

Standard Model PDEs: 1) Advection equation

prototypical ex. of a hyperbolic eqn $\rightarrow u_t + C u_x = 0$

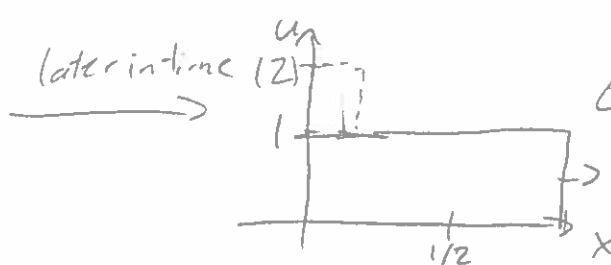
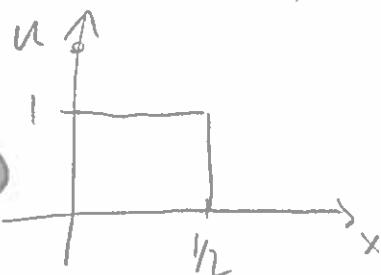
2) Diffusion eqn: $u_t = D u_{xx}$ \leftarrow prototypical parabolic eqn.
(heat)

3) Poisson eqn: $u_{xx} = f$ \leftarrow prototypical elliptic eqn.

1) $u_t + u_x = 0, x > 0$

$\begin{cases} u(0, t) = 1, \\ u(x, 0) = \begin{cases} 1, & x < 1/2 \\ 0, & x \geq 1/2 \end{cases} \end{cases}$ (2)

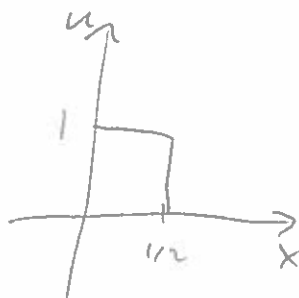
\leftarrow discontinuity hard to deal w/ numerically



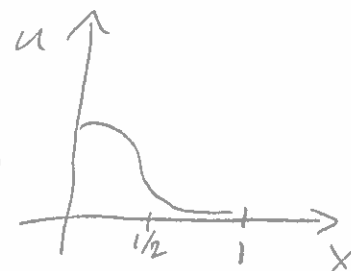
\leftarrow discontinuity exists for all time!!

2) $u_{tt} = u_{xx}, 0 \leq x \leq 1$

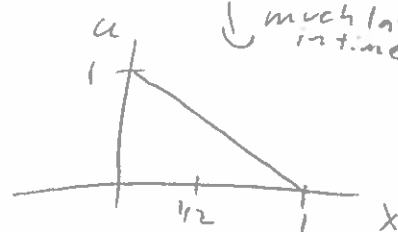
$\begin{cases} u(0, t) = 1, \\ u(1, t) = 0 \\ u(x, 0) = \begin{cases} 1, & x < 1/2 \\ 0, & x \geq 1/2 \end{cases} \end{cases}$



later in time



\downarrow much later in time



super easy to represent numerically

* Need different schemes for numerically rep'ing time & space for the different problems

In 228A, we will focus on the Poisson eqn, which is time independent

Where does the Poisson eqn. come up?

1. Steady-state diffusion problems

E.g. u is a concentration: $u_t = \underbrace{Du_{xx}}_{\text{transport by diffusion}} + \underbrace{f}_{\text{input}}$

At steady-state, $u_t \rightarrow 0 \Rightarrow -Du_{xx} = f$

Reaction-diffusion eqn: $u_t = Du_{xx} - ku + f$

Steady-state: $-Du_{xx} + ku = f$

2. Electrostatics:

E.g. \vec{E} is a steady electric field, ρ is a charge distribution, ϵ_0 is a parameter

$$\nabla \cdot \vec{E} = \rho / \epsilon_0 \quad \& \quad \nabla \times \vec{E} = 0 \rightarrow \text{there is a potential } \phi \text{ s.t.}$$
$$\vec{E} = -\nabla \phi$$

$$\text{Thus } \nabla \cdot (-\nabla \phi) = \rho / \epsilon_0 \Rightarrow -\Delta \phi = \rho / \epsilon_0$$

3. Potential Flow:



$$\nabla \cdot \vec{u} = 0, \quad \vec{u} \text{ is a velocity,}$$

divergence free / incompressible

& we also have $\nabla \times \vec{u} = 0$ no vorticity

$$\Rightarrow \vec{u} = \nabla \phi \Rightarrow \nabla \cdot \vec{u} = \boxed{\Delta \phi = 0}$$

4. Low Reynolds # Flow: very small length scales, e.g. bacteria swimming
drag-dominated, no inertia

$$\mu \Delta \vec{u} - \nabla p = 0$$
$$\nabla \cdot \vec{u} = 0$$

stoke's eqns.

take div. of first eqn

$$-\Delta p = 0$$

Poisson Eqn: $u_{xx} = f$ 1D

$$u_{xx} + u_{yy} = f \quad 2D \text{ (cartesian)}$$

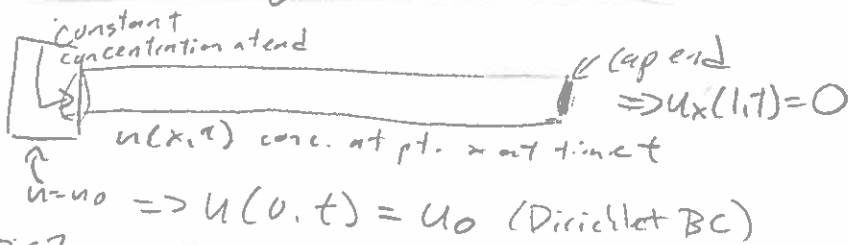
$$\Delta u = f$$

$$\text{or } \nabla^2 u = f$$

I need a domain, I need boundary conditions

1D - Steady state diffusion eqn.

$$u_t = Du_{xx} + f$$



How does capped end translate to a BC?

Given (a, b) , $q(t) = \int_a^b u(x, t) dx$ = total amt. of chemical in interval (a, b) .

Consider $u_t = Du_{xx} \rightarrow \int_a^b u_t dx = \int_a^b Du_{xx} dx$

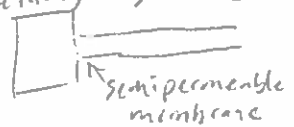
$$\frac{d}{dt} q(t) = D(u_x(b, t) - u_x(a, t))$$

transport rates on boundary of (a, b)

$J = -Du_x$ is the diffusive flux, capped end \Rightarrow no diffusive flux!
(at end) $u_x(0, t) = 0$
Neumann B.C.

What about a given flux? $-Du_x(0, t) = g$ (non-homogeneous Neumann B.C.)

something in between?



$$J = \alpha(u_r - u_l) \Rightarrow \text{flux balance:}$$

right & left of membrane

$$-Du_x(0, t) = \alpha(u(0, t) - u_0)$$

$$\Rightarrow \text{Robin BC} \quad \alpha u(0, t) + Du_x(0, t) = \alpha u_0$$

Solving Poisson eqn need Dirichlet, Neumann, or Robin BC's.



$$u_{xx} = f \quad \text{on } (0,1)$$

$$u(0) = u(1) = 0$$

Fourier series soln: What if $f = -(n\pi)^2 \sin(n\pi x)$, $n=1,2,\dots$
 Soln $u = \sin(n\pi x)$

Reframe higher level: $Lu = f$ w/ L operator $\frac{d^2}{dx^2}$ on space of fns.
 satisfying $u(0) = u(1) = 0$.

$$L(\sin(n\pi x)) = -(n\pi)^2 \sin(n\pi x)$$

$\Rightarrow u_n = \sin(n\pi x)$, $n=1,\dots$ are eigenfunctions of L w/ eigenvalues $\lambda_n = -n^2\pi^2$

These eigenfns are orthogonal in $L^2[0,1]$ \leftarrow square integrable fns.

$$\text{inner prod: } \langle f(x), g(x) \rangle = \int_0^1 f(x)g(x)dx$$

eigenfns form a complete set in $L^2[0,1]$ (w/o any BC's - FS will always converge in L^2 , def not pointwise)

$$\text{For any } f \in L^2[0,1], \quad f(x) = \sum_{n=1}^{\infty} a_n \sin(n\pi x)$$

$$\text{Known soln can be written in form } u(x) = \sum_{n=1}^{\infty} \beta_n \sin(n\pi x)$$

We can compute a_n by orthogonality: $a_n = 2 \int_0^1 f(x) \sin(n\pi x) dx = \frac{\langle f, \sin(n\pi x) \rangle}{\langle \sin(n\pi x), \sin(n\pi x) \rangle}$
 a_n is the projection of f onto the eigenspace spanned by $\sin(n\pi x)$.

$$L\left(\sum_{n=1}^{\infty} \beta_n \sin(n\pi x)\right) = \sum_{n=1}^{\infty} a_n \sin(n\pi x)$$

$$\sum_{n=1}^{\infty} \lambda_n \beta_n \sin(n\pi x) = \sum_{n=1}^{\infty} a_n \sin(n\pi x)$$

$$\text{using orthogonality } \Rightarrow \beta_n = \frac{a_n}{\lambda_n}$$

Same idea on a 2D rectangle: $\Delta u = u_{xx} + u_{yy} = f$ on $(0,1) \times (0,1)$
 $u = 0$ on the boundary

$$\text{Eigenfns: } u_{n,m} = \sin(n\pi x) \sin(m\pi y), \quad n,m=1,2,\dots$$

$$\lambda_{n,m} = -(n^2 + m^2)\pi^2$$

Finite Difference Methods

- Discretize space

Given PDE, domain, & boundary conditions



• Discretize the domain

• Represent functions by their values at the points

• Use discrete values to approximate derivatives using algebraic formulas

\Rightarrow Get an algebraic eqn $A\vec{u} = \vec{b}$

Finite Element Methods

- Discretize function space

Reformulate the problem as a variational problem

$$\Delta u = f; \quad F(u) = \int_{\Omega} \frac{1}{2} \nabla u \cdot \nabla u + u f \, dx$$

Find $u \in S$ to minimize $F(u)$, where S is the space of admissible fns.

Assume the minimization problem gives the PDE.

Now discretize the function space S , i.e. pick a finite-dim space $S_h \subset S$.

so that for $u_h \in S_h$, $u_h = \sum_{k=1}^n u_k \phi_k(x)$

Ex:



ϕ 's are locally supported, small overlap w/ other ϕ 's
these are a basis for piecewise linear fns.

Solve the minimization problem "exactly" on the finite-dim space S_h .

$$F(u_h) = \frac{1}{2} \sum_{i,j} A_{ij} u_i u_j + \sum_i b_i u_i, \quad \text{where } A_{ij} = \int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, dx \in \begin{matrix} \text{mostly} \\ 0's \text{ si.} \\ \phi's \text{ loc.} \\ \text{supp.} \end{matrix}$$
$$\& \quad b_i = \int_{\Omega} f \phi_i \, dx$$

The unknowns are the u_i 's \rightarrow coefficients for expansion of min. soln u_h

If there is a minimum, it must be true that $\frac{\partial F}{\partial u_i} = 0 \quad \forall i \Rightarrow \nabla F = 0$ at min.

computing this, get $A\vec{u} = \vec{b}$.

The structure of this is very similar to that of a Finite Difference method.

Spectral Methods

Use approx: $U_h(x) = \sum_{k=1}^n a_k \phi_k(x)$

generally, the ϕ 's are not locally supported (e.g. $\phi_k = \sin(k\pi x)$)

From here, solve the variational problem

\hookrightarrow changes structure of $A\vec{u} = \vec{b}$, structure of A is dense.

or set of points, solve $(\Delta u)(x_j) = f(x_j)$.

get $\sum_{j=1}^N a_j \phi_j''(x_i) = f(x_i)$ for $i=1, \dots, N$

\Rightarrow linear system $A_{ij} = \phi_j''(x_i)$ still dense $\rightarrow A\vec{x} = \vec{f}$.



9/29 - MAT 228A - Lecture 3

● Finite-Difference Methods: idea - approximate derivatives using function values at points

Approx. derivative w/ a difference:

Forward difference operator: $D_+ u(x) = \frac{u(x+h) - u(x)}{h}$, h fixed

Backward difference operator: $D_- u(x) = \frac{u(x) - u(x-h)}{h}$

How accurately do these approximate $\frac{d}{dx}$?

Error $E = D_+ u(x) - \frac{d}{dx} u(x)$

Expand using Taylor series as $h \rightarrow 0$

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

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

● $E = D_+ u(x) - \frac{d}{dx} u(x) = \frac{h}{2} u''(x) + \frac{h^2}{6} u'''(x) + \dots$

Assume u'' is bounded, then $E = \mathcal{O}(h)$.

Same idea for D_- : $u(x-h) = u(x) - hu'(x) + \frac{h^2}{2} u''(x) - \frac{h^3}{6} u'''(x) + \dots$

$$D_- u(x) = u'(x) - \frac{h}{2} u''(x) + \frac{h^2}{6} u'''(x) + \dots$$

$\xrightarrow{\text{error}}$ $E = \mathcal{O}(h)$. ← error is same but opp. sign as D_+ at leading order!

Consider $D_0 u(x) = \frac{1}{2} (D_+ + D_-) u(x) = \frac{1}{2} \left[\frac{u(x+h) - u(x)}{h} + \frac{u(x) - u(x-h)}{h} \right]$

centered difference operator

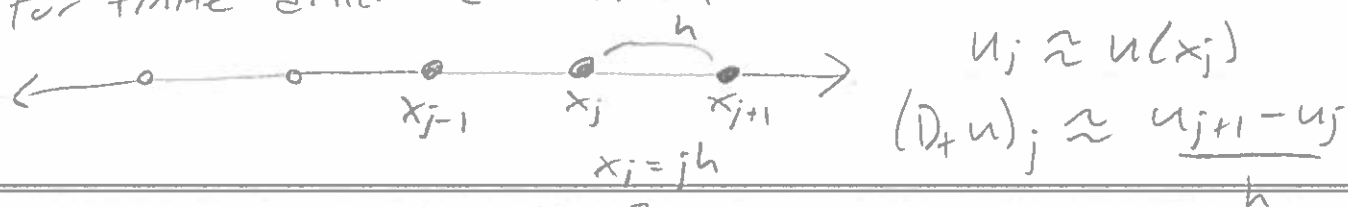
$$= \frac{u(x+h) - u(x-h)}{2h} = u'(x) + \frac{h^2}{6} u'''(x) + \dots$$

$\xrightarrow{\text{error}}$

● $D_0 u(x)$ is $\mathcal{O}(h^2)$ accurate (assuming $u'''(x)$ bounded)

Terminology: D_+ & D_- provide first-order accurate approximations to $\frac{d}{dx}$
 D_0 provides second-order " " " "

For fixed u , $D_+ u(x)$ can be evaluated everywhere.
 For finite-difference methods, start on a discrete domain.



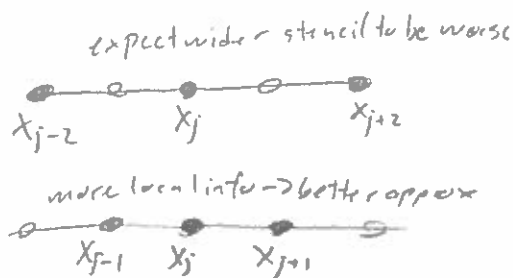
How do we get higher derivatives?

→ apply difference operators multiple times (but which ones?)
 $D_+^2, D_-^2, D_0^2, D_+ D_0, D_- D_0, D_+ D_-,$ etc. \Rightarrow all approx. $\frac{d^2}{dx^2}$

2 of these are second-order accurate: $D_0^2, D_+ D_-$

$$D_0(D_0 u)_j = D_0\left(\frac{u_{j+1} - u_{j-1}}{2h}\right) = \frac{u_{j+2} - 2u_j + u_{j-2}}{4h^2}$$

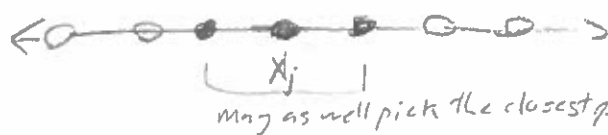
$$D_+(D_- u)_j = D_+\left(\frac{u_j - u_{j-1}}{h}\right) = \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2}$$



Different approach to approx. $u''(x_j)$:

need at least 3 pts. to det. $u''(x)$

(since 2 pts can only give you a line $\rightarrow 0^{th}$ deriv.)



Take a linear combination of these 3 values:

$$(D^2 u)_j = a u_{j-1} + b u_j + c u_{j+1} \quad a, b, c \text{ constants TBD.}$$

$$= a u(x-h) + b u(x) + c u(x+h) \quad \leftarrow \text{assume the } u_j \text{ sample a cont's fn.}$$

$$\begin{aligned} \text{TS}_2 &= a \left[u(x) - h u'(x) + \frac{h^2}{2} u''(x) - \frac{h^3}{6} u'''(x) + \frac{h^4}{24} u^{(4)}(x) + \dots \right] \\ &+ b u(x) + c \left[u(x) + h u'(x) + \frac{h^2}{2} u''(x) + \frac{h^3}{6} u'''(x) + \frac{h^4}{24} u^{(4)}(x) + \dots \right] \end{aligned}$$

\leftarrow note how again the sym. saves the day and gives us more acc. for free!

$$= \underbrace{(a+b+c)}_0 u(x) + \underbrace{h(c-a)}_0 u'(x) + \underbrace{\frac{h^2}{2}(a+c)}_1 u''(x) + \frac{h^3}{6}(c-a) u'''(x) + \frac{h^4}{24}(c+a) u^{(4)}(x)$$

$$\Rightarrow \begin{cases} a+b+c=0 \\ -ha+hc=0 \\ \frac{h^2}{2}a + \frac{h^2}{2}c = 1 \end{cases} \rightarrow \begin{cases} a = 1/h^2 \\ b = -2/h^2 \\ c = 1/h^2 \end{cases}$$

note we just derived the $D_+ D_-$ operator!

$$(D^2 u)_j = \frac{1}{h^2} u_{j-1} - \frac{2}{h^2} u_j + \frac{1}{h^2} u_{j+1}$$

$$\Rightarrow (D^2 u)_j = u''(x) + \frac{h^2}{12} u^{(4)}(x) + \dots$$

hence $D^2 = D_+ D_-$ is second-order accurate approx. of $\frac{d^2}{dx^2}$

For this problem, Absolute error = $D_+ u(z) - u'(z)$

9/29

$$\text{relative error} = \frac{D_+ u(z) - u'(z)}{u'(z)}$$

Forward differences and centered differences of $f(x) = \exp(x)$ at $x = 2$.

h	error of		ratio	error of		ratio
	forward difference			centered difference		
1.00000e+00	5.307e+00	---		1.295e+00	---	
5.00000e-01	2.198e+00	2.41		3.117e-01	4.15	
2.50000e-01	1.006e+00	2.19		7.721e-02	4.04	
1.25000e-01	4.817e-01	2.09		1.926e-02	4.01	
6.25000e-02	2.358e-01	2.04		4.812e-03	4.00	
3.12500e-02	1.167e-01	2.02		1.203e-03	4.00	
1.56250e-02	5.803e-02	2.01		3.007e-04	4.00	
7.81250e-03	2.894e-02	2.01		7.517e-05	4.00	
3.90625e-03	1.445e-02	2.00		1.879e-05	4.00	

halving
h
results in
halving
errors
(asymptotically)

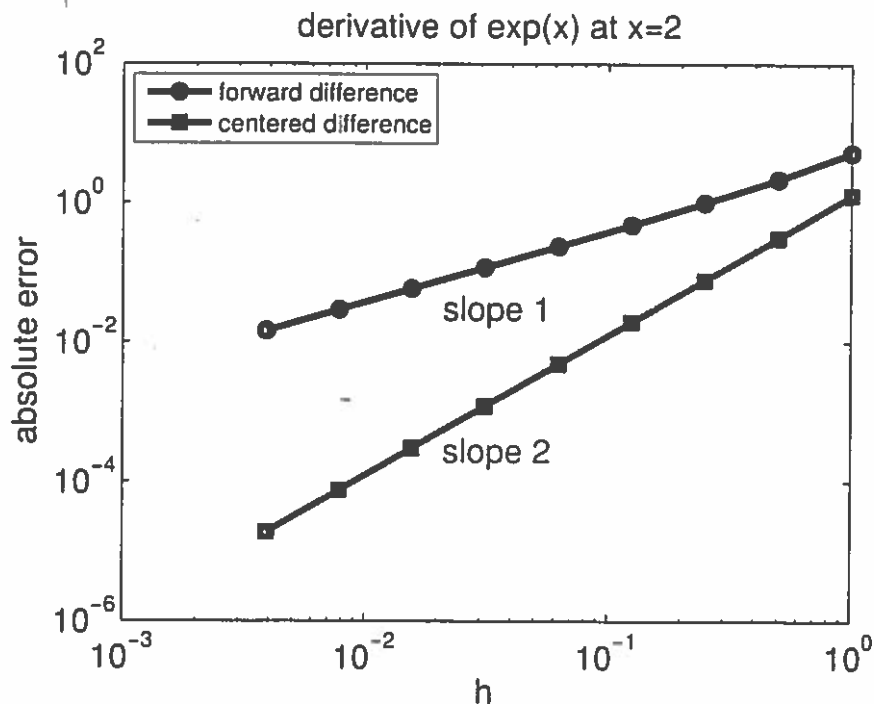
half h
grid spacing
should
quarter
the
error
(asymptotically)

$$\frac{1}{2} \rightarrow \frac{1}{2} \Rightarrow \mathcal{O}(h)$$

$$\frac{1}{2} \rightarrow \frac{1}{4} \Rightarrow \mathcal{O}(h^2)$$

Recall: D_+ is $\mathcal{O}(h)$ accurate

D_0 is $\mathcal{O}(h^2)$ accurate

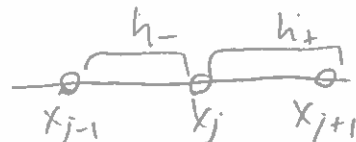


★ often demonstrate the desired order of accuracy by displaying this table or this log-log plot.

$D_+ D_+$ 

this is only first-order accurate
since it lacks the symmetry of our other
approach

Similarly

 $h_+ \neq h_-$ also would only be 1st order acc.

Expect generically that 3 point operators give only 1st order approx.
to the 2nd derivative.

\hookrightarrow can add another pt. to give another eqn/coefficient
to cancel out a leading order error.



How many points m do we need to approx. n^{th} deriv. to p^{th} order accuracy,

$\hookrightarrow w_1 u_1 + w_2 u_2 + \dots + w_m u_m$

Taylor series $\Rightarrow \underbrace{A_0 u(x) + A_1 u'(x) + \dots + A_n u^{(n)}(x)}_{n+1 \text{ constraints}} + A_{n+1} u^{(n+1)}(x) + \dots$
 \Rightarrow at least $n+1$ points to get n^{th} derivative.

\hookrightarrow solving for these constraints, we expect $w \sim \frac{1}{h^n}$

so then $A_{n+1} \sim \text{size } h$ (generically)

To get up to p^{th} order accuracy, need $m = (n+1) + (p-1)$

$$\Rightarrow \underline{m = n + p}$$

MAT228A - Lecture 4 - 10/4

One-D finite diff. soln. of Poisson eqn

$$\begin{cases} u_{xx} = f \\ u(0) = \alpha \\ u(1) = \beta \end{cases}$$

$j=0$ $\xrightarrow{\text{points } 1-N}$ $j=N+1$ $x_0=0$ $x=1$ $x_j = jh$ $x_{N+1} = (N+1)h = 1$
 Discretize interval into equally spaced points $\Rightarrow h = \frac{1}{N+1}$

$u_j \approx u(x_j) \rightarrow \frac{u_{j-1} - 2u_j + u_{j+1}}{h^2} = f_j = f(x_j) \quad (*)$
 Discretize eqns \rightarrow for $j=1, \dots, N$

Now $\underline{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{pmatrix}$ eqns at the form $A\underline{u} = \underline{b}$ where A is D_+ -op & \underline{b} is the f vector.
 j^{th} eqn $\rightarrow \begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_N \end{pmatrix}$

$(a_{j1} \dots a_{jN}) \begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix} = f_j$ should match $(*)$, hence all a_{ji} but 3 should be 0

\Rightarrow For $j=2, \dots, N-1$: $a_{jj} = -\frac{2}{h^2}$, $a_{j,j+1} = a_{j,j-1} = \frac{1}{h^2}$, $b_j = f_j = f(x_j)$

$j=1$, $\frac{u_0 - 2u_1 + u_2}{h^2} = f_1$ but u_0 is known by BC $\Rightarrow \frac{\alpha - 2u_1 + u_2}{h^2} =$

$\Rightarrow \frac{-2u_1 + u_2}{h^2} = f_1 - \frac{\alpha}{h^2} \Rightarrow a_{11} = -\frac{2}{h^2}$, $a_{1,2} = \frac{1}{h^2}$, $b_1 = f_1 - \frac{\alpha}{h^2}$

hence $A = \frac{1}{h^2} \begin{pmatrix} -2 & 1 & & & \\ 1 & -2 & & & \\ & & \ddots & & \\ & & & -2 & 1 \\ & & & 1 & -2 \end{pmatrix}$, $b = \begin{pmatrix} f_1 - \frac{\alpha}{h^2} \\ f_2 \\ \vdots \\ f_N - \frac{\beta}{h^2} \end{pmatrix}$
A is tri-diagonal

What if $\underline{u} = \begin{pmatrix} u_0 \\ u_1 \\ \vdots \\ u_N \\ u_{N+1} \end{pmatrix}$? Then

$\frac{1}{h^2} \begin{pmatrix} h^2 & & & & \\ & -2 & 1 & & \\ & 1 & -2 & & \\ & & & \ddots & \\ & & & & -2 & 1 \\ & & & & 1 & -2 \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ \vdots \\ u_N \\ u_{N+1} \end{pmatrix} = \begin{pmatrix} \alpha \\ f_1 \\ \vdots \\ f_N \\ \beta \end{pmatrix}$

this is weird but OK

Does A have an inverse? Yes.

How close is u_j to $u(x_j)$ \leftarrow soln to PDE?

Error at a point: $e_j^h = u_j^h - u(x_j) \rightarrow$ Error vector $\underline{e}^h = \underline{u}^h - \underline{u}_{\text{sol}}$
 explicit notation of grid spacing \rightarrow We'd like $\|\underline{e}^h\| \rightarrow 0$ as $h \rightarrow 0$ for some norm.

If the error goes to zero as the mesh spacing goes to zero, the method/scheme is converged.

Vector, matrix, function norms: Consider $u(x) = 1$ on $(0, 1)$

Then $\|u\|_2 = \|u\|_1 = \|u\|_\infty = 1$ fn. norms

Sample $u(x)$: $(\vec{u})_j = u(x_j)$, $\vec{u}_j = 1$ sample using N points.

$$\|\vec{u}\|_2 = \left(\sum_{j=1}^N u_j^2 \right)^{1/2} = \left(\sum_{j=1}^N 1 \right)^{1/2} = \sqrt{N} \quad \|\vec{u}\|_\infty = \max_j |u_j| = 1$$

$$\|\vec{u}\|_1 = \sum_{j=1}^N |u_j| = N$$

Vector norms \rightarrow values depend on length of vector.

We should use discretized function norms for errors! \hookrightarrow not appropriate here.

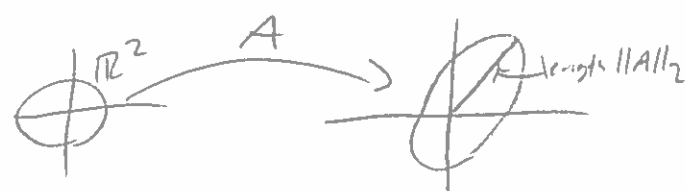
Let \vec{e}^h be a grid function: $\|\vec{e}^h\|_2 = \left(h \cdot \sum_{j=1}^N e_j^2 \right)^{1/2}$

$$\|\vec{e}^h\|_1 = h \cdot \sum_{j=1}^N |e_j|$$

define the norms based on the mesh.

Matrix norms: $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$, A is an $n \times n$ matrix

$$\|A\|_p = \max_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_p = 1}} \|Ax\|_p$$



$$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}| = \text{max row sum}$$

$$\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}| = \text{max column sum}$$

$$\|A\|_2 = \sqrt{\rho(A^*A)}$$

where ρ is the spectral radius = modulus of largest eigenvalue
= largest singular value of A

Proof of $\|A\|_\infty = \text{max row sum}$:

$$\|A\|_\infty = \max_{\|x\|_\infty = 1} \|Ax\|_\infty: \quad \|Ax\|_\infty = \max_i \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_i \sum_{j=1}^n |a_{ij}| \cdot |x_j| \stackrel{\text{by } \|x\|_\infty = 1}{\leq} \max_i \sum_{j=1}^n |a_{ij}|$$

Note the max is attained at some I : $\max_i \sum_{j=1}^n |a_{ij}| = \sum_{j=1}^n |a_{Ij}|$.

So pick $x_j = \text{sign}(a_{Ij})$, then $\|Ax\|_\infty = \max_i \left| \sum_{j=1}^n a_{ij} x_j \right| = \sum_{j=1}^n |a_{Ij}|$.

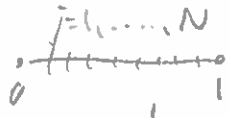
Hence max row sum $\leq \|A\|_\infty$, hence equality. \square

For any induced matrix norm,

$$\frac{\|A\vec{x}\|}{\|\vec{x}\|} \leq \max_{\vec{x} \neq 0} \frac{\|A\vec{x}\|}{\|\vec{x}\|} = \|A\| \Rightarrow \|A\vec{x}\| \leq \|A\| \cdot \|\vec{x}\|$$

228A - Lecture 5 - 10/6

$u_{xx} = f$
 $u(0) = \alpha$
 $u(1) = \beta$

$j=1, \dots, N$

 $h = \frac{1}{N+1}$

$A \vec{u}^h = \vec{b}$ finite difference approx.

Know how big the error is $\vec{e}^h = \vec{u}^h - \vec{u}_{sol}$, $(\vec{u}_{sol})_j = u(x_j)$

Hope that $\|\vec{e}^h\| = \mathcal{O}(h^2)$.

We get error by using our difference operator D_{+-} instead of $\frac{d^2}{dx^2}$ of $\mathcal{O}(h^2)$

↳ Error in equation → how does it translate to error in soln?

↳ Discretization error / local truncation error

In general, \vec{u}_{sol} is not a solution to $A\vec{u} = \vec{b}$

local truncation error $\vec{\tau}^h = A\vec{u}_{sol} - \vec{b}$

for $j=2, \dots, N-1$, $\tau_j^h = \frac{1}{h^2} (u(x_{j-1}) - 2u(x_j) + u(x_{j+1})) - f(x_j)$
 $= \frac{1}{h^2} (u(x_j - h) - 2u(x_j) + u(x_j + h)) - f(x_j)$

$\tau_j^h \rightarrow u_{xx}(x_j) + \frac{h^2}{12} u^{(4)}(x_j) + h.o.t. - f(x_j) = \mathcal{O}(h^2)$

Now $A\vec{u}^h = \vec{b}$

$A\vec{u}_{sol} = \vec{b} + \vec{\tau}^h$

expect error to stay $\mathcal{O}(h^2)$

$A(\vec{u}^h - \vec{u}_{sol}) = -\vec{\tau}^h \Rightarrow A(\vec{e}^h) = -\vec{\tau}^h \Rightarrow \boxed{\vec{e}^h = -A^{-1}\vec{\tau}^h}$

A numerical scheme is consistent if $\vec{\tau}^h \rightarrow 0$ as $h \rightarrow 0$.

A numerical scheme is convergent if $\vec{e}^h \rightarrow 0$ as $h \rightarrow 0$.

For linear schemes applied to linear PDEs:

Lax-equivalence Thm: If a scheme is consistent & stable, it is convergent.

↳ errors don't separate too rapidly in time

For a time-independent BVP, stable $A\vec{u}^h = \vec{b}^h \Rightarrow \|(A^h)^{-1}\| \leq C$ for all $h \leq h_0$.

$\|\vec{e}^h\| = \|A^{-1}\vec{\tau}^h\| \leq \|A^{-1}\| \cdot \|\vec{\tau}^h\| \leq C \cdot \|\vec{\tau}^h\| \rightarrow 0$ as $h \rightarrow 0$

↑ if scheme is stable if scheme is consistent

⊗ Analyze stability via bound on $\|A^{-1}\|$

Stability in 2-norm:

$$A = \frac{1}{h^2} \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & \ddots & \\ & \ddots & \ddots & 1 \\ & & 1 & -2 \end{bmatrix} \text{ is symmetric.}$$

$$A_{\text{sym}} \Rightarrow \|A\|_2 = \rho(A) = \max_j |\lambda_j|$$

$$A_{\text{sym}} \Rightarrow A^{-1} \text{ also sym} \Rightarrow \|A^{-1}\|_2 = \rho(A^{-1}) = \max_j \left| \frac{1}{\lambda_j} \right| = \frac{1}{\min_j |\lambda_j|}$$

(since A^{-1} has inverse eigenvalues of A)

$$\text{eigenvalues } \lambda_k = -k^2 \pi^2$$

Recall: $u^k(x) = \sin(k\pi x)$ is an eigenfn. of $\frac{\partial^2}{\partial x^2}$ on fns. zero at $x=0, 1$.

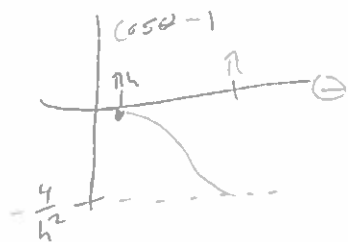
It turns out that the eigenvectors of A are the eigenfns. of $\frac{\partial^2}{\partial x^2}$ evaluated at the grid points. (boundary mismatch doesn't matter since boundary data is 0)

\Rightarrow Show $u_j^k = \sin(k\pi x_j)$ ($k=1, \dots, N$) is an eigenvector, compute the eigenvalues

$$\begin{aligned} & \frac{1}{h^2} (\sin(k\pi x_{j-1}) - 2\sin(k\pi x_j) + \sin(k\pi x_{j+1})) \\ &= \frac{1}{h^2} (\sin(k\pi(x_j-h)) - 2\sin(k\pi x_j) + \sin(k\pi(x_j+h))) \\ &= \frac{1}{h^2} (\sin(k\pi x_j) \cos(k\pi h) - \sin(k\pi h) \cos(k\pi x_j) - 2\sin(k\pi x_j) + \sin(k\pi x_j) \cos(k\pi h) + \sin(k\pi h) \cos(k\pi x_j)) \\ &= \frac{1}{h^2} \sin(k\pi x_j) [2\cos(k\pi h) - 2] = \lambda_k \sin(k\pi x_j) \end{aligned}$$

$$\begin{aligned} \Rightarrow \text{eigenvalues are } \lambda_k &= \frac{2}{h^2} [\cos(k\pi h) - 1], \quad k=1, \dots, N \\ &= -\frac{4}{h^2} \sin^2(k\pi h/2) \end{aligned}$$

For k from 1 to N , $k\pi h$ goes from πh to $N\pi h = \frac{N\pi}{N+1}$



Smallest eigenvalue λ_k occurs at $k=1$

$$\lambda_1 = \frac{2}{h^2} (\cos(\pi h) - 1)$$

$$\Rightarrow \|A^{-1}\| = 1/|\lambda_1|$$

$$\begin{aligned} & \text{also for small } h, \\ & \lambda_k = -k^2 \pi^2 + \mathcal{O}(h^2) \end{aligned}$$

$$\lambda_1 \rightarrow ? \text{ as } h \rightarrow 0: \lambda_1 = \frac{2}{h^2} (\cos(\pi h) - 1) = \frac{2}{h^2} (1 - \frac{1}{2} \pi^2 h^2 + \mathcal{O}(h^4) - 1) = -\pi^2 + \mathcal{O}(h^2)$$

$$\begin{aligned} \text{hence } \|A^{-1}\| &= \frac{1}{\pi^2} + \mathcal{O}(h^2) \Rightarrow \tilde{e}^h = -A^{-1} \tau^h \Rightarrow \|\tilde{e}^h\| \leq \|A^{-1}\|_2 \|\tau^h\| = 1 \\ &= (\frac{1}{\pi^2} + \mathcal{O}(h^2)) \cdot \mathcal{O}(h^2). \end{aligned}$$

Lecture 5 - 10/6 handout

Problem

$$u_{xx} = -\sin(3\pi x)$$

$$u(0) = 1, \quad u(1) = 0$$

Errors

$$e^h = u^h - u(x)$$

$$\text{expect } e^h \approx Ch^2 \Rightarrow e^{h/2} \approx \frac{Ch^2}{4}$$

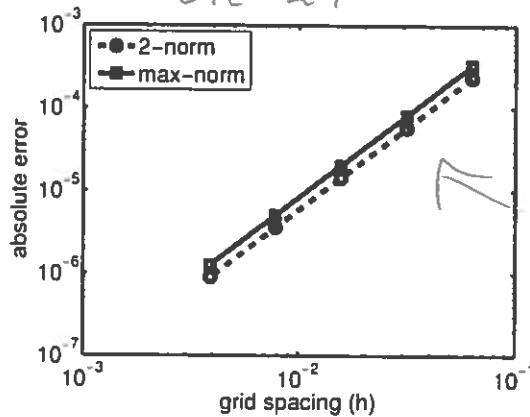
halving
grid
spacing

h	e2	ratio	emax	ratio
6.250000e-02	2.342e-04	---	3.312e-04	---
3.125000e-02	5.779e-05	4.053	8.173e-05	4.053
1.562500e-02	1.440e-05	4.013	2.037e-05	4.013
7.812500e-03	3.598e-06	4.003	5.088e-06	4.003
3.906250e-03	8.992e-07	4.001	1.272e-06	4.001

expect ratio 4

$$e^h/e^{h/2} \approx 4$$

also 2nd-order accurate
in the max norm



slope 2.

need 2 spacings
for 1 ratio.

$$e^h/e^{h/2} \approx 4$$

if we don't
know the soln?

Make your own
by plugging in &
computing f + B.C.

or

estimate the
error by
looking at
successive
approximations

Differences

$$d^h = u^h - u^{h/2}$$

h	d2	ratio	dmax	ratio
6.250000e-02	1.764e-04	---	2.495e-04	---
3.125000e-02	4.339e-05	4.066	6.137e-05	4.066
1.562500e-02	1.080e-05	4.016	1.528e-05	4.016
7.812500e-03	2.698e-06	4.004	3.816e-06	4.004

$$u^h \approx u(x) + Ch^p$$

$$\Rightarrow u^{h/2} \approx u(x) + C\left(\frac{h}{2}\right)^p$$

$$d^h = u^h - u^{h/2} \approx C\left(h^p - \left(\frac{h}{2}\right)^p\right)$$

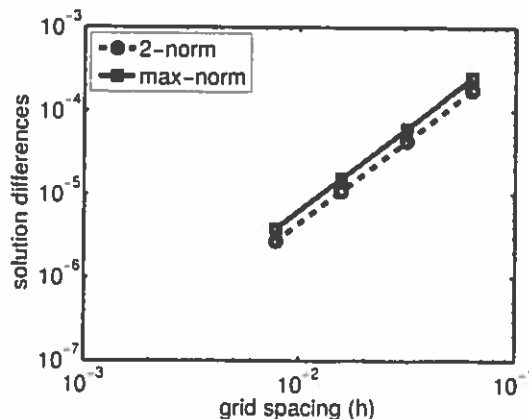
$$= C h^p \left(1 - \frac{1}{2^p}\right)$$

$$\approx C_2 h^p$$

$$d^{h/2} \approx C_2 (h/2)^p$$

need 3 spacings
for 1 ratio

$$d^h/d^{h/2} \approx 4$$





MAT 228A - Lecture 6 - 10/11

$$u_{xx} = f, \quad u(0) = \alpha, \quad u(1) = \beta \rightarrow \text{discretize to } A\vec{u} = \vec{b}$$

where $A = \frac{1}{h^2} \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots \\ & & 1 & -2 \end{bmatrix}$. How do we solve this linear system?

Solve by Gaussian elimination $A = LU = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ \text{stuff} & & 1 \end{pmatrix} \begin{pmatrix} 1 & \text{stuff} \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$

Then $A\vec{u} = \vec{b} \Rightarrow LU\vec{u} = \vec{b} \Rightarrow L\vec{v} = \vec{b}$ & then $U\vec{u} = \vec{v}$ ↑
How expensive is this?

Suppose A is $n \times n$, takes n^2 computations to eliminate one of the n rows

↳ Work is $\mathcal{O}(n^3)$ to factor A

Work to solve a triangular system is $\mathcal{O}(n^2)$

Our matrix has a special structure \rightarrow tridiagonal & diagonally-dominant

$$\begin{pmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots \\ & & 1 & -2 \end{pmatrix} = \begin{pmatrix} 1 & & & \\ -1/2 & 1 & & \\ & \ddots & \ddots & \\ & & 1 & \end{pmatrix} \begin{pmatrix} -2 & 1 & & \\ 0 & -3/2 & 1 & 0 \\ & 1 & -2 & 1 \\ & & \ddots & \ddots \end{pmatrix}$$

$$= \begin{pmatrix} 1 & & & \\ -1/2 & 1 & & \\ & -1/3 & 1 & \\ & & \ddots & \ddots \end{pmatrix} \begin{pmatrix} -2 & 1 & & \\ 0 & -1/2 & 1 & \\ & 0 & -1/3 & 1 \\ & & 1 & -2 \end{pmatrix}$$

← work on each row doesn't depend on n ! But n rows, so work is $\mathcal{O}(n)$
So work to turn A into LU is $\mathcal{O}(n)$!

For this matrix, $\begin{pmatrix} \text{super} \\ \text{diag} \\ \text{sub} \end{pmatrix} = \begin{pmatrix} 1 \\ & \ddots & \\ & & 1 \end{pmatrix} \begin{pmatrix} \text{super} \\ \text{diag} \\ \text{sub} \end{pmatrix} \leftarrow \mathcal{O}(n) \text{ to factor}$
 $\leftarrow \mathcal{O}(n) \text{ to forward \& back-sub}$

This process is the Thomas Algorithm. In Matlab, $\vec{u} = A \backslash \vec{b}$ where A is sparse
We can get the inverse analytically, but $A^{-1}\vec{b}$ takes $\mathcal{O}(n^2)$ work!!!
(but it is dense)

MATLAB: $n=100$; $e = \text{ones}(n,1)$; $A = \text{spdiags}([e \ -2*e \ e], -1:1, n, n)$;
 $Af = \text{full}(A)$; $b = \text{rand}(n,1)$;

tic; $y = A \backslash b$; toc, tic; $z = Af \backslash b$; toc
 ↳ 0.000994 sec ↳ 0.228348 sec

tic; $B = \text{inv}(A)$; toc $\rightarrow .001892$ sec
 tic; $w = B*b$; toc $\rightarrow .001858$ sec

last time, we showed $A\vec{e} = -\vec{e}$

$$\& \|A^{-1}\|_2 = \mathcal{O}(1)$$

had $\mathcal{O}(h^2)$ from
↓ discretization

$$\Rightarrow \|\vec{e}\|_2 = \|A^{-1}\vec{e}\|_2 \leq \|A^{-1}\|_2 \|\vec{e}\|_2 = \mathcal{O}(h^2)$$

~~So an example where $\|\vec{e}\|_\infty = \mathcal{O}(h^2)$~~

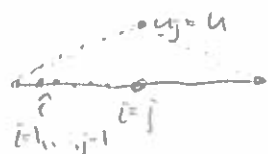
Try norm equivalence: $C\|\vec{e}\|_\infty \leq \|\vec{e}\|_2 \leq C\|\vec{e}\|_\infty$

$$\|\vec{e}\|_2 = h^{1/2} \left(\sum_{j=1}^n e_j^2 \right)^{1/2} \geq h^{1/2} |e_i| \text{ for any } i \Rightarrow \geq h^{1/2} \|\vec{e}\|_\infty$$

$$\Rightarrow h^{1/2} \|\vec{e}\|_\infty \leq \|\vec{e}\|_2 \leq Ch^2 \Rightarrow \|\vec{e}\|_\infty \leq Ch^{3/2} \leftarrow \text{we can do better.}$$

Max-norm analysis: $A\vec{u} = \vec{b}$ where $b_i = \begin{cases} 1 & \text{for } i=j \\ 0 & \text{for } i \neq j \end{cases}$ for some j

Then $A^{-1}\vec{b} = j^{\text{th}}$ column of A^{-1}



(Constructing a sort of Green's fn / inverse operator)
for our difference operator

$$\begin{cases} \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} = 0 \\ u_0 = 0 \\ u_j = u \end{cases}$$

can satisfy this
w/ a linear function
of x_i

$$\text{for } i < j, u_i = \frac{U}{x_j} x_i$$

$$\text{Similar idea for } i > j, u_i = \frac{U(1-x_i)}{1-x_j}$$

$$\text{For } i=j, \frac{u_{j-1} - 2u_j + u_{j+1}}{h^2} = 1$$

$$\frac{U x_{j-1}}{x_j} - 2U + \frac{U(1-x_{j+1})}{1-x_j} = h^2 \leftarrow \text{solve this for } U$$

$$\Rightarrow U = h(x_j - 1)x_j$$

So the soln to $\vec{u} = A^{-1}\vec{b}$ is $u_i = \begin{cases} h(x_j - 1)x_i & i \leq j \\ h(x_i - 1)x_j & i > j \end{cases} \leftarrow \text{looks like a green's fn!!}$

this is the i^{th} element of the j^{th} column of A^{-1} . $= h(x_i - 1)x_j$

$$\Rightarrow \|A^{-1}\|_\infty = \max_i \sum_{j=1}^n |A_{ij}^{-1}| \leq nh \leq 1 \text{ since } n \text{ scales like } 1/h.$$

$$\text{Hence } \|\vec{e}\|_\infty \leq \|A^{-1}\|_\infty \|\vec{e}\|_\infty = \mathcal{O}(h^2).$$

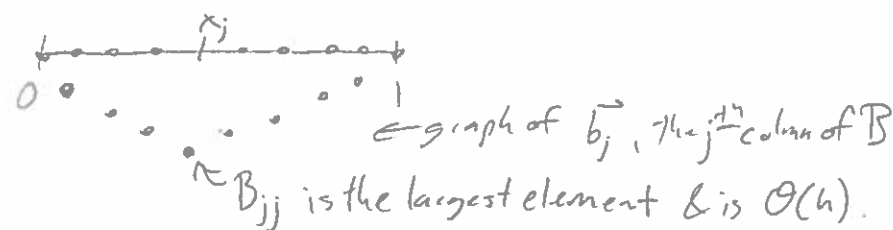
MAT228A - Lecture 7 - 10/13

last time $A = \frac{1}{h^2} \begin{pmatrix} -2 & 1 & & \\ 1 & -2 & & \\ & & \ddots & \\ & & & -2 \end{pmatrix}$ analyzed errors in the max norm by computing the inverse $A^{-1} = B$

where $B_{ij} = h(x_j - 1)x_i \rightarrow$ where $x_j = x_i$ for $j \geq i$
Looks like a Green's fn (kernel of inverse op to $\frac{\partial^2}{\partial x^2}$).

Recall: $A\vec{e} = -\vec{\tau} \Rightarrow \vec{e} = -A^{-1}\vec{\tau} = -B\vec{\tau} = -\begin{pmatrix} \vec{b}_1 & \vec{b}_2 & \dots & \vec{b}_n \end{pmatrix} \begin{pmatrix} \tau_1 \\ \tau_2 \\ \vdots \\ \tau_n \end{pmatrix}$
 $= \sum_{j=1}^n \vec{b}_j \tau_j$ known this is $\mathcal{O}(h^2)$

How much does truncation error at a point affect the error?

\rightarrow Consider an interior point x_j : 

Hence each term in the sum

$\sum_{j=1}^n \vec{b}_j \tau_j$ is $\mathcal{O}(h^3)$. But we add up $\mathcal{O}(\frac{1}{h})$ terms \Rightarrow total error still $\mathcal{O}(h^2)$.

What about truncation error at a point near the boundary?

\vec{b}_1 biggest element is $B_{11} = h(x_1 - 1)x_1 = h(h-1)h = \mathcal{O}(h^2)$.
 $\Rightarrow \vec{b}_1 \tau_1 = \mathcal{O}(h^4)$ Can get away w/ bigger truncation error near boundaries in Dirichlet problem since BC's keep you close to 0.

Neumann Boundaries

- 1) discretize
- 2) solve linear system

$x_0=0 \quad x_1=h \quad x_2=2h \dots$
 $u_x(0) = g$

Look at discretized operator at x_1 : $\frac{u_0 - 2u_1 + u_2}{h^2} = f_1$ don't know u_0 this time!

Why not discretize the BC? $u_x(0) = g \approx \frac{u_1 - u_0}{h} \Rightarrow u_0 = u_1 - gh$

$\Rightarrow \frac{1}{h^2} \begin{pmatrix} -h & h & & \\ 1 & -2 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} \vec{u} = \begin{pmatrix} g \\ f_1 \\ \vdots \\ f_n \end{pmatrix}$ or $\frac{1}{h^2} \begin{pmatrix} -1 & 1 & & \\ 1 & -2 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} = \begin{pmatrix} f_1 + \frac{g}{h} \\ f_2 \\ \vdots \\ f_n \end{pmatrix}$

Discretizing the BC introduces an $\mathcal{O}(h)$ approximation

\rightarrow unobviously, this makes the total error $\mathcal{O}(h)$.

Boundary errors didn't matter in Dirichlet problem b/c of its Green's fn. but this problem has a different Green's function!

This method is only first-order accurate.

Another way to discretize the Neumann BC:

Note that $\frac{u_1 - u_0}{h} = u_x(0) + O(h)$, but $\frac{u_1 - u_{-1}}{2h} = u_x(h/2) + O(h^2)$
(as a centered difference approx.)

Imagine extending the domain by adding ghost points $\rightarrow x_{-1} = -h, x_0, x_1 = h$

~~Finite PDE at $x_0 = 0$: $\frac{u_{-1} - 2u_0 + u_1}{h^2} = f_0$~~

Use BC $u_x(0) = g$ to define the ghost point, centered difference: $\frac{u_1 - u_{-1}}{2h} = g + O(h)$

$\Rightarrow u_{-1} = u_1 - 2hg \Rightarrow f_0 = \frac{u_{-1} - 2u_0 + u_1}{h^2} = \frac{u_1 - 2hg - 2u_0 + u_1}{h^2} \Rightarrow \frac{-2u_0 + 2u_1}{h^2} = f_0 + \frac{2g}{h}$

$\Rightarrow \frac{1}{h^2} \begin{pmatrix} -2 & 2 \\ 1 & -2 & 1 \\ \vdots & \vdots & \vdots \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ \vdots \end{pmatrix} = \begin{pmatrix} f_0 + \frac{2g}{h} \\ f_1 \\ \vdots \end{pmatrix} \leftarrow \text{this method is } O(h^2) \text{ accurate, but our matrix is not symmetric}$

can divide first row by 2 to get $\frac{1}{h^2} \begin{pmatrix} -1 & 1 \\ 1 & -2 & 1 \\ \vdots & \vdots & \vdots \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ \vdots \end{pmatrix} = \begin{pmatrix} \frac{f_0}{2} + \frac{g}{h} \\ f_1 \\ \vdots \end{pmatrix}$

Difference btwn these methods lies in the size of the matrices & misplacement of the boundary

Solvability:

$$\begin{cases} u_{xx} = f, & 0 < x < 1 \\ u_x(0) = \alpha \\ u_x(1) = \beta \end{cases}$$

$$\Rightarrow \int_0^1 u_{xx} dx = \int_0^1 f dx$$

$$u_x(1) - u_x(0) = \int_0^1 f dx$$

$$\text{solvability condition} \rightarrow \boxed{\beta - \alpha = \int_0^1 f(x) dx}$$

Suppose u is a soln to the problem. Then $u + C$ for any constant is also a sol

This doesn't matter physically since this comes from potential problems.

(who cares what actual voltage/pressure is, only care about it relative to norm)

This will cause issues in our numerical method, though:

Discretize the problem: $\frac{1}{h^2} \begin{pmatrix} -1 & 1 \\ 1 & -2 & 1 \\ \vdots & \vdots & \vdots \\ -1 & 1 \end{pmatrix} \begin{pmatrix} u_0 \\ \vdots \\ u_{n-1} \end{pmatrix} = \begin{pmatrix} f_0/2 + \alpha/h \\ f_1 \\ \vdots \\ f_{n-1}/2 - \beta/h \end{pmatrix}$

Matrix is singular!

Supp \vec{u} a soln, then $\vec{u} + C\vec{1}$ is also a soln (all row sums of our matrix are 0)
 $\Rightarrow \vec{1}$ spans the nullspace of A . $A\vec{1} = \vec{0}$.

$A\vec{u} = \vec{b}$ will have a soln only for $\vec{b} \in \text{ran } A \Rightarrow \vec{b} \perp \ker A^* \Rightarrow \vec{b} \perp \vec{1}$.

$\Rightarrow A\vec{u} = \vec{b} - \langle \vec{b}, \vec{1} \rangle \vec{1}$ has a soln.

$$\Rightarrow A\vec{u} = f - \int_0^1 f(x) dx = f - \beta + \alpha$$

MAT228A - Lecture 8 - 10/18

Finish up 1-D stuff: Poisson eqn w/ Neumann BC's:

boundary values unknown

2nd order discretization:

$$\frac{1}{h^2} \begin{pmatrix} -2 & 2 & & & \\ 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & & 1 & -2 \\ & & & & -2 & 2 \end{pmatrix} \vec{u} = \begin{pmatrix} f_0 + 2\alpha/h \\ f_1 \\ \vdots \\ f_n \\ f_{n+1} - 2\beta/h \end{pmatrix}$$

$$\begin{cases} u_{xx} = f \\ u_x(0) = \alpha \\ u_x(1) = \beta \end{cases}$$

soln is not unique &

condition for solvability: $\int_0^1 f(x) dx = \beta - \alpha$

Soln is not unique & matrix is singular.

Note: $A \cdot \vec{1} = \vec{0} \Rightarrow \vec{1} \in \text{null } A$ & w/ more work, $\vec{1}$ spans $\text{null } A$.

How to solve $A\vec{u} = \vec{b}$ w/ singular A :

Need $\vec{b} \in \text{ran } A = (\text{null } A^*)^\perp$. Luckily, $\dim(\text{null } A^*) = 1$, so only need to check orthogonality against 1 spanning vec.

$$A^* = \frac{1}{h^2} \begin{pmatrix} -2 & 1 & & & \\ 2 & -2 & 1 & & \\ & 1 & -2 & \ddots & \\ & & \ddots & \ddots & 1 \\ & & & 1 & -2 \end{pmatrix} \rightarrow \text{nullspace spanned by } \vec{v} = \begin{pmatrix} 1/2 \\ \vdots \\ 1/2 \end{pmatrix} \text{ so check } \vec{b} \cdot \vec{v} = 0 \Rightarrow \vec{b} \in \text{ran } A.$$

$$\vec{v}^T \vec{b} = \frac{1}{2} f_0 + \frac{\alpha}{h} + f_1 + \dots + f_n + \frac{1}{2} f_{n+1} - \frac{\beta}{h} = 0$$

$$\Rightarrow \frac{h}{2} f_0 + h \sum_{j=1}^n f_j + \frac{h}{2} f_{n+1} = \beta - \alpha$$

This is a 2nd order accurate approximation to $\int_0^1 f(x) dx$ using Trapezoidal rule.

$$\Rightarrow \int_0^1 f(x) dx + O(h^2) = \beta - \alpha.$$

Satisfying the conds condition doesn't guarantee discrete condition holds!

One way to solve $A\vec{u} = \vec{b}$ is w/ an iteration scheme

E.g. $\vec{u}^{(k+1)} = T\vec{u}^{(k)} + \vec{c}$ To get a soln, need $\vec{b} \in \text{ran } A$

Project \vec{b} onto $\text{ran } A$ if it's not exactly there (but $O(h^2)$ away!)

$$P\vec{b} = \vec{b} - \frac{\vec{v}^T \vec{b}}{\vec{v}^T \vec{v}} \vec{v} \text{ is guaranteed to be in } \text{ran } A$$

this is $O(h^2)$, so changing discretization but not accuracy.

Direct Solve: perturbed system: $A\vec{u} = \vec{b} - \lambda \vec{v}$ - assume don't know λ .

$$\Rightarrow A\vec{u} + \lambda \vec{v} = \vec{b} \text{ For an appropriate choice of } \lambda, \text{ have soln. } \vec{u}.$$

choose: $\vec{1}^T \vec{u} = 0$ - want soln of above to have mean zero.

Needed an extra eqn. for λ . so just picked one arbitrarily since solns. are unique.

Augmented system:

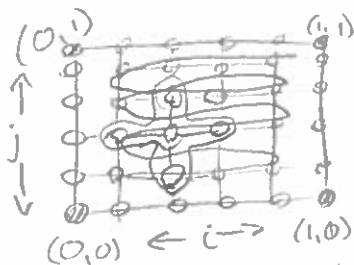
This system has a unique soln.

$$\begin{pmatrix} A & \begin{pmatrix} \bar{I} \\ \bar{V} \\ \bar{I} \\ 0 \end{pmatrix} \\ \begin{pmatrix} \bar{I}^{-1} & 0 \end{pmatrix} \end{pmatrix} \begin{pmatrix} \bar{u} \\ \bar{v} \\ \bar{z} \end{pmatrix} = \begin{pmatrix} \bar{b} \\ \bar{1} \\ \bar{0} \end{pmatrix}$$

2D Eqns

2D Poisson
w/ Dirichlet BC's

$$\begin{cases} \Delta u = f & \text{on } (0,1) \times (0,1) \\ u(0,y) = u(1,y) = u(x,0) = u(x,1) = 0 \end{cases}$$



$$\begin{aligned} x_i &= i h, i=0, \dots, n+1 \\ y_j &= j h, j=0, \dots, n+1 \\ h &= \frac{1}{n+1} \end{aligned}$$

all points including bdy

$$\begin{aligned} i=0, n+1 &\rightarrow \text{all } j \\ j=0, n+1 &\rightarrow \text{all } i \end{aligned}$$

Discrete fn. on lattice: $u_{ij} \approx u(x_i, y_j)$ 2D grid fn.

Discretize $\Delta u = u_{xx} + u_{yy} = f$

$$(u_{xx})_{ij} \approx \frac{u_{i+1,j} - 2u_{ij} + u_{i-1,j}}{h^2}, \quad (u_{yy})_{ij} \approx \frac{u_{i,j+1} - 2u_{ij} + u_{i,j-1}}{h^2}$$

$$\Rightarrow (\Delta u)_{ij} \approx \frac{u_{i-1,j} + u_{i,j-1} - 4u_{ij} + u_{i+1,j} + u_{i,j+1}}{h^2} = f_{ij} \text{ for } i,j=1, \dots, n$$

$\Rightarrow n^2$ linear eqns.

Stencil of the difference op.

$$\hookrightarrow \frac{1}{h^2} \begin{bmatrix} 1 & -4 & 1 \\ 1 & & 1 \end{bmatrix}$$

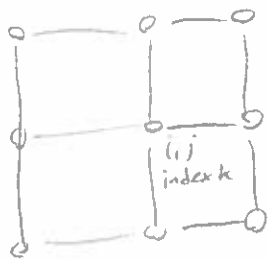
$$\hookrightarrow A \vec{u} = \vec{b}$$

need to pack grid fn. into vector

$u_{ij} \rightarrow$ single indexed vector

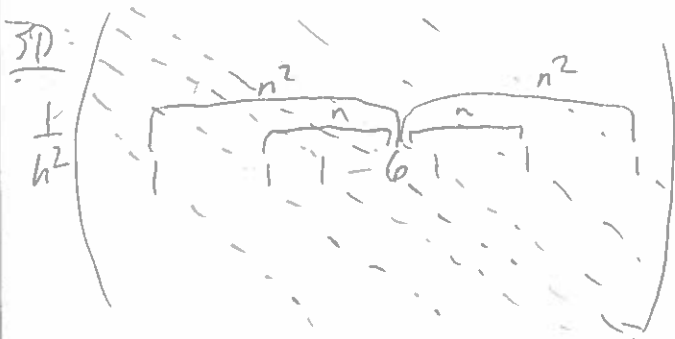
Standard: order by x , then y

$$\vec{u} = \begin{pmatrix} u_{11} \\ u_{21} \\ u_{31} \\ \vdots \\ u_{12} \\ u_{22} \\ \vdots \end{pmatrix}$$



$$\Rightarrow \frac{1}{h^2} \begin{pmatrix} \text{pentadiagonal structure} \\ 1 & -4 & 1 & & \\ & 1 & -4 & 1 & \\ & & 1 & -4 & 1 \\ & & & 1 & -4 & 1 \\ & & & & 1 & -4 & 1 \end{pmatrix} \begin{pmatrix} \vec{u} \end{pmatrix} = \begin{pmatrix} \vec{b} \end{pmatrix}$$

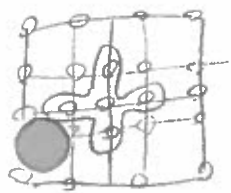
pentadiagonal structure



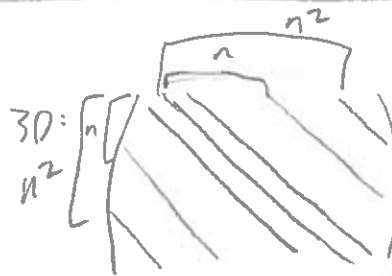
Gaussian elimination on banded matrices doesn't preserve sparsity of matrix.

\hookrightarrow factoring will fill in between bands

MAT228A - Lecture #9 - 10/20



$A\vec{u} = \vec{b}$
 2D: $n \times \begin{pmatrix} \text{banded} \end{pmatrix}$



Sparsity patterns
 not good enough
 for Gaussian elimination
 → fills in lower bands

Work to factor this banded matrix? (LU decomp.)

$b \times \begin{pmatrix} \text{banded} \end{pmatrix} \quad \mathcal{O}(abN)$
 $N \times N$

So in 2D, $N = n^2, a = b = n \Rightarrow \mathcal{O}(n^4) = \mathcal{O}(N^2)$
 better than full matrix $\rightarrow \mathcal{O}(N^3)$
 not as good as possible $\rightarrow \mathcal{O}(N)$.

backward/forward solve the factored system

work = $\mathcal{O}(Nn) = \mathcal{O}(n^3) = \mathcal{O}(N^{3/2})$.

In 3D, work to factor is $\mathcal{O}(N^3 \cdot n^2 \cdot n^2) = \mathcal{O}(n^7) = \mathcal{O}(N^{7/3})$.
 ($N = n^3, a = b = n^2$)

Work to back/forward solve factored system is $\mathcal{O}(n^2 N) = \mathcal{O}(n^5) = \mathcal{O}(N^{5/3})$.

Memory usage in 3D? $\mathcal{O}(N)$ nonzero points in original matrix b/c of banded sparse structure

Filled in factored matrix: $n^3 = N \times \begin{pmatrix} \text{banded} \end{pmatrix}$ storage scales like $\mathcal{O}(n^5) = \mathcal{O}(N^{5/3})$.
 (oldsky factored)
 ↳ symmetric matrices both factors similar

how big for $n = 100$? $100^3 = 10^6 = 1 \text{ million gridpoints} = N$
 $\frac{100^5}{100^3} = 10^2 \Rightarrow n^5 = 10^{10} = 10 \cdot 10^9 = 10 \text{ billion}$

In double-precision, 2 bytes/16 bits per number. Have ~10 gigabytes in RAM, need ~20 gigabytes for this

For systems like this, we know the eigenvalues & eigenvectors, so we can diagonalize

$A\vec{u} = \vec{b}$, Q is the matrix of eigenvectors of A , Λ diag. matrix of eigenvalues

In our problem, eigenvectors are orthogonal $\Rightarrow Q^{-1} = Q^T$

$\Rightarrow AQ = Q\Lambda \Rightarrow A = Q\Lambda Q^T \Rightarrow Q\Lambda Q^T \vec{u} = \vec{b}$

$\Rightarrow \Lambda Q^T \vec{u} = Q^T \vec{b}$

$\Rightarrow Q^T \vec{u} = \Lambda^{-1} Q^T \vec{b}$

$\Rightarrow \vec{u} = Q \Lambda^{-1} Q^T \vec{b}$

Storing Q is expensive as a dense matrix

Dealing w/ Λ is trivial

Fast Fourier Transform: Does multiplication by Q^T in $\mathcal{O}(N \ln N)$ time.

Then inverse FFT for Q for another $\mathcal{O}(N \ln N)$ time.

special solver A .

In 2D, $\begin{pmatrix} \text{banded} \end{pmatrix}$ or $\frac{1}{n^2} \begin{pmatrix} T & I & & \\ I & & & \\ & & I & \\ & & & T \end{pmatrix}$ w/ $T = \begin{pmatrix} -4 & 1 & & \\ 1 & & & \\ & & 1 & \\ & & & -4 \end{pmatrix}$ $\leftarrow n \times n$ identity

Exploit this structure to get $\mathcal{O}(n^{3/2})$.
 w/ "nested dissection"
 can be used for 3D

Convergence in 2D: $A \underline{e}^h = -\underline{\tau}^h$ w/ \underline{e}^h error & $\underline{\tau}^h$ discretization error.

For the 2-norm, same idea as in 1D.

known eigenvalues, eigenfunctions: $u_{ij}^{km} = \sin(k\pi x_i) \sin(m\pi y_j)$
 $\hookrightarrow \underline{\tau}^{km} = \frac{2}{h^2} [\cos(k\pi h) + \cos(m\pi h) - 2]$

& do similar proof as in 1D.

For the ∞ -norm: Show in 2 steps: ① Discrete Maximum Principle:

① $Lu = f$. If $Lu \geq 0$ on some region, then max. value of u is attained on the boundary of this region.
 \uparrow discrete operator

Similarly, if $Lu \leq 0$ on region, min u attained on the boundary.

② If u is a discrete fn. zero on the boundary of the discrete unit square, then $\|u\|_\infty \leq \frac{1}{8} \|Lu\|_\infty$

How do we use this? $Le = -\tau \Rightarrow$ using ②, $\|e\|_\infty \leq \frac{1}{8} \|Le\|_\infty = \frac{1}{8} \|\tau\|_\infty = O(h^2)$.

Hence convergence!

Pf. of 1: Supp. $Lu \geq 0 \Rightarrow \frac{1}{h^2} (u_{i-1,j} + u_{i,j-1} - 4u_{ij} + u_{i,j+1} + u_{i+1,j}) \geq 0$

$$\Rightarrow u_{ij} \leq \frac{1}{4} (u_{i-1,j} + u_{i,j-1} + u_{i,j+1} + u_{i+1,j})$$

Center point of stencil is \leq average of neighbors

If u_{ij} is a maximum, then all neighbors must be equal $\Rightarrow u$ constant.
Hence any max. must be on the boundary.

Pf. of 2: $Lu = f$, u is zero on the boundary

Let $w_{ij} = \frac{1}{4} ((x_i - \frac{1}{2})^2 + (y_j - \frac{1}{2})^2)$. Then $Lw = 1$.

Then $L(u + w\|f\|_\infty) = f + \|f\|_\infty \geq 0$.

So max value of $u + w\|f\|_\infty$ is on the boundary by ① (and $u=0$ on bdy)

$$\Rightarrow u \leq u + w\|f\|_\infty \leq \max(w)\|f\|_\infty = \frac{1}{8} \|f\|_\infty$$

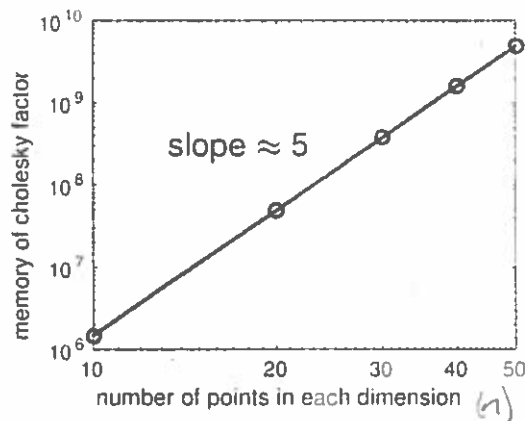
$$\text{Hence } u \leq \frac{1}{8} \|Lu\|_\infty \Rightarrow \|u\|_\infty \leq \frac{1}{8} \|Lu\|_\infty$$

$$\Rightarrow \|e\|_\infty \leq \frac{1}{8} \|Le\|_\infty = \frac{1}{8} \|\tau\|_\infty = O(h^2) \rightarrow 0.$$

□

Lecture 9 Handout - 10/20

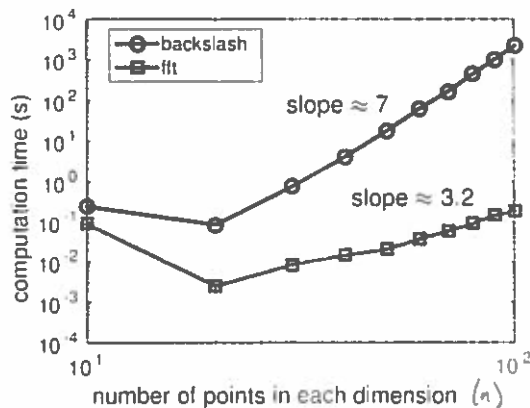
3D w/ gridsize n in each dim. Memory Usage



n	mem (B) ^{bytes}
10	1.48e+06
20	4.90e+07
30	3.77e+08
40	1.60e+09
50	4.90e+09 ~ 5 gb.

Hence memory scales like $O(n^5)$

Solution Time



n	t (GE)	t (FFT)
10	0.2415	0.0893
20	0.0834	0.0025
30	0.7502	0.0085
40	4.0408	0.0145
50	17.3654	0.0202
60	61.7943	0.0365
70	163.9611	0.0569
80	443.7283	0.0891
90	1016.0200	0.1412
100	2185.5417	0.1790

Comp. time grows like $O(n^7)$

FFT is near-optimal, time $O(n^{3.2}) \approx O(N \ln N)$



$$Au = b \quad \text{or} \quad Lu = f \rightarrow u = \underbrace{u_{\text{sol}}}_{\text{soln. to PDE}} + \mathcal{O}(h^2) \rightarrow \text{can approx. } u \text{ algebraically within } \mathcal{O}(h^2) \text{ \& remain within } \mathcal{O}(h^2) \text{ to } u_{\text{sol.}}$$

Suppose u^k is an approximate soln to $Lu = f$, where u is the exact soln. of the alg. eqn. (discrete PDE).

Algebraic error: $e^k = u - u^k$.

$$Le^k = Lu - Lu^k = f - Lu^k = r^k \leftarrow \text{residual (defect)}$$

Correction eqn: $u = u^k + e^k = u^k + L^{-1} r^k$ easy to compute

Let $B \approx L^{-1}$, then example iterative scheme: approximate this operator to get iterative scheme

$$u^{k+1} = u^k + B r^k$$

If $B \approx L^{-1}$ is good approx, hope u^{k+1} is better approx to u than u^k .

A simple method: Let $B = D^{-1}$ where D is the diagonal part of L .

$$\Rightarrow u^{k+1} = u^k + D^{-1} r^k = u^k + D^{-1} (f - Lu^k) = (I - D^{-1}L) u^k + D^{-1} f$$

For our problem, $D = -\frac{4}{h^2} I \Rightarrow D^{-1} = -\frac{h^2}{4} I$ (our problem is Dirichlet Laplacian on rectangle)

$$\Rightarrow u^{k+1} = (I + \frac{h^2}{4} L) u^k - \frac{h^2}{4} f$$

At a point: $u_{ij}^{k+1} = u_{ij}^k + \frac{h^2}{4} (u_{i-1,j}^k + u_{i,j-1}^k - 4u_{ij}^k + u_{i+1,j}^k + u_{i,j+1}^k) - \frac{h^2}{4} f_{ij}$

$$u_{ij}^{k+1} = \frac{1}{4} (u_{i-1,j}^k + u_{i,j-1}^k + u_{i+1,j}^k + u_{i,j+1}^k) - \frac{h^2}{4} f_{ij}$$

This is the Jacobi iterative scheme, and for our problem it converges (very slowly).

As we do this, we go through the grid & compute better values u_{ij} & then update at the end all at once - why not take advantage of better values "behind you" as we go?



loop some # of times

loop i

loop j

$$u_{ij} = \frac{1}{4} (u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1}) - \frac{h^2}{4} f_{ij}$$

can parallelize

Gauss-Seidel (lexicographic order)

Loop same max #

Loop i

Loop j

$$u_{ij} = \frac{1}{4} (u_{i,j-1} + u_{i-1,j} + u_{i+1,j} + u_{i,j+1}) - \frac{h^2}{4} f_{ij}$$

converges faster than Jacobi, but still slow

but not parallelizable

Think of both methods as splitting methods:

$$Au = b \rightarrow A = M - N \text{ where } M \text{ is easy to invert}$$

$$\Rightarrow Mu - Nu = b \Rightarrow u = M^{-1}Nu + M^{-1}b$$

$$\text{Iterative scheme: } u^{k+1} = M^{-1}Nu^k + M^{-1}b$$

For our problem:

$$A = D - L - U \text{ where } D \text{ is diagonal, } -L \text{ is lower triangular part, } -U \text{ is upper triangular part}$$

$$\text{Jacobi: } M = D, N = L + U \Rightarrow u^{k+1} = D^{-1}(L + U)u^k + D^{-1}b$$

$$\text{Gauss-Seidel: } M = D - L, N = U \Rightarrow u^{k+1} = (D - L)^{-1}Uu^k + (D - L)^{-1}b$$

Both methods are of the form: $u^{k+1} = Tu^k + c$ a fixed pt. iteration

$$\text{fixed point soln: } u = Tu + c$$

When does this fixed pt. iteration converge?

$$e^k = u - u^k, \quad -) u^{k+1} = Tu^k + c$$

$$\Rightarrow e^k = T^k e^0$$

$$e^{k+1} = u - u^{k+1} = T(u - u^k)$$

when does $e^k \rightarrow 0$?
only for $\rho(T) < 1$. (spectral radius i.e. max magnitude of eigenvalues)

Suppose we can diagonalize T : $TX = X\Lambda$ Λ diagonal matrix
 $T = X\Lambda X^{-1}$

$$\Rightarrow T^2 = (X\Lambda X^{-1})(X\Lambda X^{-1}) = X\Lambda^2 X^{-1} \rightarrow T^k = X\Lambda^k X^{-1}$$

$$\Lambda^k = \begin{pmatrix} \lambda_1^k & & \\ & \lambda_2^k & \\ & & \lambda_n^k \end{pmatrix} \rightarrow 0 \text{ if } |\lambda_i| < 1 \forall i$$

$$\Rightarrow T^k e^0 \rightarrow 0 \text{ for any } e^0 \text{ if } |\lambda_i| < 1 \forall i.$$

11/1

$\underline{u}^{k+1} = T\underline{u}^k + \underline{c}$ When do we stop the iteration?

1) Based on the size of the residual $r^k = f - A u^k$

How do these control error?

relative tolerance: $\|r^k\| \leq (tol) \cdot \|f\|$

Bound error based on residual bound: $Ae^k = r^k$

Calculating residuals takes extra work, but already have successive iterates!

Use $u^{k+1} - u^k$: absolute tolerance $\|u^{k+1} - u^k\| < tol$
relative tolerance $\|u^{k+1} - u^k\| < (tol) \|u^k\|$

$$\Rightarrow u^{k+1} - u^k = Br^k.$$

Note $\|e^k\| = \|A^{-1}r^k\| = \|A^{-1}B^{-1}(u^{k+1} - u^k)\| \leq \|A^{-1}B^{-1}\| \|u^{k+1} - u^k\|$.

Jacobi does not depend on the ordering of the unknowns

GS does depend on the ordering of the unknowns

Another way to order:

update
call red.

then
update
all black

Take colored points
 $i+j = \begin{cases} \text{even} - \text{red} \\ \text{odd} - \text{black} \end{cases}$

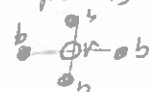
Order reds,
then order blacks
& then apply GS-RB
update all red, turn all,

update one point
after another

$$u_{i,j} = \frac{1}{4} (\underbrace{u_{i-1,j} + u_{i,j-1}}_{\text{neu.}} + \underbrace{u_{i+1,j} + u_{i,j+1}}_{\text{alt.}} - h^2 f_{i,j})$$

GS-JRB :

loop k
 loop red
 $u_{ij} = \frac{1}{4} (u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - h^2 f_{ij})$ ← note these are all black points
 loop black
 $u_{ij} = "$

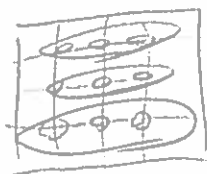


GS-JRB makes sense for our Poisson solver since a red point is updated via all black points

↳ Easier to //ize than GSlex.

& vice-versa.

Other variations → block or line relaxation methods



← could update each row at same time
 (like 1D Poisson eqn.)
 ↳ tridiagonal solve on each row of pts.

Not especially better
 except very useful
 for problems such as

weak coupling in y direction, so just do 1D solve → $u_{xx} + \epsilon u_{yy} = f$
 on the x direction & loop over y-rows. "anisotropy" / could also do blocks

SOR (successive over-relaxation)

$$GS: u_{ij} = \frac{1}{4} (u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - h^2 f_{ij})$$

$$SOR: u_{ij} = \frac{\omega}{4} (u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - h^2 f_{ij}) + (1-\omega) u_{ij}$$

Linear combination of GS update & current iterate.

$\omega < 1$ under-relaxation, $\omega > 1$ over-relaxation

*SOR requires that $0 < \omega < 2$ for convergence. ★

Recall: All these methods are "matrix-splitting" $A = M - N$

$$M u^{k+1} = N u^k + f$$

$$u^{k+1} = M^{-1} N u^k + M^{-1} f$$

$$SOR: M = \frac{1}{\omega} D - L \quad (\text{where } A = D - L - U)$$

$$N = \frac{1-\omega}{\omega} D + U$$

$$T_{SOR} = M^{-1} N = \left(\frac{1}{\omega} D - L \right)^{-1} \left(\frac{1-\omega}{\omega} D + U \right)$$

$$= \omega (D - \omega L)^{-1} \left(\frac{1-\omega}{\omega} D + U \right) = (D - \omega L)^{-1} ((1-\omega)D + \omega U)$$

$$\det(T_{SOR}) = \underbrace{\det(D - \omega L)^{-1}}_{\text{triangular}} \cdot \underbrace{\det((1-\omega)D + \omega U)}_{\text{triangular}} = \frac{\det((1-\omega)D)}{\det(D)} = (1-\omega)^N$$

of grid
 points

Require $|\det(T_{SOR})| < 1 \Rightarrow |1-\omega| < 1 \Rightarrow 0 < \omega < 2$.

□

Convergence analysis for SOR:

Lecture 11 cont# (14/1)

(use same trick as in GS convergence analysis
to compute eigenvalues of the update matrix in terms of Jacobi update)

Let μ be an eigenvalue of Jacobi update T_J .

Then $\mu = \frac{\lambda + \omega - 1}{\omega \lambda^{1/2}}$ for λ eigenvalue of T_{SOR}

Rearranged: $\lambda - \omega \mu \lambda^{1/2} + (\omega - 1) = 0$

$$\Rightarrow 2\lambda^{1/2} = \omega \mu \pm \sqrt{\omega^2 \mu^2 - 4(\omega - 1)}$$

As $\omega \rightarrow 0$, real eigenvalues $|\lambda^{1/2}| \rightarrow 1$

As $\omega \rightarrow 2$, complex eigenvalues
if λ complex, know $|\lambda^{1/2}| = |\omega - 1|$, hence $|\lambda^{1/2}| \rightarrow 1$.

Want to minimize $|\lambda^{1/2}|$ so decrease ω from 2 \rightarrow better convergence.

If λ real, to minimize $|\lambda^{1/2}|$, $\frac{\partial \lambda^{1/2}}{\partial \omega} < 0$ so increase ω from 0
for better convergence.

\Rightarrow Pick on boundary of real & complex λ :

Optimal ω^* satisfies $\omega^{*2} \mu^2 = 4(\omega^* - 1)$

Work $\Rightarrow \omega^* = \frac{2}{1 + (1 - \rho_J^2)^{1/2}}$, where $\rho_J = \rho(T_J)$
spectral radius of Jacobi matrix

$\Rightarrow \rho_{SOR} = \omega^* - 1$.

Recall: $\rho_J = \cos(\pi h) \Rightarrow \omega^* = \frac{2}{1 + (1 - \cos^2 \pi h)^{1/2}} = \frac{2}{1 + \sin \pi h}$

As $h \rightarrow 0$, $\omega^* = \frac{2}{1 + \sin \pi h} \approx \frac{2}{1 + \pi h} \approx 2(1 - \pi h) + \mathcal{O}(h^2)$

\Rightarrow As $h \rightarrow 0$, $\rho_{SOR} \approx 1 - 2\pi h + \mathcal{O}(h^2)$.

iterations to reduce error by 10^{-1} :

SOR has better
iteration count
and scaling !!

$N \times N$	GS	SOR
32 x 32	254	12
64 x 64	985	24
128 x 128	3882	47
256 x 256	14404	94



MAT228A - Lecture # 12 11/3

last time SOR: $w^* = \frac{2}{1+(1-\rho_J^2)^{1/2}}$ optimal $\rho_{SOR} = w^* - 1$

Asymptotically as $h \rightarrow 0$, $w^* = \frac{2}{1+\sin \pi h} = 2(1-\pi h) + O(h^2) \Rightarrow \rho_{SOR} = 1 - 2\pi h + O(h^2)$

compared to GS on 256^2 mesh, # iterations per digit of accuracy: GS $\rightarrow 15409$
SOR $\rightarrow 94$

Scaling on the work of SOR to converge to a $\text{tol} = Ch^2$.

$$\rho^k = Ch^2 \Rightarrow k = \frac{\ln Ch^2}{\ln \rho} \quad \ln \rho \approx \ln(1 - 2\pi h) \approx -2\pi h$$

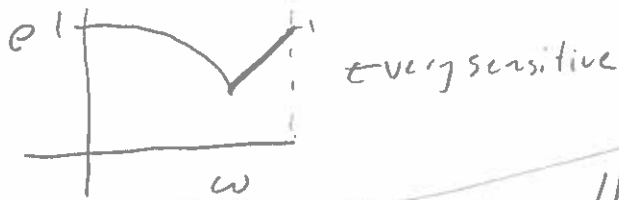
$$\approx \frac{\ln C + \ln h^2}{-2\pi h} \Rightarrow k = O(h^{-1} \ln h)$$

So in 2D, $h \sim \frac{1}{n}$, $n = \#$ of pts. in each direction, $N = n^2 = \text{total } \#$ of points

$$\Rightarrow k = O(\sqrt{N} \ln N) \text{ iteration count}$$

work per iteration is $O(N) \Rightarrow \text{work to solve is } O(N^{3/2} \ln N)$

Drawback: Need to know w^* . Generically, expect $w^* = \frac{2}{1+Ch}$



applies for large k , what happens for small k ?

All methods so far, $\frac{\|e^{k+1}\|}{\|e^k\|} \approx \rho \rightarrow 1$ as $h \rightarrow 0$

Can we find a method with $\rho < C < 1$ as $h \rightarrow 0$?

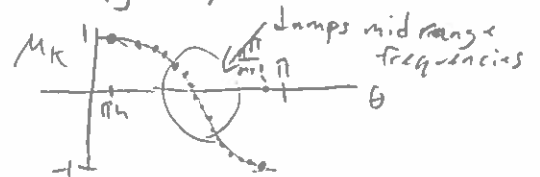
Then we can get iteration count to fixed tolerance indep. of mesh!

Analyze Jacobi in 1D: $T_J = I + \frac{h^2}{2} A$

eigvals of A are $\lambda_k = \frac{2}{h^2}(\cos(k\pi h) - 1)$, eigvalues of T_J are $\mu_k = 1 + \frac{h^2}{2} \lambda_k = \cos k\pi h$

$k=1, \dots, n$, $h = \frac{1}{n+1} \Rightarrow k\pi h = \pi h, \dots, \frac{n\pi}{n+1}$

high k & low k have same order magnitude λ_k .
high & low spatial frequencies are damped the least



Jacobi iterations result in low-freq. average w/ high freq. oscillations on top.

(can we try to improve convergence, smoothing ... for Jacobi)

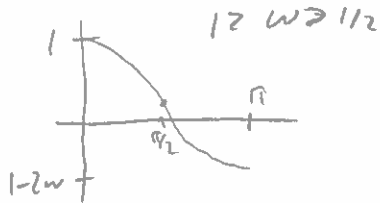
→ ω -Jacobi method: $u_j^{k+1} = \frac{\omega}{2}(u_{j-1}^k + u_{j+1}^k - h^2 f_j) + (1-\omega)u_j^k$

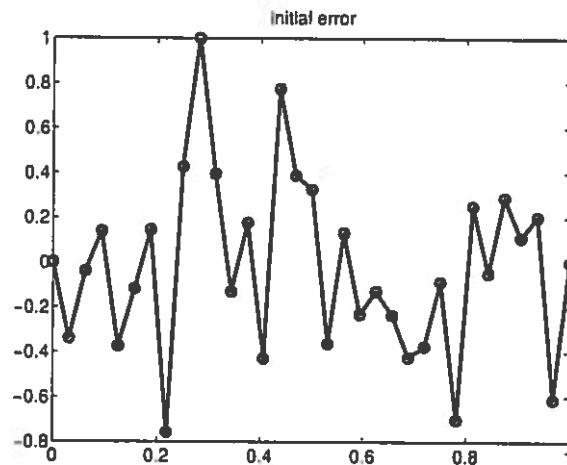
update matrix $T_{\omega J} = \omega T_J + (1-\omega)I = I + \omega \frac{h^2}{2} A = \omega(I + \frac{h^2}{2} A) + (1-\omega)I$

eigenvalues $\mu_k^\omega = \omega \cos(k\pi h) + (1-\omega)$

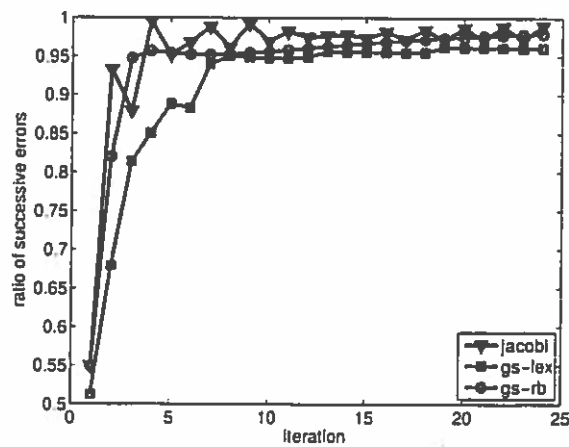
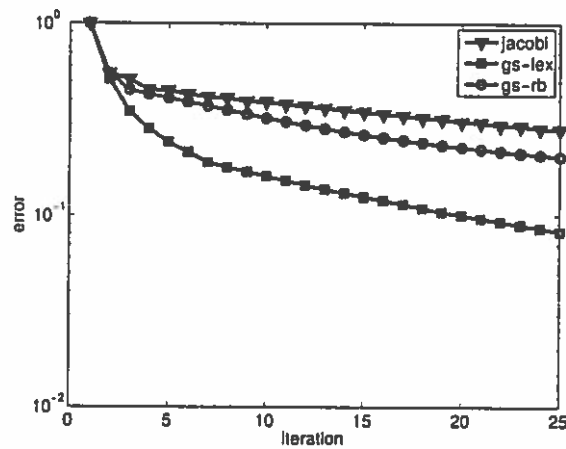
for k small, $\mu_k^\omega = \omega(1 - \frac{k^2 \pi^2 h^2}{2} + \dots) + (1-\omega) = 1 - \omega \frac{k^2 \pi^2 h^2}{2} \leq$ not much less than 1

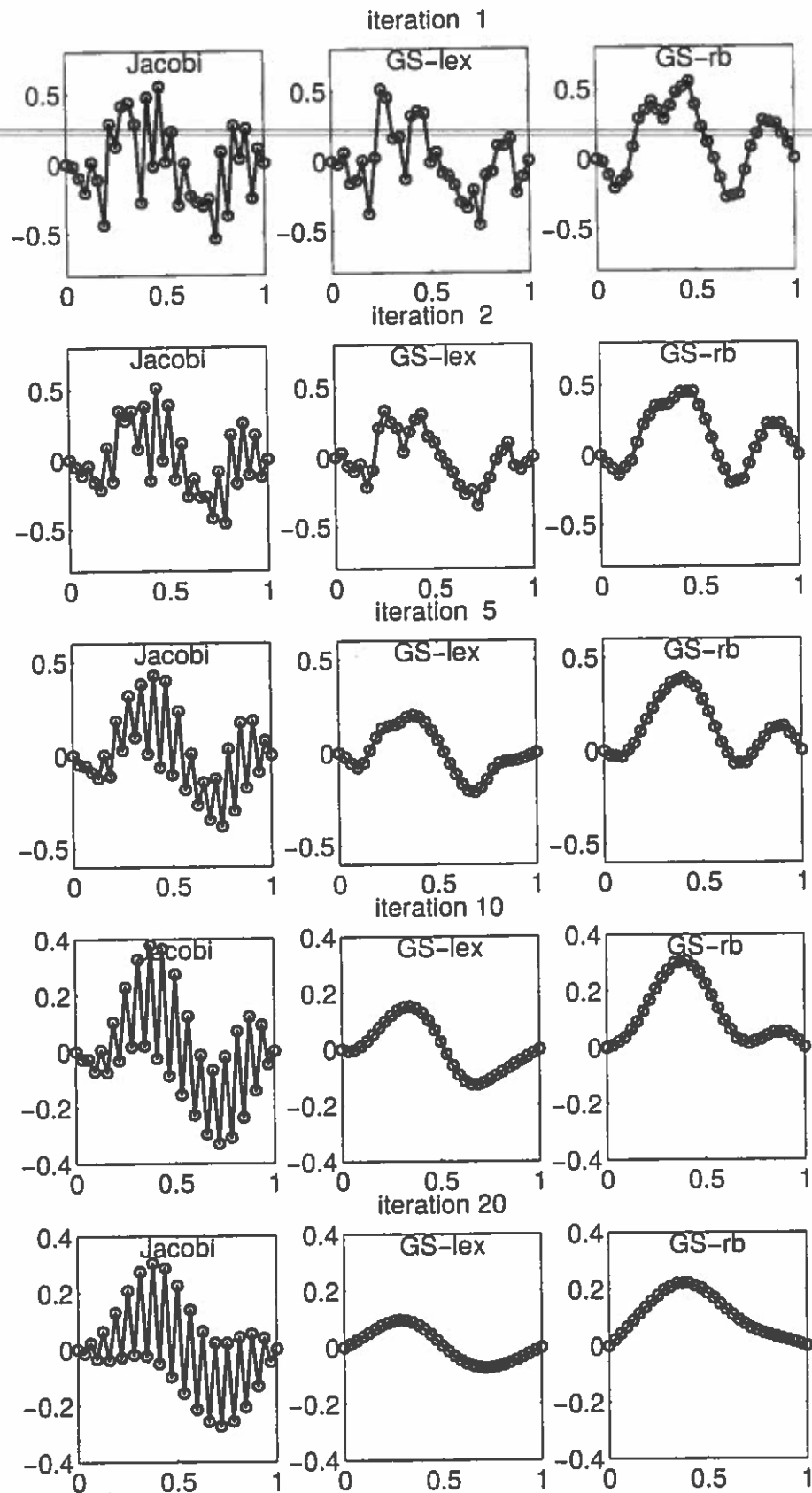
for k big, $\mu_k^\omega = 1 - 2\omega + \dots \leq$ much less than 1, will damp high frequencies





← "equal contribution"
of all frequencies





low frequencies
dominate \rightarrow
dodge
dynamics

high frequencies
make grid scale
oscillations

MAT228A - Lecture 13 - 11/8

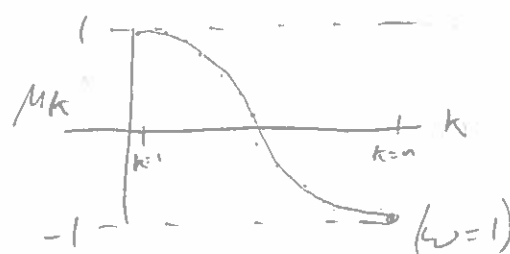
Multigrid: Couple simple iterative methods w/ more complicated

ω -Jacobi in 1D: $u_j^{k+1} = \frac{\omega}{2} (u_{j-1}^k + u_{j+1}^k - \text{RHS}) + (1-\omega)u_j^k$

update matrix: $\omega T_J + (1-\omega)I$

eigenvectors: $u_{j,k} = \sin(k\pi x_j)$

lig values $M_k = \omega \cos(k\pi h) + (1-\omega)$



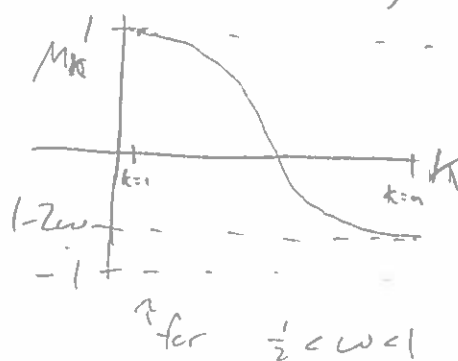
When $\omega=1$, superposition of low & high frequencies, rapid damping of mid freqs.

for $\omega \neq 1$, shift & scale this graph:

$$M_1 = \omega(1 - \frac{\pi^2 h^2}{2}) + (1-\omega)$$

$$= 1 - \omega \frac{\pi^2 h^2}{2} \leftarrow \text{looks like should make } \omega > 1 \text{ to get better convergence}$$

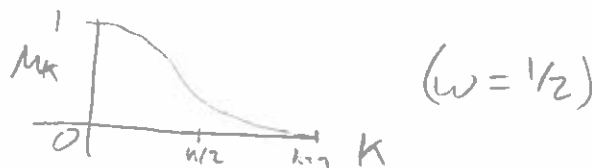
$$M_n = 1 - 2\omega + \dots \leftarrow \text{should make } 0 < \omega \leq 1 \text{ or NO convergence}$$



\hookrightarrow So only low freq. are damped slowly, higher freqs. damp rapidly but convergence for the low freqs. is now slowed for $\omega \neq 1$.

Using GS/SOR yields smooth errors \rightarrow can represent on a coarse grid!
Smooth using GS/SOR iterative methods, then use coarse-correction!

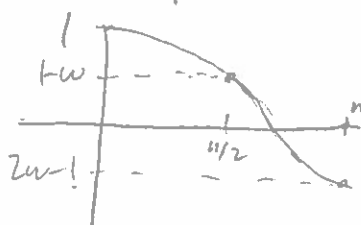
Pick ω to eliminate highest frequencies:
 $\omega = 1/2$ most effective on highest freqs.



Want to reduce frequencies w/ wavenumbers $k \geq n/2$

(since coarsening mesh by throwing out every other point makes representing freqs of $k \geq n/2$ impossible!)

reduces these by a factor of 3 every iteration!



min in abs. value = min. of both

optimizing problem in 1D gets worse in 2D better as ω is

optimal ω^* : $1 - \omega^* = 2\omega^* - 1 \Rightarrow \omega^* = 2/3$

More formally, $\mu = \max_{\frac{n}{2} \leq k \leq n} |\mu_k|$ is the smoothing factor, i.e. the largest factor by which high freqs ($k \geq n/2$) are reduced w/ one application of the smoother w-Jac. $\mu/\omega = 2/3$

Smoothing factors

w-Jacobi:	ω^*	μ	GS lex μ	GS RB μ	ω -GS RB
1D	2/3	1/3	0.45	0.125	
2D	4/5	3/5	0.5	0.25	
3D	6/7	5/7	0.567	0.445	0.23

winner winner

GS has nonlinear eigenfns, not important how these are eliminated
 How does GS eliminate high freq. sine waves?

GS lex analysis of smoothing: Local Fourier analysis:

ignore boundaries & analyze the problem on the infinite domain ($x_j = jh, j \in \mathbb{Z}$)

Now the eigenfns are $u_k^h = e^{ikx} = e^{ikh} \leftarrow$ Fourier modes

What values does k take?

Highest wavenumber k corresponds to $\Rightarrow k(2h) = 2\pi \Rightarrow k = \pi/h$



\Rightarrow All k must be in $-\pi/h \leq k \leq \pi/h \Rightarrow -\pi \leq kh \leq \pi$
 $-\pi \leq \theta \leq \pi$

θ is "continuous" wave number

eigenfns: $u_k(\theta) = \exp(i\theta)$ high freq. $|\theta| \geq \pi/2$

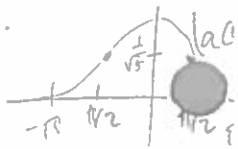


GS-lex:

$e^{kh} = T e^k \xrightarrow{\text{on infinite lattice}} e_k^{kh} = \frac{1}{2}(e_{k-1}^{kh} + e_{k+1}^{kh})$
 error at $k+1$ iterate $\Rightarrow e_k^{kh} = T e_k^k = a e_k^k$ where a is "amplification factor"

Assume $e_k^k = \exp(i\theta)$
 error in one eigenspace, fix θ

$\Rightarrow a e_k^k = \frac{1}{2}(a e_{k-1}^k + e_{k+1}^k)$ $\left\{ \begin{array}{l} e_{k-1}^k = e_k \exp(-i\theta) \\ e_{k+1}^k = e_k \exp(i\theta) \end{array} \right.$ Use sub.
 $\Rightarrow a = \frac{1}{2}(a \exp(-i\theta) + \exp(i\theta))$



$\Rightarrow a = \frac{\exp(i\theta)}{2 - \exp(i\theta)}$ smoothing factor $\Rightarrow \mu = \max_{|\theta| \geq \pi/2} |a(\theta)| = \frac{1}{\sqrt{5}} \approx 0.45$

MAT228A - Lecture 14 - 11/10

In 3D. $GE \sim N^{2/3}$ scaling, $FFT \sim N^{2/3}$ scaling, $SOR \sim N^{1/3}$ MATLAB

● Vectorize code for speedup!

$\sim N^{4.4/3}$ Fortran
 $\sim N^{1.5}$

Recall: Big Idea of Multigrid is to use smoothers (Jacobi, GS) to get rid of high spatial freq. error & then use a coarser mesh to eliminate low freq. errors

Coarse Grid Correction

Let u_h be the algebraic soln to $L_h u_h = f_h$

u_h^k is the approximate soln after k -iterations (i.e. through Jacobi or GS)

$e_h^k = u_h - u_h^k$ algebraic error $r_h^k = f_h - L_h u_h^k$ residual

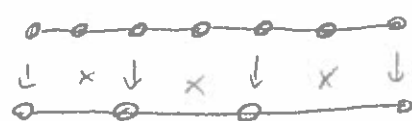
$$\hookrightarrow u_h = u_h^k + e_h^k = u_h^k + \underbrace{L_h^{-1} r_h^k}_{\text{approximate } L_h^{-1} \text{ to gen. iterative scheme}}$$

● a coarse mesh to "solve" $L_h e_h^k = r_h^k$ for e_h^k (since we wiped out high freq. error!)

Notation:

Let Ω_h = mesh w/ spacing h (original fine grid)

then Ω_{2h} is a coarser grid \rightarrow half as many pts. in 1D



$G(\Omega_h)$ set of grid functions on Ω_h (i.e., f, u, L, \dots)

Transfer operators

Restriction operator: $I_h^{2h} : G(\Omega_{2h}) \rightarrow G(\Omega_h)$

Interpolation/prolongation: $I_{2h}^h : G(\Omega_h) \rightarrow G(\Omega_{2h})$

Have u_h^k , compute fine grid residual $r_h^k = f_h - L_h u_h^k$
restrict residual $r_{2h}^k = I_h^{2h} r_h^k$

Solve for error $e_{2h}^k = L_{2h}^{-1} r_{2h}^k$ & does it matter how we solve since may be cheaper

Interpolate to fine grid: $\tilde{e}_h^k = I_{2h}^h e_{2h}^k$

Correct the approx soln:

$$u_h^{k+1} = u_h^k + \tilde{e}_h^k$$

Iterative Scheme:

$$u_h^{k+1} = u_h^k + \tilde{e}_h^k$$

$$= u_h^k + I_{2h}^h e_{2h}^k$$

$$= u_h^k + I_{2h}^h L_{2h}^{-1} r_{2h}^k$$

$$= u_h^k + I_{2h}^h L_{2h}^{-1} I_h^{2h} r_h^k$$

$$= u_h^k + I_{2h}^h L_{2h}^{-1} I_h^{2h} (f_h - L_h u_h^k)$$

$$= \left(I - \underbrace{I_{2h}^h L_{2h}^{-1} I_h^{2h} L_h}_{\text{hope this is close to } L_h^{-1}} \right) u_h^k + I_{2h}^h L_{2h}^{-1} I_h^{2h} f_h$$

hope this is close to $L_h^{-1} \Rightarrow$ product close to I

Still a f.p. scheme!

$$u_h^{k+1} = K u_h^k + c$$

will converge quickly when $\rho(K) \approx 0$.

Perform smoothing to converge!

Let S = smoothing operator

Two-Grid Iteration:

1. Pre-smoothing \Rightarrow apply S ν_1 times

2. Apply Coarse Grid Correction \rightarrow compute residual
 \rightarrow restrict residual
 \rightarrow solve for coarse grid error \star add more grids here for multigrid
 \rightarrow interpolate error
 \rightarrow correct (add error back in)

3. Post-smooth \Rightarrow apply S ν_2 times

\hookrightarrow multigrid iteration Matrix $M = S^{\nu_2} (I - I_{2h}^h L_{2h}^{-1} I_h^{2h} L_h) S^{\nu_1}$

Questions: 1. How to pick ν_1 & ν_2 ?

2. What are the transfer operators

3. How to solve for L_{2h}^{-1} ? What is this coarse grid operator?

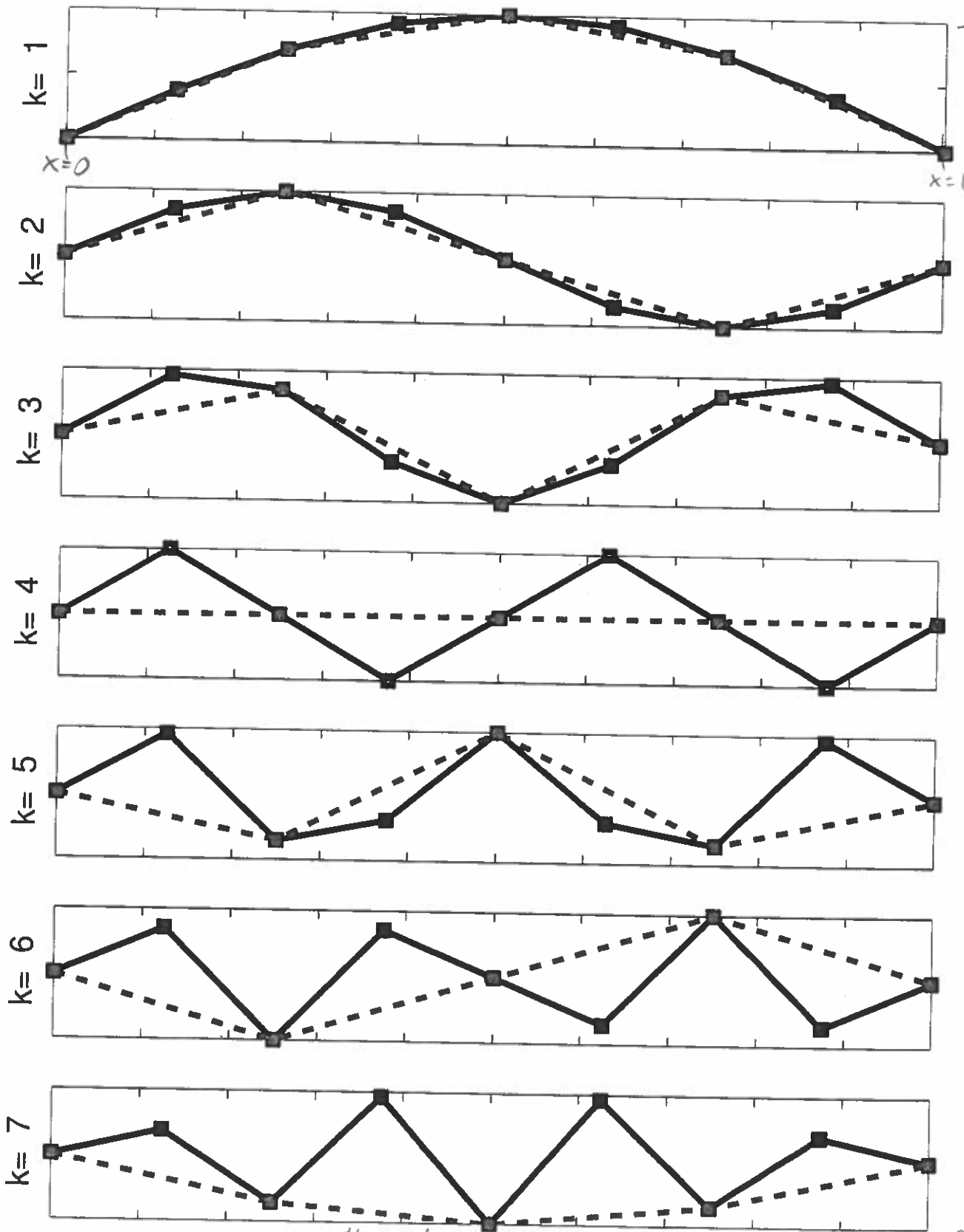
4. Which smoothing operator?

5. How efficient is this? $\Rightarrow O(N \ln N)$ work

What happens to sine waves
as we coarsen mesh?

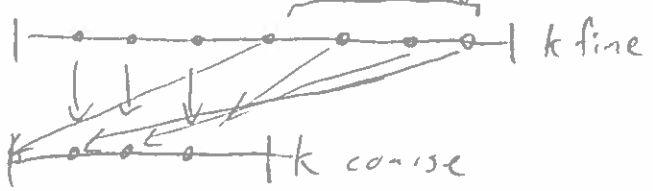
Lecture 17 today
11/10

— $\Rightarrow h=1$
--- $\Rightarrow h=\frac{1}{4}$



these look like $k=3, 2, 1$ on coarse mesh

use smoothing to erase these first!



(Aliasing error)

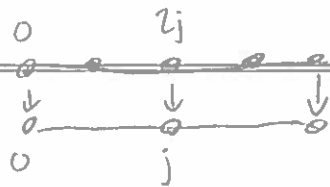
coarse grid correction alone is terrible!

creates extra low freq. error upon coarsening

5. Multigrid will be a $\mathcal{O}(N \ln N)$ algorithm

2. Transfer operators:

Restriction



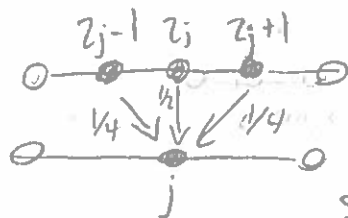
Injection

$$(u_h)_j = (u_h)_{z_j}$$

turns out to be not very good

Full-weighting

Wipes out high freqs. very well



$$(u_h)_j = \frac{1}{4}(u_h)_{z_{j-1}} + \frac{1}{2}(u_h)_{z_j} + \frac{1}{4}(u_h)_{z_{j+1}}$$

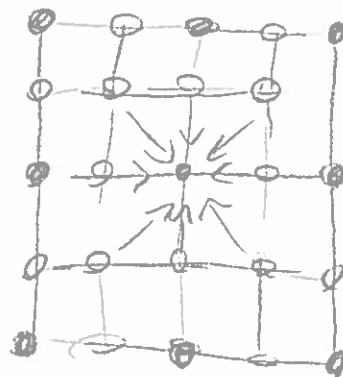
Adjoint of the linear interpolation op.

Stencil for Full-weighting:

1D: $I_h^{2h} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$

2D: $I_h^{2h} = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$
 $= \frac{1}{4} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$

$= (I_h^{2h})^T I_h^{2h} \leftarrow \text{outer/tensor product}$



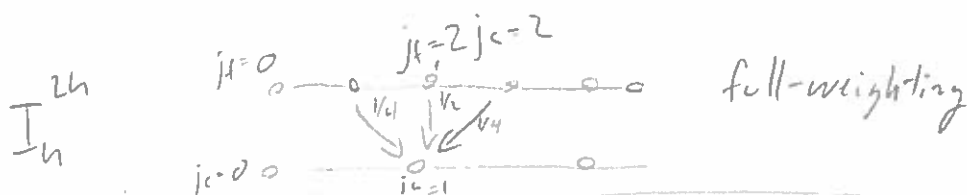
MAT228A - Lecture 16 - 11/15

Transfer operators

Restriction Operator

$$1D: I_h^{2h} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$$

$$2D: I_h^{2h} = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



pseudocode for 1D:

loop over coarse mesh j_c

compute j_f

$$u_c(j_c) = \frac{1}{4} (u_f(j_c-1) + 2u_f(j_c) + u_f(j_c+1))$$

(to vectorize, add 3 shifted arrays together)

Half-weighting

$$I_h^{2h} = \frac{1}{8} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 4 & 1 \\ 0 & 1 & 0 \end{bmatrix} \leftarrow \text{less information, so need better smoothing}$$

Not as good as full-weighting

Interpolation/Prolongation Operator

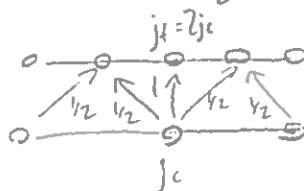
pseudocode for 1D: $u_f = \vec{0}$

loop over coarse mesh j_c

$$u_f(2j_c) = u_c(j_c)$$

$$u_f(2j_c-1) += \frac{1}{2} u_c(j_c)$$

$$u_f(2j_c+1) += \frac{1}{2} u_c(j_c)$$



$$(u_h)_{2j_c} = (u_{2h})_{j_c}$$

$$(u_h)_{2j_c-1} = \frac{1}{2} ((u_{2h})_{j_c-1} + (u_{2h})_{j_c})$$

Bilinear interp.

$$I_h^{1D}: I_h^{2h} = \frac{1}{2} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$$

$$I_h^{2D}: I_h^{2h} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



Matrix for Full-weighting restriction:

I_h^{2h} is a 3×7 matrix

$$I_h^{2h} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 2 & 1 \end{bmatrix}$$

Matrix for Interpolation:

I_h^{2h} is 7×3 matrix

$$I_h^{2h} = \frac{1}{2} \begin{bmatrix} 1 & 0 & 0 \\ 2 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\text{Clearly, } I_h^{2h} = \frac{1}{2} (I_h^{1D})^T$$

With the inner product $\langle u, v \rangle_h = h u^T v$,

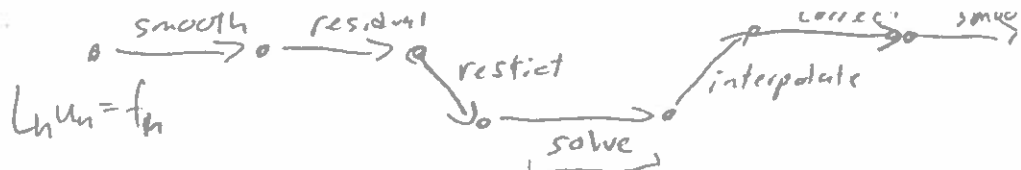
$$\langle u_h, I_h^{2h} v_{2h} \rangle_h = h u_h^T I_h^{2h} v_{2h}$$

$$= h (I_h^{2h})^T u_h^T v_{2h} = h (2 I_h^{1D} u_h)^T v_{2h}$$

$$= \langle T_h^{2h} u_h, v_{2h} \rangle$$

Full-weighting is the adjoint of interpolation

Multigrid/Z-grid overview



One way to define L_h :

Galerkin coarse grid op.

$$L_h = I_h^T L_h I_h$$

$$L_{2h} e_{2h} = r_{2h}$$

In 2D, stencil $L_h = \frac{1}{h^2} \begin{bmatrix} 1 & -4 & 1 \\ -4 & 16 & -4 \\ 1 & -4 & 1 \end{bmatrix} \Rightarrow L_{2h} = \frac{1}{(2h)^2} \begin{bmatrix} 1 & 2 & 1 \\ 2 & -12 & 2 \\ 1 & 2 & 1 \end{bmatrix}$, $L_{4h} = \dots$ else $L_{4h} = \dots$

This came from a PDE though! So just rediscretize the PDE!

Simpler approach:

$$L_{2h} = \frac{1}{(2h)^2} \begin{bmatrix} 1 & -4 & 1 \\ -4 & 16 & -4 \\ 1 & -4 & 1 \end{bmatrix} \leftarrow \text{rediscretization}$$

(Galerkin used in Algebraic multigrid \rightarrow more general, but slower!)

How well does this work?

Let $v_1 = \#$ of presmooth steps, $v_2 = \#$ of postsmooth steps, $V = v_1 + v_2$

\leftarrow away from start, v only # that matters in Z-grid.

Spectral radii for Z-grid, 2D, fullweighting, bilinear interpol, rediscretization

v	ω -Jac	GS-RB
1	0.6	0.25
2	0.360	0.074
3	0.216	0.053
4	0.137	0.041

\star These spectral radii are

independent of mesh spacing!

(before: $p = f(h)$ & $p \rightarrow 1$ as $h \rightarrow 0$)

$$\star \rho < C < 1$$

Recall: To det. iteration counts to reduce error by ϵ : $\rho^k = \epsilon$

$$\Rightarrow k = \ln \epsilon / \ln \rho$$

Q:

How much smoothing should we do?

$$k \propto -1 / \log_{10}(\rho)$$

A: Balance smoothing work w/ iteration work

$$\text{Work} = \underbrace{v}_{\text{smoothing}} + \underbrace{w}_{\text{everything else}}$$

Work per digit of accuracy:

total work

comparable work to 5 smooths

$$\frac{(v+w)}{-\log_{10}(\rho)} = (\text{work}) (\# \text{its. per digit})$$

$w=1$ is comparable work to 1 smooth

v	$w=0$	$w=1$	$w=5$	$w=6$
1	1.66	3.32	9.97	11.63
2	1.77	2.65	6.19	7.07
3	2.35	3.14	6.27	7.05
4	2.88	3.60	6.49	7.21

$$\text{total work} \propto \frac{(v+w)}{\log_{10}(\rho)}$$

Suggests $v=2$ gives lowest work!

Last time, we saw that $\boxed{r=2 \text{ smoothing steps}}$ gives the lowest work for GS-RB, full-weighting, bilinear interpolation, rediscritization 2-grid approx.

Now to pick r_1 & r_2 (The pre-smooth step # & post-smooth step #)?

Let M be the multigrid iteration operator. Then error at iteration k is

$$e^k = M e^{k-1} \Rightarrow \|e^k\|_2 \leq \|M\|_2 \|e^{k-1}\|_2 \Rightarrow \frac{\|e^k\|_2}{\|e^{k-1}\|_2} \leq \|M\|_2$$

for small k , $\|M\|_2$ is a better error approx. than ρ^k .

r_1	r_2	$\ M\ _2$
1	0	0.559
0	1	1.414
2	0	0.200
1	1	0.141
0	2	1.414 \leftarrow \text{bad!}

even the same amt. of work!

better to mix!

r_1	r_2	$\ M\ _2$
3	0	0.137
2	1	0.081
1	2	0.081
0	3	1.414

In general, common choices: $r_1=1, r_2=1$
(or $r_1=2, r_2=1$) since $r=2$ optimal
if $r=3$ optimal

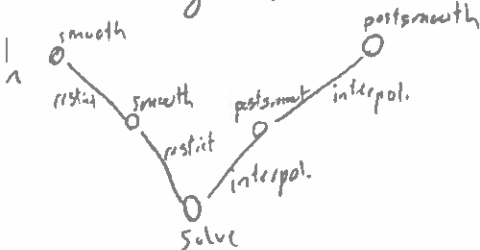
Moving from 2-grid to Multigrid:

Why use direct solve on the coarse grid?

Use another 2-grid as the solver!

Multigrid iteration:

Solving $L_h u_h = f_h$



This is one iteration of a 3-grid algorithm using a "V-cycle"

→ other forms of cycles:

γ iterations before returning to fine grid

V-cycle $\Rightarrow \gamma=1$.

smooth r_1 times & compute r_h

↳ restrict to r_h , solve $L_{2h} e_{2h} = r_{2h}$

↳ smooth $L_{2h} u_{2h} = f_{2h}$ r_1 times (initial guess $e_{2h}^0 = 0$)

↳ actually error & residual

↳ compute residual to error equation " r_{2h} "

↳ restrict to $r_{4h} \rightarrow$ solve $L_{4h} e_{4h} = r_{4h}$

rename & solve $L_{4h} u_{4h} = f_{4h}$ (for u_{4h})

interpolate to correct $u_{2h} = u_{2h} + I_{4h}^{2h} u_{4h}$

smooth r_2 times

interpolate & correct $u_h = u_h + I_{2h}^h u_{2h}$

smooth r_2 times

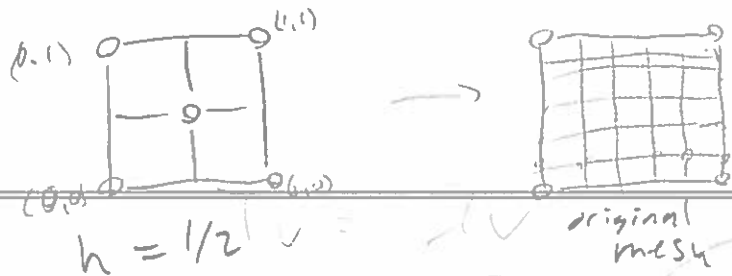
→ e.g., $\gamma=2$ on 3-grid



$\gamma=2$ on 4 grid, W-cycle



In practice, $\nu=1$ V-cycles & $\nu=2$ W-cycles are sufficient.
 How deep should we go? As far as we can! Until direct solve is trivial.



Pick powers of 2 as mesh spacing
 Go down until mesh $h=1/2$.

coarsest
 one-unknown \rightarrow trivial solve

total work is $C \cdot N$

for fine mesh

All fine
 mesh ops.
 are same
 order of
 work

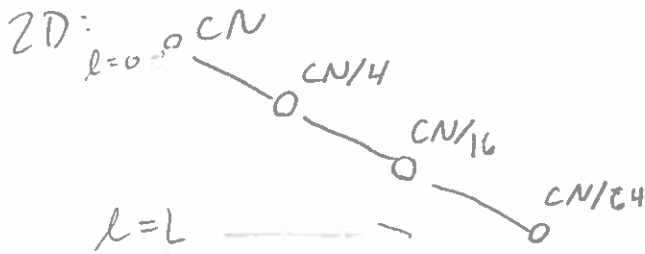
How much work in a V-cycle?

On the fine mesh, have to smooth ν times.
 $\rightarrow \mathcal{O}(N)$ work

\rightarrow compute residual $\rightarrow \mathcal{O}(N)$ work

\rightarrow restrict to coarse mesh $\rightarrow \mathcal{O}(N)$ work

\rightarrow correct/interp $\rightarrow \mathcal{O}(N)$ work



$$\Rightarrow \text{total work: } \sum_{l=0}^L CN \left(\frac{1}{4}\right)^l \approx \text{limit in large } L \quad CN \left(\frac{1}{1-\frac{1}{4}}\right) = \frac{4}{3} CN$$

Hence V-cycle work is not very different
 from work on fine level

GS-TB, FW, bilin. int.

2D: $\nu=2$

two-grid $\rho \approx 0.074$

V-cycle $\rho \approx 0.10$

256x256 mesh

tol = 10^{-7}

64x64 mesh

tol = 10^{-6}

	Iterations	work/it.	total work	
SOR	144	1	144	} 3x faster when
MG	6	$\frac{4}{3}(2+4)=8$	48	

	# its	work/it	total work
SOR	658	1	658
MG	7	8	56

$\sim 13x$ faster!
 order of magnitude faster!

$$L\vec{u} = \vec{f}$$

$$-\frac{4}{h^2}u = f$$

$$u = -f \cdot h^2/4$$

Conjugate Gradient Method: $Au = f$

Can solve using matrix-vector product
given u , compute Au
don't need explicit matrix.

where A is symmetric positive-definite
 $A^T = A \iff y^T A y > 0 \quad \forall y \neq 0$
 \Rightarrow all eigenvalues of A are positive & real
& complete set of orthogonal eigenvectors

CG is an example of a Krylov method \rightarrow others: minres - symmetric indefinite &
gmres - not sym, indefinite

CG is related to minimization

\hookrightarrow Suppose A is sym. pos. def., $N \times N$

Define functional $\phi: \mathbb{R}^N \rightarrow \mathbb{R} : \phi(u) = \frac{1}{2} u^T A u - u^T f$

Soln to $Au = f$ is the minimizer of ϕ

Pf. $\nabla \phi = \frac{1}{2} A u + \frac{1}{2} A^T u - f = A u - f = 0$ for $Au = f$

$\nabla \nabla \phi = A \leftarrow A$ pos. def., so ϕ is convex. $\Rightarrow Au = f$ soln minimizes ϕ

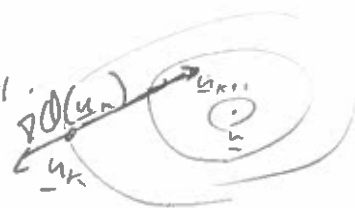
Method of Steepest Descents.

Have u_k , k^{th} iterate, want method for generating u_{k+1} .

$$\Rightarrow u_{k+1} = u_k + \alpha r_k$$

Pick α to minimize $\phi(u_{k+1}) = \min_{\alpha} \phi(u_k + \alpha r_k)$

$$\Rightarrow \alpha = \arg \min_{\alpha} \phi(u_k + \alpha r_k)$$



$$\nabla \phi(u_k) = A u_k - f = -r_k$$

compute:
 $\frac{d}{d\alpha} \phi(u + \alpha r) = 0$. $\phi(u + \alpha r) = \frac{1}{2} (u + \alpha r)^T A (u + \alpha r) - (u + \alpha r)^T f$
 $= \frac{1}{2} u^T A u - u^T f + \alpha (\frac{1}{2} r^T A u + \frac{1}{2} u^T A r - r^T f)$
 $+ \alpha^2 \frac{1}{2} r^T A r$ \swarrow same since scalar & A symm.

$$\Rightarrow \frac{d}{d\alpha} \phi(u + \alpha r) = (r^T A u - r^T f) + \alpha r^T A r$$

$$0 = r^T (A u - f) + \alpha r^T A r = -r^T r + \alpha r^T A r$$

\nwarrow note $\frac{d^2}{d\alpha^2} \phi(u + \alpha r) = r^T A r > 0$
so convex!

$$\Rightarrow \alpha = \frac{r^T r}{r^T A r}$$

Steepest Descents Algorithm

initialize u_0

loop in k

$$r_k = f - Au_k$$

check $\|r_k\|$ for stopping

$$\alpha = \frac{r_k^T r_k}{r_k^T A r_k}$$

$$u_{k+1} = u_k + \alpha r_k$$

work:

mainly 2 matrix-vector products

$$\text{faster: } r_{k+1} = f - Au_{k+1}$$

$$= f - A(u_k + \alpha r_k)$$

$$= r_k - \alpha A r_k$$

which we already compute in finding A .

More efficient method

initialize u_0 & $r_0 = f - Au_0$

loop k

check $\|r_k\|$

compute $w = Ar_k$

$$\alpha = \frac{r_k^T r_k}{r_k^T w}$$

$$u_{k+1} = u_k + \alpha r_k$$

$$r_{k+1} = r_k - \alpha w$$

only one matrix-vector product in $w = Ar_k$ step!

★ residual used as stopping criterion
b/c it is available & generally more robust \rightarrow doesn't dep. on method

Level curves of ϕ are ellipses
the shape is related to the eigenvalues

$$r = \lambda(u - v)$$

$$f - Av = \lambda(u - v)$$

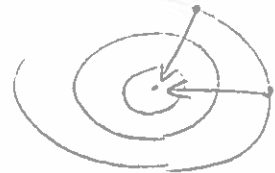
$$Au - Av = \lambda(u - v)$$

$$A(u - v) = \lambda(u - v)$$

Hence $u - v$ is an eigenvector of A
so residual at these points (in pic) are



If $\lambda_1 = \lambda_2$



any start. pt. will do
1 it. to find sol.

If $\frac{\lambda_2}{\lambda_1} \approx 1$

takes like 2 its.



method
converges
quickly

If $\frac{\lambda_2}{\lambda_1} \gg 1$



lots of little steps
slow convergence!

The convergence rate of steepest descents

depends on $K = \|A\|_2 \|A^{-1}\|_2$ the condition #.

For symm. matrix, $K = \frac{\max_k |\lambda_k|}{\min_k |\lambda_k|}$

MAT228A - Lecture 19 - 11/29

Recall: Have problem $A\vec{u} = \vec{f}$ where A is symm. pos. definite

- Only need to be able to compute matrix-vector products!

Steepest Descents guaranteed to converge!

↳ $\phi(u) = u^T A u - u^T f$ & minimize this functional to solve $A\vec{u} = \vec{f}$.

Speed depends on condition #, $\kappa = \|A\|_2 \cdot \|A^{-1}\|_2$

For symm. matrix, $\kappa = \max_k |\lambda_k| / \min_k |\lambda_k|$

Conjugate Gradient Descent - search direction is different than residual

Let p_k be the vector to follow/search: $u_{k+1} = u_k + \alpha p_k$.

Follow this search path until ϕ increases, i.e. the minimizer of ϕ on a 1D space

Find $\alpha = \frac{p_k^T r_k}{p_k^T A p_k}$

Start w/ Conj. grad. in 2D:

- initial guess u_0

↳ compute initial residual r_0

↳ initialize $p_0 = r_0$

→ one step of steepest descent to get u_1 : $u_1 = u_0 + \alpha p_0$, $\alpha = \frac{p_0^T r_0}{p_0^T A p_0}$

- pick p_1 so that $p_0^T A p_1 = 0$.

↳ new direction not orthogonal to old, but A -conjugate,
i.e. orthogonal in the A -inner product $\langle p_0, p_1 \rangle_A = p_0^T A p_1 = 0$
since A is symm. pos. def.

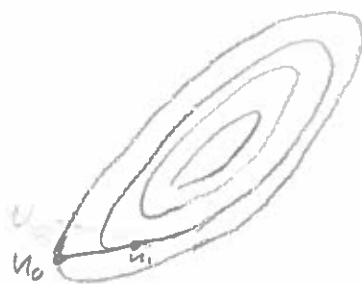
Why?

↳ p_0 is tangent to level set of ϕ at $\phi(u_1)$, i.e. $p_0^T r_1 = 0$

$\Rightarrow p_0^T (f - A u_1) = 0 \Rightarrow p_0^T (A u - A u_1) = 0 \Rightarrow p_0^T A (u - u_1) = 0$

So p_1 will be in the direction of $u - u_1$.

$\Rightarrow u_2$ is the solution (in 2D).



Conj. Grad Desc. in 3D:

• Initialize u_0 , compute r_0 , set $p_0 = r_0$.

• get $u_1 = u_0 + \alpha p_0$, $\alpha = \frac{p_0^T r_0}{p_0^T A p_0}$.

~~Pick p_1 to be A -conjugate to p_0 (2D space to pick from)~~

~~\hookrightarrow pick some direction in that space (projection of residual into A -conj. space)~~

~~p_1 & r_1 define new α , $u_2 = u_1 + \alpha p_1$.~~

~~p_0 & p_1 span a plane $u_2 + c_0 p_0 + c_1 p_1$ which is tangent to level surface of ϕ at $\phi(u_2)$.~~

~~Now pick p_2 A -conjugate to both p_0 & p_1 .~~

~~\hookrightarrow Minimize ϕ along p_2 : $u_3 = u_2 + \alpha p_2$ (solved).~~

Done b/c $(c_0 p_0 + c_1 p_1)^T r_2 = 0$

$$(c_0 p_0 + c_1 p_1)^T (f - A u_2) = 0$$

$$(c_0 p_0 + c_1 p_1)^T (A u - A u_2) = 0$$

$$(c_0 p_0 + c_1 p_1)^T A (u - u_2) = 0$$

\hookrightarrow pick p_2 in this direction!

★ CGD guaranteed to converge (w/ infinite arithmetic precision, no round-off) to the exact solution in N steps ($N = \#$ of unknowns).

Usually get "close" in far fewer steps.

Pseudo code for CGD algorithm: • initialize: $u_0, r_0 = f - A u_0, p_0 = r_0$

• loop k (still in loop)

$$\hookrightarrow w = A p_k$$

$$\hookrightarrow \alpha = \frac{r_k^T r_k}{w^T p_k}$$

$$\hookrightarrow u_{k+1} = u_k + \alpha p_k$$

$$\hookrightarrow r_{k+1} = r_k - \alpha w$$

\hookrightarrow check $\|r_{k+1}\|$ for stopping

end loop

$$\hookrightarrow \beta = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

$$\hookrightarrow p_{k+1} = r_{k+1} + \beta p_k$$

\hookrightarrow this is r_{k+1} with stuff in the (p_0, p_1, \dots, p_k) -space projected off in A -inner prod

★ save $r_k^T r_k$ since we use it several times.

Sketch of analysis of convergence of CGD:

- B/c A is sym. pos. def., can define inner product $\langle u, v \rangle_A = u^T A v$ & a norm $\|u\|_A^2 = u^T A u$.

Can show that error after n iterations of CGD is

$$\|e_n\|_A^2 \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^n \|e_0\|_A^2$$

If κ is close to 1, error converges rapidly!

If κ large, error converges very slowly.

Preconditioning - idea is to change condition number & solve unrelated problem.

Thm (from some book) The vectors in the CG algorithm have the following properties provided $r_k \neq 0$:

- 1) p_k is A -conjugate to all previous search directions p_0, p_1, \dots, p_{k-1}
- 2) r_k is orthogonal to all previous residuals r_0, r_1, \dots, r_{k-1} .
- 3) The following subspaces of \mathbb{R}^N are identical:
 $\text{span}(p_0, p_1, \dots, p_{k-1})$, $\text{span}(r_0, A r_0, A^2 r_0, \dots, A^{k-1} r_0)$,
 $\text{span}(A e_0, A^2 e_0, \dots, A^k e_0)$

Define $K_n = \text{span}(r_0, A r_0, \dots, A^{n-1} r_0)$ n -dim'l Krylov space associated w/ r_0 .

Krylov methods perform minimizations on successive/nested Krylov spaces. ^{*}

Note $u_n \in u_0 + K_n$ (affine space) & u_n minimizes ϕ in this space.

→ Minimizing ϕ is equivalent to minimizing $\|e_j\|_A^2$.

$$\begin{aligned} \|e_j\|_A^2 &= e_j^T A e_j = (u_j - u)^T A (u_j - u) = u_j^T A u_j - 2 u_j^T A u + u^T A u \\ &= u_j^T A u_j - 2 u_j^T f + u^T A u \\ &= 2 \phi(u_j) + C \end{aligned}$$



Recall: Minimizing $\phi(u)$ is equivalent to $\min \|e_k\|_A^2$ for CG.

Key thm: $\text{span}(p_0, p_1, \dots, p_{k-1}) = \text{span}(Ae_0, A^2e_0, \dots, A^ke_0)$

$$u_k = u_0 + \alpha_0 p_0 + \dots + \alpha_{k-1} p_{k-1}$$

$$u_k = u_0 + c_1 Ae_0 + c_2 A^2e_0 + \dots + c_k A^ke_0$$

subtract solution!

$$e_k = e_0 + c_1 Ae_0 + c_2 A^2e_0 + \dots + c_k A^ke_0 = q(A) \cdot e_0 \quad \leftarrow \text{matrix polynomial of } A.$$

CG picks $q \in \Pi_k = \text{space of polynomials of deg. at most } k \text{ with } q(0)=1$ to minimize $\|e_k\|_A = \|q(A)e_0\|_A$.

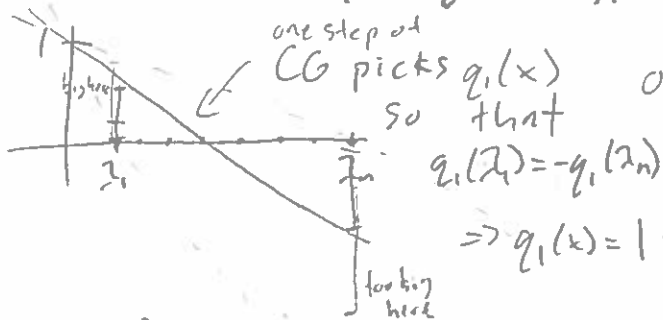
Note: A is diagonalizable: $A = Q \Lambda Q^T$ and $Q^T = Q^{-1}$.

$$\text{so } A^2 = Q \Lambda Q^{-1} Q \Lambda Q^{-1} = Q \Lambda^2 Q^{-1}$$

$$\text{thus } A^k = Q \Lambda^k Q^{-1}$$

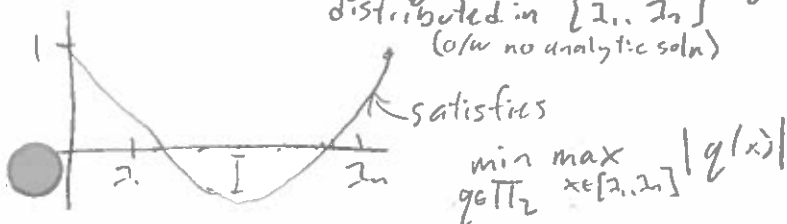
Fact: $q(A) = Q q(\Lambda) Q^{-1} = Q \begin{pmatrix} q(\lambda_1) & & \\ & \ddots & \\ & & q(\lambda_n) \end{pmatrix} Q^{-1}$

can show $\|e_k\|_A^2 = \|q(A)e_0\|_A^2 \leq \max_j (q(\lambda_j))^2 \|e_0\|_A^2$



our error bound comes from minimizing this (q on the spectrum of A).

2 steps of CG: assuming eigenvalues uniformly distributed in $[\lambda_1, \lambda_n]$ (or no analytic soln)



In general, can use scaled & shifted Chebyshev polys. to get the bound

$$\|e_k\|_A = \|q_k(A)e_0\|_A \leq 2 \left(\frac{\sqrt{K}-1}{\sqrt{K}+1} \right)^k \|e_0\|_A \quad \leftarrow \text{condition \#}$$

How well does this work? For Discrete Laplacian, $K = \Theta(h^{-2})$ since smallest $\lambda = \Theta(1)$ largest $\lambda = \Theta(h^{-2})$. For large K , #iterations to converge to a rel. tol $\frac{\|e_k\|_A}{\|e_0\|_A} < \epsilon$ is $\sqrt{K} \log \epsilon$. So for its $\Theta(h^{-1})$ it.

For discrete Laplacian, # its. to converge to rel. tol. of ϵ ($\frac{\|e_k\|_A}{\|e_0\|_A} < \epsilon$) is $\sqrt{K} \log \epsilon = \mathcal{O}(h^{-1}) = \mathcal{O}(n) = \mathcal{O}(\sqrt{N})$.
 each iteration costs $\mathcal{O}(N)$, so in 2D, total work = $\mathcal{O}(N^{3/2})$
 (same scaling as SOR).

~~What if we cluster the eigenvalues into groups? Then the polynomials we're minimizing over the spectrum can be lower order & CG converges way faster! This is Preconditioning~~

Preconditioning CG: $Au = f$ \leftarrow slow convergence if K is big for A .
 Change the problem $\Rightarrow M^{-1}Au = M^{-1}f$ \leftarrow same solution u , but spectrum of $M^{-1}A$ is diff. than.

Pick M^{-1} s.t. $M^{-1}A$ has a smaller K & clustered eigenvalues.

For CG, the pre-conditioner M should be: 1) symmetric, positive-definite

2) $M^{-1}A$ better conditioned
 3) M^{-1} easy to apply, i.e., $Mx = b$ easy to solve

Naively, M^{-1} should approx. A^{-1} so that $K \approx 1$.

Problem: $M^{-1}A$ is not symmetric!

$$Au = f \Rightarrow B^{-1}Au = B^{-1}f \Rightarrow \underbrace{(B^{-1}AB^{-T})}_{\text{left \& right preconditioning}} \underbrace{(B^T u)}_{\tilde{u}} = \underbrace{(B^{-1}f)}_{\tilde{f}}$$

$$\text{Define } \tilde{A} = B^{-1}AB^{-T}, \tilde{u} = B^T u, \tilde{f} = B^{-1}f \Rightarrow \tilde{A}\tilde{u} = \tilde{f}.$$

\hookrightarrow symmetric, pos. def. since $y^T B^{-1}AB^{-T}y = (B^{-T}y)^T AB^{-T}y = z^T A z > 0$ for $y \neq 0$.

$$\text{Note: } B^{-T}B^{-1}AB^{-T}B^T = B^{-T}B^{-1}A = (BB^T)^{-1}A = M^{-1}A \Rightarrow M = BB^T$$

\tilde{A} has the same eigenvalues as $M^{-1}A$ since $B^{-T}\tilde{A}B^T$ is just a change of basis.

Write CG in \sim variables & then transform back to original variables.

$$u_k = B^{-T}\tilde{u}_k, p_k = B^{-T}\tilde{p}_k, r_k = B\tilde{r}_k, \text{ then } B, B^T \text{ drop out of algorithm.}$$

$$\text{One step of CG: } \tilde{p}_{k+1} = \tilde{r}_{k+1} + \tilde{\rho}_k \tilde{p}_k$$

$$B^T p_{k+1} = B^T \tilde{r}_{k+1} + \tilde{\rho}_k B^T p_k \Rightarrow p_{k+1} = B^{-T}B^T r_{k+1} + \tilde{\rho}_k p_k = M^{-1}r_{k+1} + \tilde{\rho}_k p_k$$

Pre-conditioned Conjugate Gradient Descent Algorithm

initialize $r_0 = f - Au_0$

1) solve $Mz_0 = r_0$ or compute $z_0 = M^{-1}r_0$ (if have action apply M^{-1})

initialize $p_0 = z_0$

loop k

$$W_k = Ap_k$$

$$\alpha = \frac{z_k^T r_k}{p_k^T W_k}$$

$$u_{k+1} = u_k + \alpha p_k$$

$$r_{k+1} = r_k - \alpha W_k$$

check $\|r_{k+1}\|$ for stopping criterion

$$\text{compute } z_{k+1} = M^{-1}r_{k+1}$$

$$\beta = \frac{z_{k+1}^T r_{k+1}}{z_k^T r_k}$$

$$p_{k+1} = z_{k+1} + \beta p_k$$

end

How do we pick the pre-conditioner M^{-1} ?

• Easy/naive implementation: $M^{-1} = D^{-1}$, Dis diag. of $A \leftarrow$ Does nothing for constant coefficient problem (ie. Poisson eqn) since D^{-1} is a scalar.

• Use other iteration schemes, but symmetrized:

• SSOR \leftarrow symmetric SOR, loop down mesh once, loop back up (2x total)

• MG \leftarrow Using a symmetric smoother (RB-BR, etc.)

• Approximate factorizations - incomplete LU / cholesky factorization

\hookrightarrow know nothing, just algebraic.

