

(a) If  $v \in \text{span}\{w_1, \dots, w_n\}$ , Then there exists  
 $\{\alpha_i\}_{j=1}^n$ , s.t.  $v = \sum_{j=1}^n \alpha_j w_j$

Then,

$$\frac{\langle v, w_j \rangle}{\|w_j\|^2} = \frac{\alpha_j \langle w_j, w_j \rangle}{\|w_j\|^2} = \alpha_j$$

Hence

$$v = \sum_{j=1}^n \alpha_j w_j = \sum_{j=1}^n \frac{\langle v, w_j \rangle}{\|w_j\|^2} w_j$$

(b) i) The CG Algorithm can be interpreted as minimizing function  $\phi(u)$  over the space

$U_0 + \text{span}\{p_0, p_1, \dots, p_{k-1}\}$  in the  $k$ th iteration.

As long as the optimizer  $u^*$  is in the subspace  $U_0 + \text{span}\{p_0, p_1, \dots, p_{n^*-1}\}$

Then the algorithm will converge in  $n^*$  steps.

The subspace  $U_0 + \text{span}\{p_0, p_1, \dots, p_{n^*-1}\}$

not necessarily has to be the whole space  $\mathbb{R}^N$  in order to contain optimizer  $u^*$  in it.

ii) We firstly prove  $\langle p_1, p_0 \rangle_A = 0$

$$\begin{aligned}\langle P_1, P_0 \rangle_A &= \left\langle r_1 - \frac{\langle r_1, P_0 \rangle_A}{\|r_0\|_A^2} P_0, P_0 \right\rangle_A \\ &= \langle r_1, P_0 \rangle_A - \frac{\langle r_1, P_0 \rangle_A}{\|r_0\|_A^2} \|P_0\|_A^2 = \langle r_1, P_0 \rangle_A - \langle r_1, P_0 \rangle_A \\ &= 0\end{aligned}$$

suppose when  $n=k$

$$\langle P_n, P_j \rangle_A = 0 \quad \forall \quad 0 \leq j < n \leq n^* - 1$$

when  $n=k+1$

$$\langle P_n, P_j \rangle_A = \left\langle r_n - \sum_{i=0}^{n-1} \frac{\langle r_n, P_i \rangle_A}{\|P_i\|_A^2} P_i, P_j \right\rangle_A$$

$$= \langle r_n, P_j \rangle_A - \sum_{i=0}^{n-1} \frac{\langle r_n, P_i \rangle_A}{\|P_i\|_A^2} \langle P_i, P_j \rangle_A$$

By induction assumption

$$= \langle r_n, P_j \rangle_A - \frac{\langle r_n, P_j \rangle_A}{\|P_j\|_A^2} \langle P_j, P_j \rangle_A$$

$$= \langle r_n, P_j \rangle_A - \langle r_n, P_j \rangle_A = 0$$

(c) (i) Using linear algebra results, we have

$$A = \sum_{n=1}^N \lambda_n \phi_n \phi_n^T$$

$$\text{Then } \langle Av, w \rangle = \left\langle \sum_{n=1}^N \lambda_n \phi_n \phi_n^T v, w \right\rangle$$

$$= \sum_{n=1}^N \langle \lambda_n \phi_n, w \rangle \langle \phi_n, v \rangle$$

$$= \sum_{n=1}^N \lambda_n \langle v, \phi_n \rangle \langle \phi_n, w \rangle$$

ii) Since  $A$  is a positive definite matrix

$$\text{Then } \phi_n^T A \phi_n > 0 \quad 1 \leq n \leq N$$

which implies

$$\phi_n^T \lambda_n \phi_n = \lambda_n > 0 \quad \forall n$$

iii) We only need to prove

$$\lambda_1 \leq \frac{\langle Av, v \rangle}{\|v\|^2} \leq \lambda_N$$

Without loss of Generality  $\|v\|^2 = 1$

$$\begin{aligned} \langle Av, v \rangle &= v^T A v = v^T \sum_{n=1}^N \lambda_n \phi_n \phi_n^T v \\ &= \sum_{n=1}^N \lambda_n \langle v, \phi_n \rangle \langle v, \phi_n \rangle \end{aligned}$$

Also notice that

$$\sum_{n=1}^N \langle v, \phi_n \rangle \langle v, \phi_n \rangle = v^T [\phi_1, \dots, \phi_N] \begin{bmatrix} \phi_1^T \\ \vdots \\ \phi_N^T \end{bmatrix} v$$

$$= v^T v = 1$$

Hence we have

$$\begin{aligned} \lambda_1 \sum_{n=1}^N \langle v, \phi_n \rangle \langle v, \phi_n \rangle &\leq \sum_{n=1}^N \lambda_n \langle v, \phi_n \rangle \langle v, \phi_n \rangle \\ &\leq \lambda_N \sum_{n=1}^N \langle v, \phi_n \rangle \langle v, \phi_n \rangle \end{aligned}$$

$$\Rightarrow \lambda_1 \leq \sum_{n=1}^N \lambda_n \langle v, \phi_n \rangle \langle v, \phi_n \rangle \leq \lambda_N$$

$$\Rightarrow \lambda_1 \leq \langle Av, v \rangle \leq \lambda_N \quad \text{when } \|v\| = 1$$

iv) We only need to show  $\|Av\|^2 \leq \lambda_N^2 \|v\|^2$   
 notice  $A^2$  has the same eigenvectors as  $A$   
 and corresponding eigenvalues  $\lambda_N^2$

$$\|Av\|^2 = v^T A^2 v = \langle A^2 v, v \rangle \leq \lambda_N^2 \|v\|^2$$

(d)  $P_{n+1} = r_{n+1} + \beta_n P_n$

$$r_{n+1} = r_n - \alpha_n A P_n$$

$$r_n = P_n - \beta_{n-1} P_{n-1}$$

Add together, we have

$$P_{n+1} = (1 + \beta_n) P_n - \alpha_n A P_n - \beta_{n-1} P_{n-1}$$

(e) Define characteristic polynomial of  $A$  as

$$P(\lambda) = \det(\lambda I_n - A)$$

$P(\lambda)$  is a  $N$  degree polynomial

$$P(\lambda) = \lambda^N + \alpha_{N-1} \lambda^{N-1} + \dots + \alpha_1 \lambda + \alpha_0$$

By Cayley - Hamilton theorem,

$$P(A) = A^N + \alpha_{N-1} A^{N-1} + \dots + \alpha_1 A + \alpha_0 I = 0$$

Hence,

$A^N$  is a linear combination of  $I, A, A^2, \dots, A^{N-1}$

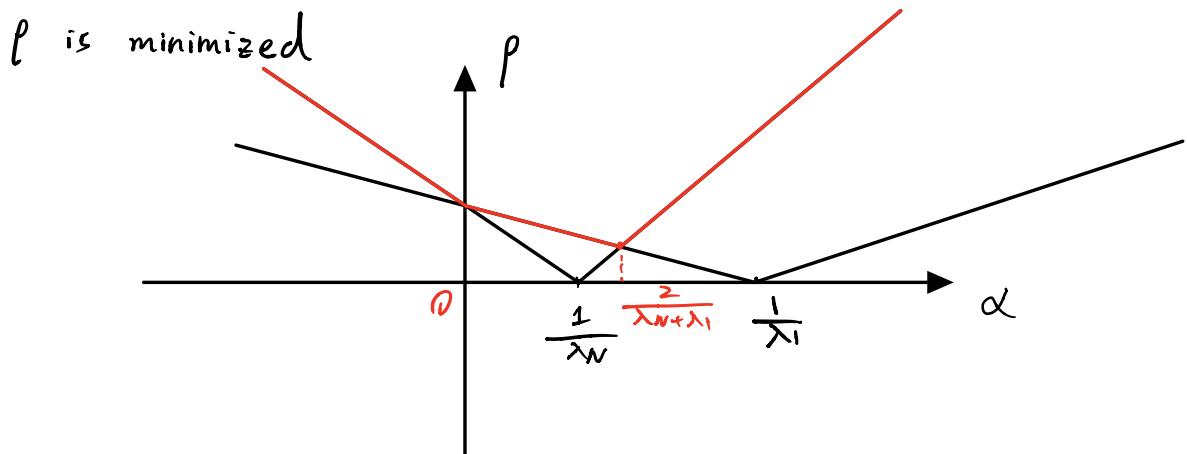
$$\begin{aligned}
 (f) \text{ i) } u_{n+1} - \mu &= u_n - \mu + \alpha (Au - A\mu) \\
 &= u_n - \mu - \alpha A(u_n - \mu) \\
 &= (I - \alpha A)(u_n - \mu) \\
 &= (I - \alpha A)e_n
 \end{aligned}$$

ii) Notice that the eigenvalues of  $I - \alpha A$  should be  $\{1 - \alpha \lambda_j\}_{j=1}^N$

$$\begin{aligned}
 \|e_{n+1}\| &= \sqrt{e_n^\top (I - \alpha A)^2 e_n} \leq \sqrt{\max_{1 \leq j \leq N} (1 - \alpha \lambda_j)^2 \|e_n\|^2} \\
 &= \max_{1 \leq j \leq N} |1 - \alpha \lambda_j| \|e_n\| = p \|e_n\|
 \end{aligned}$$

$$\text{iii) } p = \max_{1 \leq j \leq N} |1 - \alpha \lambda_j| = \max\{|1 - \alpha \lambda_1|, |1 - \alpha \lambda_N|\}$$

$$\text{when } \alpha = \frac{2}{\lambda_1 + \lambda_N} \quad |1 - \alpha \lambda_1| = |1 - \alpha \lambda_N|$$



$$\text{at } \frac{2}{\lambda_N + \lambda_1}$$

iv) by iii), it is obvious that

$$\rho \leq \frac{C - c}{C + c} = \frac{k' - 1}{k' + 1} < 1$$

$$\text{since } k' = \frac{C}{c} > 1$$

(g) i) This is obvious

$$\begin{aligned} \text{ii}) \quad r_{n+1} &= r_n - \alpha_n A P_n \\ &= r_n - \alpha_n A (r_n + \beta_{n-1} p_{n-1}) \\ &= r_n - \alpha_n A r_n - \alpha_n \beta_{n-1} A P_{n-1} \\ &= r_n - \alpha_n A r_n + \alpha_n \beta_{n-1} \left( \frac{r_n - r_{n-1}}{\alpha_{n-1}} \right) \\ &= r_n - \alpha_n A r_n + \frac{\alpha_n \beta_{n-1}}{\alpha_{n-1}} (r_n - r_{n-1}) \end{aligned}$$

$$\text{iii}) \quad r_1 = r_0 - \alpha_0 A r_0$$

$$\frac{r_1}{\|r_0\|} = q_0 - \alpha_0 A q_0$$

$$\frac{\|r_1\|}{\|r_0\|} = \sqrt{\beta_0}$$

$$q_1 \sqrt{\beta_0} = q_0 - \alpha_0 A q_0$$

$$\Rightarrow q_1 \frac{\sqrt{\beta_0}}{\alpha_0} = \frac{1}{\alpha_0} q_0 - A q_0$$

$$\Rightarrow A q_0 = \gamma_0 q_0 - \delta_0 q_1$$

$$\begin{aligned}
 r_{n+1} &= r_n - \alpha_n A r_n + \frac{\alpha_n \beta_{n-1}}{\alpha_{n-1}} (r_n - r_{n-1}) \\
 \Rightarrow \frac{r_{n+1}}{\|r_n\|} &= q_n - \alpha_n A q_n + \frac{\alpha_n \beta_{n-1}}{\alpha_{n-1}} \left( q_n - \frac{r_{n-1}}{\|r_n\|} \right) \\
 \Rightarrow q_{n+1} \sqrt{\beta_n} &= q_n - \alpha_n A q_n + \frac{\alpha_n \beta_{n-1}}{\alpha_{n-1}} \left( q_n - q_{n-1} / \sqrt{\beta_{n-1}} \right) \\
 \Rightarrow q_{n+1} \delta_n &= p_n q_n - A q_n - \delta_{n-1} q_{n-1} \\
 \Rightarrow A q_n &= -\delta_{n-1} q_{n-1} + p_n q_n - \delta_n q_{n+1}
 \end{aligned}$$

iv)  $A q_n = [q_{n-1} \ q_n \ q_{n+1}] \begin{bmatrix} -\delta_{n-1} \\ p_n \\ -\delta_n \end{bmatrix}$

Hence

$$A q_j = Q_n \cdot \begin{bmatrix} 0 \\ \vdots \\ -\delta_{j-1} \\ q_j \\ \vdots \\ -\delta_j \\ \vdots \\ 0 \end{bmatrix}$$

and  $A q_{n-1} = Q_n \begin{bmatrix} 0 \\ \vdots \\ -\delta_{n-2} \\ p_{n-1} \end{bmatrix} - \delta_n q_n$

Combine these equations, we have

$$A Q_n = Q_n T_n - \delta_{n-1} q_n e_n^T$$

$$\text{v) } Q_n^T A Q_n = Q_n^T Q_n T_n - Q_n^T S_{n-1} q_n e_n^T \\ = T_n$$

$Q_n^T q_n = 0$  since  $q_n$  orthogonal to  
 $\text{span}\{q_0, q_1, \dots, q_{n-1}\}$

## 1 Problem B:

Here we use the interp1 function in Matlab to find the threshold numerically. Below are some selected output of uniform norm difference.

Sample size N	Difference in the uniform norm
90	0.0118
91	0.0120
92	0.0113
93	0.0115
94	0.0109
95	0.0110
96	0.0104
97	0.0106
98	0.0100
99	0.0102
100	0.0097
101	0.0098
102	0.0093

From the above numerical output, we can clearly see that,  $N = 100$  is the smallest value of  $N$  such that  $f$  differs from its linear interpolant by at most  $10^{-2}$  in the uniform norm.

## 2 Problem C:

(a) In the coding session, we set  $\Delta t = 1/2\Delta x$ . The time axes range from 0 to 0.5 i.e.  $T = 0.5$ . We solve the 2D-wave equation in space  $[0, 1] \times [0, 1] \times [0, T]$ . In order for comparing the error, the "true solution" is computed numerically with grid spacing  $\Delta x = 1/1024$ .

The updated formula should be

$$\begin{aligned} \frac{U_{i,j}^{n+1} - 2U_{i,j}^n + U_{i,j}^{n-1}}{\Delta t} &= \frac{1}{\Delta x^2} [U_{i+1,j}^n + U_{i,j+1}^n - 4U_{i,j}^n + U_{i-1,j}^n + U_{i,j-1}^n] \\ U_{i,j}^0 &= 0 \\ \frac{U_{i,j}^1 - U_{i,j}^0}{\Delta t} &= f(x_i)f(y_j) \end{aligned} \tag{1}$$

The final plot is given below

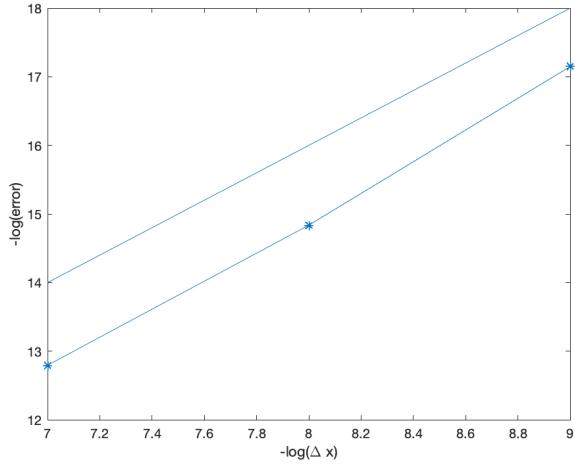


Figure 1: error vs grid spacing

We only plot three points here, grid spacing  $\Delta x = 1/128$ ,  $\Delta x = 1/256$  and  $\Delta x = 512$ . The blue straight line above is  $y = 2x$  plotted for comparison.

(b) Here we firstly obtain the following equation

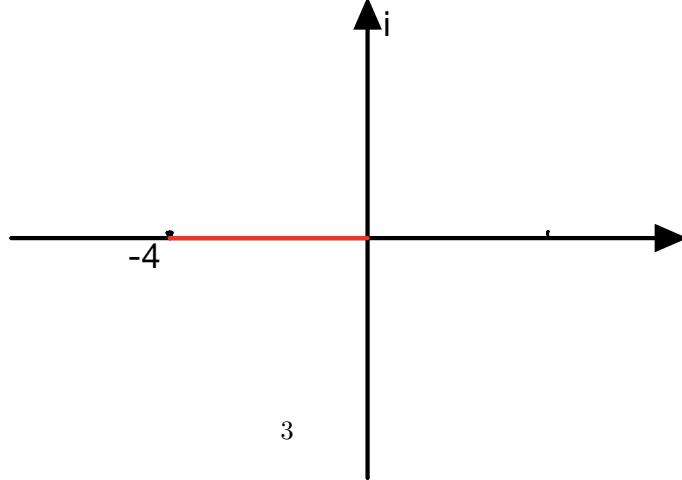
$$\frac{y^{n+1} - 2y^n + y^{n-1}}{\Delta t^2} = \lambda y^n$$

set  $y^n = \rho^n$  and rearrange, we will have

$$\rho^2 - (\lambda(\Delta t)^2 + 2)\rho + 1 = 0$$

We need the  $|\rho| \leq 1$ , which will result in a stability region for  $\lambda(\Delta t)^2$  as

$$\begin{aligned} -4 &\leq \operatorname{Re}(\lambda(\Delta t)^2) \leq 0 \\ \operatorname{Im}(\lambda(\Delta t)^2) &= 0 \end{aligned}$$



(c) Just like what we did in the previous HW, we rewrite the problem using Kronecker product.

$$U'' = \frac{1}{\Delta^2 x} (A \otimes I_N + I_N \otimes A) U$$

Here we have

$$A = \begin{pmatrix} -2 & 1 & 0 & 0 & \cdots & 0 \\ 1 & -2 & 1 & 0 & \cdots & 0 \\ 0 & 1 & -2 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & 1 & -2 \end{pmatrix} \quad (2)$$

By the previous HW, the eigenvalues of  $A \otimes I_N + I_N \otimes A$  should be  $\lambda_i + \lambda_j$  where  $\lambda_i$  and  $\lambda_j$  are eigenvalues of  $A$ .

Hence the eigenvalue of matrix  $A \otimes I_N + I_N \otimes A$  is bounded by -8. Combine this with the result in b, we obtain

$$\frac{\Delta^2 t}{\Delta^2 x} \leq 1/2$$

(d) We use  $\exp(ikj_1 h)\exp(ikj_2 h)$  to replace  $U_{j_1, j_2}$  in the updated equation and after some rearrangement, we finally observe the following

$$\begin{aligned} g(k) - 2 + \frac{1}{g(k)} \\ = \frac{\Delta^2 t}{\Delta^2 x} [\exp(ikh) + \exp(ikh) + \exp(-ikh) + \exp(-ikh) - 4] \end{aligned}$$

Hence

$$\begin{aligned} g(k) &= \frac{\Delta^2 t}{\Delta^2 x} (4 \cos kh - 4) - 1/g(k) + 2 \\ &= -8 \frac{\Delta^2 t}{\Delta^2 x} \sin^2 \frac{kh}{2} - \frac{1}{g(k)} + 2 \end{aligned}$$

We need  $|g(k)| \leq 1$ , and thus we still can derive that

$$\frac{\Delta^2 t}{\Delta^2 x} \leq \frac{1}{2}$$

Consistent with which we obtained in (c).

(e) The Fourier series is given below

$$u(x, y, t) = \sum_{n=1}^{\infty} \sum_{m=1}^{\infty} \sin m\pi x \sin n\pi y (B_{mn} \cos \lambda_{mn} t + B_{mn}^* \sin \lambda_{mn} t) \quad (3)$$

Here  $\lambda_{mn} = \sqrt{(m\pi)^2 + (n\pi)^2}$

The  $B_{mn}$  are determined by the boundary conditions. The physics extra terms are dispersive.

### **3 Github Link**

[https://github.com/BZHANG327/HW2\\_Numerical\\_Analysis.git](https://github.com/BZHANG327/HW2_Numerical_Analysis.git)