

Math 6630: Numerical Solutions of Partial Differential Equations

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

January 18, 2023



Finite difference methods

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.

Finite difference methods

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.

We will need some notation to understand how finite difference methods work: 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.,

$$u_j \approx u(x_j)$$

More notation: difference operators

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 .

More notation: difference operators

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:

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

We can chain these operators to approximate higher order derivatives:

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

1D steady-state diffusion

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.

1D steady-state diffusion

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:

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

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 $f_j = f(x_j)$.

The scheme

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

$$u_0 = g_0,$$

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

If we define vectors,

$$\mathbf{u} = (u_1, \dots, u_M)^T, \quad \mathbf{f} = (f_1, \dots, f_M)^T,$$

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:

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

Goal: compute the vector \mathbf{u} .

Numerical considerations

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

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

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.

Numerical considerations

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 $h \downarrow 0$?

Consistency, I

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

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

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 .

$$\tau_j := (D_+ D_- u(x_j) - f(x_j))$$

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

Consistency, I

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

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

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 .

$$\tau_j := (D_+ D_- u(x_j) - f(x_j))$$

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

An exercise shows that,

$$\tau_j = ch^2 u^{(4)}(x_j) + \mathcal{O}(h^4), \quad (c = 1/12)$$

where c is an absolute constant.

Consistency, I

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

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

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 .

$$\tau_j := (D_+ D_- u(x_j) - f(x_j))$$

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. Our estimates will use the norm of the LTE:

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

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:

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

Consistency, II

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.

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

Consistency, II

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

Definition

We say that a numerical scheme is *consistent* if

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

In our particular case, we have,

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

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.

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

Consistency, II

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.

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

Stability, I

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

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

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,

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

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

Definition

We say that our scheme is *stable* if

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

where C is independent of h .

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

Stability, I

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

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

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

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

where C is independent of h .

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

Stability, I

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

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

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

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

where C is independent of h .

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

Stability, II

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:

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

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 $u \mapsto u''$.

Stability, II

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:

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

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 $u \mapsto u''$.

A linear algebraic interlude

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{R}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{R}^{M \times M}$ and $\mathbf{V} \in \mathbb{R}^{N \times N}$, 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 $\sigma_1 \geq \sigma_2 \geq \dots$, 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 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

A linear algebraic interlude

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{R}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{R}^{M \times M}$ and $\mathbf{V} \in \mathbb{R}^{N \times N}$, 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 $\sigma_1 \geq \sigma_2 \geq \dots$, 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 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

A linear algebraic interlude

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{R}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{R}^{M \times M}$ and $\mathbf{V} \in \mathbb{R}^{N \times N}$, 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 $\sigma_1 \geq \sigma_2 \geq \dots$, 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 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

A linear algebraic interlude

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{R}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{R}^{M \times M}$ and $\mathbf{V} \in \mathbb{R}^{N \times N}$, 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 $\sigma_1 \geq \sigma_2 \geq \dots$, 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 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

A linear algebraic interlude

We now review the following facts from linear algebra:

- If $\mathbf{A} \in \mathbb{R}^{M \times N}$ is any matrix, then it admits a singular value decomposition: there exist two unitary matrices $\mathbf{U} \in \mathbb{R}^{M \times M}$ and $\mathbf{V} \in \mathbb{R}^{N \times N}$, 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 $\sigma_1 \geq \sigma_2 \geq \dots$, 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 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_M|$.

Stability, III

We can now finish our stability verification:

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

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 .

Stability, III

We can now finish our stability verification:

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

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 .

Convergence, I

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:

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

Now note that,

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

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

Convergence, I

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:

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

Now note that,

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

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

Convergence, II

Therefore:

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

and so,

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

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.

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

Convergence, II

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.

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

Convergence, II

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.

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

Convergence for linear FD methods

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 small relative to mesh spacing h .
- *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.

Convergence for linear FD methods

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 small relative to mesh spacing h .
- *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.

Convergence for linear FD methods

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 small relative to mesh spacing h .
- *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.

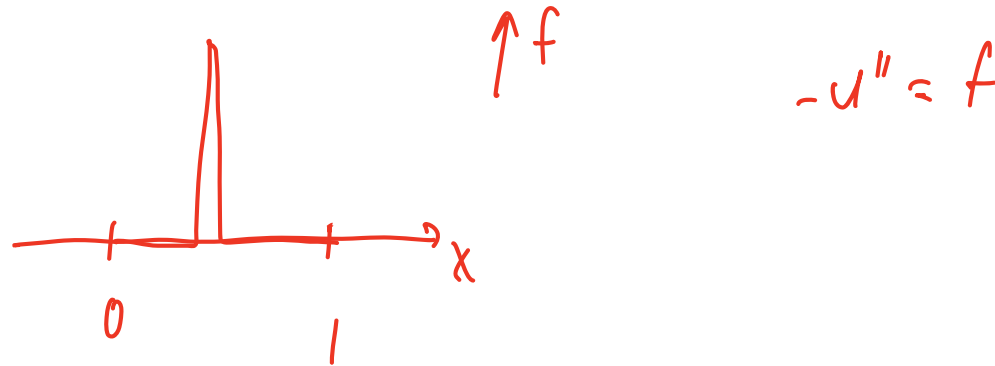
Thus, the following idea is true for linear FD schemes:

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

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.

Some generalizations



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: $(\kappa(x)u'(x))' = f(x)$

$$u(0) = g_0 \quad (\text{Dirichlet})$$

$$u'(0) = h_0 \quad (\text{Neumann})$$

$$\alpha u(0) + \beta u'(0) = j_0 \quad (\text{Robin})$$

References I

-  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. URL: <https://onlinelibrary.wiley.com/doi/abs/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.