

Numerical Analysis - Part II

Anders C. Hansen

Lecture 4

Solving PDEs with finite difference methods

Special structure of 5-point equations

Observation 1 (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

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

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).

Special structure of 5-point equations

Observation 2 (Special structure of 5-point equations)

By Exercise 4, its eigenvalues and orthonormal eigenvectors are given as

$$B\mathbf{q}_\ell = \lambda_\ell \mathbf{q}_\ell, \quad \lambda_\ell = -4 + 2 \cos \frac{\ell\pi}{m+1},$$

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

where $\gamma_m = \sqrt{\frac{2}{m+1}}$ is the normalization factor. Hence

$B = QDQ^{-1} = QDQ$, where $D = \text{diag}(\lambda_\ell)$ and $Q = Q^T = (\mathbf{q}_{j\ell})$.

Note that all $m \times m$ TST matrices share the same full set of eigenvectors, hence they all commute!

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.

The Hockney method

This results in new system

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

where $k = 1 \dots m$.

The Hockney method

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 $\mathbf{c}_k = \mathbf{Q}\mathbf{b}_k$, $k = 1 \dots m$ $\mathcal{O}(m^3)$
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 $\mathbf{u}_k = \mathbf{Q}\mathbf{v}_k$, $k = 1 \dots m$ $\mathcal{O}(m^3)$

The improved Hockney method

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

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

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

The discrete Fourier transform (DFT)

Definition 3 (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

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

where $j = 0 \dots n-1$.

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

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

The discrete Fourier transform (DFT)

An important observation: Thus, multiplication by Q in (1) 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!

The fast Fourier transform (FFT)

We assume that n is a power of 2, i.e. $n = 2m = 2^p$, and for $\mathbf{y} \in \Pi_{2m}$, denote by

$$\mathbf{y}^{(E)} = \{y_{2j}\}_{j \in \mathbb{Z}} \quad \text{and} \quad \mathbf{y}^{(O)} = \{y_{2j+1}\}_{j \in \mathbb{Z}}$$

the even and odd portions of \mathbf{y} , respectively. Note that $\mathbf{y}^{(E)}, \mathbf{y}^{(O)} \in \Pi_m$.

Suppose that we already know the inverse DFT of both 'short' sequences,

$$\mathbf{x}^{(E)} = \mathcal{F}_m^{-1} \mathbf{y}^{(E)}, \quad \mathbf{x}^{(O)} = \mathcal{F}_m^{-1} \mathbf{y}^{(O)}.$$

Computing the fast Fourier transform (FFT)

It is then possible to assemble $\mathbf{x} = \mathcal{F}_{2m}^{-1}\mathbf{y}$ in a small number of operations. Since $\omega_{2m}^{2m} = 1$, we obtain $\omega_{2m}^2 = \omega_m$, and

$$\begin{aligned}x_\ell &= \sum_{j=0}^{2m-1} \omega_{2m}^{j\ell} y_j = \sum_{j=0}^{m-1} \omega_{2m}^{2j\ell} y_{2j} + \sum_{j=0}^{m-1} \omega_{2m}^{(2j+1)\ell} y_{2j+1} \\ &= \sum_{j=0}^{m-1} \omega_m^{j\ell} y_j^{(E)} + \omega_{2m}^\ell \sum_{j=0}^{m-1} \omega_m^{j\ell} y_j^{(O)} = x_\ell^{(E)} + \omega_{2m}^\ell x_\ell^{(O)},\end{aligned}$$

where $\ell = 0, \dots, m-1$.

Computing the fast Fourier transform (FFT)

Therefore, it costs just m products to evaluate the first half of \mathbf{x} , provided that $\mathbf{x}^{(E)}$ and $\mathbf{x}^{(O)}$ are known. It actually costs nothing to evaluate the second half, since

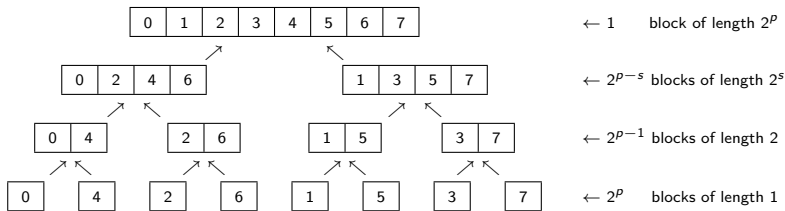
$$\omega_m^{j(m+\ell)} = \omega_m^{j\ell}, \quad \omega_{2m}^{m+\ell} = -\omega_{2m}^{\ell} \quad \Rightarrow \quad x_{m+\ell} = x_{\ell}^{(E)} - \omega_{2m}^{\ell} x_{\ell}^{(O)},$$

where $\ell = 0, \dots, m-1$.

Note: To execute FFT, we start from vectors of unit length and in each s -th stage, $s = 1 \dots p$, assemble 2^{p-s} vectors of length 2^s from vectors of length 2^{s-1} : this costs $2^{p-s} 2^{s-1} = 2^{p-1}$ products.

Computing the fast Fourier transform (FFT)

Altogether, the cost of FFT is $p2^{p-1} = \frac{1}{2}n \log_2 n$ products.



For $n = 1024 = 2^{10}$, say, the cost is $\approx 5 \times 10^3$ products, compared to $\approx 10^6$ for naive matrix multiplication! For $n = 2^{20}$ the respective numbers are $\approx 1.05 \times 10^7$ and $\approx 1.1 \times 10^{12}$, which represents a saving by a factor of more than 10^5 .

Numerical analysis vs Foundations of computational mathematics

Numerical analysis is mostly concerned with providing upper bounds.

Foundations of computational mathematics is concerned with determining the boundaries of what computers can achieve. That means also lower bounds.

Is $n \log(n)$ the best one can do when computing the discrete Fourier transform? Does there exist an algorithm that is faster?

Partial differential equations of evolution

Recall the Poisson equation

Recall the *Poisson equation*

$$\nabla^2 u = f \quad (x, y) \in \Omega, \quad (2)$$

where $\nabla^2 = \Delta = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$ is the Laplace operator and Ω is an open connected domain of \mathbb{R}^2 with a Jordan boundary, specified together with the *Dirichlet boundary condition*

$$u(x, y) = \phi(x, y) \quad (x, y) \in \partial\Omega. \quad (3)$$

(You may assume that $f \in C(\Omega)$, $\phi \in C^2(\partial\Omega)$, but this can be relaxed by an approach outside the scope of this course.)

Solving the diffusion equation

We consider the solution of the *diffusion equation*

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, \quad 0 \leq x \leq 1, \quad t \geq 0,$$

with *initial conditions* $u(x, 0) = u_0(x)$ for $t = 0$ and *Dirichlet boundary conditions* $u(0, t) = \phi_0(t)$ at $x = 0$ and $u(1, t) = \phi_1(t)$ at $x = 1$. By Taylor's expansion

$$\begin{aligned} \frac{\partial u(x, t)}{\partial t} &= \frac{1}{k} [u(x, t+k) - u(x, t)] + \mathcal{O}(k), & k = \Delta t, \\ \frac{\partial^2 u(x, t)}{\partial x^2} &= \frac{1}{h^2} [u(x-h, t) - 2u(x, t) + u(x+h, t)] + \mathcal{O}(h^2), & h = \Delta x, \end{aligned}$$

so that, for the true solution, we obtain

$$u(x, t+k) = u(x, t) + \frac{k}{h^2} [u(x-h, t) - 2u(x, t) + u(x+h, t)] + \mathcal{O}(k^2 + kh^2). \quad (4)$$

Numerical scheme for the diffusion equation

That motivates the numerical scheme for approximation $u_m^n \approx u(x_m, t_n)$ on the rectangular mesh $(x_m, t_n) = (mh, nk)$:

$$u_m^{n+1} = u_m^n + \mu (u_{m-1}^n - 2u_m^n + u_{m+1}^n), \quad m = 1 \dots M. \quad (5)$$

Here $h = \frac{1}{M+1}$ and $\mu = \frac{k}{h^2} = \frac{\Delta t}{(\Delta x)^2}$ is the so-called *Courant number*. With μ being fixed, we have $k = \mu h^2$, so that the local truncation error of the scheme is $\mathcal{O}(h^4)$. Substituting whenever necessary initial conditions u_m^0 and boundary conditions u_0^n and u_{M+1}^n , we possess enough information to advance in (5) from $\mathbf{u}^n := [u_1^n, \dots, u_M^n]$ to $\mathbf{u}^{n+1} := [u_1^{n+1}, \dots, u_M^{n+1}]$.

Convergence

Similarly to ODEs or Poisson equation, we say that the method is *convergent* if, for a fixed μ , and for every $T > 0$, we have

$$\lim_{h \rightarrow 0} |u_m^n - u(x_m, t_n)| = 0 \quad \text{uniformly for } (x_m, t_n) \in [0, 1] \times [0, T].$$

In other words, if $e_m^n := u_m^n - u(mh, nk)$ is the error of approximation, and $\mathbf{e}^n = [e_1^n, \dots, e_M^n]$ with $\|\mathbf{e}^n\| := \max_m |e_m^n|$, then convergence is equivalent to

$$\lim_{h \rightarrow 0} \max_{1 \leq n \leq T/k} \|\mathbf{e}^n\| = 0.$$

Note: In the present case, however, a method has an extra parameter μ , and it is entirely possible for a method to converge for some choice of μ and diverge otherwise.

Proving convergence

Theorem 4

If $\mu \leq \frac{1}{2}$, then method (5) converges.

Proof. Let $e_m^n := u_m^n - u(mh, nk)$ be the error of approximation, and let $\mathbf{e}^n = [e_1^n, \dots, e_M^n]$ with $\|\mathbf{e}^n\| := \max_m |e_m^n|$. Convergence is equivalent to

$$\lim_{h \rightarrow 0} \max_{1 \leq n \leq T/k} \|\mathbf{e}^n\| = 0$$

for every constant $T > 0$. Subtracting (4) from (5), we obtain

$$\begin{aligned} e_m^{n+1} &= e_m^n + \mu(e_{m-1}^n - 2e_m^n + e_{m+1}^n) + \mathcal{O}(h^4) \\ &= \mu e_{m-1}^n + (1 - 2\mu)e_m^n + \mu e_{m+1}^n + \mathcal{O}(h^4). \end{aligned}$$

Then

$$\|\mathbf{e}^{n+1}\| = \max_m |e_m^{n+1}| \leq (2\mu + |1 - 2\mu|) \|\mathbf{e}^n\| + ch^4 = \|\mathbf{e}^n\| + ch^4,$$

by virtue of $\mu \leq \frac{1}{2}$. Since $\|\mathbf{e}^0\| = 0$, induction yields

$$\|\mathbf{e}^n\| \leq cnh^4 \leq \frac{cT}{k} h^4 = \frac{cT}{\mu} h^2 \rightarrow 0 \quad (h \rightarrow 0)$$

□

In practice we wish to choose h and k of comparable size, therefore $\mu = k/h^2$ is likely to be large. Consequently, the restriction of the last theorem is disappointing: unless we are willing to advance with tiny time step k , the method (5) is of limited practical interest. The situation is similar to stiff ODEs: like the Euler method, the scheme (5) is simple, plausible, explicit, easy to execute and analyse – but of very limited utility. . . .