# Convergence of Finite Difference

## Methods for BVPs

A. Donev, Courant

We showed last class how to use the FD method to convert the BVP

$$u''(x) = f(x) \qquad 0 < x < 1$$

$$u(0) = \alpha \qquad\qquad u(1) = \beta$$

to the system $\qquad AU = F$

$$\frac{1}{h^2}\begin{bmatrix} -2 & 1 & & & \\ 1 & & & 1 & \\ & & & & 1 \\ & 1 & & & -2 \end{bmatrix} \qquad F = \begin{bmatrix} f_1 - \alpha/h^2 \\ f_2 \\ \vdots \\ f_m - \beta/h^2 \end{bmatrix}$$

Let the true solution evaluated pointwise be

$$\hat{u} = \left[ u(x_1), \ldots, u(x_m) \right]$$

Then the <u>global error</u> is

$$E = u - \hat{u}$$

<u>Question</u> : Does $\| E \| \to 0$ as $h \to 0$ ??? 
$$m \to \infty$$

Put the exact solution in the finite difference to get the <u>local truncation error</u> (*LTE*)

$$\bar{\tau}_j = \frac{1}{h^2} \left[ u(x_{j-1}) - 2u(x_j) + u(x_j) \right] - \underbrace{f(x_j)}_{u''(x_j) \text{ from PDE}}$$

$$\overline{\tau}_j = \frac{h^2}{12} u''''(x_j) + O(h^4)$$

$$\|\overline{\tau}\| = O(h^2)$$

$$\Rightarrow \quad \overline{\tau} = A\hat{U} - F = O(h^2) \quad \Rightarrow$$

$$\begin{cases} A\hat{U} = F + \overline{\tau} \\ A U = F \end{cases} \Big\} \text{ take difference}$$

$$\Rightarrow \quad \boxed{A E = -\overline{\tau}} + \begin{cases} E_0 = 0 \\ E_m = 0 \end{cases} \quad \text{(homogeneous)}$$

Back to this over & over again:

Error satisfies same equation as solution but with LTE as a source term on the R.H.S.

Now we can guess that the global error is also $O(h^2)$

=second order accurate by the following argument:

Since $AE = -\tau$ approximates

$$e''(x) = -\tau(x) \quad , \quad e(0) = e(1) = 0$$

$$e''(x) = -\frac{h^2}{12} u^{(4)}(x) + O(h^4)$$

$$\Rightarrow \quad e(x) \approx -\frac{h^2}{12} u''(x) + \text{boundary terms } O(h^2)$$

$$\boxed{e(x) \approx -\frac{h^2}{12} f(x)} \quad \begin{array}{l}\text{continuum} \\ \text{estimate}\end{array}$$

But this is not a discrete proof !

Discretely, $E = -A^{-1}\bar{z}$ $\Rightarrow$

$$\|E\| \leq \|A^{-1}\| \|\bar{z}\|$$

so if we want both $\|E\|$ and $\|\bar{z}\|$ to be $O(h^2)$ independent of $h$ then

we want

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

A method to solve ; a linear BVP is __stable__ if $A^{-1}$ exists and

$$\|A^{-1}\| \leq C \quad \text{for all } h < h_0$$

A method is <u>consistent</u> if

$$\|\tau\| \to 0 \quad \text{as} \quad h \to 0$$

| stability + consistency $\Rightarrow$ convergence |

| $O(h^p)$ LTE + stability $\Rightarrow$ $O(h^p)$ error |

So now we just need to prove

$$\|A^{-1}\| \leq C$$

and we have proven second order convergence.
For this we need to <u>pick our norm</u>!
And now all norms are <u>NOT equivalent</u>
because infinite dimensional as $h \to 0$

# Stability m $L_2$

$$\| A \|_2 = \underset{\uparrow}{S}(A) = \max_{p} | \lambda_p |$$

spectral
radius

Here A is symmetric so all $\lambda$'s
are real, and same for inverse, so

$$\boxed{\| A^{-1} \|_2 = \left( \min_{p} | \lambda_p | \right)^{-1}}$$

We need to find the eigenvalues.
Now go back to the <u>continuum</u>
PDE and recall that the eigen functions
are just sin functions. So <u>guess</u>
(anzatz)

$$u_j^p = \sin(p\pi j h)$$

$p'$th eigen_vector_

Plug into $A u^p = \lambda^p u^p$ and get

$$\lambda_p = \frac{2}{h^2}\left[\cos(p\pi h) - 1\right]$$

Observe that for small wave_index_
(wavenumber) $p$ we have

$$\lambda_p \approx -\pi^2 p^2 + \frac{1}{12}\pi^4 p^4 h^2$$

also shows second-order_
convergence

Now smallest eigenvalue corresponds
to the longest wavelength ⑨

$$\lambda_1 = \frac{2}{h^2}\left(\cos(\pi h) - 1\right) \approx -\pi^2 + O(h^2)$$

$$= \underline{\underline{const}}$$

$$\Rightarrow \quad \|E\| \leq \|A^{-1}\| \|\tau\| \approx \frac{\|\tau\|}{\pi^2}$$

and we have second order accuracy.

Also recall $\tau \approx \frac{h^2}{12} u^{(4)}(x) = \frac{h^2}{12} f^{(2)}(x)$

So the smoother $f(x)$ is, the less
points we need, which makes sense
physically.

# Max norm stability in $L_\infty$

If we just use finite-dimensional linear algebra, we would bound

$$\| E \|_\infty \leq \frac{1}{\sqrt{h}} \| E \|_2 \implies$$

$$\| E \|_\infty = O(h^{3/2})$$

But this is overly pessimistic, in fact,

$$\| E \|_\infty = O(h^2) \quad \text{as well.}$$

So we need to show

$$\| A^{-1} \|_\infty \leq C$$

For this we go back to Green's functions

What is the $j$'th column
of $A^{-1}$?

Remember

$$\tilde{G}_j = A^{-1} e_j$$

$j$'th column
of $A^{-1}$

zeros but $1$ in entry $j$

$\Rightarrow$ $\boxed{A \tilde{G}_j = e_j}$

But A is a discretitation of <u>Laplacian</u>
with homogeneous BCs (Dirichlet), so
this is a discretitation of

$$\begin{cases} u''(x) = \delta(x - x_j) \cdot h \\ u(0) = u(1) = 0 \end{cases}$$

↑ grid spacing

That is

$\widetilde{G}_j$ is a discrete Green's function

and tells us how an error (local truncation error, roundoff, etc.) at node $j$ translates to error at other nodes.

We can compute $\widetilde{G}_j$ explicitly in this simple case:

$$(\tilde{G}_j)_i = h \begin{cases} (x_j - 1) x_i, & i = 1, \dots, j \quad \text{⑭} \\ (x_i - 1) x_j, & i = j, j+1, \dots, m \end{cases}$$

which is exactly what we expect,

$$(\tilde{G}_j)_i = \boxed{h\, G(x_i, x_j) = A^{-1}_{ij}}$$

↑
Actual Green's function
of PDE

Note: Inhomogeneous Dirichlet BCs can be easily handled, see 2.11 in le Veque.

Now go back to

$$AU = F \implies AE = \tau$$

We need to bound $\|A^{-1}\|$

$$\|A^{-1}\|_\infty = \max_{1 < i \leq m} \sum_{j=1}^{m} (\tilde{G}_j)_i$$

Note that $\left| (\tilde{G}_j)_i \right| \leq h \implies$

since $\quad x(1-x) \leq 1$

$$\|A^{-1}\|_\infty \leq mh = 1$$

so indeed we have <u>stability</u> !

Observe that $\bar{\tau}_j = O(h^2)$ and so if we only made a localized error $\bar{\tau}_j = O(h^2)$ and a much smaller error at other points, then the global error

$$E_i = h \cdot O(h^2) \cdot G(x_i ; x_j)$$

$$E_i = O(h^3)$$

This shows an important fact:

We can make an error of order $p$ at a few points (e.g., boundaries) and still get global order $q > p$
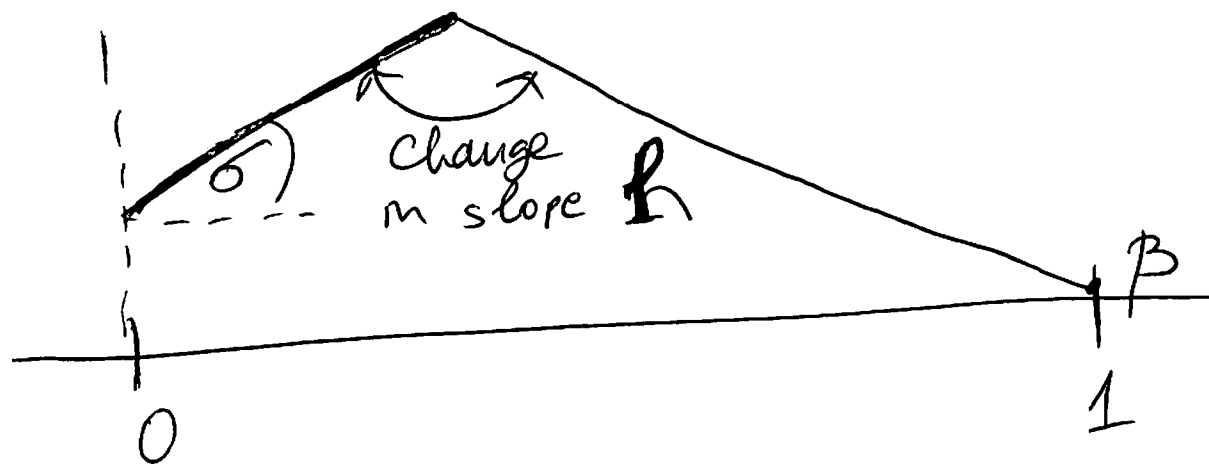
_sometimes_

But <u>not always</u>.

E.g. Consider Neumann BCs

$$u'(0) = \delta \qquad u(1) = \beta$$

(see 2.12 in LeVeque)

What is Green's function now?



If we introduce an error $O(h)$
at the boundary when we impose BC

e.g. $\dfrac{u_1 - u_0}{h} = \sigma$

$$\frac{1}{h^2}\begin{bmatrix} -h & h & & & \\ 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & & \ddots \end{bmatrix}\begin{bmatrix} E_0 \\ E_1 \\ E_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} O(h) \\ O(h^2) \\ O(h^2 \\ \vdots \end{bmatrix}$$

Now we need the first column of $A^{-1}$ — this corresponds to the BVP

$$\begin{cases} u''(x) = 0 \\ u'(0) = 1 \end{cases} \qquad u(1) = 0$$

↑ no factor of $h$ to help us

Now $O(h)$ error on boundary causes $O(h)$ error everywhere

Note that we can compute the Green's function here by hand algebraically. or geometrically:

$$\frac{u_1 - u_0}{h} = 1 \quad \Rightarrow \quad u_1 = u_0 + h$$

After that, slopes are the same because discrete Laplacian is zero

$$u_{j+1} - u_j = u_j - u_{j-1}$$



$$\Rightarrow \quad u_0 = - mh = 1$$

For interior points the difference is that we have a factor of (20)

$$\frac{\frac{u_{j+1} - u_j}{h} - \frac{u_j - u_{j-1}}{h}}{h} = 1$$

So difference in __slopes__ is $\underline{\underline{h}}$, which gave an extra factor of $h$. Therefore for elliptic PDEs it seems OK to make $O(h)$ error at a few interior points and still be second-order, but __NOT__ at boundaries.

[NOTE: PARABOLIC is EASIER for B.C.s]

We can of course use Richardson extrapolation to get higher-order

$$U_j \approx u(jh) \quad , \quad j = 1, \ldots, m$$

$$V_j \approx u(jh/2) \quad , \quad j = 1, \ldots, 2m+1$$

$$U_j \approx V_{2j}$$

Note:

$$e_c = U_j - u(jh) = c_2 h^2 + c_4 h^4 + O(h^6)$$

$$e_f = V_{2j} - u(jh) = c_2 \left(\frac{h}{2}\right)^2 + c_4 \left(\frac{h}{2}\right)^4 + O(h^6)$$

therefore

$$\boxed{\overline{u}_j = \frac{1}{3}\left(4 V_{2j} - V_j\right)} \text{ is } \underline{\text{fourth order}}$$

$$= u(jh) + \frac{1}{3}\left(\frac{1}{4} - 1\right) C_4 \, h^4$$

But for the Poisson equation there
is an even simpler trick:

$$AE = -\overline{\tau} = -\frac{h^2}{12} u^{(4)} + O(h^4)$$

$$= -\frac{h^2}{12} f'' + O(h^4)$$

$$AU = F$$

$$AE \approx -\frac{h^2}{12} \cdot (D^2 f) + O(h^4)$$

numerical second derivative

$$\hat{U} = U - E$$

"true" solution

$$\boxed{A\hat{U} = F + \frac{h^2}{12} D^2 f}$$

is a <u>fourth-order</u> discretization that costs almost the same as the second-order discretization!

$$A\hat{U} = \left(I + \frac{h^2}{12}D^2\right)f$$

$$\left[\left(I + \frac{h^2}{12}D^2\right)^{-1}D^2\right]\hat{U} = f$$

$$\hat{A} = \left(I + \frac{h^2}{12}D^2\right)^{-1}D^2$$

must be a fourth-order approximation of the Laplacian (?) — this is called a "compact difference"

How to show this?

For Periodic BCs, use Fourier

We know the Fourier basis

diagonalizes $D^2$ (in ~~fact~~, any ~~finite~~ difference in a periodic domain).

So let's work in the Fourier basis.

$$D^2 e^{ikx} = \left( \frac{e^{ikh} - 2 + e^{-ikh}}{h^2} \right) e^{ikx}$$

$\underbrace{\phantom{\frac{e^{ikh} - 2 + e^{-ikh}}{h^2}}}$

$\lambda_k = \underline{\text{symbol}}$ of $D^2$

(eigenvalue)

$$\lambda_k = -\frac{\sin^2(kh/2)}{(h/2)^2} = -k^2 + \frac{h^2 k^4}{12} + O(h^4)$$

$\uparrow$ second order

So in Fourier space

$$\hat{D}^2 = -\frac{\sin^2(kh/2)}{(h/2)}$$

So

$$\left(I + \frac{h^2}{12}\hat{D}^2\right)^{-1}\hat{D}^2 = -k^2 + \frac{k^6 h^4}{240} + O(h^6)$$
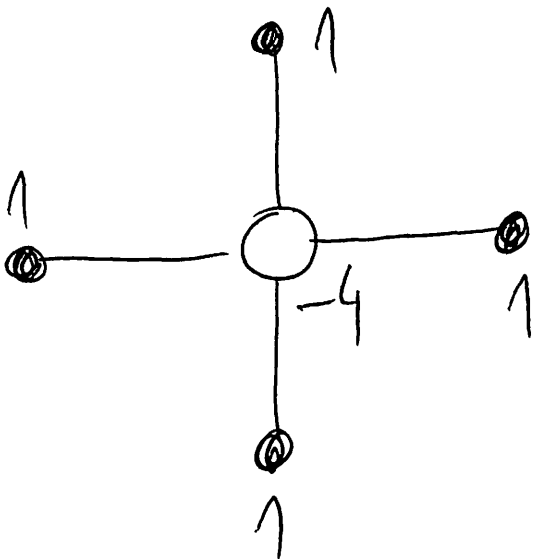
Fourth order!

$\nearrow$ (I used Maple)

So indeed the compact difference
is fourth order

Generalizing to <u>Higher Dimensions</u>
is in <u>principle</u> straightforward

$$\nabla^2 u = U_{xx} + U_{yy} = f(x,y)$$

$$\frac{1}{h^2}\left( U_{i-1,j} + U_{i+1,j} + U_{i,j-1} + U_{i,j+1} - 4U_{ij} \right)$$

$$= f_{ij}$$

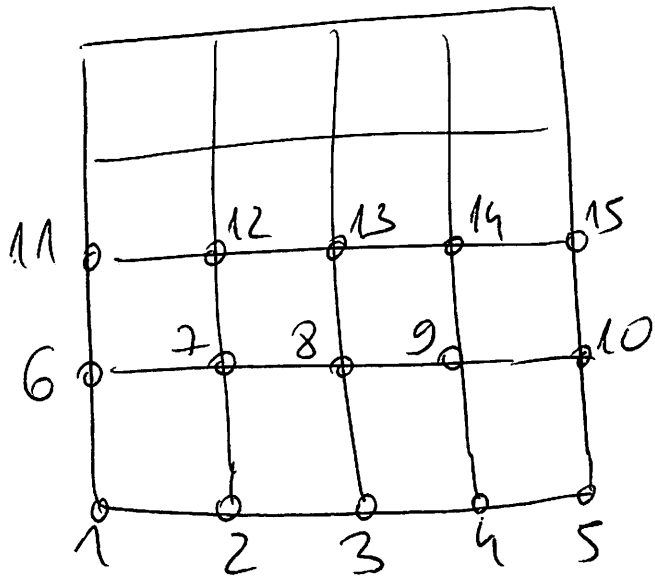$5^{pt}$ Laplacian stencil

Now take a box of $m \times m$ nodes. ㉘

The linear system

$$AU = F \text{ has } m^2 \text{ variables}$$

so $A$ is $[m^2 \times m^2]$ & <u>very sparse</u>

Number of non-zeros is $5m^2$ only



$$A = \frac{1}{h^2}$$

What is the <u>local truncation error</u>?

Since x and y directions are 1D:

$$\tau_{ij} = \frac{1}{12} h^2 \left( u_{xxxx} + u_{yyyy} \right) + O(h^4)$$

To prove second order accuracy we need to bound the norm of $\|A^{-1}\|$

With <u>Dirichlet</u> BCs on all sides the <u>eigen vectors</u> of A are

$$u_{j,j'}^{p,q} = \sin(p\pi j h) \sin(q\pi j' h)$$

or $e^{i(2\pi(p j h + q j' h))}$ for <u>periodic</u> BCs

Since $\quad D^2_{2D} = D^2_x + D^2_y \Rightarrow$

$$\lambda_{2D} = \lambda_x + \lambda_y$$

$$\boxed{\lambda_{p,q} = -\frac{\sin^2(k_x h/2)}{(h/2)^2} - \frac{\sin^2(k_y h/2)}{(h/2)^2}}$$

where $\quad k_x = \frac{2\pi}{L} p \quad$ and $\quad k_y = \frac{2\pi}{L} q$

So all of the properties / analysis from 1D carries through directly

The conditioning number of A is

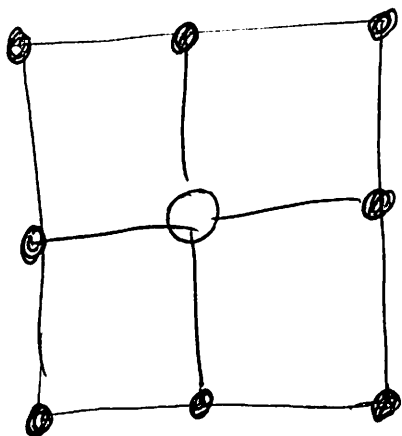$$K_2(A) = \frac{|\lambda_{max}|}{|\lambda_{min}|} = \frac{\frac{8}{h^2}}{2\pi^2} = \frac{4}{\pi^2 h^2}$$

$$\boxed{K_2(A) \sim \frac{1}{h^2} \quad \Big| \quad \sim \frac{1}{m^2}}$$

grows very fast with the resolution $m$

So iterative methods for solving

$AU = F$  will not converge

fast unless we use a preconditioner

Consider also the

$$\underline{\underline{9^{pt} \text{ Laplacian}}} \text{ with stencil}$$

that includes the next-nearest

neighbors:



$$\left(\nabla_9^2 u\right)_{i,j} = \frac{1}{6h^2}\left[4\left(u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1}\right)\right.$$

$$+ \left(u_{i-1,j-1} + u_{i-1,j+1} + u_{i+1,j+1} + u_{i+1,j-1}\right)$$

$$\left. - 20\, u_{ij}\right]$$

$$= \nabla^2 u(x_i, y_j) + \frac{h^2}{12}\left(u_{xxxx} + u_{yyyy} + 2u_{xxyy}\right)$$

$$+ O(h^4)$$

$$(\nabla_g^2 u)_{i,j} = \nabla^2 u(x_i, y_j) + \frac{h^2}{12} \underbrace{(\nabla^2)^2 u}_{\text{isotropic}} + O(h^4)$$

this Laplacian is <u>isotropic</u> to $O(h^4)$, meaning, it is <u>rotationally-invariant</u>. this can be important to preserve physics and symmetries

Furthermore,

$$(\nabla^2)^2 u = \nabla^2 f$$

from the PDE

$$\nabla_g^2 u = \nabla^2 u + \frac{h^2}{12} \nabla^2 f + O(h^4)$$

So we can get a __fourth-order__ discretization of $\nabla^2$ in 2D, i.e., a __compact finite difference__, from

$$\left(\nabla_g^2 u\right)_{i,j} = f_{i,j} + \frac{h^2}{12} \nabla_{i,j}^2 f \quad , \text{e.g.}$$

$$\nabla_g^2 U = F + \frac{h^2}{12} \nabla_5^2 F = \left(I + \frac{h^2}{12} \nabla_5^2\right) F$$

or

$$\boxed{\nabla_g^2 U = \left(I + \frac{h^2}{12} \nabla_g^2\right) F} \quad \underline{\text{fourth order}}$$

Which is related to

$$\nabla^2_{compact} = \left(I + \frac{h^2}{12}\nabla^2_g\right)^{-1}\nabla^2_g$$

Observe that if $F = 0$, i.e., we are solving the Laplace equation, then

$$\nabla^2_g U = 0$$ is a $4^{th}$ order scheme!

So in this sense $\nabla^2_g$ is a <u>better</u> discretization than $\nabla^2_5$. However, it leads to a <u>denser more coupled</u> linear system that is harder to solve efficiently. So often we use $\nabla^2_5$

How expensive is it to solve

$$AU = F$$

in 1D, 2D & 3D ?

$\rightarrow$ <u>Discuss on board</u>

(P)CG note:

$$\|e_k\|_A \leq 2\left(\frac{\sqrt{K_2}-1}{\sqrt{K_2}+1}\right)^k \|e_0\|_A$$

$$\approx 2\left(1-\frac{2}{\sqrt{K_2}}\right)^k \|e_0\|_A$$

<u>Theorem</u> (George)

$\left\{\begin{array}{l} \text{Any } \underline{direct} \text{ method for } AU=F \text{ in 2D requires at} \\ \text{least } \underline{O(m^3) \text{ operations}} \end{array}\right.$

this bound is achieved by a <u>nested-dissection</u> algorithm

We will show in future that <u>multigrid</u> can solve in time $m^2 \log(m)$ in 2D and $m^3 \log(m)$ in 3D — same as FFT & <u>optimal</u>

$\boxed{\text{Factorizing } (a+b) \text{ diagonals matrix takes } N \cdot a \cdot b \text{ operations}}$

Conjugate gradients :

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq 2 \left( \frac{\sqrt{K}-1}{\sqrt{K}+1} \right)^k \simeq 2 \left( 1 - \frac{2}{\sqrt{K}} \right)^k$$

$$\simeq 2 e^{-2k/\sqrt{K}}$$

so $\log \frac{\|e_k\|_A}{\|e_0\|_A} \approx - \frac{2k}{\sqrt{K}}$

So to get a fixed number of <u>accurate digits</u>
we need a number of iterations $\sim \sqrt{K}$

$\#$ iterations $\sim \sqrt{\frac{1}{h^2}} \sim m$

Each iteration costs $O(m)$ in 1D,

$O(m^2)$ in 2D or $O(m^3)$ in 3D