

Math 6620: Analysis of Numerical Methods, II

Solvers for initial value problems, Part I

See Ascher and Petzold 1998, Chapters 1-5

Akil Narayan¹

¹Department of Mathematics, and Scientific Computing and Imaging (SCI) Institute
University of Utah



We've discussed how to solve simple stationary problems.

Before delving into time-dependent problems, we'll spend some time discussing numerical methods for solving [initial value problems](#).

We've discussed how to solve simple stationary problems.

Before delving into time-dependent problems, we'll spend some time discussing numerical methods for solving [initial value problems](#).

Our focal problem will be the ordinary differential equation:

$$\mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0,$$

where the initial condition \mathbf{u}_0 is given.

The assumption is that we seek to compute $\mathbf{u}(t)$ for $t \in [0, T]$.

(Or possibly only at the endpoint, or at some discrete values in the interval)

The function \mathbf{f} can be nonlinear and/or the ODE can be autonomous, and higher-order ODE's can be written as first-order ODE's with an expanded state vector.

Initial value problems appear ubiquitously in numerical solutions of PDE's, typically through [method of lines](#) (MOL) discretizations.

In MOL discretizations, one discretizes *all but one* variable, and the one remaining variable is typically time:

$$u_t = u_{xx} + u u_x \longrightarrow \mathbf{u}'(t) = \mathbf{D}^2 \mathbf{u} + \mathbf{u} \circ (\mathbf{D} \mathbf{u}) =: \mathbf{f}(\mathbf{u}).$$

Above, \mathbf{D} is some matrix discretization of the differentiation operation, and \circ is the Hadamard (elementwise) product.

I.e., first we have discretized in x (physical) space, but left t (time) alone.

Initial value problems appear ubiquitously in numerical solutions of PDE's, typically through [method of lines](#) (MOL) discretizations.

In MOL discretizations, one discretizes *all but one* variable, and the one remaining variable is typically time:

$$u_t = u_{xx} + u u_x \longrightarrow \mathbf{u}'(t) = \mathbf{D}^2 \mathbf{u} + \mathbf{u} \circ (\mathbf{D} \mathbf{u}) =: \mathbf{f}(\mathbf{u}).$$

Above, \mathbf{D} is some matrix discretization of the differentiation operation, and \circ is the Hadamard (elementwise) product.

I.e., first we have discretized in x (physical) space, but left t (time) alone.

MOL discretizations are not the only way to develop PDE solvers, but they appear so frequently that understanding how to solve the resulting IVP's is quite important.

In these slides we'll focus exclusively on numerically solving the IVP, assuming \mathbf{f} is given.

Well-posedness for IVP's is a quite well-studied topic.

Theorem (Picard-Lindelöf)

Suppose that $\mathbf{f} : \mathbb{R} \times \mathbb{R}^M \rightarrow \mathbb{R}^M$ is continuous in some open ball around $(0, \mathbf{u}_0) \in \mathbb{R} \times \mathbb{R}^M$, and further is Lipschitz continuous in the \mathbf{u} variable in this ball. Consider the initial value problem,

$$\mathbf{u}'(t) = \mathbf{f}(t; \mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0.$$

Then there exists some $\epsilon > 0$ such that there is a unique solution $\mathbf{u}(t)$ to the above problem for $t \in [-\epsilon, \epsilon]$.

Well-posedness for IVP's is a quite well-studied topic.

Theorem (Picard-Lindelöf)

Suppose that $\mathbf{f} : \mathbb{R} \times \mathbb{R}^M \rightarrow \mathbb{R}^M$ is continuous in some open ball around $(0, \mathbf{u}_0) \in \mathbb{R} \times \mathbb{R}^M$, and further is Lipschitz continuous in the \mathbf{u} variable in this ball. Consider the initial value problem,

$$\mathbf{u}'(t) = \mathbf{f}(t; \mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0.$$

Then there exists some $\epsilon > 0$ such that there is a unique solution $\mathbf{u}(t)$ to the above problem for $t \in [-\epsilon, \epsilon]$.

In the proof of this theorem, one constructs \mathbf{u} such that,

$$\mathbf{u}(t) = \mathbf{u}(0) + \int_0^t \mathbf{f}(t; \mathbf{u}(t)) dt.$$

This formula is the starting point for many numerical schemes.

$$\mathbf{u}'(t) = \mathbf{f}(t; \mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0.$$

An easy strategy is to discretize time with equispaced values up to a terminal time T :

$$t_0 := 0, \quad t_n := nk = n\Delta t, \quad T = Nk.$$

I.e., we either choose $\Delta t = k$, or N .

$$\mathbf{u}'(t) = \mathbf{f}(t; \mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0.$$

An easy strategy is to discretize time with equispaced values up to a terminal time T :

$$t_0 := 0, \quad t_n := nk = n\Delta t, \quad T = Nk.$$

I.e., we either choose $\Delta t = k$, or N .

The exact solution satisfies

$$\mathbf{u}(t_{n+1}) = \mathbf{u}(t_n) + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt,$$

and therefore we seek to compute numerical approximations $u_n \approx u(t_n)$ that would satisfy,

$$\mathbf{u}_{n+1} \approx \mathbf{u}_n + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt$$

$$\mathbf{u}'(t) = \mathbf{f}(t; \mathbf{u}),$$

$$\mathbf{u}_n \approx \mathbf{u}(t_n)$$

$$\mathbf{u}_{n+1} \approx \mathbf{u}_n + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt$$

$$\mathbf{u}(0) = \mathbf{u}_0.$$

A great many schemes result from discretization of the integral via quadrature.

$$\begin{aligned} \mathbf{u}'(t) &= \mathbf{f}(t; \mathbf{u}), & \mathbf{u}(0) &= \mathbf{u}_0. \\ \mathbf{u}_n &\approx \mathbf{u}(t_n) \\ \mathbf{u}_{n+1} &\approx \mathbf{u}_n + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt \end{aligned}$$

A great many schemes result from discretization of the integral via quadrature.

The [Forward Euler](#) method uses the quadrature rule,

$$\int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt \approx (t_{n+1} - t_n) \mathbf{f}(t, \mathbf{u}(t)) \Big|_{t=t_n} \approx k \mathbf{f}(t_n, \mathbf{u}_n) =: k \mathbf{f}_n,$$

leading to the scheme,

$$\mathbf{u}_{n+1} = \mathbf{u}_n + k \mathbf{f}_n.$$

This scheme is [explicit](#).

$$\begin{aligned} \mathbf{u}'(t) &= \mathbf{f}(t; \mathbf{u}), & \mathbf{u}(0) &= \mathbf{u}_0. \\ \mathbf{u}_n &\approx \mathbf{u}(t_n) \\ \mathbf{u}_{n+1} &\approx \mathbf{u}_n + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt \end{aligned}$$

A great many schemes result from discretization of the integral via quadrature.

The [Backward Euler](#) method uses the quadrature rule,

$$\int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt \approx (t_{n+1} - t_n) \mathbf{f}(t, \mathbf{u}(t)) \Big|_{t=t_{n+1}} \approx k \mathbf{f}(t_{n+1}, \mathbf{u}_{n+1}) =: k \mathbf{f}_{n+1},$$

leading to the scheme,

$$\mathbf{u}_{n+1} = \mathbf{u}_n + k \mathbf{f}_{n+1}.$$

This scheme is [implicit](#).

$$\begin{aligned} \mathbf{u}'(t) &= \mathbf{f}(t; \mathbf{u}), & \mathbf{u}(0) &= \mathbf{u}_0. \\ \mathbf{u}_n &\approx \mathbf{u}(t_n) \\ \mathbf{u}_{n+1} &\approx \mathbf{u}_n + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt \end{aligned}$$

A great many schemes result from discretization of the integral via quadrature.

The [Crank-Nicolson](#) method uses the Trapezoid rule approximation,

$$\begin{aligned} \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{u}(t)) dt &\approx \frac{(t_{n+1} - t_n)}{2} \left(\mathbf{f}(t, \mathbf{u}(t)) \Big|_{t=t_n} + \mathbf{f}(t, \mathbf{u}(t)) \Big|_{t=t_{n+1}} \right) \\ &= \frac{k}{2} (\mathbf{f}_n + \mathbf{f}_{n+1}). \end{aligned}$$

leading to the scheme,

$$\mathbf{u}_{n+1} = \mathbf{u}_n + \frac{k}{2} (\mathbf{f}_n + \mathbf{f}_{n+1}).$$

This scheme is also implicit.

We might reasonably define convergence as the following:

Definition

A scheme for an ODE is convergent to order p if

$$\max_{n \in [N]} \|e_n\| = \mathcal{O}(k^p),$$

$$e_n := \mathbf{u}(t_n) - \mathbf{u}_n.$$

for some choice of norm $\|\cdot\|$ on M -dimensional vectors.

We might reasonably define convergence as the following:

Definition

A scheme for an ODE is convergent to order p if

$$\max_{n \in [N]} \|e_n\| = \mathcal{O}(k^p), \quad e_n := \mathbf{u}(t_n) - \mathbf{u}_n.$$

for some choice of norm $\|\cdot\|$ on M -dimensional vectors.

As you might expect, we cannot tackle convergence directly without some setup involving consistency and stability.

For ODE IVP's, our **local truncation error** is again defined as the residual of the scheme when the *exact* solution is inserted, and is **consistent** if $k \downarrow 0$ causes the LTE to vanish.

To fix details, let's consider forward Euler:

$$\mathbf{u}_{n+1} = \mathbf{u}_n + k\mathbf{f}_n.$$

First, we rewrite the scheme to match the units of the original ODE:

$$\frac{\mathbf{u}_{n+1} - \mathbf{u}_n}{k} = \mathbf{f}_n \quad \longrightarrow \quad D^+ \mathbf{u}_n - \mathbf{f}_n = \mathbf{0},$$

where we have introduced the difference operator $D^+ \mathbf{u}_n := 1/k(\mathbf{u}_{n+1} - \mathbf{u}_n)$.

Note the *superscript* for +.

For ODE IVP's, our **local truncation error** is again defined as the residual of the scheme when the *exact* solution is inserted, and is **consistent** if $k \downarrow 0$ causes the LTE to vanish.

To fix details, let's consider forward Euler:

$$\mathbf{u}_{n+1} = \mathbf{u}_n + k\mathbf{f}_n.$$

First, we rewrite the scheme to match the units of the original ODE:

$$\frac{\mathbf{u}_{n+1} - \mathbf{u}_n}{k} = \mathbf{f}_n \quad \longrightarrow \quad D^+ \mathbf{u}_n - \mathbf{f}_n = \mathbf{0},$$

where we have introduced the difference operator $D^+ \mathbf{u}_n := 1/k(\mathbf{u}_{n+1} - \mathbf{u}_n)$.

Note the *superscript* for $+$.

After a short computation, we conclude:

$$\text{LTE}_n := \|D^+ \mathbf{u}(t_n) - \mathbf{f}(t_n, \mathbf{u}(t_n))\| \simeq k\|\mathbf{u}''(t_n)\| = Ck,$$

i.e., our scheme is consistent to first order.

This *suggests* what convergence we should expect.

There are several notions of stability for numerical methods for IVP's.

Our previous definition of stability concerned the map from scheme inputs to outputs. Here, that is the map from initial data to, say, \mathbf{u}_n for any/all $n \in [N]$.

There are several notions of stability for numerical methods for IVP's.

Our previous definition of stability concerned the map from scheme inputs to outputs. Here, that is the map from initial data to, say, \mathbf{u}_n for any/all $n \in [N]$.

Unfortunately, this map is *much* more complicated than our previous example. (FD methods for $u'' = f$)

To see why, note that

$$\frac{\mathbf{u}_{n+1} - \mathbf{u}_n}{k} = \mathbf{f}_n = \mathbf{f}(t_n, \mathbf{u}_n) \neq \mathbf{f}(t_n, \mathbf{u}(t_n)),$$

i.e., at every step, we are actually solving a *perturbed* ODE system that accumulates errors from previous steps.

It turns out that one way to successfully control these perturbations is to assume that the scheme residual operating on the exact solution does not grow out of control.

For Forward Euler, this is:

$$R_n \mathbf{u}(t_n) := D^+ \mathbf{u}(t_n) - \mathbf{f}(t_n, u(t_n)).$$

Definition

A numerical scheme is **0-stable** if there is some constant C such that for all h sufficiently small,

$$\max_{n \in [N]} \|\mathbf{e}_n\| \leq C \left(\|\mathbf{e}_0\| + \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\| \right).$$

(Above, R_n is the residual of the corresponding time integration scheme.)

$$\max_{n \in [N]} \|\mathbf{e}_n\| \leq C \left(\|\mathbf{e}_0\| + \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\| \right).$$

It's not too difficult to see how 0-stability and consistency yield convergence:

- In most practical cases, $\mathbf{e}_0 = \mathbf{0}$.
- We have $\|R_n \mathbf{u}(t_n)\| = \text{LTE}_n$.

$$\max_{n \in [N]} \|\mathbf{e}_n\| \leq C \left(\|\mathbf{e}_0\| + \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\| \right).$$

It's not too difficult to see how 0-stability and consistency yield convergence:

- In most practical cases, $\mathbf{e}_0 = \mathbf{0}$.
- We have $\|R_n \mathbf{u}(t_n)\| = \text{LTE}_n$.

i.e.,

$$\text{Consistency} + 0\text{-stability} \implies \text{Convergence}$$

It remains to establish that schemes of interest are consistent and 0-stable.

A technical but somewhat straightforward computation shows that Forward Euler is 0-stable.
With,

$$R_n = \frac{u(t_n) - u(t_{n-1})}{k} - k f(t_{n-1}, u(t_{n-1}))$$

$$\omega := \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\|,$$

$$\underline{u_n = u_{n-1} + k f_{n-1}}$$

then

$$\begin{aligned} \mathbf{e}_n = \mathbf{u}(t_n) - \mathbf{u}_n &= \mathbf{u}(t_n) - \mathbf{u}(t_{n-1}) + \mathbf{u}(t_{n-1}) - \mathbf{u}_{n-1} + \mathbf{u}_{n-1} - \mathbf{u}_n \\ &= \mathbf{e}_{n-1} + k R_{n-1} \mathbf{u}(t_n) + k \mathbf{f}(t_{n-1}, \mathbf{u}(t_{n-1})) - k \mathbf{f}_{n-1} \\ &= \mathbf{e}_{n-1} + k [\mathbf{f}(t_{n-1}, \mathbf{u}(t_{n-1})) - \mathbf{f}(t_{n-1}, \mathbf{u}_{n-1})] + k R_{n-1} \mathbf{u}(t_n) \end{aligned}$$

A technical but somewhat straightforward computation shows that Forward Euler is 0-stable. With,

$$\omega := \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\|,$$

then

$$\begin{aligned} \mathbf{e}_n &= \mathbf{u}(t_n) - \mathbf{u}_n = \mathbf{u}(t_n) - \mathbf{u}(t_{n-1}) + \mathbf{u}(t_{n-1}) - \mathbf{u}_{n-1} + \mathbf{u}_{n-1} - \mathbf{u}_n \\ &= \mathbf{e}_{n-1} + kR_n \mathbf{u}(t_n) + k\mathbf{f}(t_{n-1}, \mathbf{u}(t_{n-1})) - k\mathbf{f}_{n-1} \\ &= \mathbf{e}_{n-1} + k[\mathbf{f}(t_{n-1}, \mathbf{u}(t_{n-1})) - \mathbf{f}(t_{n-1}, \mathbf{u}_{n-1})] + kR_n \mathbf{u}(t_n) \end{aligned}$$

Therefore, we have,

$$\begin{aligned} \|\mathbf{e}_n\| &\leq \|\mathbf{e}_{n-1}\| + k\|\mathbf{f}(t_{n-1}, \mathbf{u}(t_{n-1})) - \mathbf{f}(t_{n-1}, \mathbf{u}_{n-1})\| + k\omega \\ &\leq \|\mathbf{e}_{n-1}\| + kL\|\mathbf{u}(t_{n-1}) - \mathbf{u}_{n-1}\| + k\omega \\ &= \|\mathbf{e}_{n-1}\| (1 + kL) + k\omega, \end{aligned}$$

where the last inequality uses (assumed!) Lipschitz continuity of \mathbf{f} :

$$\|\mathbf{f}(t, \mathbf{x}) - \mathbf{f}(t, \mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

$$\omega = \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\|,$$

$$\|\mathbf{e}_n\| \leq \|\mathbf{e}_{n-1}\| (1 + kL) + k\omega,$$

Iterating on the inequality, we have,

$$\begin{aligned} \|\mathbf{e}_n\| &\leq \|\mathbf{e}_{n-1}\| (1 + kL) + k\omega \\ &\leq (1 + kL)^2 \|\mathbf{e}_{n-2}\| + k\omega (1 + (1 + kL)) \end{aligned}$$

⋮

$$\leq (1 + kL)^n \|\mathbf{e}_0\| + k\omega \sum_{j=0}^{n-1} (1 + kL)^j$$

$$\|\mathbf{e}_0\| = 0 \quad k\omega \sum_{j=0}^{n-1} (1 + kL)^j$$

$$(1 + a/N)^N \leq \exp a \quad k\omega \sum_{j=0}^{n-1} e^{Lt_j}$$

$$x \mapsto e^x \text{ is convex} \quad \omega \int_0^T e^{Lt} dt \leq \frac{e^{LT}}{L} \omega = C \max_{n \in [N]} \|R_n \mathbf{u}(t_n)\|.$$



Hence, we have that Forward Euler is convergent:

$$\text{LTE}_n = \mathcal{O}(k) \xrightarrow{\text{0-stability}} \max_{n \in [N]} \|e_n\| = \mathcal{O}(k).$$

Hence, we have that Forward Euler is convergent:

$$\text{LTE}_n = \mathcal{O}(k) \xrightarrow{\text{0-stability}} \max_{n \in [N]} \|e_n\| = \mathcal{O}(k).$$

This proves strict first-order convergence, but typically the hidden constants are *very* large, e.g., behave like e^{LT} .

Note that the statements above are asymptotic in k as $k \downarrow 0$. More refined (and in some sense useful) analysis can be achieved if we consider finite, nonzero k .



Ascher, Uri M. and Linda R. Petzold (1998). *Computer Methods for Ordinary Differential Equations and Differential-Algebraic Equations*. SIAM. ISBN: 978-1-61197-139-2.