

Mathematical Tripos Part II: Michaelmas Term 2015

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 $e = (e_{i,j}) \in \mathbb{R}^{m^2}$.

Theorem 1.14 *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*

$$\|e\| \leq ch.$$

Proof. 1) We already know (having constructed the 5-point formula by matching Taylor expansions) that

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

Subtracting this from (1.5), 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

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

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

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

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,

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

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

Observation 1.15 (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

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

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

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

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

Method 1.16 (The 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

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

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

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

Method 1.17 (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 $q_{\ell j} = \gamma_m \sin \frac{\pi \ell j}{m+1}$. 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

$$(Q\mathbf{y})_\ell = \sum_{j=1}^m \sin \frac{\pi \ell j}{m+1} y_j = \text{Im} \sum_{j=0}^m \exp \frac{i\pi \ell j}{m+1} y_j = \text{Im} \sum_{j=0}^{2m+1} \exp \frac{2i\pi \ell j}{2m+2} y_j, \quad \ell = 1, \dots, m, \quad (1.8)$$

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

Problem 1.18 (The discrete Fourier transform) 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

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

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

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

An important observation: Thus, multiplication by Q in (1.8) 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 $\mathcal{O}(n^2)$ operations. This, however, is suboptimal and the cost of calculation can be lowered a great deal!