions and the analytical solution, as well as the error. We then compare the efficiency of all three algorithms, and finally, in the discussion section we will consider the advantages and disadvantages of each algorithm, and discuss their uses.

## 2 Methods

We would like to solve eq. 1 numerically. Generalising the equation, we get

$$-\frac{d^2u}{dx^2} = f(x), \quad (2)$$

where we have assumed that  $\rho \propto \frac{1}{r}e^{-r}$  and let  $r \rightarrow x$ ,  $\phi \rightarrow u$ . Summing the backward and forward Taylor expansions of  $u(x)$  and discretising the equation for  $n$  integration points, we get:

$$\frac{v_{i+1} - 2v_i + v_{i-1}}{h^2} + \mathcal{O}(h^2) = f_i, \quad (3)$$

where  $h = \frac{1}{n+1}$ . Using the Dirichlet boundary conditions  $v_0 = v_{n+1} = 0$  and only considering  $x \in (0, 1)$ , we can rewrite the equation as a set of linear equations, represented as the following matrix equation:

$$\mathbf{A}\mathbf{v} = \mathbf{d}, \quad (4)$$

where

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 & \dots & \dots & 0 \\ -1 & 2 & -1 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & -1 & 2 & -1 \\ 0 & \dots & \dots & 0 & -1 & 2 \end{bmatrix},$$

is a tridiagonal matrix and  $d_i = h^2 f_i$ . Having the equation in this form, we can use linear algebra to solve the set of linear equations, and obtain the 2nd derivative of  $u(x)$ .

## 2.1 Gaussian elimination

The fact that  $\mathbf{A}$  is tridiagonal, means that we can develop a general algorithm to solve the equations by the method of Gaussian elimination. If we generalise our matrix

$$\mathbf{A} = \begin{bmatrix} b_1 & c_1 & 0 & \dots & \dots & 0 \\ a_1 & b_2 & c_2 & 0 & \dots & 0 \\ 0 & a_2 & b_3 & c_3 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & a_{n-2} & b_{n-1} & c_{n-1} \\ 0 & \dots & \dots & 0 & a_{n-1} & b_n \end{bmatrix},$$

and row reduce  $\mathbf{A}$  to become upper-triangular (see appendix A.) This yields us the following relations:

$$\tilde{b}_i = b_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}c_{i-1}, \quad \tilde{d}_i = d_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}\tilde{d}_{i-1}, \quad (5)$$

where  $\tilde{b}_1 \equiv b_1$  and  $\tilde{d}_1 \equiv d_1$ . Performing the matrix-vector multiplication in eq. 4, we get the following relations from the second-to-last and last row:

$$v_i = \frac{d_i - c_i v_{i+1}}{\tilde{b}_i}, \quad v_n = \frac{\tilde{d}_n}{\tilde{b}_n}. \quad (6)$$

With these expressions, we can construct our general algorithm as follows:

initialise  $\mathbf{v}$ ,  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{d}$ ;

**for**  $i = 2, \dots, n$  **do**

$$| \quad b_i = b_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}c_{i-1};$$

$$| \quad d_i = d_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}\tilde{d}_{i-1};$$

$$v_n = \frac{\tilde{d}_n}{\tilde{b}_n};$$

**for**  $i = n-1, \dots, 1$  **do**

$$| \quad v_i = \frac{d_i - c_i v_{i+1}}{\tilde{b}_i};$$

This algorithm has a total of  $8n$  FLOPS.

If we now take into account that our matrix has only twos along the diagonal, and -1 above and below the diagonal, we can make a specialised algorithm. Plugging in these values in the expression for  $\tilde{b}_i$  in equation 5 gives us the following:

$$\tilde{b}_i = 2 - \frac{1}{\tilde{b}_{i-1}}, \quad \tilde{b}_1 = 2 \quad (7)$$

We can prove by induction that this is the same as:

$$\tilde{b}_i = \frac{i+1}{i}, \quad i \geq 1 \quad (8)$$

Since  $\tilde{b}_i$  is only a function of  $i$  we can now set up  $\tilde{\mathbf{b}}$  in the initialisation. Plugging in the values in the remaining expressions gives:

$$\tilde{d}_i = d_i + \frac{\tilde{d}_{i-1}}{\tilde{b}_{i-1}}, \quad v_i = \frac{\tilde{d}_i + u_{i+1}}{\tilde{b}_i} \quad (9)$$

where  $v_n = \frac{\tilde{d}_n}{\tilde{b}_n}$ . This gives us the following specialised algorithm, with a total of  $4n$  FLOPS (not counting the initialisation of  $\tilde{\mathbf{b}}$ ):

```

initialise  $\mathbf{v}$ ,  $\mathbf{b}$ ,  $\mathbf{d}$ ;
for  $i = 2, \dots, n$  do
     $d_i = d_i + \frac{d_{i-1}}{b_{i-1}};$ 
 $v_n = \frac{d_n}{b_n};$ 
for  $i = n - 1, \dots, 1$  do
     $v_i = \frac{d_i + v_{i+1}}{b_i};$ 

```

## 2.2 LU-decomposition

Let  $A$  be a square matrix, then we can decompose it in the following form

$$A = LU, \quad (10)$$

where  $L$  is a lower triangular matrix and  $U$  is an upper triangular matrix. For further conditions on  $A$ , see [1]. This lets us rewrite eq. 4 in the form:

$$(LU)v = d, \quad (11)$$

where we let

$$y = Uv, \quad (12)$$

and rewrite eq. 11 to

$$Ly = d. \quad (13)$$

To solve for  $v$  then we need to first solve eqs. 13 for  $y$  and then substitute  $y$  into eq. 12 where we can then solve for  $v$ , yielding us the solution to eq. 4.

This then is our general LU-decomposition method, and it has around  $\frac{2}{3}n^3$  flops. We could elaborate further on the finite difference details hidden away in the notation, but this algorithm is implemented in armadillo (see the Programming technicalities section), and so most of these details aren't handled by us.

## 2.3 Error

We define the error for each mesh point as:

$$\epsilon_i = \log_{10} \left( \left| \frac{v_i - u_i}{u_i} \right| \right) \quad (14)$$

where  $u_i = u(x_i)$  is the analytical solution given by  $u(x) = 1 - (1 - e^{-10})x - e^{-10x}$ .

## 2.4 Programming technicalities

We implement our algorithms and methods in C++. For our Gaussian elimination algorithm we use only the standard library to declare and operate arrays. For the LU decomposition algorithm we use the armadillo library with LAPACK to declare matrices and perform the LU-decomposition as well as solving equations 13 and 12.

## 3 Results

Let's start with the solution to equation 4. Our general matrix solver gave the results shown in Figure 1 for  $n=10$ ,  $n=100$  and  $n=1000$ . The analytic solution is also shown. For  $n=1000$  the numerical and analytical solutions are so close to each other that we can not distinguish them in the plot, even after zooming in with a factor of 100. Both the specialised algorithm and the LU decomposition algorithm gave the same result as shown in figure 1.

We wanted to get an idea about the efficiency of the different algorithms and how they compare. Table 1 shows the ratio between the CPU time for the general algorithm and the specialised algorithm, as well as the ratio between the LU decomposition algorithm and the specialised algorithm for different values of  $n$ . We see that the specialised algorithm is between 40% to 100% faster than the general algorithm, where the difference appears to decrease as  $n$  increases. On the other hand, the ratio between the LU decomposition algorithm

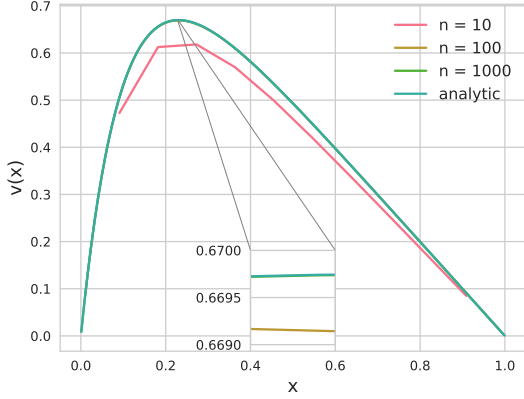


Figure 1: The numeric solution using different numbers of steps and the analytic solution. The window shows a section of the plot zoomed in with a factor of 100.

and the specialised algorithm increases rapidly with  $n$ . For  $n = 10^4$  the specialised algorithm is more than a million times faster than the LU decomposition. As we were not able to run the program for  $n > 10^4$ , we have no values to present in this domain.

$n$	$t_g/t_s$	$t_{LU}/t_s$
10	2.08	3.70
$10^2$	1.89	$1.00 \cdot 10^2$
$10^3$	1.48	$1.05 \cdot 10^4$
$10^4$	1.43	$1.18 \cdot 10^6$
$10^5$	1.39	-
$10^6$	1.41	-
$10^7$	1.39	-

Table 1: Ratio between CPU time for the general algorithm ( $t_g$ ), the special algorithm ( $t_s$ ) and the LU decomposition algorithm ( $t_{LU}$ ) for different matrix sizes ( $n$ ). The LU decomposition crashed for  $n$  greater than  $10^4$ .

Lastly we wanted to assert the error in our numerical solutions. Figure 2 shows the maximum relative error in our specialised algorithm

as a function of  $n$  with a logarithmic scale. We see that the error decreases with  $n$  until  $n = 10^6$ , and then increases.

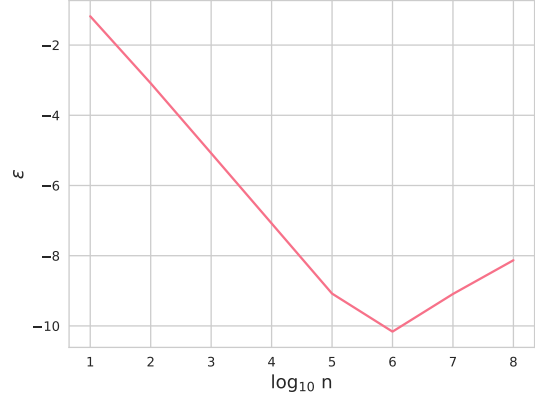


Figure 2: The maximum error in our specialised matrix solver as a function of the number of steps/matrix size on a logarithmic scale.

## 4 Discussion

Having presented the different methods and their relative efficiency and accuracy, we would like to discuss their uses as well as their boundaries. We see that for our problem of solving eq. 2, we could construct a very specialised and efficient algorithm because we used a 2nd order central approximation to the 2nd derivative. This let us reduce the number of flops by a half, and our computer should, with an optimised CPU, only have to do half the work. In terms of accuracy, we found no significant difference between the algorithms, and so this also means that we get the most accuracy per CPU time with our specialised algorithm up until  $n=10^{-6}$  (see Figure 2). Thus for our problem of a linear second order differential equation with Dirichlet boundaries, this is clearly the most advantageous algorithm.

However, while the specialised algorithm is

good at solving eq. 2, it is, as our naming suggests, not very general. There exists a wide variety of problems that can be expressed in the form of our matrix equation, eq. 4, where the matrix  $A$  takes different forms. What we have called the general algorithm, which is almost as efficient as the specialised algorithm (see Table 1), can be used if the matrix  $A$  is tridiagonal. If matrix  $A$  can take *any* form neither our specialised nor our general algorithm can be used. In the latter case, the LU-decomposition is useful, and

## References

- [1] Pavel Okunev and Charles R Johnson. Necessary and sufficient conditions for existence of the lu factorization of an arbitrary matrix. *arXiv preprint math/0506382*, 2005.

## A

### Derivation of Gaussian elimination algorithm terms

Considering our generalised matrix

$$A = \begin{bmatrix} b_1 & c_1 & 0 & \dots & \dots & 0 \\ a_1 & b_2 & c_2 & 0 & \dots & 0 \\ 0 & a_2 & b_3 & c_3 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & a_{n-2} & b_{n-1} & c_{n-1} \\ 0 & \dots & \dots & 0 & a_{n-1} & b_n \end{bmatrix},$$

we will obtain our algorithm from row reduction by subtracting the first row times  $\frac{a_1}{b_1}$  from the second row ( $I - \frac{a_1}{b_1}II$ ), and we see that the new diagonal element is given by

$$\tilde{b}_2 \equiv b_2 - \frac{a_1}{b_1}c_1.$$

Considering eq. 4 then, the corresponding change on the left hand side is given by

$$\tilde{d}_2 \equiv d_2 - \frac{a_1}{b_1}d_1.$$

Performing the same operation on the next row ( $III - \frac{a_2}{b_2}II$ ), we get

$$\tilde{b}_3 \equiv b_3 - \frac{a_2}{b_2}c_2, \quad \tilde{d}_3 \equiv d_3 - \frac{a_2}{b_2}\tilde{d}_2.$$

Letting  $\tilde{b}_1 \equiv b_1$ ,  $\tilde{d}_1 \equiv d_1$ , and generalising, we get:

$$\tilde{b}_i = b_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}c_{i-1}, \quad \tilde{d}_i = d_i - \frac{a_{i-1}}{\tilde{b}_{i-1}}\widetilde{d_{i-1}}. \quad (4.15)$$