

Mathematical Tripos Part II: Michaelmas Term 2020

Numerical Analysis – Lecture 10

Technique 2.31 (Splitting) We would like to find a fast solver to the system (2.16):

$$(I - \frac{1}{2}\mu A_*) \mathbf{u}^{n+1} = (I + \frac{1}{2}\mu A_*) \mathbf{u}^n, \quad (2.17)$$

The matrix $B = I - \frac{1}{2}\mu A_*$ has a structure similar to that of A_* (for the 5-point scheme), so we may apply the Hockney method for solving $B\mathbf{u}^{n+1} = \mathbf{v}^n = C\mathbf{u}^n$. However, since the method (2.17) has a local truncation error $\mathcal{O}(k^3 + kh^2)$, we don't need the exact solution of the system: it would be enough to have one within the error.

Let us employ the notation

$$\Delta_x^2 u_{\ell,m} = u_{\ell-1,m} - 2u_{\ell,m} + u_{\ell+1,m}, \quad \Delta_y^2 u_{\ell,m} = u_{\ell,m-1} - 2u_{\ell,m} + u_{\ell,m+1}.$$

Then the Crank-Nicolson method calculates \mathbf{u}^{n+1} by solving the system

$$[I - \frac{1}{2}\mu(\Delta_x^2 + \Delta_y^2)] u_{\ell,m}^{n+1} = [I + \frac{1}{2}\mu(\Delta_x^2 + \Delta_y^2)] u_{\ell,m}^n, \quad \ell, m = 1 \dots M. \quad (2.18)$$

The local error is however preserved if we replace this formula by the difference equation

$$[I - \frac{1}{2}\mu\Delta_x^2][I - \frac{1}{2}\mu\Delta_y^2] u_{\ell,m}^{n+1} = [I + \frac{1}{2}\mu\Delta_x^2][I + \frac{1}{2}\mu\Delta_y^2] u_{\ell,m}^n, \quad (2.19)$$

which is called the *split version of Crank-Nicolson*. Indeed, the difference between two schemes is equal to

$$\begin{aligned} \frac{1}{4}\mu^2 \Delta_x^2 \Delta_y^2 (u_{\ell,m}^{n+1} - u_{\ell,m}^n) &= \frac{k^2}{4} \frac{1}{h^2} \Delta_x^2 \frac{1}{h^2} \Delta_y^2 \left(k \frac{\partial}{\partial t} u_{\ell,m}^n + \mathcal{O}(k^2) \right) \\ &= \frac{k^3}{4} \left(\frac{\partial^2}{\partial x^2} \frac{\partial^2}{\partial y^2} \frac{\partial}{\partial t} u_{\ell,m}^n + \mathcal{O}(k + h^2) \right) = \mathcal{O}(k^3), \end{aligned}$$

the same magnitude as of the local error. In the matrix form, (2.19) is equivalent to splitting the matrix A_* into the sum of two matrices A_x and A_y as

$$A_* = A_x + A_y, \quad A_x = \begin{bmatrix} -2I & I & & & \\ & I & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & I \\ I & & & & -2I \end{bmatrix}, \quad A_y = \begin{bmatrix} H & & & & \\ & H & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & H \end{bmatrix}, \quad H = \begin{bmatrix} -2 & 1 & & & \\ & 1 & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & 1 & -2 \end{bmatrix},$$

and solving the uncoupled system

$$[I - \frac{1}{2}\mu A_x][I - \frac{1}{2}\mu A_y] \mathbf{u}^{n+1} = [I + \frac{1}{2}\mu A_x][I + \frac{1}{2}\mu A_y] \mathbf{u}^n$$

as

$$B_x \mathbf{u}^{n+1/2} = C_x C_y \mathbf{u}^n, \quad B_y \mathbf{u}^{n+1} = \mathbf{u}^{n+1/2}.$$

Matrix $B_y = I - \frac{1}{2}\mu A_y$ is block diagonal, and solving $B_y \mathbf{u} = \mathbf{v}$ is just solving one and the same tridiagonal system $Bu_i = v_i$ with different right-hand sides. Matrix $B_x = I - \frac{1}{2}\mu A_x$ is of the same form up to a permutation (reordering of the grid), so solving $B_x \mathbf{v} = \mathbf{b}$ is again a fast procedure.

Example 2.32 Consider the general diffusion equation

$$\frac{\partial u}{\partial t} = \nabla^\top (a(x, y) \nabla u) + f(x, y) = \frac{\partial}{\partial x} \left(a(x, y) \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left(a(x, y) \frac{\partial u}{\partial y} \right) + f(x, y), \quad (2.20)$$

where $a(x, y) > \alpha > 0$ and $f(x, y)$ are given, together with initial conditions on $[0, 1]^2$ and Dirichlet boundary conditions along $\partial[0, 1]^2 \times [0, \infty)$. Replace each space derivative by *central differences* at midpoints,

$$\frac{dg(\xi)}{d\xi} \approx \frac{g(\xi + \frac{1}{2}h) - g(\xi - \frac{1}{2}h)}{h},$$

resulting in the ODE system

$$u'_{\ell,m} = \frac{1}{h^2} \left[a_{\ell-\frac{1}{2},m} u_{\ell-1,m} + a_{\ell+\frac{1}{2},m} u_{\ell+1,m} + a_{\ell,m-\frac{1}{2}} u_{\ell,m-1} + a_{\ell,m+\frac{1}{2}} u_{\ell,m+1} - \left(a_{\ell-\frac{1}{2},m} + a_{\ell+\frac{1}{2},m} + a_{\ell,m-\frac{1}{2}} + a_{\ell,m+\frac{1}{2}} \right) u_{\ell,m} \right] + f_{\ell,m}. \quad (2.21)$$

Assuming zero boundary conditions and $f \equiv 0$, we have a system $\mathbf{u}' = A\mathbf{u}$, and we may solve it again by Crank–Nicolson, and apply the split $A = A_x + A_y$. Here, A_x and A_y are again constructed from the contribution of discretizations in the x - and y -directions respectively, namely A_x includes all the $a_{\ell \pm \frac{1}{2},m}$ terms, and A_y consists of the remaining $a_{\ell,m \pm \frac{1}{2}}$ components. Arguments similar to what we used in moving from (2.18) to (2.19) justify the use of the split version in this general case as well.

Intermezzo 2.33 (Linear systems of ODEs) With greater generality, let us consider the ODE system

$$\mathbf{y}' = A\mathbf{y}, \quad \mathbf{y}(0) = \mathbf{y}_0. \quad (2.22)$$

We define formally a *matrix exponential* by Taylor series, $e^B := \sum_{k=0}^{\infty} \frac{1}{k!} B^k$, and easily verify by formal differentiation that $\text{de}^{tA}/\text{dt} = Ae^{tA}$, therefore $\mathbf{y}(t) = e^{tA}\mathbf{y}_0$ is a solution.

One observes that one-step methods for solving (2.22) are approximating a matrix exponential. Thus, with $k = \Delta t$,

$$\begin{aligned} \text{Euler: } \mathbf{y}_n &= (I + kA)^n \mathbf{y}_0, & 1 + z &= e^z + \mathcal{O}(z^2); \\ \text{TR: } \mathbf{y}_n &= \left[\left(I - \frac{1}{2}kA \right)^{-1} \left(I + \frac{1}{2}kA \right) \right]^n \mathbf{y}_0, & \frac{1+\frac{1}{2}z}{1-\frac{1}{2}z} &= e^z + \mathcal{O}(z^3). \end{aligned}$$

Technique 2.34 (Splitting methods) Recall that, for $z_1, z_2 \in \mathbb{C}$, we have $e^{z_1+z_2} = e^{z_1}e^{z_2}$ and had this been true for the matrices, i.e. that $e^{tA} = e^{t(B+C)} = e^{tB}e^{tC}$, we could have approximated each component of the exponent of $A = A_x + A_y$ with the trapezoidal rule, say, to produce

$$\mathbf{u}^{n+1} = \left(I - \frac{1}{2}\mu A_x \right)^{-1} \left(I + \frac{1}{2}\mu A_x \right) \left(I - \frac{1}{2}\mu A_y \right)^{-1} \left(I + \frac{1}{2}\mu A_y \right) \mathbf{u}^n, \quad \mu = k/h^2, \quad (2.23)$$

and since both $I - \frac{1}{2}\mu A_x$ and $I - \frac{1}{2}\mu A_y$ are tridiagonal, this system can be solved very cheaply.

Unfortunately, the assumption that $e^{t(B+C)} = e^{tB}e^{tC}$ is, in general, false. Not all hope is lost, though, and we will demonstrate that, suitably implemented, splitting is a powerful technique to reduce drastically the expense of numerical solution.

Method 2.35 (Splitting) Comparing the Taylor expansions of $e^{t(B+C)}$ with $e^{tB}e^{tC}$ we obtain

$$e^{tB}e^{tC} = e^{t(B+C)} + \frac{1}{2}t^2(BC - CB) + \mathcal{O}(t^3). \quad (2.24)$$

In particular, $e^{tB}e^{tC} = e^{t(B+C)}$ for all $t \geq 0$ if and only if B and C commute. The good news is, however, that approximating $e^{\Delta t(B+C)}$ with $e^{\Delta t B}e^{\Delta t C}$ incurs an error of $\mathcal{O}((\Delta t)^2)$. So, if r is a rational function such that $r(z) = e^z + \mathcal{O}(z^2)$, then

$$\mathbf{u}^{n+1} = r(\mu A_x)r(\mu A_y)\mathbf{u}^n \quad (2.25)$$

produces an error of $\mathcal{O}((\Delta t)^2)$. The choice $r(z) = (1 + \frac{1}{2}z)/(1 - \frac{1}{2}z)$ results in a *split Crank–Nicolson* scheme, whose implementation reduces to a solution of tridiagonal algebraic linear systems.

It is easy to prove that

$$e^{t(B+C)} = \frac{1}{2} \left(e^{tB}e^{tC} + e^{tC}e^{tB} \right) + \mathcal{O}(t^3), \quad e^{t(B+C)} = e^{\frac{1}{2}tB}e^{tC}e^{\frac{1}{2}tB} + \mathcal{O}(t^3),$$

the second formula is called the *Strang splitting*. Thus, as long as $r(z) = e^z + \mathcal{O}(z^3)$, the time-stepping formula $\mathbf{u}^{n+1} = r\left(\frac{1}{2}\mu A_x\right)r(\mu A_y)r\left(\frac{1}{2}\mu A_x\right)\mathbf{u}^n$ carries a local error of $\mathcal{O}((\Delta t)^3)$.

As far as stability is concerned, we observe that both A_x and A_y are symmetric, hence normal, therefore so are $r(\mu A_x)$ and $r(\mu A_y)$. Then Euclidean ℓ_2 -norm equals the spectral radius, therefore for the splitting (2.25), we have

$$\|\mathbf{u}^{n+1}\| \leq \|r(\mu A_x)\| \cdot \|r(\mu A_y)\| \cdot \|\mathbf{u}^n\| = \rho[r(\mu A_x)] \cdot \rho[r(\mu A_y)] \cdot \|\mathbf{u}^n\|.$$

It is easy to verify by Gershgorin theorem that the eigenvalues of the matrices A_x and A_y are nonpositive, hence provided that r fulfils $|r(z)| < 1$ for $z \in \mathbb{C}$ with $\text{Re } z < 0$, it is then true that $\rho[r(\mu A_x)], \rho[r(\mu A_y)] \leq 1$. This proves $\|\mathbf{u}^{n+1}\| \leq \|\mathbf{u}^n\| \leq \dots \leq \|\mathbf{u}^0\|$, hence stability.

Method 2.36 (Splitting of inhomogeneous systems) Recall our goal, namely fast methods for the two-dimensional diffusion equation. Our exposition so far has been contrived, because of the assumption that the boundary conditions are zero. In general, the linear ODE system is of the form

$$\mathbf{u}' = A\mathbf{u} + \mathbf{b}, \quad \mathbf{u}(0) = \mathbf{u}^0, \quad (2.26)$$

where \mathbf{b} originates in boundary conditions (and in a forcing term $f(x, y)$ in the original PDE (2.20)). Note that our analysis should accommodate $\mathbf{b} = \mathbf{b}(t)$, since boundary conditions might vary in time! The *exact* solution of (2.26) is provided by the *variation of constants* formula

$$\mathbf{u}(t) = e^{tA}\mathbf{u}(0) + \int_0^t e^{(t-s)A}\mathbf{b}(s) \, ds, \quad t \geq 0,$$

therefore

$$\mathbf{u}(t_{n+1}) = e^{\Delta t A}\mathbf{u}(t_n) + \int_{t_n}^{t_{n+1}} e^{(t_{n+1}-s)A}\mathbf{b}(s) \, ds.$$

The integral can be frequently evaluated explicitly, e.g. when \mathbf{b} is a linear combination of polynomial and exponential terms. For example, $\mathbf{b}(t) \equiv \mathbf{b} = \text{const}$ yields

$$\mathbf{u}(t_{n+1}) = e^{\Delta t A}\mathbf{u}(t_n) + A^{-1} \left(e^{\Delta t A} - I \right) \mathbf{b}.$$

This, unfortunately, is not a helpful observation, since, even if we split the exponential e^{tA} , how are we supposed to split $A^{-1} = (B + C)^{-1}$? The remedy is not to evaluate the integral explicitly but, instead, to use quadrature. For example, the trapezoidal rule $\int_0^k g(\tau) \, d\tau = \frac{1}{2}k[g(0) + g(k)] + \mathcal{O}(k^3)$ gives

$$\mathbf{u}(t_{n+1}) \approx e^{\Delta t A}\mathbf{u}(t_n) + \frac{1}{2}\Delta t [e^{\Delta t A}\mathbf{b}(t_n) + \mathbf{b}(t_{n+1})],$$

with a local error of $\mathcal{O}((\Delta t)^3)$. We can now replace exponentials with their splittings. For example, Strang's splitting results in

$$\mathbf{u}^{n+1} = r\left(\frac{1}{2}\Delta t B\right) r\left(\Delta t C\right) r\left(\frac{1}{2}\Delta t B\right) \left[\mathbf{u}^n + \frac{1}{2}\Delta t \mathbf{b}^n \right] + \frac{1}{2}\Delta t \mathbf{b}^{n+1}.$$

As before, everything reduces to (inexpensive) solution of tridiagonal systems!