Mathematical Tripos Part II: Michaelmas Term 2021

Numerical Analysis – Lecture 3

Let $\hat{u}_{i,j} = u(ih, jh)$ be the grid values of the exact solution of the Poisson equation, and let $e_{i,j} = u_{i,j} - \hat{u}_{i,j}$ be the pointwise error of the 5-point formula. Set $e = (e_{i,j}) \in \mathbb{R}^n$ where $n = m^2$, and for $x \in \mathbb{R}^n$ let $\|x\| = \|x\|_{\ell_2}$ be the Euclidean norm of the vector x:

$$\|\boldsymbol{x}\|^2 = \sum_{k=1}^n |x_k|^2 = \sum_{i=1}^m \sum_{j=1}^m |x_{i,j}|^2$$

Theorem 1.11 Subject to sufficient smoothness of the function f and of the boundary conditions, there exists a number c > 0, independent of $h = \frac{1}{m+1}$, such that

$$\|\boldsymbol{e}\| \leq ch$$
.

Proof. 1) We already know (having constructed the 5-point formula by matching Taylor expansions) that, for the exact solution, we have

$$\widehat{u}_{i-1,j} + \widehat{u}_{i+1,j} + \widehat{u}_{i,j-1} + \widehat{u}_{i,j+1} - 4\widehat{u}_{i,j} = h^2 f_{i,j} + \eta_{i,j}, \qquad \eta_{i,j} = \mathcal{O}(h^4).$$

Subtracting this from numerical approximation (1.6), 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, $Ae = \eta$, where A is symmetric (negative definite). It follows that

$$Ae = \eta \Rightarrow e = A^{-1}\eta \Rightarrow ||e|| \le ||A^{-1}|| ||\eta||.$$

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

$$\|m{\eta}\|^2 = \sum_{i=1}^m \sum_{j=1}^m |\eta_{i,j}|^2 \le c^2 m^2 h^8 < c^2 \frac{1}{h^2} h^8 = c^2 h^6 \quad \Rightarrow \quad \|m{\eta}\| \le ch^3.$$

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,

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

Therefore $\|\boldsymbol{e}\| \leq \|A^{-1}\| \|\boldsymbol{\eta}\| \leq ch$ for some constant c > 0.

Observation 1.12 (Special structure of 5-point equations) We wish to motivate and introduce a family of efficient solution methods for the 5-point equations: the *fast Poisson solvers*. Thus, suppose that we are solving $\nabla^2 u = f$ in a square $m \times m$ grid with the 5-point formula (all this can be generalized a great deal, e.g. to the nine-point formula). Let the grid be enumerated in *natural ordering*, i.e. by columns. Thus, the linear system Au = b can be written explicitly in the block form

$$\underbrace{\begin{bmatrix} B & I \\ I & B & \ddots \\ & \ddots & \ddots & I \\ & & I & B \end{bmatrix}}_{A} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}, \qquad B = \begin{bmatrix} -4 & 1 \\ 1 & -4 & \ddots \\ & \ddots & \ddots & 1 \\ & & 1 & -4 \end{bmatrix}_{m \times m},$$

where $u_k, b_k \in \mathbb{R}^m$ are portions of u and 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

$$B\boldsymbol{q}_{\ell} = \lambda_{\ell} \boldsymbol{q}_{\ell}, \qquad \lambda_{\ell} = -4 + 2\cos\frac{\ell\pi}{m+1}, \qquad \boldsymbol{q}_{\ell} = \gamma_m \left(\sin\frac{j\ell\pi}{m+1}\right)_{j=1}^m, \qquad \ell = 1..m,$$

where $\gamma_m = \sqrt{\frac{2}{m+1}}$ is the normalization factor. Hence $B = QDQ^{-1} = QDQ$, 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!

Method 1.13 (The Hockney method) Set $v_k = Qu_k$, $c_k = Qb_k$, therefore our system becomes

	$oldsymbol{v}_1$		$igcap_1$	
$I D$ $\cdot \cdot$	$oldsymbol{v}_2$		$oldsymbol{c}_2$	
$\vdots \vdots I$	÷	=	÷	•
	$oldsymbol{v}_m$		c_m	

Let us by this stage reorder the grid by rows, instead of by columns.. In other words, we permute $v \mapsto \hat{v} = Pv$, $c \mapsto \hat{c} = Pc$, so that the portion \hat{c}_1 is made out of the first components of the portions c_1, \ldots, c_m , the portion \hat{c}_2 out of the second components and so on. This results in new system

$$\begin{bmatrix} \Lambda_1 \\ & \Lambda_2 \\ & & \ddots \\ & & & \Lambda_m \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{v}}_1 \\ \hat{\boldsymbol{v}}_2 \\ \vdots \\ \hat{\boldsymbol{v}}_m \end{bmatrix} = \begin{bmatrix} \hat{\boldsymbol{c}}_1 \\ \hat{\boldsymbol{c}}_2 \\ \vdots \\ \hat{\boldsymbol{c}}_m \end{bmatrix}, \qquad \Lambda_k = \begin{bmatrix} \lambda_k & 1 \\ 1 & \lambda_k & 1 \\ & \ddots & \ddots \\ & & 1 & \lambda_k \end{bmatrix}_{m \times m}, \quad k = 1 \dots m.$$

These are *m* uncoupled systems, $\Lambda_k \hat{v}_k = \hat{c}_k$ 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 $c_k = Qb_k$, k = 1...m $\mathcal{O}(m^3)$
- 2. Solve $m \times m$ tridiagonal systems $\Lambda_k \hat{v}_k = \hat{c}_k$, k = 1...m $\mathcal{O}(m^2)$
- 3. Form the products $\boldsymbol{u}_k = Q \boldsymbol{v}_k$, k = 1...m $\mathcal{O}(m^3)$

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

Method 1.14 (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

$$(Q\boldsymbol{y})_{\ell} = \sum_{j=1}^{m} \sin\frac{\pi j\ell}{m+1} y_j = \operatorname{Im} \sum_{j=0}^{m} \exp\frac{i\pi j\ell}{m+1} y_j = \operatorname{Im} \sum_{j=0}^{2m+1} \exp\frac{2i\pi j\ell}{2m+2} y_j, \quad \ell = 1...m,$$
(1.7)

where $y_{m+1} = \cdots = y_{2m+1} = 0$.

Definition 1.15 (The discrete Fourier transform (DFT)) Let Π_n be the space of all *bi-infinite complex n-periodic sequences* $\mathbf{x} = \{x_\ell\}_{\ell \in \mathbb{Z}}$ (such that $x_{\ell+n} = x_\ell$). Set $\omega_n = \exp \frac{2\pi i}{n}$, the primitive root of unity of degree *n*. The *discrete Fourier transform (DFT)* of \mathbf{x} is

$$\mathcal{F}_n: \Pi_n \to \Pi_n$$
 such that $\boldsymbol{y} = \mathcal{F}_n \boldsymbol{x}$, where $y_j = \frac{1}{n} \sum_{\ell=0}^{n-1} \omega_n^{-j\ell} x_\ell$, $j = 0...n-1$.

Trivial exercise: You can easily prove that \mathcal{F}_n is an isomorphism of Π_n onto itself and that

$$\boldsymbol{x} = \mathcal{F}_n^{-1} \boldsymbol{y}, \quad \text{where} \quad x_\ell = \sum_{j=0}^{n-1} \omega_n^{j\ell} y_j, \quad \ell = 0...n-1.$$

An important observation: Thus, multiplication by Q in (1.7) can be reduced to calculating an inverse of DFT.

Since we need to evaluate DFT (or its inverse) only in a single period, we can do so by multiplying a vector by a matrix, at the cost of $O(n^2)$ operations. This, however, is suboptimal and the cost of calculation can be lowered a great deal!