

Mathematical Tripos Part II: Michaelmas Term 2021

Numerical Analysis – Lecture 10

Linear systems of ODEs In all the examples of semi-discretization we have seen so far, we always reach a linear system of ODE of the form:

$$\mathbf{u}' = A\mathbf{u}, \quad \mathbf{u}(0) = \mathbf{u}_0. \quad (2.17)$$

The solution of this linear system of ODE is given by

$$\mathbf{u}(t) = e^{tA}\mathbf{u}_0 \quad (2.18)$$

where the *matrix exponential* function is defined by $e^B := \sum_{k=0}^{\infty} \frac{1}{k!} B^k$. It is easily verified that $\frac{d}{dt} e^{tA} = A e^{tA}$, therefore (2.18) is indeed a solution of (2.17).

If A can be diagonalized $A = V D V^{-1}$, then $e^{tA} = V e^{tD} V^{-1}$ where e^{tD} is the diagonal matrix consisting $\text{diag}(e^{tD_{ii}})$. As such one can compute the solution of (2.17) exactly. However computing an eigenvalue decomposition can be costly, and so one would like to consider more efficient methods.

Observe that one-step methods for solving (2.17) are approximating a matrix exponential. Indeed, with $k = \Delta t$, we have:

$$\begin{aligned} \text{Euler:} \quad & \mathbf{u}^{n+1} = (I + kA)\mathbf{u}^n, & e^z &= 1 + z + \mathcal{O}(z^2); \\ \text{Implicit Euler:} \quad & \mathbf{u}^{n+1} = (I - kA)^{-1}\mathbf{u}^n, & e^z &= (1 - z)^{-1} + \mathcal{O}(z^2); \\ \text{Trapezoidal Rule:} \quad & \mathbf{u}^{n+1} = (I - \tfrac{1}{2}kA)^{-1} (I + \tfrac{1}{2}kA)\mathbf{u}^n, & e^z &= \frac{1+\frac{1}{2}z}{1-\frac{1}{2}z} + \mathcal{O}(z^3). \end{aligned}$$

In practice the matrix A is very sparse, and this can be exploited when solving linear systems e.g., for the implicit Euler or Trapezoidal Rule.

Splitting In many cases, the matrix A is naturally expressed as a sum of two matrices, $A = B + C$. For example, when discretizing the diffusion equation in 2D with zero boundary conditions, we have $A = \frac{1}{h^2}(A_x + A_y)$ where $A_x \in \mathbb{R}^{M^2 \times M^2}$ corresponds to the 3-point discretization of $\frac{\partial^2}{\partial x^2}$, and $A_y \in \mathbb{R}^{M^2 \times M^2}$ corresponds to the 3-point discretization of $\frac{\partial^2}{\partial y^2}$. In matrix notations, if the grid points are ordered by columns, then we have:

$$A_x = \begin{bmatrix} -2I & I & & \\ & \ddots & \ddots & \\ & & \ddots & I \\ I & & & -2I \end{bmatrix}, \quad A_y = \begin{bmatrix} G & & \\ & G & \\ & & \ddots \\ & & & G \end{bmatrix}, \quad G = \begin{bmatrix} -2 & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & -2 \end{bmatrix} \in \mathbb{R}^{M \times M}. \quad (2.19)$$

When the matrices B and C commute, we can use the following fact about the matrix exponential.

Proposition 2.31 If B and C commute, then $e^{B+C} = e^B e^C$.

Proof. We have

$$e^{B+C} = \sum_{k=0}^{\infty} \frac{1}{k!} (B+C)^k = \sum_{k=0}^{\infty} \frac{1}{k!} \sum_{i+j=k} \binom{k}{i} B^i C^j = \sum_{i,j=0}^{\infty} \frac{1}{i!j!} B^i C^j = e^B e^C$$

where in the second step we used the fact that B and C commute.

The matrices A_x and A_y in (2.19) happen to commute (easy to check), and so $e^{\Delta t A} = e^{\frac{\Delta t}{h^2} A_x} e^{\frac{\Delta t}{h^2} A_y}$. This means that the solution of the semi-discretized diffusion equation in 2D, with zero boundary conditions, satisfies

$$\mathbf{u}^{n+1} = e^{\mu A_x} e^{\mu A_y} \mathbf{u}^n. \quad (2.20)$$

Split Crank-Nicolson: In the split Crank-Nicolson scheme, we approximate each exponential map in (2.20) by the rational function $r(z) = (1 + z/2)(1 - z/2)^{-1}$, which leads to

$$\mathbf{u}^{n+1} = (I + \frac{\mu}{2}A_x)(I - \frac{\mu}{2}A_x)^{-1}(I + \frac{\mu}{2}A_y)(I - \frac{\mu}{2}A_y)^{-1}\mathbf{u}^n. \quad (2.21)$$

Note that computing $\mathbf{u}^{n+1/2} = (I + \frac{\mu}{2}A_y)(I - \frac{\mu}{2}A_y)^{-1}\mathbf{u}^n$ can be done efficiently in $\mathcal{O}(M^2)$ time as A_y is block-diagonal, and the matrices G are tridiagonal (each tridiagonal solve requires $\mathcal{O}(M)$ time, and we have M of these). Computing $\mathbf{u}^{n+1} = (I + \frac{\mu}{2}A_x)(I - \frac{\mu}{2}A_x)^{-1}\mathbf{u}^{n+1/2}$ can also be done in $\mathcal{O}(M^2)$ time, since A_x is also block-diagonal provided we appropriately permute the rows and columns so that the grid ordering is by rows instead of columns. This means that the update step (2.21) of Split-Crank-Nicolson can be performed in time $\mathcal{O}(M^2)$ and only requires tridiagonal matrix solves (no FFT needed).

In general, however, the matrices B and C in $A = B + C$ do not have to commute, as in the following example:

Example 2.32 The general diffusion equation with a diffusion coefficient $a(x, y) > 0$ is given by:

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(a(x, y) \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left(a(x, y) \frac{\partial u}{\partial y} \right), \quad (2.22)$$

together with initial conditions on $[0, 1]^2$ and Dirichlet boundary conditions along $\partial[0, 1]^2 \times [0, \infty)$. We 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

$$\begin{aligned} u'_{\ell, m} = & \frac{1}{h^2} \left[a_{\ell-\frac{1}{2}, m} u_{\ell-1, m} + a_{\ell+\frac{1}{2}, m} u_{\ell+1, m} + a_{\ell, m-\frac{1}{2}} u_{\ell, m-1} + a_{\ell, m+\frac{1}{2}} u_{\ell, m+1} \right. \\ & \left. - (a_{\ell-\frac{1}{2}, m} + a_{\ell+\frac{1}{2}, m} + a_{\ell, m-\frac{1}{2}} + a_{\ell, m+\frac{1}{2}}) u_{\ell, m} \right]. \end{aligned} \quad (2.23)$$

Assuming zero boundary conditions, we have a system $\mathbf{u}' = A\mathbf{u}$, and the matrix A can be split as $A = \frac{1}{h^2}(A_x + A_y)$. 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.

In this case the matrices A_x and A_y do not necessarily commute. The next proposition tells us that approximating $e^{t(B+C)}$ by $e^{tB}e^{tC}$ results in an error of $\mathcal{O}(t^2)$.

Proposition 2.33 For any matrices B, C ,

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

Proof. We Taylor-expand both expressions $e^{tB}e^{tC}$ and $e^{t(B+C)}$:

$$\begin{aligned} e^{tB}e^{tC} &= (I + tB + t^2B^2/2 + \mathcal{O}(t^3))(I + tC + t^2C^2/2 + \mathcal{O}(t^3)) \\ &= I + t(B + C) + \frac{t^2}{2}(B^2 + C^2 + 2BC) + \mathcal{O}(t^3) \end{aligned}$$

and

$$\begin{aligned} e^{t(B+C)} &= I + t(B + C) + \frac{t^2}{2}(B + C)^2 + \mathcal{O}(t^3) \\ &= I + t(B + C) + \frac{t^2}{2}(B^2 + C^2 + BC + CB) + \mathcal{O}(t^3). \end{aligned}$$

The result follows.

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 $r(z) = (1 + \frac{1}{2}z)/(1 - \frac{1}{2}z) = e^z + \mathcal{O}(z^3)$ results in a *split Crank–Nicolson* scheme, whose implementation reduces to a solution of tridiagonal algebraic linear systems.

Strang splitting: One can obtain better splitting approximations of $e^{t(B+C)}$. For example it is not hard to prove that $e^{\frac{1}{2}tB}e^{tC}e^{\frac{1}{2}tB}$ gives a $\mathcal{O}(t^3)$ approximation of $e^{t(B+C)}$, i.e.,

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

Thus, as long as $r(z) = e^z + \mathcal{O}(z^3)$, the time-stepping formula

$$\mathbf{u}^{n+1} = r(\tfrac{1}{2}\mu A_x) r(\mu A_y) r(\tfrac{1}{2}\mu A_x) \mathbf{u}^n$$

carries a local error of $\mathcal{O}((\Delta t)^3)$.

Stability: Consider the general diffusion equation with the splitting scheme (2.25). 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 we have

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

The function $r(z) = (1 + \frac{1}{2}z)(1 - \frac{1}{2}z)^{-1}$ satisfies $|r(z)| \leq 1$ for $z \in \mathbb{C}$ with $\operatorname{Re} z \leq 0$. By the Gersgorin theorem, we see that the eigenvalues of A_x and A_y are nonpositive. Then it is 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.

Remark 2.34 (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, possibly, in a forcing term $f(x, y)$ in the original PDE (2.22)). 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

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

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

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

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