

Mathematical Tripos Part II: Michaelmas Term 2023

Numerical Analysis – Lecture 3

Let \hat{u} be the exact solution of the Poisson equation, and let $\hat{u}_{i,j} = \hat{u}(ih, jh)$ be its values on the grid. Let

$$e_{i,j} = \hat{u}_{i,j} - u_{i,j} \quad (1.7)$$

be the pointwise error of the 5-point formula. Set $\mathbf{e} = (e_{i,j}) \in \mathbb{R}^N$ where $N = M^2$, and for $\mathbf{x} \in \mathbb{R}^N$ let $\|\mathbf{x}\|_2$ be the Euclidean norm of the vector \mathbf{x} :

$$\|\mathbf{x}\|_2^2 = \sum_{k=1}^N |x_k|^2 = \sum_{i=1}^M \sum_{j=1}^M |x_{i,j}|^2.$$

Theorem 1.11 Assume the solution \hat{u} of Poisson's equation is C^4 and let

$$c = \frac{1}{12} \max_{0 < x, y < 1} \left| \frac{\partial^4 \hat{u}}{\partial x^4}(x, y) \right| + \left| \frac{\partial^4 \hat{u}}{\partial y^4}(x, y) \right| > 0. \quad (1.8)$$

Then the error vector \mathbf{e} defined in (1.7) satisfies

$$\|\mathbf{e}\|_2 \leq (c/8)h.$$

Proof. For a C^4 univariate function $g : (a, b) \rightarrow \mathbb{R}$, the finite-difference approximation of $g''(x)$ for $x \in (a+h, b-h)$ satisfies

$$|g''(x) - (g(x+h) + g(x-h) - 2g(x))/h^2| \leq \frac{h^2}{12} \max_{\xi \in (x-h, x+h)} |g^{(iv)}(\xi)|.$$

Applied to the Laplacian of a C^4 bivariate function $u(x, y)$ we get

$$\begin{aligned} & |\nabla^2 u(x, y) - (u(x+h, y) + u(x-h, y) + u(x, y+h) + u(x, y-h) - 4u(x, y))/h^2| \\ & \leq \frac{h^2}{12} \max_{\substack{\xi \in (x-h, x+h) \\ \kappa \in (y-h, y+h)}} \left| \frac{\partial^4 u}{\partial x^4}(\xi, \kappa) \right| + \left| \frac{\partial^4 u}{\partial y^4}(\xi, \kappa) \right|. \end{aligned}$$

1) Since \hat{u} is the exact solution of Poisson's equation, we know that $\nabla^2 \hat{u}(ih, jh) = f_{i,j}$ for all $1 \leq i, j \leq M$. Replacing the left-hand side with the five-point approximation, and using the error bound above we can write:

$$\hat{u}_{i-1,j} + \hat{u}_{i+1,j} + \hat{u}_{i,j-1} + \hat{u}_{i,j+1} - 4\hat{u}_{i,j} = h^2 f_{i,j} + \eta_{i,j}, \quad |\eta_{i,j}| \leq ch^4 \quad (1.9)$$

where c is as defined in (1.8).

The solution of the five-point method u satisfies, for all $1 \leq i, j \leq M$:

$$u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1} - 4u_{i,j} = h^2 f_{i,j}. \quad (1.10)$$

Subtracting (1.10) from (1.9), we obtain

$$e_{i-1,j} + e_{i+1,j} + e_{i,j-1} + e_{i,j+1} - 4e_{i,j} = \eta_{i,j}$$

or, in the matrix form, $A\mathbf{e} = \boldsymbol{\eta}$, where A is symmetric (negative definite). It follows that

$$A\mathbf{e} = \boldsymbol{\eta} \Rightarrow \mathbf{e} = A^{-1}\boldsymbol{\eta} \Rightarrow \|\mathbf{e}\|_2 \leq \|A^{-1}\| \|\boldsymbol{\eta}\|_2,$$

where $\|A^{-1}\|$ is operator norm (also known as the spectral norm) of A^{-1} defined as $\|A^{-1}\| = \max_{\mathbf{x} \neq 0} \|A^{-1}\mathbf{x}\|_2 / \|\mathbf{x}\|_2$.

2) Since every component of $\boldsymbol{\eta}$ satisfies $|\eta_{i,j}|^2 \leq c^2 h^8$, where $h = \frac{1}{M+1}$, and there are M^2 components, we have

$$\|\boldsymbol{\eta}\|_2^2 = \sum_{i=1}^M \sum_{j=1}^M |\eta_{i,j}|^2 \leq c^2 M^2 h^8 < c^2 \frac{1}{h^2} h^8 = c^2 h^6 \Rightarrow \|\boldsymbol{\eta}\|_2 \leq ch^3.$$

3) The matrix A is symmetric, hence so is A^{-1} and therefore $\|A^{-1}\| = \rho(A^{-1})$. Here $\rho(A^{-1})$ is the spectral radius of A^{-1} , that is $\rho(A^{-1}) = \max_i |\lambda_i|$, where λ_i are the eigenvalues of A^{-1} . The eigenvalues of A^{-1} are the reciprocals of the eigenvalues of A , and the latter are given by Proposition 1.12. Thus, using the fact that $\sin(\pi h/2) \geq 1$ for $h \leq 1$ we get

$$\|A^{-1}\| = \frac{1}{4} \max_{k,\ell=1\dots m} \left(\sin^2 \frac{k\pi h}{2} + \sin^2 \frac{\ell\pi h}{2} \right)^{-1} = \frac{1}{8 \sin^2(\frac{1}{2}\pi h)} \leq \frac{1}{8h^2}.$$

Therefore $\|e\|_2 \leq \|A^{-1}\| \|\boldsymbol{\eta}\|_2 \leq (c/8)h$ as desired. \square

Fast Poisson solvers Suppose that we are solving $\nabla^2 u = f$ in a square $M \times M$ grid with the 5-point formula. Let the grid be enumerated in as before, i.e., by columns. Thus, the linear system $A\mathbf{u} = \mathbf{b}$ can be written explicitly in the block form

$$\underbrace{\begin{bmatrix} H & I & & & \\ I & H & \ddots & & \\ & \ddots & \ddots & I & \\ & & & I & H \end{bmatrix}}_A \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_M \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_M \end{bmatrix}, \quad H = \begin{bmatrix} -4 & 1 & & & \\ & 1 & -4 & \ddots & \\ & & \ddots & \ddots & 1 \\ & & & 1 & -4 \end{bmatrix}_{M \times M},$$

where $\mathbf{u}_k, \mathbf{b}_k \in \mathbb{R}^M$ are portions of \mathbf{u} and \mathbf{b} , respectively, and B is a TST-matrix which means *tridiagonal*, *symmetric* and *Toeplitz* (i.e., constant along diagonals). By Exercise 4, its eigenvalues and orthonormal eigenvectors are given as

$$H\mathbf{q}_\ell = \lambda_\ell \mathbf{q}_\ell, \quad \lambda_\ell = -4 + 2 \cos \frac{\ell\pi}{M+1}, \quad \mathbf{q}_\ell = \gamma_M \left(\sin \frac{j\ell\pi}{M+1} \right)_{j=1}^M, \quad \ell = 1..M,$$

where $\gamma_M = \sqrt{\frac{2}{M+1}}$ is the normalization factor. Hence $H = QDQ^{-1} = QDQ$, where $D = \text{diag}(\lambda_\ell)$ and $Q = Q^T = (q_{j\ell})$. Note that all $M \times M$ TST matrices share the same full set of eigenvectors, hence they all commute!

Hockney method Set $\mathbf{v}_k = Q\mathbf{u}_k, \mathbf{c}_k = Q\mathbf{b}_k$, therefore our system becomes

$$\begin{bmatrix} D & I & & & \\ I & D & \ddots & & \\ & \ddots & \ddots & I & \\ & & & I & D \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_M \end{bmatrix} = \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_M \end{bmatrix}.$$

Let us by this stage reorder the grid *by rows, instead of by columns*. In other words, we permute $\mathbf{v} \mapsto \hat{\mathbf{v}} = P\mathbf{v}$, $\mathbf{c} \mapsto \hat{\mathbf{c}} = P\mathbf{c}$, so that the portion $\hat{\mathbf{c}}_1$ is made out of the first components of the portions $\mathbf{c}_1, \dots, \mathbf{c}_M$, the portion $\hat{\mathbf{c}}_2$ out of the second components and so on. This results in new system

$$\begin{bmatrix} \Lambda_1 & & & & \\ & \Lambda_2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \Lambda_M \end{bmatrix} \begin{bmatrix} \hat{\mathbf{v}}_1 \\ \hat{\mathbf{v}}_2 \\ \vdots \\ \hat{\mathbf{v}}_M \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{c}}_1 \\ \hat{\mathbf{c}}_2 \\ \vdots \\ \hat{\mathbf{c}}_M \end{bmatrix}, \quad \Lambda_k = \begin{bmatrix} \lambda_k & 1 & & & \\ 1 & \lambda_k & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & & 1 & \lambda_k \end{bmatrix}_{M \times M}, \quad k = 1..M.$$

These are M *uncoupled* systems, $\Lambda_k \hat{\mathbf{v}}_k = \hat{\mathbf{c}}_k$ for $k = 1..M$. Being *tridiagonal*, each such system can be solved fast, at the cost of $\mathcal{O}(M)$. Thus, the steps of the algorithm and their computational cost are as follows.

1. Form the products $\mathbf{c}_k = Q\mathbf{b}_k, \quad k = 1 \dots M$ $\mathcal{O}(M^3)$
2. Solve $M \times M$ tridiagonal systems $\Lambda_k \hat{\mathbf{v}}_k = \hat{\mathbf{c}}_k, \quad k = 1 \dots M$ $\mathcal{O}(M^2)$
3. Form the products $\mathbf{u}_k = Q\mathbf{v}_k, \quad k = 1 \dots M$ $\mathcal{O}(M^3)$

(Permutations $\mathbf{c} \mapsto \hat{\mathbf{c}}$ and $\hat{\mathbf{v}} \mapsto \mathbf{v}$ are basically free.)

Improved Hockney algorithm We observe that the computational bottleneck is to be found in the $2M$ matrix-vector products by the matrix Q . Recall further that the elements of Q are $q_{j\ell} = \gamma_M \sin \frac{\pi j\ell}{M+1}$. This special form lends itself to a considerable speedup in matrix multiplication. Before making the problem simpler, however, let us make it more complicated! We write a typical product in the form

$$(Q\mathbf{y})_\ell = \sum_{j=1}^M \sin \frac{\pi j\ell}{M+1} y_j = \text{Im} \sum_{j=0}^M \exp \frac{i\pi j\ell}{M+1} y_j = \text{Im} \sum_{j=0}^{2M+1} \exp \frac{2i\pi j\ell}{2M+2} y_j, \quad \ell = 1 \dots M, \quad (1.11)$$

where $y_{M+1} = \dots = y_{2M+1} = 0$.

The discrete Fourier transform (DFT) The *discrete Fourier transform* of a vector $\mathbf{y} \in \mathbb{C}^n$ is $\mathbf{x} = \mathcal{F}_n \mathbf{y}$ defined by

$$x_\ell = \sum_{j=0}^{n-1} \omega_n^{j\ell} y_j \quad \ell = 0, \dots, n-1$$

where $\omega_n = \exp(2i\pi/n)$. (We assume in the above that vectors are indexed from 0 to $n-1$.) Thus, we see that multiplication by Q in (1.11) can be reduced to calculating a DFT. In the next lecture, we see how to compute the DFT of a vector \mathbf{y} in $\mathcal{O}(n \log n)$ operations, instead of $\mathcal{O}(n^2)$.