

Math 6620: Analysis of Numerical Methods, II

Finite difference methods for 1D stationary problems

See LeVeque 2007, Chapter 2

Akil Narayan¹

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



One of our main strategies for (numerically) solving differential equations will be to replace derivatives at locations with (secant, “finite”) differences of function evaluations.

Finite difference methods are a good starting point to understand numerical methods: they are simple, easy to understand, and (typically) easy to implement.

One of our main strategies for (numerically) solving differential equations will be to replace derivatives at locations with (secant, "finite") differences of function evaluations.

Finite difference methods are a good starting point to understand numerical methods: they are simple, easy to understand, and (typically) easy to implement.

The basic idea is to approximate derivatives in DE's with *finite* difference approximations:

$$u'(x) \approx \frac{u(x+h) - u(x)}{h} \quad \text{or} \quad \frac{u(x) - u(x-h)}{h} \quad \text{or} \quad \frac{u(x+h) - u(x-h)}{2h}$$

The above are only examples, but are conceptually useful to understand the overall picture.

How accurate are these approximations?

$$\frac{u(x+h) - u(x)}{h} = u'(x) + \text{"error"}$$

$$\text{Note: } u(x+h) = \sum_{j=0}^{\infty} \frac{u^{(j)}(x)}{j!} h^j = u(x) + hu'(x) + \frac{h^2}{2} u''(x) + \dots$$

$$\frac{u(x+h)-u(x)}{h} = u'(x) + \frac{h}{2} u''(x) + \dots$$
$$= u'(x) + O(h) \quad \left(\lim_{h \rightarrow 0} \frac{O(h)}{h} = C \in \mathbb{R} \right)$$

$$\frac{u(x+h)-u(x-h)}{2h} = u'(x) + O(h^2)$$

$\frac{u(x+h)-u(x)}{h}$ is "first-order accurate", "stencil" is $\{x, x+h\}$

$\frac{u(x+h)-u(x-h)}{2h}$ is "second-order accurate", "stencil" is $\{x, x-h, x+h\}$

One of our main strategies for (numerically) solving differential equations will be to replace derivatives at locations with (secant, “finite”) differences of function evaluations.

Finite difference methods are a good starting point to understand numerical methods: they are simple, easy to understand, and (typically) easy to implement.

The basic idea is to approximate derivatives in DE's with *finite* difference approximations:

$$u'(x) \approx \frac{u(x+h) - u(x)}{h} \quad \text{or} \quad \frac{u(x) - u(x-h)}{h} \quad \text{or} \quad \frac{u(x+h) - u(x-h)}{2h}$$

The above are only examples, but are conceptually useful to understand the overall picture.

The above differences have an *order* of accuracy, which corresponds to the homogeneous degree of asymptotic behavior of the error:

$$u'(x) = \frac{u(x+h) - u(x)}{h} + \mathcal{O}(h)$$

$$u'(x) = \frac{u(x) - u(x-h)}{h} + \mathcal{O}(h)$$

$$u'(x) = \frac{u(x+h) - u(x-h)}{2h} + \mathcal{O}(h^2)$$

One can similarly generate approximations in very general scenarios:

$$u'(x) = Au(x+2h) + Bu(x+h) + Cu(x-h) + \mathcal{O}(h^p), \quad p=?$$

$$A, B, C = ?$$

$$u(x \pm h) = u(x) \pm hu'(x) + \frac{h^2}{2} u''(x) \pm \frac{h^3}{6} u'''(x) + \dots$$

$$u(x+2h) = u(x) + 2hu'(x) + 2h^2 u''(x) + \dots$$

$$A u(x+2h) + B u(x+h) + C u(x-h)$$

$$= u(x) [A+B+C] + u'(x) [2hA+Bh-Ch]$$

$$+ u''(x) [2h^2 A + \frac{h^2}{2} B + \frac{h^2}{2} C] + \dots$$

$$A+B+C=0$$

$$2hA+Bh-Ch=1$$

$$2h^2 A + \frac{h^2}{2} B + \frac{h^2}{2} C = 0$$

$$\Rightarrow A=0, B=-\frac{1}{2h}, C=\frac{1}{2h}, p=2$$

With $u(x)$ an unknown function for $x \in [0, 1]$, we will *discretize* u by considering its value at $M + 2$ equispaced points on $[0, 1]$:

$$h := \frac{1}{M + 1}, \quad x_j := jh, \quad j = 0, \dots, M + 1.$$

We will use u_j to denote our computational approximation to $u(x_j)$, i.e.,

$$u_j \approx u(x_j)$$

It will be convenient to use some shorthand for finite difference operators.

With $u(x)$ an unknown function for $x \in [0, 1]$, we will *discretize* u by considering its value at $M + 2$ equispaced points on $[0, 1]$:

$$h := \frac{1}{M + 1}, \quad x_j := jh, \quad j = 0, \dots, M + 1.$$

We will use u_j to denote our computational approximation to $u(x_j)$, i.e.,

$$u_j \approx u(x_j)$$

It will be convenient to use some shorthand for finite difference operators.

With $u_j \approx u(x_j)$ and x_j an equidistant grid of spacing h , we define:

$$\begin{aligned} D_+ u(x) &:= \frac{u(x+h) - u(x)}{h}, & D_- u(x) &:= \frac{u(x) - u(x-h)}{h}, & D_0 u(x) &:= \frac{u(x+h) - u(x-h)}{2h}, \\ D_+ u_j &:= \frac{u_{j+1} - u_j}{h}, & D_- u_j &:= \frac{u_j - u_{j-1}}{h}, & D_0 u_j &:= \frac{u_{j+1} - u_{j-1}}{2h}. \end{aligned}$$

Note that $D_{\pm,0}$ apply in conceptually similar ways to functions $u(x)$ as to discrete values u_j .

With $u(x)$ an unknown function for $x \in [0, 1]$, we will *discretize* u by considering its value at $M + 2$ equispaced points on $[0, 1]$:

$$h := \frac{1}{M + 1}, \quad x_j := jh, \quad j = 0, \dots, M + 1.$$

We will use u_j to denote our computational approximation to $u(x_j)$, i.e.,

$$u_j \approx u(x_j)$$

It will be convenient to use some shorthand for finite difference operators.

With $u_j \approx u(x_j)$ and x_j an equidistant grid of spacing h , we define:

$$\begin{aligned} D_+ u(x) &:= \frac{u(x+h) - u(x)}{h}, & D_- u(x) &:= \frac{u(x) - u(x-h)}{h}, & D_0 u(x) &:= \frac{u(x+h) - u(x-h)}{2h}, \\ D_+ u_j &:= \frac{u_{j+1} - u_j}{h}, & D_- u_j &:= \frac{u_j - u_{j-1}}{h}, & D_0 u_j &:= \frac{u_{j+1} - u_{j-1}}{2h}. \end{aligned}$$

Note that $D_{\pm,0}$ apply in conceptually similar ways to functions $u(x)$ as to discrete values u_j .

These difference operators are convenient for shorthand. For example:

$$D_+ u(x) = u'(x) + \mathcal{O}(h), \quad D_- u(x) = u'(x) + \mathcal{O}(h), \quad D_0 u(x) = u'(x) + \mathcal{O}(h^2).$$

We can chain these operators to approximate higher order derivatives:

$$D_+ D_- u(x) = D_- D_+ u(x) = u''(x) + \mathcal{O}(h^2)$$

Our prototypical equation is an ODE describing the steady-state temperature distribution on a one-dimensional domain:

$$\begin{aligned} -u''(x) &= f(x), & x \in (0, 1) \\ u(0) &= g_0, \\ u(1) &= g_1. \end{aligned}$$

where f , g_0 , and g_1 are presumed given and known.

This is a model for steady-state heat diffusion:

- The ODE models homogeneous, isotropic heat diffusion in an environment.
- $u(x)$ is the temperature at location x .
- The boundary conditions correspond to pinning the temperature value.

Our prototypical equation is an ODE describing the steady-state temperature distribution on a one-dimensional domain:

$$\begin{aligned} -u''(x) &= f(x), & x \in (0, 1) \\ u(0) &= g_0, \\ u(1) &= g_1. \end{aligned}$$

where f , g_0 , and g_1 are presumed given and known.

We construct a finite-difference (FD) scheme for this equation as follows:

$$u''(x) \longrightarrow D_+ D_- u_j = \frac{1}{h^2} (u_{j-1} - 2u_j + u_{j+1})$$

$$[M] = \{1, 2, \dots, M\}$$

The scheme is to assert that $-u''(x_j) = f(x_j)$ for $j \in [M]$. Thus, we have:

$$\begin{aligned} -u_{j-1} + 2u_j - u_{j+1} &= h^2 f_j, & j = 1, \dots, M, \\ u_0 &= g_0, \\ u_{M+1} &= g_1, \end{aligned}$$

where we define $f_j = f(x_j)$, and the last two equalities explicitly enforce boundary conditions at the discrete level.

$$-u_{j-1} + 2u_j - u_{j+1} = h^2 f_j, \quad j = 1, \dots, M,$$

$$u_0 = g_0,$$

$$u_{M+1} = g_1,$$

If we define vectors,

$$\mathbf{u} = (u_1, \dots, u_M)^T,$$

$$\mathbf{f} = (f_1, \dots, f_M)^T,$$

where the vector \mathbf{u} contains our unknowns, we have the linear system,

$$\mathbf{A}\mathbf{u} = \mathbf{f} + \frac{g_0}{h^2} \mathbf{e}_1 + \frac{g_1}{h^2} \mathbf{e}_M,$$

and the matrix \mathbf{A} is symmetric:

$$\mathbf{A} = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & 2 & -1 & & \\ & & & \ddots & \ddots & \ddots \\ & & & & & -1 & 2 \end{pmatrix}$$

$$e_j \in \mathbb{R}^M$$

$$(e_j)_i = 1$$

$$(e_j)_k = 0, k \neq j$$

Goal: compute the vector \mathbf{u} .

To summarize: we have discretized the ODE

$$\begin{aligned} -u''(x) &= f(x), & x \in (0, 1) \\ u(0) &= g_0, \\ u(1) &= g_1 \end{aligned}$$

to obtain the linear system

$$\mathbf{A}\mathbf{u} = \mathbf{f} + \frac{g_0}{h^2}\mathbf{e}_1 + \frac{g_1}{h^2}\mathbf{e}_M,$$

Some observations worth noting:

- $u \mapsto u''$ is a *symmetric* operator, and \mathbf{A} is a symmetric matrix.
- \mathbf{A} is invertible (actually, its spectrum is explicitly computable)
- \mathbf{A} is *sparse*, having only $3M - 2$ nonzero entries. (\mathbf{A} is tridiagonal)
- The naive computational cost of this approach is $\mathcal{O}(M^3)$, as that is the brute-force cost to invert an $M \times M$ matrix.
- For this particular problem, there are $\mathcal{O}(M)$ algorithms to solve the linear system.

To summarize: we have discretized the ODE

$$\begin{aligned} -u''(x) &= f(x), & x \in (0, 1) \\ u(0) &= g_0, \\ u(1) &= g_1 \end{aligned}$$

to obtain the linear system

$$\mathbf{A}\mathbf{u} = \mathbf{f} + \frac{g_0}{h^2}\mathbf{e}_1 + \frac{g_1}{h^2}\mathbf{e}_M,$$

Some observations worth noting:

- $u \mapsto u''$ is a *symmetric* operator, and \mathbf{A} is a symmetric matrix.
- \mathbf{A} is invertible (actually, its spectrum is explicitly computable)
- \mathbf{A} is *sparse*, having only $3M - 2$ nonzero entries.
- The naive computational cost of this approach is $\mathcal{O}(M^3)$, as that is the brute-force cost to invert an $M \times M$ matrix.
- For this particular problem, there are $\mathcal{O}(M)$ algorithms to solve the linear system.

We can compute \mathbf{u} . But is it true that $u_j \approx u(x_j)$?

Does the numerical solution become more accurate as $h \downarrow 0$?

Our high-level questions regard *convergence* of the scheme.

Before addressing these, consider a simpler question about “consistency”.

Definition

The **Local Truncation Error** (LTE) τ for a scheme is the residual of the scheme when the *exact* solution $u(x_j)$ is inserted in place of u_j .

$$(-u'' = f)$$

Our high-level questions regard *convergence* of the scheme.

Before addressing these, consider a simpler question about “consistency”.

Definition

The **Local Truncation Error** (LTE) τ for a scheme is the residual of the scheme when the exact solution $u(x_j)$ is inserted in place of u_j .

$$\tau_j := \overbrace{(-D_+ D_- u(x_j) - f(x_j))}^{(-D_+ D_- u(x_j) - f(x_j))}$$

This is the error in the ODE statement at x_j due to our discretization.

An exercise shows that,

$$\tau_j = ch^2 u^{(4)}(x_j) + \mathcal{O}(h^4),$$

where c is an absolute constant.

$$\begin{aligned} \tau_j &= \frac{1}{h^2} (-u(x_{j-1}) + 2u(x_j) - u(x_{j+1})) - f(x_j) \\ &= -u''(x_j) + \mathcal{O}(h^2) - f(x_j) = \mathcal{O}(h^2) \end{aligned}$$

The LTE, τ_j , is actually a function on the grid.

A 'good' outcome would be that this grid function decays to 0 as $h \downarrow 0$.

Hence, our estimates will use the norm of the LTE:

$$\boldsymbol{\tau} = (\tau_1, \dots, \tau_M)^T,$$

$$\|\boldsymbol{\tau}\|_{2,h} = \sqrt{h} \|\boldsymbol{\tau}\|_2$$
$$\|\boldsymbol{\tau}\|_2^2 := h \sum_{j=1}^M |\tau_j|^2.$$

The LTE, τ_j , is actually a function on the grid.

A ‘good’ outcome would be that this grid function decays to 0 as $h \downarrow 0$.

Hence, our estimates will use the norm of the LTE:

$$\boldsymbol{\tau} = (\tau_1, \dots, \tau_M)^T, \quad \|\boldsymbol{\tau}\|_2^2 := h \sum_{j=1}^M |\tau_j|^2.$$

Note that M (the size of $\boldsymbol{\tau}$) scales like $1/h$, and that the 2-norm is scaled by h . This scaling factor is sensible,

$$\int_0^1 \tau^2(x) dx \approx h \sum_{j=1}^M \tau^2(x_j),$$

where $\tau(\cdot)$ is a putative function representing the LTE function on the continuum $[0, 1]$.

Clearly a “small” LTE is desirable – a particular notion of “small” is called consistency.

Definition

We say that a numerical scheme is **consistent** if

$$\lim_{h \downarrow 0} \|\tau\|_2 = 0.$$

Clearly a “small” LTE is desirable – a particular notion of “small” is called consistency.

Definition

We say that a numerical scheme is **consistent** if

$$\lim_{h \downarrow 0} \|\tau\|_2 = 0.$$

In our particular case, we have,

$$\|\tau\|_2 = \mathcal{O}(h^2) \xrightarrow{h \downarrow 0} 0,$$

hence our discretization is consistent.

Clearly a “small” LTE is desirable – a particular notion of “small” is called consistency.

Definition

We say that a numerical scheme is **consistent** if

$$\lim_{h \downarrow 0} \|\tau\|_2 = 0.$$

In our particular case, we have,

$$\|\tau\|_2 = \mathcal{O}(h^2) \xrightarrow{h \downarrow 0} 0,$$

hence our discretization is consistent.

Because we know that the LTE is $\mathcal{O}(h^2)$, we might also say that the scheme is consistent to second order.

Our use of the 2-norm in the definition of consistency is a choice – other norms are/can be useful in other situations.

Note that consistency does *not* immediately translate into accuracy of the computed numerical solution, though it does suggest what we should expect.

In order to translate consistency into scheme accuracy, we will need the scheme to “behave well” for small h . This is stability.

In order to translate consistency into scheme accuracy, we will need the scheme to “behave well” for small h . This is stability.

Recall our scheme is

$$\mathbf{A}\mathbf{u} = \mathbf{f} + \frac{g_0}{h^2}\mathbf{e}_1 + \frac{g_1}{h^2}\mathbf{e}_M,$$

and that everything on the right hand side is an input parameter (\mathbf{f}, g_0, g_1) .

Thus, abstractly we can view our scheme as the input-to-output map,

$$\mathbf{f}, g_0, g_1 \xrightarrow{\mathbf{A}^{-1}} \mathbf{u},$$

and hence we need \mathbf{A}^{-1} to behave well.

In order to translate consistency into scheme accuracy, we will need the scheme to “behave well” for small h . This is stability.

Recall our scheme is

$$\mathbf{A}u = \mathbf{f} + \frac{g_0}{h^2}e_1 + \frac{g_1}{h^2}e_M,$$

and that everything on the right hand side is an input parameter (f, g_0, g_1) .

Thus, abstractly we can view our scheme as the input-to-output map,

$$\mathbf{f}, g_0, g_1 \xrightarrow{\mathbf{A}^{-1}} u,$$

and hence we need \mathbf{A}^{-1} to behave well.

Definition

We say that our scheme is **stable** if

$$\|\mathbf{A}^{-1}\|_2 \leq C \quad \text{for all } h \text{ sufficiently small,}$$

$\|A\|_2 := \sup_{\|x\|_2=1} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{\|x\|_2=1} \|Ax\|_2$

where C is independent of h . “Sufficiently small” means $\exists h_0 > 0$ so that the inequality holds $\forall h \in (0, h_0)$.

Note that the size of \mathbf{A} depends on h , and in particular goes to infinity as h goes to 0.

We can verify stability for our scheme. Recall that,

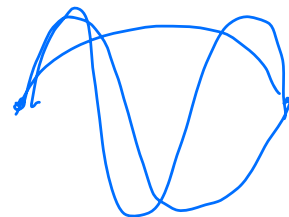
$$\mathbf{A} = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & -1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 \end{pmatrix}$$

One can explicitly compute the spectrum of this matrix:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}, \quad (\lambda_j, \mathbf{v}_j) = \left(\frac{4}{h^2} \sin^2(\pi h j / 2), \sin(\mathbf{x} j \pi) \right), \quad j \in [M]$$

where $\mathbf{x} = (x_1, \dots, x_M)^T$.

$$\mathbf{v}_j = \begin{pmatrix} \sin(x_{1j}\pi) \\ \vdots \\ \sin(x_{mj}\pi) \end{pmatrix}$$



$$-u'' = \lambda u \\ u(0) = u(1) = 0$$

\underline{v}_j is an eigenvector of \underline{A} : $(\underline{A}\underline{v}_j)_k = \frac{1}{h^2} \left[-\sin(x_{k-1}j\pi) \right.$
 $\left. + 2\sin(x_kj\pi) \right.$
 $\left. - \sin(x_{k+1}j\pi) \right]$
 $2 \leq k \leq M-1$

$$x_{k \pm 1} = x_k \pm h$$

$$\sin(x_{k-1}j\pi) + \sin(x_{k+1}j\pi) = 2\sin(x_kj\pi) \cos(j\pi h)$$

$$\Rightarrow (\underline{A}\underline{v}_j)_k = \frac{1}{h^2} \left[2\sin(x_kj\pi) - 2\sin(x_kj\pi) \cos(j\pi h) \right]$$

$$\equiv 2 \frac{1 - \cos(j\pi h)}{h^2} \sin(x_kj\pi) = 2 \frac{1 - \cos(j\pi h)}{h^2} (\underline{v}_j)_k$$

We can verify stability for our scheme. Recall that,

$$\mathbf{A} = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & -1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 \end{pmatrix}$$

One can explicitly compute the spectrum of this matrix:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}, \quad (\lambda_j, \mathbf{v}_j) = \left(\frac{4}{h^2} \sin^2(\pi h j / 2), \sin(\mathbf{x} j \pi) \right), \quad j \in [M]$$

where $\mathbf{x} = (x_1, \dots, x_M)^T$.

Of particular interest is that the eigenvectors of \mathbf{A} “look” similar to those of the second derivative operator $u \mapsto -u''$. (This suggests we’re not doing anything too crazy.)

How does this help with stability? We need to identify asymptotic behavior of $\|\mathbf{A}\|^{-1}$.

$$\|\mathbf{A}^{-1}\|.$$

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{C}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{C}^{M \times M}$ and $\mathbf{V} \in \mathbb{C}^{N \times N}$, and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{M \times N}$, such that $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The diagonal elements of $\mathbf{\Sigma}$ are $\{\sigma_j\}_j$, ordered such that $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, are called the singular values of \mathbf{A} .

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{C}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{C}^{M \times M}$ and $\mathbf{V} \in \mathbb{C}^{N \times N}$, and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{M \times N}$, such that $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The diagonal elements of $\mathbf{\Sigma}$ are $\{\sigma_j\}_j$, ordered such that $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, are called the singular values of \mathbf{A} .
- If \mathbf{A} is an $M \times M$ matrix, then $\|\mathbf{A}\|_2 = \sigma_1$, where σ_1 is the largest singular value of \mathbf{A} .

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{C}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{C}^{M \times M}$ and $\mathbf{V} \in \mathbb{C}^{N \times N}$, and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{M \times N}$, such that $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The diagonal elements of $\mathbf{\Sigma}$ are $\{\sigma_j\}_j$, ordered such that $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, are called the singular values of \mathbf{A} .
- If \mathbf{A} is an $M \times M$ matrix, then $\|\mathbf{A}\|_2 = \sigma_1$, where σ_1 is the largest singular value of \mathbf{A} .
- If \mathbf{A} is symmetric, then there exists an orthogonal matrix \mathbf{U} and a diagonal matrix \mathbf{D} , both real-valued, such that

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T.$$

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{C}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{C}^{M \times M}$ and $\mathbf{V} \in \mathbb{C}^{N \times N}$, and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{M \times N}$, such that $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The diagonal elements of $\mathbf{\Sigma}$ are $\{\sigma_j\}_j$, ordered such that $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, are called the singular values of \mathbf{A} .
- If \mathbf{A} is an $M \times M$ matrix, then $\|\mathbf{A}\|_2 = \sigma_1$, where σ_1 is the largest singular value of \mathbf{A} .
- If \mathbf{A} is symmetric, then there exists an orthogonal matrix \mathbf{U} and a diagonal matrix \mathbf{D} , both real-valued, such that

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T.$$

- If \mathbf{A} is both symmetric and invertible, then the diagonal elements of \mathbf{D} are non-zero, and

$$\mathbf{A}^{-1} = \mathbf{U}\mathbf{D}^{-1}\mathbf{U}^T.$$

$$[\sigma(A)]^2 = \lambda(A^*A)$$

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{C}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{C}^{M \times M}$ and $\mathbf{V} \in \mathbb{C}^{N \times N}$, and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{M \times N}$, such that $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The diagonal elements of $\mathbf{\Sigma}$ are $\{\sigma_j\}_j$, ordered such that $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$, are called the singular values of \mathbf{A} .
- If \mathbf{A} is an $M \times M$ matrix, then $\|\mathbf{A}\|_2 = \sigma_1$, where σ_1 is the largest singular value of \mathbf{A} .
- If \mathbf{A} is symmetric, then there exists an orthogonal matrix \mathbf{U} and a diagonal matrix \mathbf{D} , both real-valued, such that

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T.$$

- If \mathbf{A} is both symmetric and invertible, then the diagonal elements of \mathbf{D} are non-zero, and

$$\mathbf{A}^{-1} = \mathbf{U}\mathbf{D}^{-1}\mathbf{U}^T.$$

- If \mathbf{A} is symmetric, then

$$|\lambda_j| = \sigma_j,$$

where $\{\sigma_j\}_j$ are the non-decreasing singular values of \mathbf{A} , and $\{\lambda_j\}_j$ are the eigenvalues of \mathbf{A} , ordered such that $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

We can now finish our stability verification:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}, \quad (\lambda_j, \mathbf{v}_j) = \left(\frac{4}{h^2} \sin^2(\pi j h/2), \sin(\mathbf{x} j \pi) \right), \quad j \in [M]$$

where $\mathbf{x} = (x_1, \dots, x_M)^T$.

Since \mathbf{A} is invertible and symmetric:

$$\begin{aligned} \|\mathbf{A}^{-1}\|_2 &= \sigma_1(\mathbf{A}^{-1}) = \max_j |\lambda_j(\mathbf{A}^{-1})| = \max_j \frac{1}{|\lambda_j(\mathbf{A})|} = \frac{1}{\min_j |\lambda_j(\mathbf{A})|} \\ &= \frac{1}{\frac{4}{h^2} \sin^2(h\pi/2)} = \frac{h^2}{4 \sin^2(h\pi/2)} \end{aligned}$$

We can now finish our stability verification:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}, \quad (\lambda_j, \mathbf{v}_j) = \left(\frac{4}{h^2} \sin^2(\pi j h/2), \sin(\mathbf{x} j \pi) \right), \quad j \in [M]$$

where $\mathbf{x} = (x_1, \dots, x_M)^T$.

Since \mathbf{A} is invertible and symmetric:

$$\begin{aligned} \|\mathbf{A}^{-1}\|_2 &= \sigma_1(\mathbf{A}^{-1}) = \max_j |\lambda_j(\mathbf{A}^{-1})| = \max_j \frac{1}{|\lambda_j(\mathbf{A})|} = \frac{1}{\min_j |\lambda_j(\mathbf{A})|} \\ &= \frac{1}{\frac{4}{h^2} \sin^2(h\pi/2)} = \frac{h^2}{4 \sin^2(h\pi/2)} \end{aligned}$$

We are interested in the $h \downarrow 0$ behavior of this quantity. Since,

$$\sin(x) \approx x \quad \text{as } x \downarrow 0,$$

we conclude that

$$\|\mathbf{A}^{-1}\|_2 \sim \frac{h^2}{4h^2\pi^2/4} = \frac{1}{\pi^2}$$

hence our scheme is stable since $\|\mathbf{A}^{-1}\|_2 \leq C$ for small h .

We are finally in a position to consider our original question: is our scheme accurate? The answer to this question will quantify how large the error e is:

$$e = (e_1, \dots, e_M)^T, \quad e_j := u_j - u(x_j).$$

Definition

A scheme is **convergent** if $\lim_{h \downarrow 0} \|e\|_2 = 0$.

This is a rather strong statement, since as $h \downarrow 0$ we require small error at a larger number of spatial points.

We are finally in a position to consider our original question: is our scheme accurate? The answer to this question will quantify how large the error e is:

$$\mathbf{e} = (e_1, \dots, e_M)^T, \quad e_j := u_j - u(x_j).$$

Definition

A scheme is **convergent** if $\lim_{h \downarrow 0} \|\mathbf{e}\|_2 = 0$.

This is a rather strong statement, since as $h \downarrow 0$ we require small error at a larger number of spatial points.

We can now show the power of *linearity* for this problem. Define a vector containing evaluations of the exact solution:

$$\mathbf{U} = (u(x_1), \dots, u(x_M))^T.$$

Generally, $\mathbf{U} \neq \mathbf{u}$. Now note that,

$$\mathbf{A}\mathbf{u} = \mathbf{f} + \frac{g_0}{h^2}\mathbf{e}_1 + \frac{g_1}{h^2}\mathbf{e}_M \quad (\text{Definition of the scheme})$$

$$\mathbf{A}\mathbf{U} = \mathbf{f} + \frac{g_0}{h^2}\mathbf{e}_1 + \frac{g_1}{h^2}\mathbf{e}_M + \boldsymbol{\tau}, \quad (\text{Consistency})$$

Therefore:

$$\mathbf{A}e = \mathbf{A}(u - U) = -\tau,$$

and so,

$$\|e\|_2 = \|\mathbf{A}^{-1}\tau\| \leq \|\mathbf{A}^{-1}\| \|\tau\| \leq C\mathcal{O}(h^2),$$

where the last inequality uses both stability and consistency.

Therefore:

$$\mathbf{A}e = \mathbf{A}(u - U) = -\tau,$$

and so,

$$\|e\|_2 = \|\mathbf{A}^{-1}\tau\| \leq \|\mathbf{A}^{-1}\| \|\tau\| \leq C\mathcal{O}(h^2),$$

where the last inequality uses both stability and consistency.

We have just proven the following:

Theorem

The second-order difference scheme is convergent, and in particular is second-order convergent.

The “second-order convergent” part means $\|e\|_2 = \mathcal{O}(h^2)$.

Therefore:

$$\mathbf{A}e = \mathbf{A}(u - U) = -\tau,$$

and so,

$$\|e\|_2 = \|\mathbf{A}^{-1}\tau\| \leq \|\mathbf{A}^{-1}\| \|\tau\| \leq C\mathcal{O}(h^2),$$

where the last inequality uses both stability and consistency.

We have just proven the following:

Theorem

The second-order difference scheme is convergent, and in particular is second-order convergent.

The “second-order convergent” part means $\|e\|_2 = \mathcal{O}(h^2)$.

In this particular case, the order of the LTE coincides with the order of convergence. This is not always the case.

We have drawn an outline for how to establish convergence for FD schemes.

Many details are specific to the problem + discretization at hand, but the broad strokes are somewhat general:

- *Consistency*: The local truncation error is related mesh spacing h , ideally polynomially related.
- *Stability*: The scheme behaves in a well-behaved way for small mesh spacing h .
- *Linearity*: The scheme residual when the global error is plugged in is equal to the local truncation error.

We have drawn an outline for how to establish convergence for FD schemes.

Many details are specific to the problem + discretization at hand, but the broad strokes are somewhat general:

- *Consistency*: The local truncation error is related mesh spacing h , ideally polynomially related.
- *Stability*: The scheme behaves in a well-behaved way for small mesh spacing h .
- *Linearity*: The scheme residual when the global error is plugged in is equal to the local truncation error.

Thus, the following idea is true for linear FD schemes:

$$\text{Stability} + \text{Consistency} = \text{Convergence}$$

This is called the *Lax Equivalence Theorem*, or the *Lax-Richtmyer Theorem*.

We have drawn an outline for how to establish convergence for FD schemes.

Many details are specific to the problem + discretization at hand, but the broad strokes are somewhat general:

- *Consistency*: The local truncation error is related mesh spacing h , ideally polynomially related.
- *Stability*: The scheme behaves in a well-behaved way for small mesh spacing h .
- *Linearity*: The scheme residual when the global error is plugged in is equal to the local truncation error.

Thus, the following idea is true for linear FD schemes:

$$\text{Stability} + \text{Consistency} = \text{Convergence}$$

This is called the *Lax Equivalence Theorem*, or the *Lax-Richtmyer Theorem*.

One might really consider this a “meta-theorem”, as the practitioner must decide on the precise definition of what consistency and stability mean.

(Recall: we had to choose norms for stability and consistency definitions, and our stability definition involved a matrix that explicitly depends on the differential equation and the scheme.)

Much of previous technique can be generalized to more complicated setups in 1D:

- Non-uniform grids (derive non-uniform versions of $D_{\pm,0}$)
- Neumann/Robin boundary conditions (discretization of boundary conditions)
- Different error norms (e.g., L^∞ norm error)
- Non-homogeneous diffusion: $(\kappa(x)u'(x))' = f(x)$



Lax, P. D. and R. D. Richtmyer (1956). “Survey of the Stability of Linear Finite Difference Equations”. In: *Communications on Pure and Applied Mathematics* 9.2, pp. 267–293. ISSN: 1097-0312. DOI: 10.1002/cpa.3160090206.



LeVeque, Randall J. (2007). *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-Dependent Problems*. SIAM. ISBN: 978-0-89871-783-9.