

## Problem 1

### a. Discretization of the 2-dimensional Poisson's equation

Starting off with given equation of

$$-\nabla \cdot (k \nabla u) = f, \quad (1)$$

with

$$\begin{aligned} k(x, y) &= 1 + 4x + 6y, \quad (x, y) \in \bar{\Omega} \\ f(x, y) &= e^{\alpha(x-1)^2 + \alpha(y-1)^2} + e^{\alpha(x-3)^2 + \alpha(y-1)^2} \\ &\quad + e^{\alpha(x-5)^2 + \alpha(y-1)^2} + e^{\alpha(x-7)^2 + \alpha(y-1)^2} \\ &\quad + e^{\alpha(x-1)^2 + \alpha(y-3)^2} + e^{\alpha(x-3)^2 + \alpha(y-3)^2} \\ &\quad + e^{\alpha(x-5)^2 + \alpha(y-3)^2} + e^{\alpha(x-7)^2 + \alpha(y-3)^2} \\ &\text{with } \alpha = -5, \quad (x, y) \in \bar{\Omega} \end{aligned} \quad (2)$$

we integrate over cell which is as such in the figure below

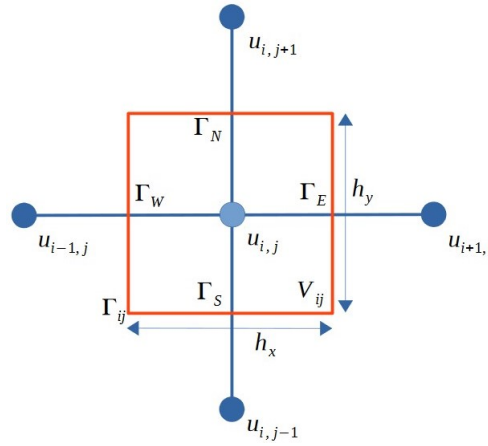


Figure 1: 2 dimensional cell (taken from lecture slide)

resulting in

$$\iint_{V_{ij}} [-\nabla \cdot (k \nabla u)] dV = \iint_{V_{ij}} f dV \quad (3)$$

and applying the 2 dimensional Divergence Theorem, we get

$$\oint_{\Gamma_{ij}} (-k \nabla u \cdot \mathbf{n}) d\Gamma = \iint_{V_{ij}} f dV \quad (4)$$

and we rewrite it as

$$\oint_{\Gamma_{ij}} \left( -k \frac{\partial u}{\partial \mathbf{n}} \right) d\Gamma = \iint_{V_{ij}} f dV \quad (5)$$

The right hand side can then be approximated as follows

$$\iint_{V_{ij}} f dV \approx f_{ij} h_x h_y \quad (6)$$

whereas the left hand side, with reference to the figure 1, can be expanded into

$$\begin{aligned} \oint_{\Gamma_{ij}} \left( -k \frac{\partial u}{\partial \mathbf{n}} \right) d\Gamma &= \int_{\Gamma_W} \left( -k \left( -\frac{\partial u}{\partial x} \right) \right) dy + \int_{\Gamma_E} \left( -k \frac{\partial u}{\partial x} \right) dy \\ &\quad + \int_{\Gamma_S} \left( -k \left( -\frac{\partial u}{\partial y} \right) \right) dx + \int_{\Gamma_N} \left( -k \frac{\partial u}{\partial y} \right) dx \end{aligned} \quad (7)$$

Combining the above two, the equation 5 can be rewritten as

$$\int_{\Gamma_W} \left( -k \left( -\frac{\partial u}{\partial x} \right) \right) dy + \int_{\Gamma_E} \left( -k \frac{\partial u}{\partial x} \right) dy + \int_{\Gamma_S} \left( -k \left( -\frac{\partial u}{\partial y} \right) \right) dx + \int_{\Gamma_N} \left( -k \frac{\partial u}{\partial y} \right) dx = f_{ij} h_x h_y \quad (8)$$

Using central difference and mid-point approximation, we can simplify the elements of the left hand side as follows

$$\begin{aligned} \int_{\Gamma_W} k \frac{\partial u}{\partial x} dy &\approx h_y k \frac{\partial u}{\partial x} \Big|_{(x_{i-1/2}, y_j)} \approx h_y k_{i-1/2, j} \frac{u_{i,j} - u_{i-1,j}}{h_x} \\ \int_{\Gamma_E} \left( -k \frac{\partial u}{\partial x} \right) dy &\approx -h_y k \frac{\partial u}{\partial x} \Big|_{(x_{i+1/2}, y_j)} \approx -h_y k_{i+1/2, j} \frac{u_{i+1,j} - u_{i,j}}{h_x} \\ \int_{\Gamma_S} k \frac{\partial u}{\partial y} dx &\approx h_x k \frac{\partial u}{\partial y} \Big|_{(x_i, y_{j-1/2})} \approx h_x k_{i, j-1/2} \frac{u_{i,j} - u_{i,j-1}}{h_y} \\ \int_{\Gamma_N} \left( -k \frac{\partial u}{\partial y} \right) dx &\approx -h_x k \frac{\partial u}{\partial y} \Big|_{(x_i, y_{j+1/2})} \approx -h_x k_{i, j+1/2} \frac{u_{i,j+1} - u_{i,j}}{h_y} \end{aligned} \quad (9)$$

Substituting the above equations into the equation 8, we get

$$\begin{aligned} &h_y k_{i-1/2, j} \frac{u_{i,j} - u_{i-1,j}}{h_x} - h_y k_{i+1/2, j} \frac{u_{i+1,j} - u_{i,j}}{h_x} \\ &+ h_x k_{i, j-1/2} \frac{u_{i,j} - u_{i,j-1}}{h_y} - h_x k_{i, j+1/2} \frac{u_{i,j+1} - u_{i,j}}{h_y} \\ &= h_x h_y f_{ij} \end{aligned} \quad (10)$$

Rearranging and dividing by  $h_x h_y$ ,

$$\begin{aligned} f_{ij} &= -\frac{k_{i-1/2, j}}{h_x^2} u_{i-1,j} - \frac{k_{i, j-1/2}}{h_y^2} u_{i,j-1} \\ &+ \left( \frac{k_{i-1/2, j}}{h_x^2} + \frac{k_{i, j-1/2}}{h_y^2} + \frac{k_{i+1/2, j}}{h_x^2} + \frac{k_{i, j+1/2}}{h_y^2} \right) u_{i,j} \\ &- \frac{k_{i+1/2, j}}{h_x^2} u_{i+1,j} - \frac{k_{i, j+1/2}}{h_y^2} u_{i,j+1} \end{aligned} \quad (11)$$

With the above equation and using  $\bar{\Omega} = [0, 8] * [0, 4]$ , the following details can be obtained:

- Dimension of the  $A$  matrix =  $((N_x - 1) * (N_y - 1), (N_x - 1) * (N_y - 1))$  where  $N_x = \frac{8}{h}$  and  $N_y = \frac{4}{h}$  for a given  $h$
- There are 5 non-zero diagonals in the matrix and letting the main diagonal be the  $0^{th}$  diagonal and rightwards is positive, the diagonals are located at diagonal numbers:  $-(N_x - 1), -1, 0, +1, +(N_x - 1)$ .
- Yes. There are zero elements for diagonal number  $-1$  and  $+1$ . Every  $N_x^{th}$  element of these two diagonals are zero.
- The elements are computed simply by identifying their location in the global  $A$  matrix to obtain the  $i$  and  $j$  values and substituting into the equation 2 for their  $k$  values.

## b. Results

The following plots were obtained by solving the system

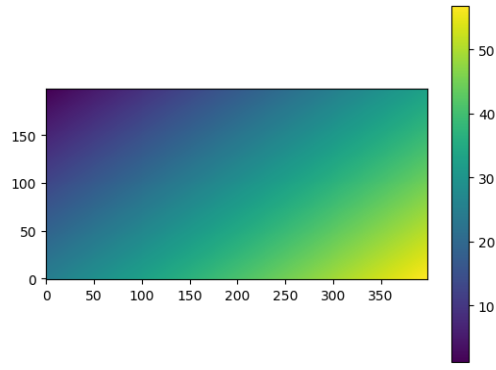


Figure 2:  $k(x, y)$

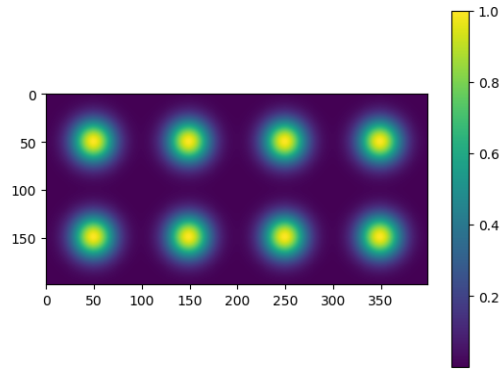


Figure 3:  $f(x, y)$

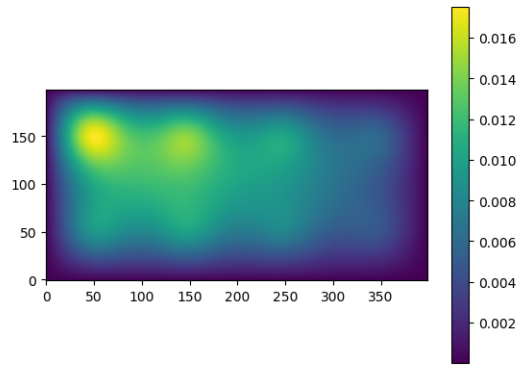


Figure 4:  $u(x, y)$

The solution shown in the figure 4 does make sense as starting from source term as shown in the figure 3, the deduction of  $k$  shown in the figure 2 has to result in top left corner being the least deducted area whereas the bottom right corner has to be the most deducted area and this is clearly shown in the figure 4.

## Problem 2

### a. Forward Euler and Backward Euler algorithms

Given

$$\frac{\partial u}{\partial t} - \Delta u = 0 \quad (12)$$

where

$$\Delta u = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \quad (13)$$

which can also be expressed as

$$\frac{\partial \mathbf{u}}{\partial t} = \mathbf{u}' = -A\mathbf{u} \quad (14)$$

for  $A$  is a negative Laplacian. Starting off with the Forward Euler algorithm, we first have

$$\mathbf{u}' = \frac{\mathbf{u}_h(t_{k+1}) - \mathbf{u}_h(t_k)}{h} = -A\mathbf{u}_h(t_k) \quad (15)$$

where  $\mathbf{u}_h$  is the iterations produced with time step  $h$  and  $k$  is the iteration number. The above equation is then rearranged into

$$\mathbf{u}_h(t_{k+1}) = \mathbf{u}_h(t_k) - hA\mathbf{u}_h(t_k) \quad (16)$$

Factorizing and simplifying,

$$\mathbf{u}_h(t_k) - hA\mathbf{u}_h(t_k) = [I - hA]\mathbf{u}_h(t_k) \quad (17)$$

for  $I$  is the identity matrix of the same size as  $A$ . Ultimately, we have

$$\mathbf{u}_h(t_{k+1}) = [I - hA]\mathbf{u}_h(t_k) \quad (18)$$

Next, for Backward Euler algorithm, instead of equation 15, we have

$$\mathbf{u}' = \frac{\mathbf{u}_h(t_k) - \mathbf{u}_h(t_{k-1})}{h} = -A\mathbf{u}_h(t_k) \quad (19)$$

Going through the same steps,

$$\mathbf{u}_h(t_k) = [I + hA]^{-1}\mathbf{u}_h(t_{k-1}) \quad (20)$$

### b. Time evolution of the heat equation

with  $h = 0.08$ , the following plots were produced with top rows representing time evolution results produced by the Forward Euler (FE) algorithm and bottom rows representing results produced by the Backward Euler (BE) algorithm.

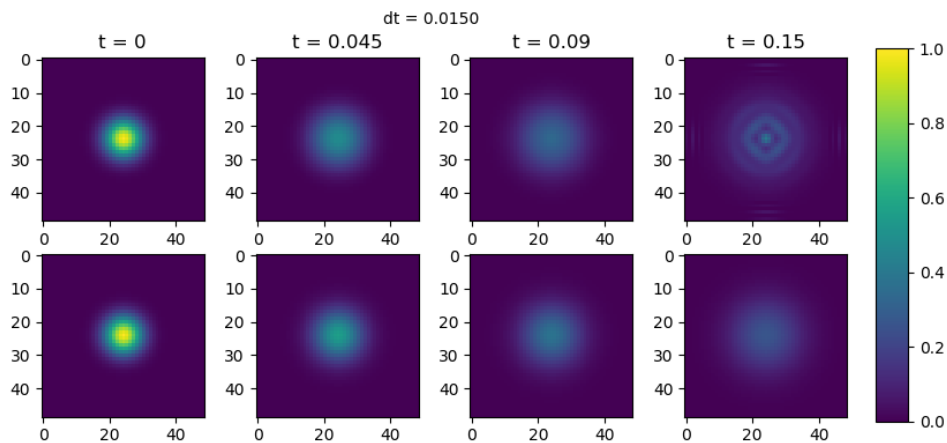


Figure 5:  $h = 0.08, dt = 0.015$

In the above figure, it can be observed that the FE becomes unstable at  $t = 0.15$  since not all points reach 0.

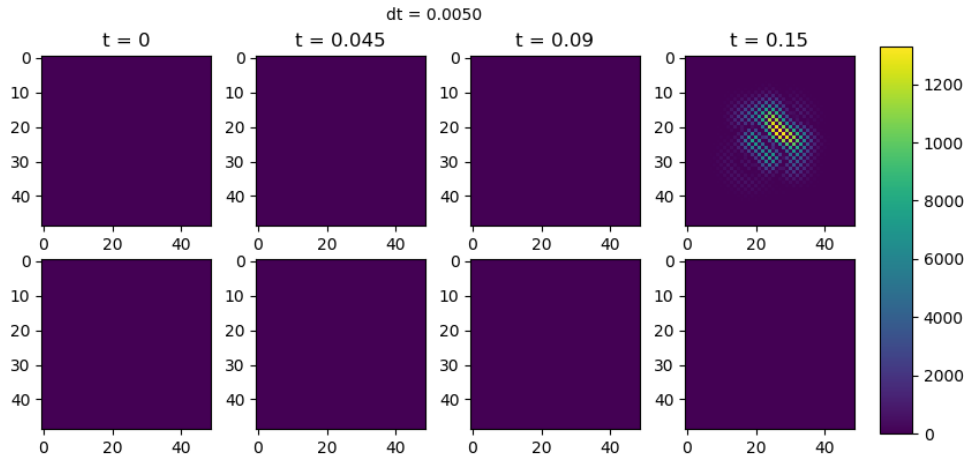


Figure 6:  $h = 0.08, dt = 0.005$

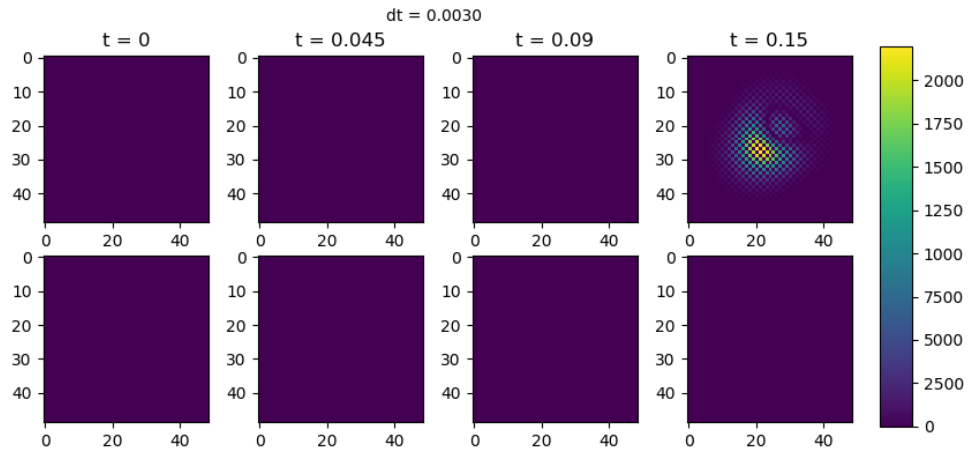


Figure 7:  $h = 0.08, dt = 0.003$

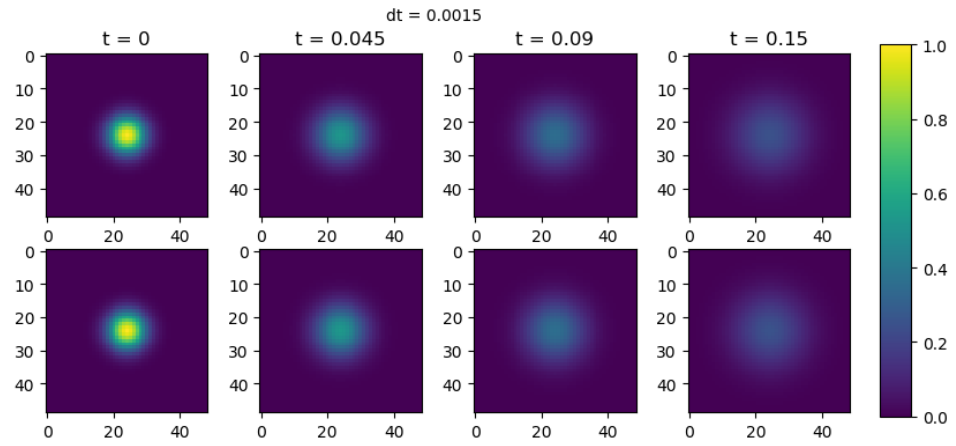


Figure 8:  $h = 0.08, dt = 0.0015$

Comparing the above 3 figures, the FE for  $dt = 0.005$  and  $dt = 0.003$  go unstable at  $t = 0.15$  whereas the FE for  $dt = 0.0015$  does not.

### c. Stability region of FE

With the equation 18 of FE, it can be seen that the absolute value of  $[I - h_t A]$  has to be less than or equal to 1 in order for the time evaluation to be stable. Thus,

$$|I - h_t A| \leq 1 \quad (21)$$

which is similar as the following expression

$$|1 - dt\lambda| \leq 1 \quad (22)$$

where  $dt = h_t$  and

$$\lambda_{k_x, k_y} = \frac{4}{h^2} \left[ \sin^2 \left( \frac{\pi k_x}{2N_x} \right) + \sin^2 \left( \frac{\pi k_y}{2N_y} \right) \right] \quad (23)$$

Since

$$0 \leq \sin^2 \left( \frac{\pi k_x}{2N_x} \right) + \sin^2 \left( \frac{\pi k_y}{2N_y} \right) \leq 2 \quad (24)$$

$$0 \leq \lambda \leq \frac{8}{h^2} \quad (25)$$

Substitution of 0 is redundant thus using  $\lambda = \frac{8}{h^2}$ , we have

$$|1 - \frac{8}{h^2} dt| \leq 1 \quad (26)$$

Splitting the above inequality into 2 inequality equations,

$$\begin{aligned} -1 &\leq 1 - \frac{8}{h^2} dt \\ dt &\leq \frac{h^2}{4} \end{aligned} \quad (27)$$

and

$$\begin{aligned} 1 - \frac{8}{h^2} dt &\leq 1 \\ dt &\geq 0 \end{aligned} \quad (28)$$

which is also redundant as long as  $dt$  is positive. Since for  $h = 0.08$  in equation 27,

$$dt \leq \frac{h^2}{4} = \frac{0.08^2}{4} = 0.0016 \quad (29)$$

$dt = 0.0015$  was clearly the only time step that was stable which was clearly depicted in the plots.

### d. Execution times

Below are comparison of the execution times (in seconds) between Forward Euler and Backward Euler algorithms. As for the time steps, stable  $dt$  was used for FE and 0.015 was used for BE.

|                     | h = 0.08            | h = 0.04            | h = 0.02          |
|---------------------|---------------------|---------------------|-------------------|
| FE (stable $dt$ )   | 0.03124833106994629 | 0.28354620933532715 | 4.802863836288452 |
| BE ( $dt = 0.015$ ) | 0.07811903953552246 | 0.42824387550354004 | 2.109236240386963 |

## Problem 3

### a. Plots for wavespeed of 1

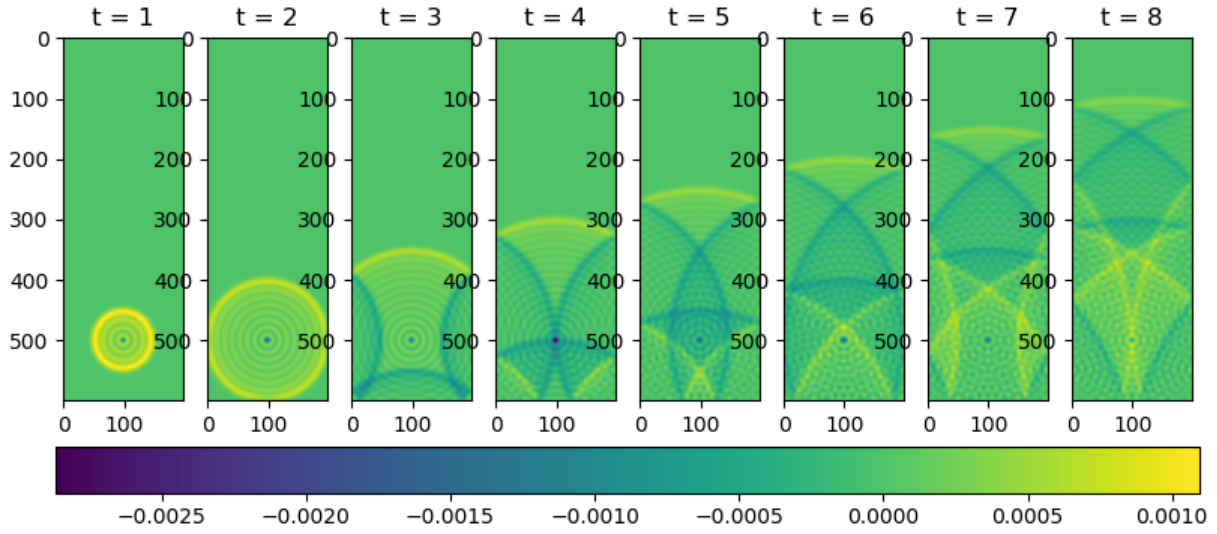


Figure 9: Plots of wave equation at 8 different timings for wavespeed,  $c = 1$

### b. Discussion

As it can be observed from the difference between the timesteps of the figure 9, the frequency of the wave is approximately  $4 \text{ units/s}$  which has already been mentioned in the problem as  $v$ . The wavelength can also be observed to be  $\frac{50 \cdot 10^{-3}}{2} = 0.25$  since there are 2 full wavelengths between the center and the circumference of  $t = 1$  plot and the distance covering them is  $50 \cdot 10^{-3}$ . With the formula

$$v = f\lambda \quad (30)$$

where  $v$  is the wavespeed,  $f$  is the frequency and  $\lambda$  is the wavelength, we get wavespeed of

$$\begin{aligned} v &= 4 * 0.25 \\ &= 1 \end{aligned} \quad (31)$$

Thus, we can say that the plots are correctly representative of wavespeed  $c = 1$

### c. Plots for wavespeed of 2

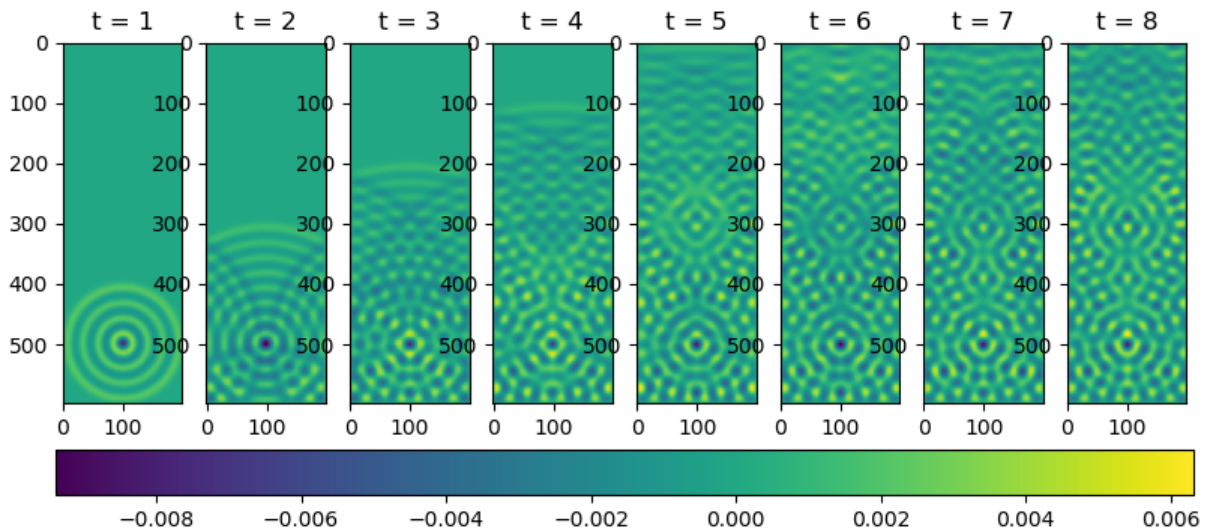


Figure 10: Plots of wave equation at 8 different timings for wavespeed,  $c = 2$

Comparing the figure 10 to the figure 9, it can be observed that for wavespeed  $c = 2$ , the magnitude of both the wave peak and trough have decreased, shown by the colour scheme. Also, standing waves can be observed from  $t = 3$  where the parts of the reflected waves cancel the progressive waves. Thus, the clear pattern shown by the figure 9 has either disappeared or become extremely vague due to the larger wavelengths.

## Problem 4

### a. Forward Euler method usage

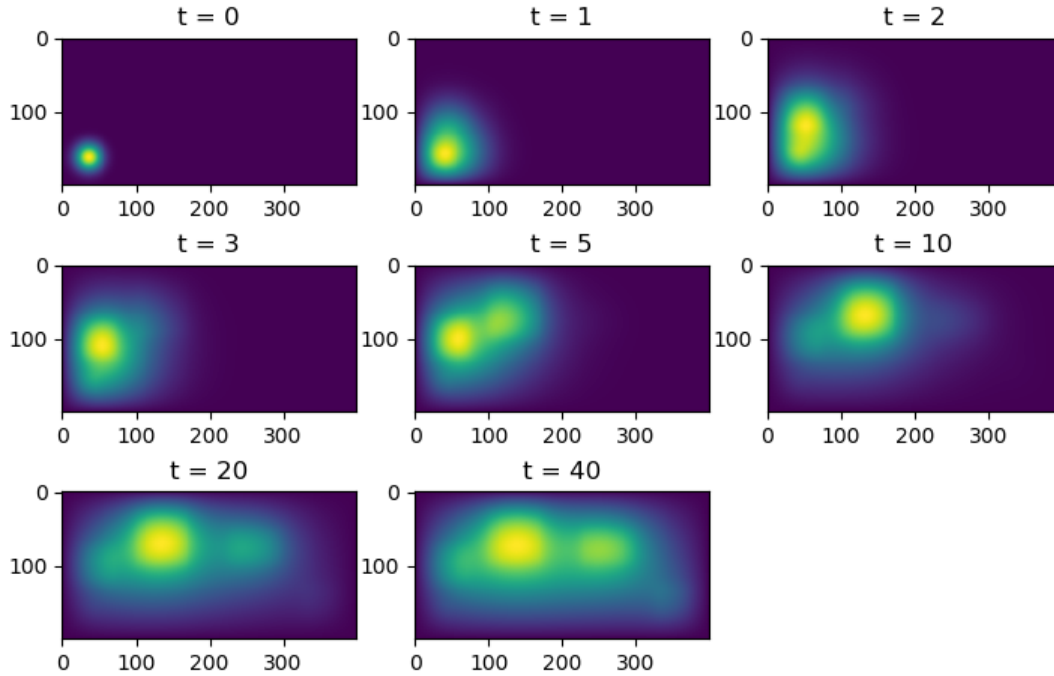


Figure 11: Solution using FE with  $h = 0.04$

### b. Convergence

The following plot describes the convergence of the Backward-Euler time-stepping method with the Newton-Raphson (NR) method with chosen time step of  $dt = 0.4$  which was found to be converging the most quickly.

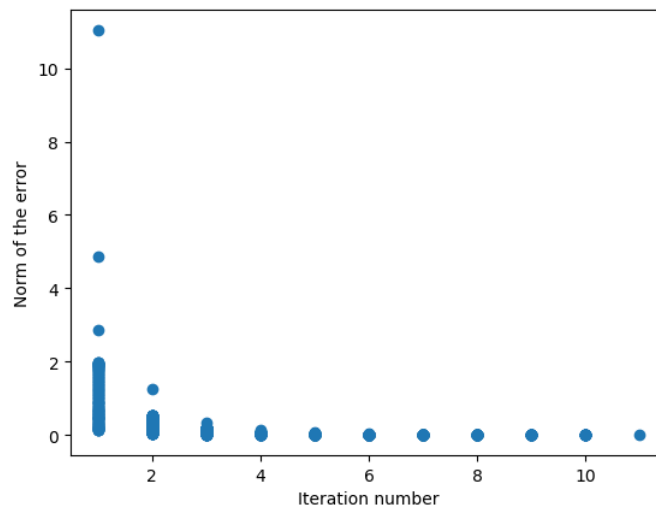


Figure 12: Convergence history of NR method for inner iterations with  $h = 0.04$  and  $dt = 0.4$

The following table compares the CPU execution times of Forward Euler method and Backward Euler method with Newton-Raphson method for the different spatial grid steps



|         | $h = 0.1$          | $h = 0.08$        | $h = 0.04$        |
|---------|--------------------|-------------------|-------------------|
| FE      | 2.302716016769409  | 5.479241847991943 | 96.46961212158203 |
| BE (NR) | 33.844868183135986 | 68.51465439796448 | 536.7286372184753 |

It can be concluded that the Backward Euler method with Newton-Raphson method integrated takes much longer than the Forward Euler counterpart. This can be due to the procedure of Newton-Raphson which involves solving of large sparse matrices for every inner iteration and this is time consuming whereas the Forward Euler method simply do not have such calculation.

### c. Change in $\alpha$

Stepping up the value of  $\alpha$  from 1 to 10, the following plots describing the solution were formed.

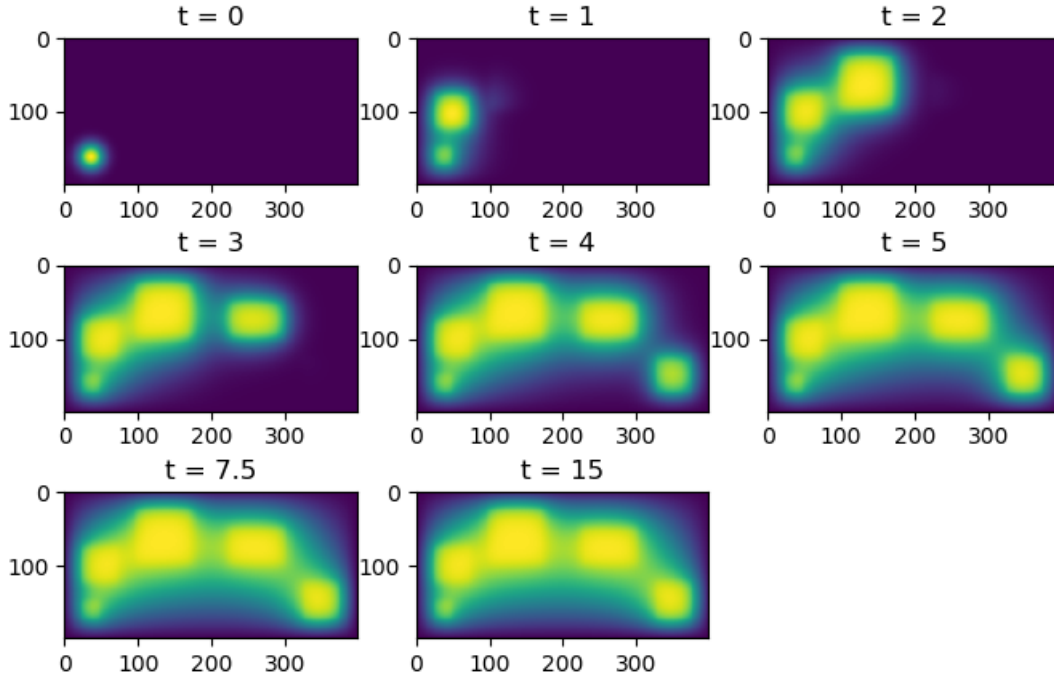


Figure 13:  $\alpha = 10$ ,  $h = 0.04$

It can be observed that the solution converged much faster than that of  $\alpha = 1$  since at  $t = 7.5s$ , it has already almost reached the final solution which is surely achieved at the shown time of  $t = 15s$ . Hence, it can be said that the speed of convergence increases with the magnitude of  $\alpha$

### d. Backward-Euler method with Picard's inner iterations

Regarding Picard's method, the iteration quickly runs into instability right at the beginning of the calculation and explodes, which concludes the program with an overflow.

## Problem 5

### a. Verification of GMRES solution

|                     | $\gamma = -40$        | $\gamma = 0$           | $\gamma = 40$          |
|---------------------|-----------------------|------------------------|------------------------|
| Absolute Difference | 4.037342676152687e-18 | 6.9510596825355345e-15 | 1.2066416308153462e-14 |

The table above records the absolute difference between the GMRES solution,  $x_k$ , and the one obtained by the direct method. It can be seen that for all three values of  $\gamma$ , the magnitudes of the absolute difference is extremely close to 0, verifying that the GMRES solution is legitimate.

## b. Plots

The following figure includes 4 plots which consist of the source function and 3 solutions, each with different values of  $\gamma$ .

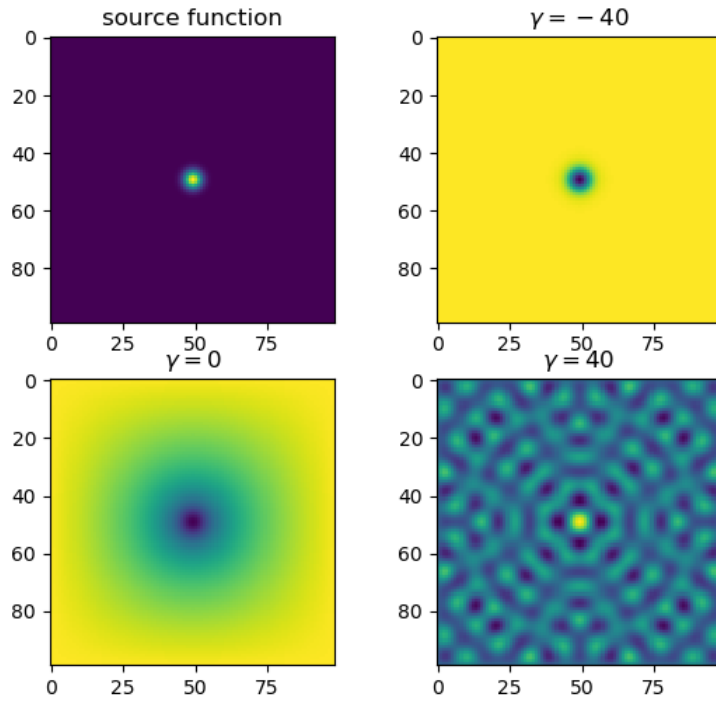


Figure 14: Source function and solutions of three  $\gamma$  values

## c. Residuals

The following figure displays the log of norm of the residuals produced by the GMRES algorithm

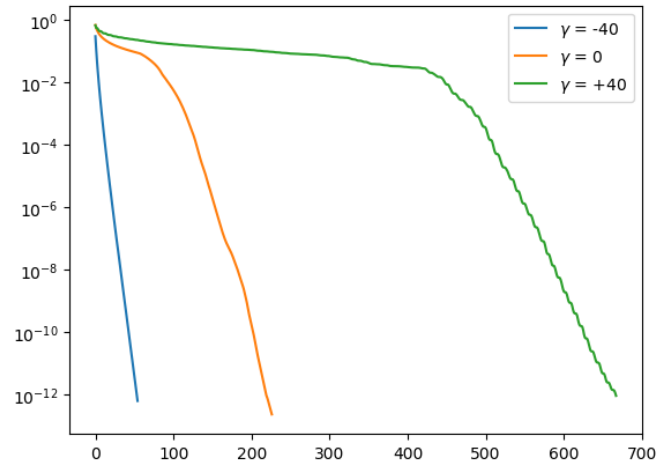


Figure 15: Norms of the residuals

## d. Matrix spectrum

Starting with the unaltered  $A$  matrix, its eigenvalues can be found with the following steps of calculation

$$\det(A - \lambda I) = 0 \quad (32)$$

where calculated values of  $\lambda$  are the eigenvalues of the A matrix. The above equation can also be expressed as such

$$(A - \lambda I)x = 0 \quad (33)$$

where vectors of  $x$  are the eigenvectors of the A matrix. On the other hand, for altered A matrix, the eigenvalues of the newly formed matrix can be found similarly as follows

$$\det(A + \gamma I - \bar{\lambda} I) = 0 \quad (34)$$

where  $\gamma$  is the alteration made to the A matrix, either  $-40$  for the Shielded Poisson equation or  $40$  for the Helmholtz equation, and the values of  $\bar{\lambda}$  are the newly derived eigenvalues. Just like the previous equation this can also be expressed as

$$(A + \gamma I - \bar{\lambda} I)x = 0 \quad (35)$$

Due to the fact that a mere shift in the values of the main diagonal of a matrix do not change its eigenvectors, we can then say that

$$\begin{aligned} A - \lambda I &= A + \gamma I - \bar{\lambda} I \\ \lambda I &= \gamma I - \bar{\lambda} I \\ \lambda &= \gamma - \bar{\lambda} \\ \bar{\lambda} &= \lambda + \gamma \end{aligned} \quad (36)$$

Hence, the new set of eigenvalues merely shift according to the value of  $\gamma$  implemented into the system and so does the matrix spectrum.

The following figure shows the values of the eigenvalues for different values of  $\gamma$ .

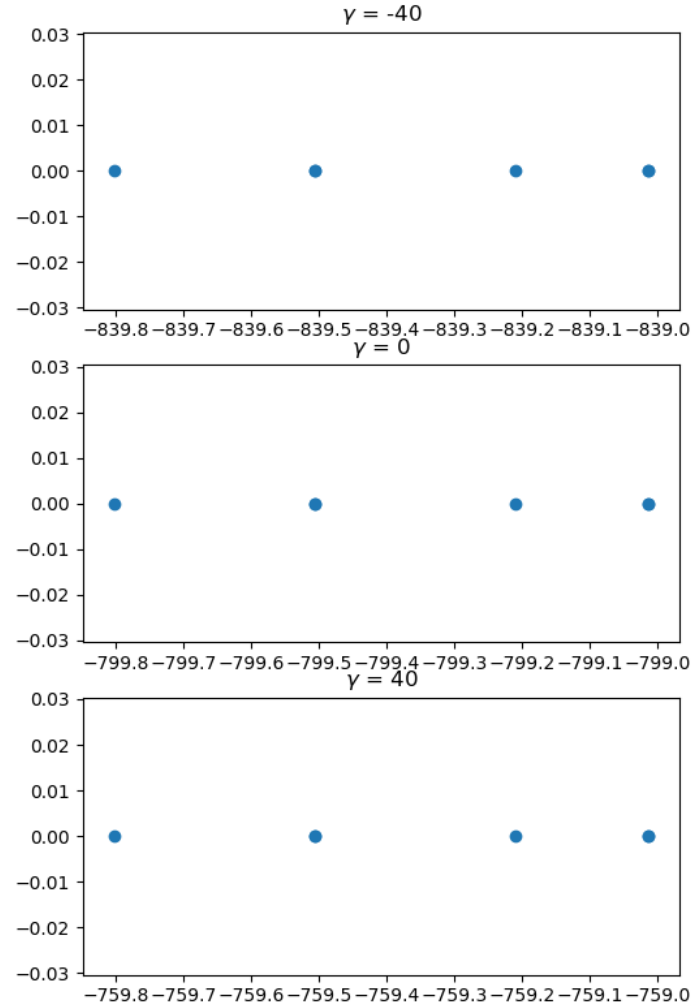


Figure 16: Eigenvalues for different values of  $\gamma$

To have a conclusion about the convergence rate, we can use the figure 15 and compare it against the figure above. It is observed that with an increase in the magnitude of the eigenvalues, the convergence rate exponentially decreases (converges more slowly).