

AMath 586

Vivek Gandhi

April 21, 2025

T1

T1 (a)

Using Taylor series to find $f(x+p)$ and $f(x-p)$, centered around x :

$$\begin{aligned}f(x+p) &= f(x) + pf'(x) + \frac{p^2}{2}f''(x) + \frac{p^3}{6}f'''(x) + \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5) \\f(x-p) &= f(x) - pf'(x) + \frac{p^2}{2}f''(x) - \frac{p^3}{6}f'''(x) + \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5)\end{aligned}$$

Rearranging the terms and adding them:

$$\begin{aligned}pf'(x) &= f(x+p) - f(x) - \frac{p^2}{2}f''(x) - \frac{p^3}{6}f'''(x) - \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5) \\pf'(x) &= f(x) - f(x-p) + \frac{p^2}{2}f''(x) - \frac{p^3}{6}f'''(x) + \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5) \\2pf'(x) &= f(x+p) - f(x-p) - \frac{p^3}{3}f'''(x) + \mathcal{O}(p^5) \\f'(x) &= \frac{f(x+p) - f(x-p)}{2p} - \frac{p^2}{6}f'''(x) + \mathcal{O}(p^4)\end{aligned}$$

Substituting $p = h/2$:

$$f'(x) = \frac{f(x+h/2) - f(x-h/2)}{h} - \frac{h^2}{24}f'''(x) + \mathcal{O}(h^4) \quad (1)$$

Using equation (1) with $f(x) = p(x_i)u'(x_i)$ to find $(p(x_i)u'(x_i))'$:

$$(p(x_i)u'(x_i))' = \frac{p(x_{i+1/2})u'(x_{i+1/2}) - p(x_{i-1/2})u'(x_{i-1/2}))}{h} - \frac{h^2}{24}(p(x_i)u'(x_i))''' + \mathcal{O}(h^4) \quad (2)$$

We know that $p \in C^4([a, b])$. But for the $\mathcal{O}(h^2)$ error term to be valid, we now require $u \in C^4([a, b])$ to also be true.

Using equation (1) again to expand $u'(x_{i+1/2})$ and $u'(x_{i-1/2})$:

$$\begin{aligned} u'(x_{i+1/2}) &= \frac{u(x_{i+1}) - u(x_i)}{h} - \frac{h^2}{24} u'''(x_{i+1/2}) + \mathcal{O}(h^4) \\ u'(x_{i-1/2}) &= \frac{u(x_i) - u(x_{i-1})}{h} - \frac{h^2}{24} u'''(x_{i-1/2}) + \mathcal{O}(h^4) \end{aligned}$$

Combining these like in equation (2):

$$\begin{aligned} & p(x_{i+1/2})u'(x_{i+1/2}) - p(x_{i-1/2})u'(x_{i-1/2}) \\ &= \frac{p(x_{i+1/2})u(x_{i+1}) - (p(x_{i+1/2}) + p(x_{i-1/2}))u(x_i) + p(x_{i-1/2})u(x_{i-1})}{h} \\ & \quad - \frac{h^2}{24} [p(x_{i+1/2})u'''(x_{i+1/2}) - p(x_{i-1/2})u'''(x_{i-1/2})] + \mathcal{O}(h^4) \end{aligned}$$

Here, the $\mathcal{O}(h^2)$ error term can be combined into one $\mathcal{O}(h^3)$ term.

Using the Taylor series approximation from before, with $f(x) = p(x)u'''(x)$:

$$\begin{aligned} f'(x) &= \frac{f(x+h/2) - f(x-h/2)}{h} - \frac{h^2}{24} f'''(x) + \mathcal{O}(h^4) \\ f(x+h/2) - f(x-h/2) &= hf'(x) + \frac{h^3}{24} f'''(x) + \mathcal{O}(h^5) \\ p(x_{i+1/2})u'''(x_{i+1/2}) - p(x_{i-1/2})u'''(x_{i-1/2}) &= h(p(x_i)u'''(x_i))' + \mathcal{O}(h^3) \end{aligned}$$

The error terms then become:

$$-\frac{h^2}{24} [p(x_{i+1/2})u'''(x_{i+1/2}) - p(x_{i-1/2})u'''(x_{i-1/2})] + \mathcal{O}(h^4) = -\frac{h^3}{24} (p(x_i)u'''(x_i))' + \mathcal{O}(h^4)$$

Plugging this back into (2):

$$\begin{aligned} (p(x_i)u'(x_i))' &= \frac{p(x_{i+1/2})u(x_{i+1}) - (p(x_{i+1/2}) + p(x_{i-1/2}))u(x_i) + p(x_{i-1/2})u(x_{i-1})}{h^2} \\ & \quad - \frac{h^2}{24} [(p(x_i)u'''(x_i))' + (p(x_i)u'(x_i))'''] + \mathcal{O}(h^3) \end{aligned}$$

Again, we see that $u \in C^4[a, b]$ needs to be true for the $\mathcal{O}(h^4)$ error term to be valid.

Using the following scheme requires dropping the $\mathcal{O}(h^2)$ and higher terms. Therefore, it has a consistency error (LTE) of $\mathcal{O}(h^2)$. As the other terms in the equation are exact, they don't create consistency error.

$$(p(x_i)u'(x_i))' \approx \frac{p(x_{i+1/2})U_{i+1} - (p(x_{i+1/2}) + p(x_{i-1/2}))U_i + p(x_{i-1/2})U_{i-1}}{h^2}$$

T1 (b)

The resulting system can be written as a linear system of the form $AU = F$:

$$A \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_N \end{bmatrix} = \begin{bmatrix} f(x_1) + \frac{p(x_{1/2})}{h^2} \alpha \\ f(x_2) \\ \vdots \\ f(x_{n-1}) \\ f(x_n) + p(x_{n+1/2}) \frac{\beta}{h^2} \end{bmatrix}$$

$$A = \frac{1}{h^2} R = \frac{1}{h^2} \begin{bmatrix} a_1 & b_1 & & & \\ b_1 & a_2 & & \ddots & \\ & \ddots & \ddots & \ddots & b_{n-1} \\ & & b_{n-1} & a_n & \end{bmatrix}$$

$$a_i = p(x_{i-1/2}) + p(x_{i+1/2}) + q(x_i)h^2$$

$$b_i = -p(x_{i+1/2})$$

My claim is that with a matrix R of size $n \times n$, the determinant $\det(R) = D_n$ is greater than or equal to $(n+1)\tilde{c}^n$, which is greater than 0.

Note that:

$$\begin{aligned} a_i &= p(x_{i-1/2}) + p(x_{i+1/2}) + q(x_i)h^2 \\ &\geq 2\tilde{c} + 0 = 2\tilde{c} \\ b_i &= -p(x_{i+1/2}) \leq \tilde{c} \end{aligned} \quad (p(x) \geq \tilde{c} > 0, \text{ and } q(x) \geq 0)$$

The n -th dimension determinant D_n can be computed recursively:

$$\begin{aligned} D_n &= a_n D_{n-1} - b_{n-1}^2 D_{n-2} \\ D_1 &= a_1 \geq 2\tilde{c} = (1+1)\tilde{c}^2 \\ D_2 &= a_1 a_2 - b_1^2 \geq 4\tilde{c}^2 - \tilde{c}^2 = (2+1)\tilde{c}^2 \end{aligned}$$

The 2 base cases D_1 and D_2 satisfy this inequality $D_n \geq (n+1)\tilde{c}^n$.

For induction, assume that this is true for the D_{n-1} and D_{n-2} cases for some $n \geq 3$.

$$\begin{aligned} a_n D_{n-1} &\geq 2\tilde{c}(n)\tilde{c}^{n-1} = 2n\tilde{c}^n \\ b_n^2 D_{n-2} &\leq \tilde{c}^2(n-1)\tilde{c}^{n-2} = (n-1)\tilde{c}^n \\ a_n D_{n-1} - b_n^2 D_{n-2} &\geq (2n - (n-1))\tilde{c}^n \\ D_n &\geq (n+1)\tilde{c}^n \end{aligned} \quad (\text{Subtracting the two})$$

As the $n-1$ and $n-2$ cases imply the n -th case, by induction, the $D_n \geq (n+1)\tilde{c}^n$ for all $n \geq 1$.

As $A = \frac{1}{h^2}R$ and $D_n = \det(R)$,

$$\begin{aligned}\det(A) &= \left(\frac{1}{h^2}\right)^n \det(R) \\ &\geq \frac{(n+1)}{h^{2n}} \tilde{c}^n\end{aligned}$$

As $n \geq 1$, $h > 0$, and $\tilde{c} > 0$, $\det(A) > 0$ must be true.

This implies that A is invertible, and the system $AU = F$ must have a unique solution.

T2

T2 (a)

$$\begin{aligned}-p(x_i) \frac{U_{i+1} - 2U_i + U_{i-1}}{h^2} - p'(x_i) \frac{U_{i+1} - U_{i-1}}{2h} + q(x_i)U_i &= f(x_i) \\ U_{i+1} \left(\frac{-p(x_i)}{h^2} - \frac{p'(x_i)}{2h} \right) + U_i \left(\frac{2p(x_i)}{h^2} + q(x_i) \right) + U_{i-1} \left(\frac{-p(x_i)}{h^2} + \frac{p'(x_i)}{2h} \right) &= f(x_i)\end{aligned}$$

Assuming 0 boundary conditions - this finite difference scheme can then be written as the system

$$\begin{aligned}A \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_N \end{bmatrix} &= \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix}, \quad A = \begin{bmatrix} b_1 & c_1 & & \\ a_2 & b_2 & \ddots & \\ & \ddots & \ddots & c_{n-1} \\ & & a_n & b_n \end{bmatrix} \\ a_i U_{i-1} + b_i U_i + c_i U_{i+1} &= f(x_i) \\ a_i &= \frac{-p(x_i)}{h^2} + \frac{p'(x_i)}{2h} \\ b_i &= \frac{2p(x_i)}{h^2} + q(x_i) \\ c_i &= \frac{-p(x_i)}{h^2} - \frac{p'(x_i)}{2h}\end{aligned}$$

As the finite difference scheme only uses U_i, U_{i+1} , and U_{i-1} , we get a tridiagonal matrix.

But as $c_i \neq a_{i+1}$, this matrix is not symmetric.

$p'(x_i)$ can be approximated using a similar finite difference scheme if necessary, but this doesn't change the symmetry of the matrix as $p(x_{i+1}) \neq p(x_i)$:

$$p'(x_i) \approx \frac{p(x_{i+1/2}) - p(x_{i-1/2})}{h}$$

T2 (b)

$$-u''(x) + a(x)u'(x) + b(x)u(x) = f(x), \quad a < x < b, \quad u(a) = \alpha, \quad u(b) = \beta$$

From **T1**, we have these:

$$f(x+p) = f(x) + pf'(x) + \frac{p^2}{2}f''(x) + \frac{p^3}{6}f'''(x) + \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5) \quad (1)$$

$$f(x-p) = f(x) - pf'(x) + \frac{p^2}{2}f''(x) - \frac{p^3}{6}f'''(x) + \frac{p^4}{24}f^{(4)}(x) + \mathcal{O}(p^5) \quad (2)$$

$$\begin{aligned} f'(x) &= \frac{f(x+p) - f(x-p)}{2p} - \frac{p^2}{6}f'''(x) + \mathcal{O}(p^4) \\ f'(x) &\approx \frac{f(x+p) - f(x-p)}{2p}, \quad \text{LTE} = \mathcal{O}(p^2) \end{aligned} \quad (3)$$

To derive an approximation of $f''(x)$, we add (1) and (2):

$$\begin{aligned} f(x+p) + f(x-p) &= 2f(x) + p^2f''(x) + \frac{p^4}{12}f^{(4)}(x) + \mathcal{O}(p^5) \\ p^2f''(x) &= f(x+p) - 2f(x) + f(x-p) - \frac{p^4}{12}f^{(4)}(x) + \mathcal{O}(p^5) \\ f''(x) &= \frac{f(x+p) - 2f(x) + f(x-p)}{p^2} - \frac{p^2}{12}f^{(4)}(x) + \mathcal{O}(p^3) \\ f''(x) &\approx \frac{f(x+p) - 2f(x) + f(x-p)}{p^2}, \quad \text{LTE} = \mathcal{O}(p^2) \end{aligned} \quad (4)$$

Using (3) to approximate $a(x)u'(x)$:

$$a(x)u'(x) \approx \frac{a(x_i)U_{i+1} - a(x_i)U_{i-1}}{2h}$$

Using (4) to approximate $-u''(x)$:

$$-u''(x) \approx \frac{-U_{i-1} + 2U_i - U_{i+1}}{h^2}$$

Combining them:

$$\begin{aligned} \frac{-U_{i-1} + 2U_i - U_{i+1}}{h^2} + \frac{a(x_i)U_{i+1} - a(x_i)U_{i-1}}{2h} + b(x_i)U_i &= f(x_i) \\ U_{i-1} \left(\frac{-1}{h^2} - \frac{a(x_i)}{2h} \right) + U_i \left(\frac{2}{h^2} + b(x_i) \right) + U_{i+1} \left(\frac{-1}{h^2} + \frac{a(x_i)}{2h} \right) &= f(x_i) \end{aligned}$$

This finite difference scheme has a consistency error of $\mathcal{O}(h^2)$.

It can be written as the following linear system:

$$A \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_N \end{bmatrix} = \begin{bmatrix} f(x_1) + \alpha \left(\frac{1}{h^2} + \frac{a(x_1)}{2h} \right) \\ f(x_2) \\ \vdots \\ f(x_{n-1}) \\ f(x_n) + \beta \left(\frac{1}{h^2} - \frac{a(x_n)}{2h} \right) \end{bmatrix}, \quad A = \begin{bmatrix} b_1 & c_1 & & \\ a_2 & b_2 & \ddots & \\ & \ddots & \ddots & c_{n-1} \\ & & a_n & b_n \end{bmatrix}$$

$$\begin{aligned} a_i U_{i-1} + b_i U_i + c_i U_{i+1} &= f(x_i) \\ a_i &= \frac{-1}{h^2} - \frac{a(x_i)}{2h} \\ b_i &= \frac{2}{h^2} + b(x_i) \\ c_i &= \frac{-1}{h^2} + \frac{a(x_i)}{2h} \end{aligned}$$

For this system to have a unique solution, the matrix A has to be non-singular (no $AU = 0$ solutions), or equivalently, if all its eigenvalues are non-zero.

We can find a sufficient condition to ensure this using Gershgorin's theorem.

For any matrix $A \in \mathbb{C}^{n \times n}$, define the disk:

$$G_i(A) = \{z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}|\}, \quad i = 1, \dots, n$$

Gerschgorin's theorem says that $\Lambda(A) \subset \bigcup_{i=1}^m G_i(A)$.

As our matrix only has 3 elements per row, the disks can be reduced to:

$$G_i(A) = \{z \in \mathbb{C} : |z - b_i| \leq |a_i| + |c_i|\}, \quad i = 1, \dots, n$$

This forms a circle around b_i with radius $|a_i| + |c_i|$. So for 0 to not be within this ring, we require $|b_i| > |a_i| + |c_i|$.

$$\begin{aligned} |b_i| &> |a_i| + |c_i| \\ \left| \frac{2}{h^2} + b(x_i) \right| &> \left| \frac{-1}{h^2} - \frac{a(x_i)}{2h} \right| + \left| \frac{-1}{h^2} + \frac{a(x_i)}{2h} \right| \\ \frac{2}{h^2} + b(x_i) &> \left| \frac{1}{h^2} + \frac{a(x_i)}{2h} \right| + \left| \frac{1}{h^2} - \frac{a(x_i)}{2h} \right| & (b(x_i) > 0) \\ &= A + B \end{aligned}$$

We can analyze the 4 cases separately based on the sign of A and B .

Case 1: $A \geq 0, B \geq 0$:

This happens when $\left| \frac{a(x_i)}{2h} \right| \leq \frac{1}{h^2}$, or $|a(x_i)| \leq \frac{2}{h}$.

$$\begin{aligned} \frac{2}{h^2} + b(x_i) &> \frac{1}{h^2} + \frac{a(x_i)}{2h} + \frac{1}{h^2} - \frac{a(x_i)}{2h} \\ &= \frac{2}{h^2} \end{aligned}$$

As $b(x_i) > 0$, this is always true.

Case 2: $A > 0$ and $B < 0$:

This happens when $\frac{a(x_i)}{2h} > \frac{1}{h^2}$, or $a(x_i) > \frac{2}{h}$.

$$\begin{aligned} \frac{2}{h^2} + b(x_i) &> \frac{1}{h^2} + \frac{a(x_i)}{2h} - \frac{1}{h^2} + \frac{a(x_i)}{2h} \\ &= \frac{a(x_i)}{h} \end{aligned}$$

As $a(x_i)/h > 2/h^2$, there's no condition that will ensure that $2/h^2 + b(x_i) > a(x_i)/h$ for all b .

Case 3: $A < 0$ and $B < 0$:

This happens when $\frac{a(x_i)}{2h} < -\frac{1}{h^2}$ and $\frac{a(x_i)}{2h} > \frac{1}{h^2}$. As $1/h^2 > 0$, these conditions can never be met.

Case 4: $A < 0$ and $B > 0$:

This happens when $\frac{a(x_i)}{2h} < -\frac{1}{h^2}$, or $a(x_i) < -\frac{2}{h}$.

$$\begin{aligned} \frac{2}{h^2} + b(x_i) &> -\frac{1}{h^2} - \frac{a(x_i)}{2h} + \frac{1}{h^2} - \frac{a(x_i)}{2h} \\ &= \frac{-a(x_i)}{h} \end{aligned}$$

As $-a(x_i)/h > 2/h^2$, there's no condition that will ensure that $2/h^2 + b(x_i) > -a(x_i)/h$ for all b .

Considering the left and right boundaries:

As a_1 and c_n don't exist, the terms A and B are omitted from the RHS. As they only add absolute values, to the lesser side of the inequality, omitting them only relaxes the condition further.

Conclusion:

The system always has a solution (A has no 0 eigenvalues) when $|a(x_i)| \leq \frac{2}{h}$.

So by choosing $h \leq \frac{2}{\max |a(x)|}$, we can ensure that the system has a consistent solution.

T3

T3 (a)

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

Consider $v(x) = u'(x)$. The problem becomes:

$$v' = -f(x), \quad x \in (0, 1), \quad v(0) = \alpha, \quad v(1) = \beta$$

From the fundamental theorem of calculus:

$$\begin{aligned} v(1) &= v(0) + \int_0^1 v'(x) dx \\ \beta &= \alpha - \int_0^1 f(x) dx \end{aligned}$$

So, a solution only exists when $\int_0^1 f(x) dx = \alpha - \beta$.

With a solution to $v(x) = u'(x)$, we still need a condition to ensure the uniqueness of the solution.

Given a solution $u(x)$, $w(x) = u(x) + C$ for any constant C will also be a solution, as $w'(x) = u'(x)$.

T3 (b)

From **T2**, we have:

$$f''(x) \approx \frac{f(x+p) - 2f(x) + f(x-p)}{p^2}, \quad \text{LTE} = \mathcal{O}(p^2)$$

This gives us the finite difference scheme:

$$\frac{-U_{i-1} + 2U_i - U_{i+1}}{h^2} = f(x_i)$$

To apply the left boundary condition, we create a ghost point U_{-1} , use it to derive formulas for $u''(0)$ and $u'(0)$, and combine them to get rid of the ghost point.

$$\frac{-U_{-1} + 2U_0 - U_1}{h^2} = -u''(0) = f(0) \tag{1}$$

We use a central difference formula as it is $\mathcal{O}(h^2)$

$$\begin{aligned} \frac{U_1 - U_{-1}}{2h} &= u'(0) = \alpha \\ U_{-1} &= U_1 - 2h\alpha \end{aligned} \tag{2}$$

Plugging (2) into (1)

$$\begin{aligned} \frac{2U_0 - 2U_1}{h^2} + \frac{2h\alpha}{h^2} &= f(0) \\ \frac{2U_0 - 2U_1}{h^2} &= f(0) - \frac{2\alpha}{h} \end{aligned}$$

Similarly for the right boundary condition, using ghost point U_{n+2} :

$$\frac{-U_{n+2} + 2U_{n+1} - U_n}{h^2} = -u''(1) = f(1) \quad (3)$$

$$\begin{aligned} \frac{U_{n+2} - U_n}{2h} &= u'(1) = \beta \\ U_{n+2} &= U_n + 2h\beta \end{aligned} \quad (4)$$

Plugging (4) into (3)

$$\begin{aligned} \frac{2U_{n+1} - 2U_n}{h^2} - \frac{2h\beta}{h^2} &= f(1) \\ \frac{2U_{n+1} - 2U_n}{h^2} &= f(1) + \frac{2\beta}{h} \end{aligned}$$

This can be combined into the following linear system:

$$\begin{array}{ccc} A & U & = F \\ \frac{1}{h^2} \begin{bmatrix} 2 & -2 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -2 & 2 \end{bmatrix} & \begin{bmatrix} U_0 \\ U_1 \\ \vdots \\ U_n \\ U_{n+1} \end{bmatrix} & = \begin{bmatrix} f(0) - 2\alpha/h \\ f(x_1) \\ \vdots \\ f(x_n) \\ f(1) + 2\beta/h \end{bmatrix} \end{array}$$

The matrix here is singular, so we don't have a guaranteed or unique solution.

It's easy to see why it's singular. Multiply the top and bottom rows by $1/2$ to get $[1, -1, \dots]$ and $[\dots, -1, 1]$.

Starting from the top, adding the i -th row to the $(i+1)$ -th row recursively leaves us with $[\dots, 1, -1, \dots]$ in each $(i+1)$ -th row, which eventually cancels out the last row.

We can fix this by introducing a constraint on $u(0)$, say $u(0) = U_0 = 0$. As for any solution $u(x)$, $u(x) + C$ is also a solution, the value of this constraint doesn't matter. We get a general family of solutions either way.

$$\begin{array}{ccc} A & U & = F \\ \frac{1}{h^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -2 & 2 \end{bmatrix} & \begin{bmatrix} U_1 \\ \vdots \\ U_n \\ U_{n+1} \end{bmatrix} & = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \\ f(1) + 2\beta/h \end{bmatrix} \end{array}$$

A is no longer singular. Recall how solutions only exist if $\beta = \alpha - \int_0^1 f(x)dx$. By removing our dependence on α and creating a consistent system, we ensure that the solution converges with $u'(0) \approx \beta + \int_0^1 f(x)dx$.

As the system is only defined on $(0, 1]$, we also manage to avoid using $f(0)$ or approximations of $u'(0), u''(0)$.

T4

T4 (a)

Showing that $L_h \Phi_{i,j} = -1$:

$$L_h U_{i,j} := -\frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} - \frac{U_{i,j-1} - 2U_{i,j} + U_{i,j+1}}{h^2}$$

The L_h operator is linear in $U_{i-1,j}, U_{i+1,j}, U_{i,j}, U_{i,j+1}, U_{i,j-1}$.

$$\begin{aligned} \Phi_{i,j} &= \frac{1}{4}((x_i - 1/2)^2 + (y_j - 1/2)^2) \\ L_h \Phi_{i,j} &= -\frac{\Phi_{i-1,j} - 2\Phi_{i,j} + \Phi_{i+1,j}}{h^2} - \frac{\Phi_{i,j-1} - 2\Phi_{i,j} + \Phi_{i,j+1}}{h^2} \\ &= \frac{4\Phi_{i,j}}{h^2} - \frac{1}{h^2} (\Phi_{i-1,j} + \Phi_{i+1,j} + \Phi_{i,j} + \Phi_{i,j+1} + \Phi_{i,j-1}) \end{aligned} \quad (1)$$

Note that

$$\begin{aligned} (a+h)^2 + (a-h)^2 &= a^2 + 2ah + h^2 + a^2 - 2ah + h^2 \\ &= 2(a^2 + h^2) \end{aligned} \quad (2)$$

Choosing some terms from the right term in (1):

$$\begin{aligned} \Phi_{i-1,j} + \Phi_{i+1,j} &= \frac{1}{4}((x_{i-1} - 1/2)^2 + (y_j - 1/2)^2) + \frac{1}{4}((x_{i+1} - 1/2)^2 + (y_j - 1/2)^2) \\ &= \frac{(y_j - 1/2)^2}{2} + \frac{(x_i - 1/2 - h)^2 + (x_i - 1/2 + h)^2}{4} \\ &= \frac{(y_j - 1/2)^2}{2} + \frac{(x_i - 1/2)^2 + h^2}{2} = 2\Phi_{i,j} + \frac{h^2}{2} \end{aligned} \quad (\text{Using (2)})$$

Choosing the other terms:

$$\begin{aligned} \Phi_{i,j+1} + \Phi_{i,j-1} &= \frac{1}{4}((x_i - 1/2)^2 + (y_{j-1} - 1/2)^2) + \frac{1}{4}((x_i - 1/2)^2 + (y_{j+1} - 1/2)^2) \\ &= \frac{(x_i - 1/2)^2}{2} + \frac{(y_j - 1/2 - h)^2 + (y_j - 1/2 + h)^2}{4} \\ &= \frac{(x_i - 1/2)^2}{2} + \frac{(y_j - 1/2)^2 + h^2}{2} = 2\Phi_{i,j} + \frac{h^2}{2} \end{aligned} \quad (\text{Using (2)})$$

Combining these back into (1):

$$\begin{aligned} L_h \Phi_{i,j} &= \frac{4\Phi_{i,j}}{h^2} - \frac{1}{h^2} (4\Phi_{i,j} + h^2) \\ &= 0 - 1 = -1 \end{aligned}$$

Therefore, $L_h \Phi_{i,j} = -1$.

Showing that $L_h\phi_{i,j} \leq 0$:

$$\phi_{i,j} := e_{i,j} + T\Phi_{i,j}$$

As the L_h operator is a linear operator,

$$\begin{aligned} L_h\phi_{i,j} &= L_h e_{i,j} + T(L_h\Phi_{i,j}) \\ &= T_{i,j} + T(L_h\Phi_{i,j}) \\ &\leq T + T(L_h\Phi_{i,j}) && (\text{As } 0 \leq |T_{i,j}| \leq T, T_{i,j} \leq T) \\ &= T(1 - 1) = 0 && (L_h\Phi_{i,j} = -1) \end{aligned}$$

Therefore, $L_h\phi_{i,j} \leq 0$.

Using the Discrete Maximum Principle:

Consider the system $-\Delta\phi(x, y) = p(x, y)$, defined on the same domain $\Omega = (0, 1)^2$, with numerical solution $\phi_{i,j}$.

Consider the value of $\Phi_{i,j}$ on the domain $[0, 1]^2$.

As Φ is a sum of squares in two different variables, minimizing or maximizing Φ can be split into two separate optimization problems.

$$\begin{aligned} \Phi_{i,j} &= \frac{(x_i - 1/2)^2}{4} + \frac{(y_j - 1/2)^2}{4}, \quad x_i, y_j \in [0, 1] \\ \max \Phi_{i,j} &= \max_{x_i} \frac{(x_i - 1/2)^2}{4} + \max_{y_j} \frac{(y_j - 1/2)^2}{4} \\ &= \frac{(1/2)^2}{4} + \frac{(1/2)^2}{4} = \frac{1}{8} && (\text{At } 0 \text{ or } 1) \\ \min \Phi_{i,j} &= \min_{x_i} \frac{(x_i - 1/2)^2}{4} + \min_{y_j} \frac{(y_j - 1/2)^2}{4} \\ &= 0 && (\text{At } (1/2, 1/2)) \end{aligned}$$

Our finite difference scheme assumes $L_h\phi_{i,j} = p(x_i, y_j) \leq 0$ (as shown above).

As $p(x, y) \leq 0$, we can use the Discrete Maximum Principle to assume that the maximum value of $\phi_{i,j}$ occurs on the boundaries.

$$\begin{aligned} \max_{\Omega} \phi_{i,j} &= \max_{\partial\Omega} (e_{i,j} + T\Phi_{i,j}) \\ &= T \max_{\partial\Omega} \Phi_{i,j} = \frac{T}{8} && (\text{The error } e_{ij} = 0 \text{ on the boundaries}) \\ \phi_{i,j} &\leq \max_{\Omega} \phi_{i,j} = \frac{T}{8} \\ e_{i,j} + T\Phi_{i,j} &\leq \frac{T}{8} \end{aligned}$$

As $\min \Phi_{i,j} = 0$

$$e_{i,j} \leq \frac{T}{8}$$

T4 (b)

Showing that $L_h \bar{\phi}_{i,j} \geq 0$:

$$\bar{\phi}_{i,j} := e_{i,j} - T\Phi_{i,j}$$

As the L_h operator is a linear operator,

$$\begin{aligned} L_h \bar{\phi}_{i,j} &= L_h e_{i,j} - T(L_h \Phi_{i,j}) \\ &= T_{i,j} - T(L_h \Phi_{i,j}) \\ &\geq -T - T(L_h \Phi_{i,j}) && (\text{As } 0 \leq |T_{i,j}| \leq T, -T_{i,j} \geq -T) \\ &= T(-1 + 1) = 0 && (L_h \Phi_{i,j} = -1) \end{aligned}$$

Therefore, $L_h \bar{\phi}_{i,j} \geq 0$.

Using the Discrete Minimum Principle:

Consider the system $-\Delta \bar{\phi}(x, y) = p(x, y)$, defined on the same domain $\Omega = (0, 1)^2$, with numerical solution $\bar{\phi}_{i,j}$.

Our finite difference scheme assumes $L_h \bar{\phi}_{i,j} = p(x_i, y_j) \geq 0$ (as shown above). As $p(x, y) \geq 0$, we can use the Discrete Minimum Principle to assume that the minimum value of $\bar{\phi}_{i,j}$ occurs on the boundaries.

$$\begin{aligned} \min_{\Omega} \bar{\phi}_{i,j} &= \min_{\partial\Omega} (e_{i,j} + T\Phi_{i,j}) \\ &= T \min_{\partial\Omega} \Phi_{i,j} = 0 && (\text{The error } e_{ij} = 0 \text{ on the boundaries}) \\ \bar{\phi}_{i,j} &\geq \min_{\Omega} \bar{\phi}_{i,j} = 0 \\ e_{i,j} - T\Phi_{i,j} &\geq 0 \end{aligned}$$

As $\max \Phi_{i,j} = \frac{T}{8}$

$$e_{i,j} \geq \frac{T}{8}$$

C1

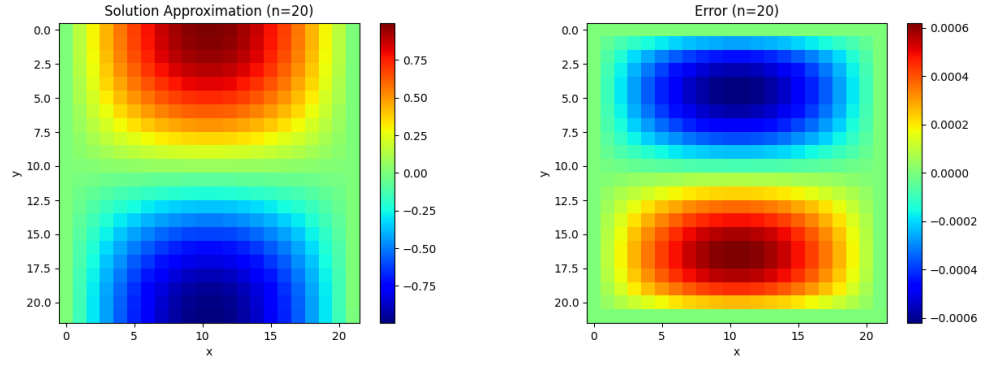


Figure 1: Numerical solution and error against analytic solution ($n = 20$)

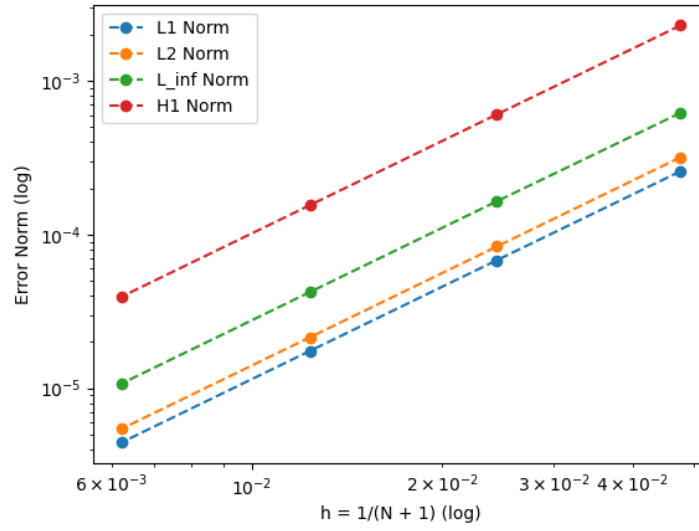


Figure 2: Solution errors with different values of h in different norms on a log-log plot

Norms	Slope
L_1	1.9971
L_2	1.9999
L_∞	1.9943
H^1	1.9979

Table 1: I used polyfit to find the slopes of each line on the log-log plot. This shows that the solution does indeed converge quadratically ($\mathcal{O}(h^2)$) in each norm.

```

1 import numpy as np
2
3 def NormL1(v, h):
4     return (h**2)*np.sum(np.abs(v[1:-1, 1:-1]))
5
6 def NormLInf(v, h):
7     return np.max(np.abs(v[1:-1, 1:-1]))
8
9 def NormL2(v, h):
10    v = v[1:-1, 1:-1]
11    return h * np.sqrt(np.sum(v**2))
12
13 def NormH1(v, h):
14    return np.sqrt(NormL2(v, h)**2 + DxV_xhNorm2(v, h) + DyV_yhNorm2(v, h))
15
16 def DxV_xhNorm2(v, h):
17    base = np.copy(v[1:-1, 1:])
18    base = base - v[1:-1, :-1]
19    return np.sum(base**2)
20
21 def DyV_yhNorm2(v, h):
22    return DxV_xhNorm2(v.transpose(), h)

```

Listing 1: 'Norms.py'

```

1
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # For A
6 from scipy.sparse import spdiags
7 # To solve AL = U
8 from scipy.sparse.linalg import gmres, spsolve
9 from scipy.linalg import solve
10 from Norms import *
11
12 # -----
13
14 do_matrix_viz = False # True to see matrix A's structure
15 do_viz = False       # True to plot the solution and error
16 mat_viz_mode = 0
17
18 # =====
19
20 def main():
21    Ns = [20, 40, 80, 160]

```

```

22 Hs = [1/(n+1) for n in Ns]
23
24 L1s = []
25 L2s = []
26 Linfs = []
27 H1s = []
28
29 for n in Ns:
30     data = gather_norms(n)
31
32     L1s.append(data[0])
33     L2s.append(data[1])
34     Linfs.append(data[2])
35     H1s.append(data[3])
36
37 get_slope = lambda errors: np.polyfit(np.log10(Hs), np.log10(errors), 1)[0]
38
39 plt.loglog(Hs, L1s, label='L1 Norm', linestyle='--', marker='o')
40 plt.loglog(Hs, L2s, label='L2 Norm', linestyle='--', marker='o')
41 plt.loglog(Hs, Linfs, label='L_inf Norm', linestyle='--', marker='o')
42 plt.loglog(Hs, H1s, label='H1 Norm', linestyle='--', marker='o')
43 plt.legend()
44 plt.xlabel('h = 1/(N + 1) (log)')
45 plt.ylabel('Error Norm (log)')
46 plt.show()
47
48 print(f'L1 slope - {get_slope(L1s)}')
49 print(f'L2 slope - {get_slope(L2s)}')
50 print(f'L_inf slope - {get_slope(Linfs)}')
51 print(f'H1 slope - {get_slope(H1s)}')
52
53 # =====
54
55 def gather_norms(n = 200):
56
57     h = 1/(n+1)
58     h2_n = (n+1)**2
59
60     # ----- Setup -----
61
62     xRange = np.linspace(0, 1, n+2)
63     yRange = np.linspace(0, 1, n+2)
64
65     # ----- Produce A -----
66
67     A = buildA(n, h2_n)
68     vizMatrix(A)
69
70     # ----- Produce F -----
71
72     BC_0 = np.sin(np.pi*xRange[1:-1])
73     BC_1 = -np.sin(np.pi*xRange[1:-1])
74     BC = np.concatenate((BC_0, np.zeros(n*(n-2)), BC_1))
75
76     initializer = lambda y, x: 2*(np.pi**2)*np.sin(np.pi*x)*np.cos(np.pi*y)
77     w0 = generate_w0(xRange[1:-1], yRange[1:-1], initializer)
78

```

```

79 F = deform(w0) + h2_n*deform(BC)
80
81 # ----- Solve for U -----
82
83 U_sol = spsolve(A, F)          # Sparse solve
84 # U_sol = solve(A.toarray(), F) # Dense solve
85 # U_sol = gmres(A, F)[0]       # GMRES
86
87 U_sol = reform(U_sol)
88
89 # ----- Append BCs -----
90
91 U = np.zeros((n+2, n+2))
92 U[1:1+U_sol.shape[0], 1:1+U_sol.shape[1]] = U_sol
93 U[0,:] = np.sin(np.pi*xRange)
94 U[-1,:] = -np.sin(np.pi*xRange)
95
96 viz(U, f'Solution (n={n})')
97
98 # ----- Compare against solution -----
99
100 initializer = lambda y, x: np.sin(np.pi*x)*np.cos(np.pi*y)
101 real = generate_w0(xRange, yRange, initializer)
102
103 error = real-U
104 viz(error, f'Error (n={n})')
105
106 return [NormL1(error, h), NormL2(error, h),
107         NormLInf(error, h), NormH1(error, h)]
108
109 # =====
110 # Helpers
111 # =====
112
113 # Unroll an nxn matrix into an mx1 one
114 def deform(A):
115     return A.reshape((A.size, 1))
116
117 # -----
118 # Reform an mx1 matrix into an nxn one
119 def reform(A):
120     n = np.sqrt(A.size).astype(int)
121     return A.reshape((n, n))
122
123 # -----
124 # Generate an nxn matrix across ranges with an initializer func
125 def generate_w0(xRange, yRange, func):
126     w0 = np.zeros((yRange.size, xRange.size))
127     for x in range(xRange.size):
128         for y in range(yRange.size):
129             w0[x,y] = func(xRange[x], yRange[y])
130     return w0
131
132 # -----
133 # Vizualize phi or w
134 def viz(A, title):
135     if do_viz:

```



```

136     plt.imshow(A, cmap='jet')
137     plt.colorbar()
138     plt.title(title)
139     plt.show()
140
141 # -----
142 # Visualize matrix A using either [0] plt.spy() or [1] plt.imshow()
143 def vizMatrix(mat):
144     global mat_viz_mode
145     if not do_matrix_viz:
146         return
147     if mat_viz_mode == 0:
148         plt.figure(5)
149         plt.spy(mat)
150     else:
151         plt.imshow(mat.toarray(), interpolation='none', cmap='binary')
152         plt.colorbar()
153         plt.title('Matrix Structure')
154         plt.show()
155
156 # =====
157 # A = -(dxx + dyy)
158 # =====
159
160 def buildA(n, h2_n):
161     m = n**2          # Size of A
162
163     e1 = np.ones(m)    # vector of 1s
164     e4 = 4*np.copy(e1) # vector of 4s
165
166     e_n = np.ones(n)
167     e_n[-1] = 0
168     e_n = -1*np.tile(e_n, n) # vector of -1s, with 0 at the end
169
170     e_m = np.roll(e_n, 1)
171
172     diagonals = [-1*e1, e_n, e4, e_m, -1*e1]
173     offsets = [-n, -1, 0, 1, n]
174
175     matA = spdiags(diagonals, offsets, m, m, format = 'csr')
176     return h2_n * matA
177
178 # =====
179
180 main()

```

Listing 2: 'HW1-C1.py